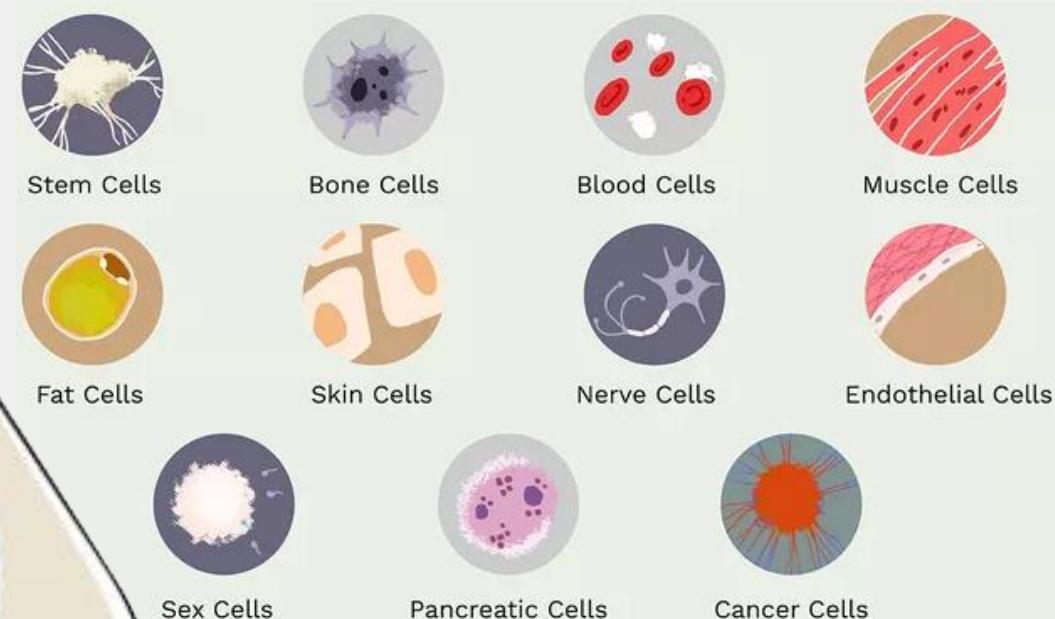
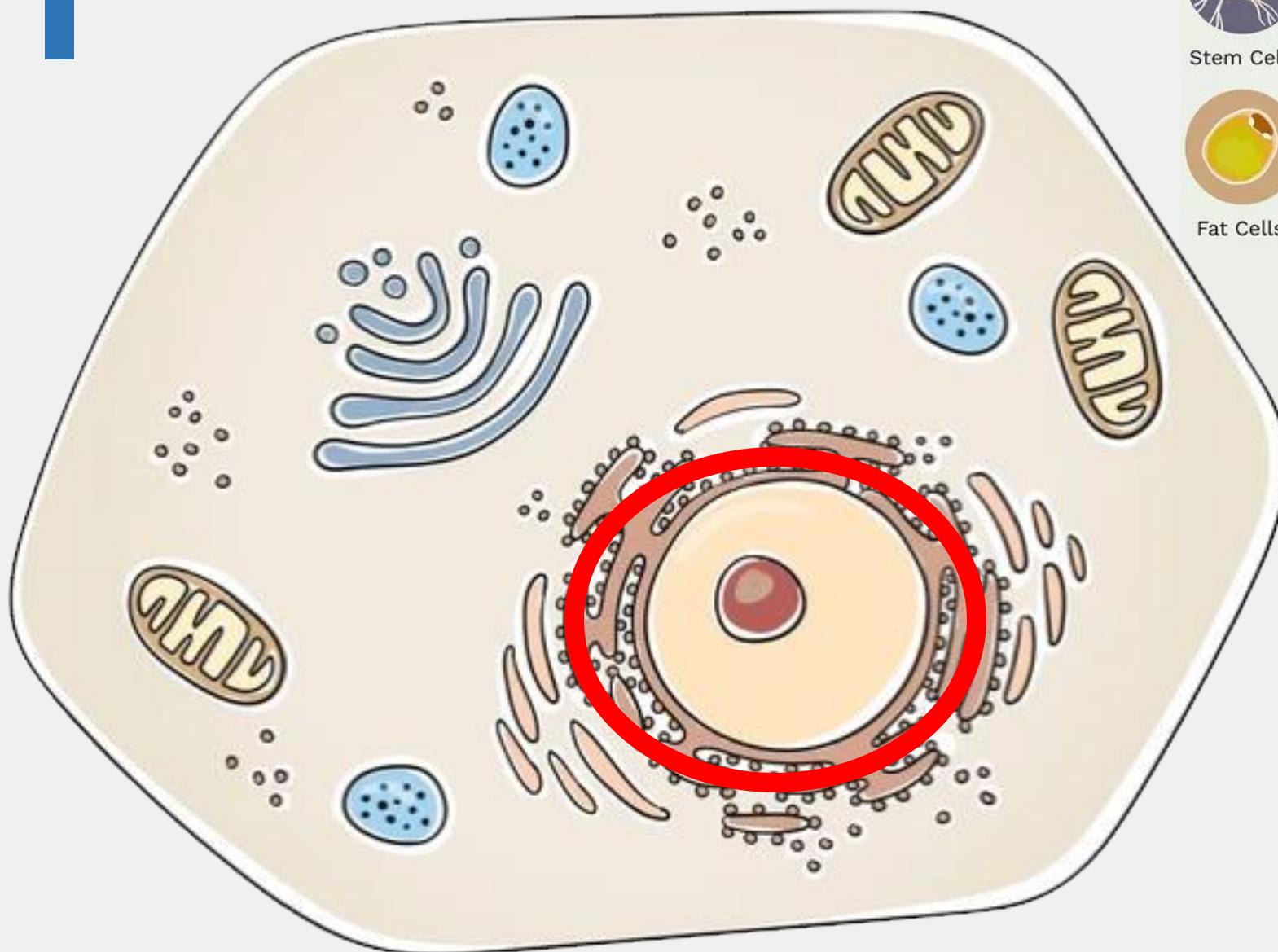


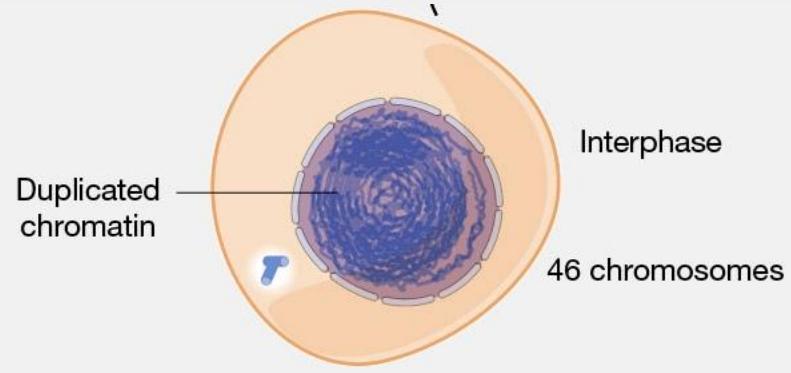


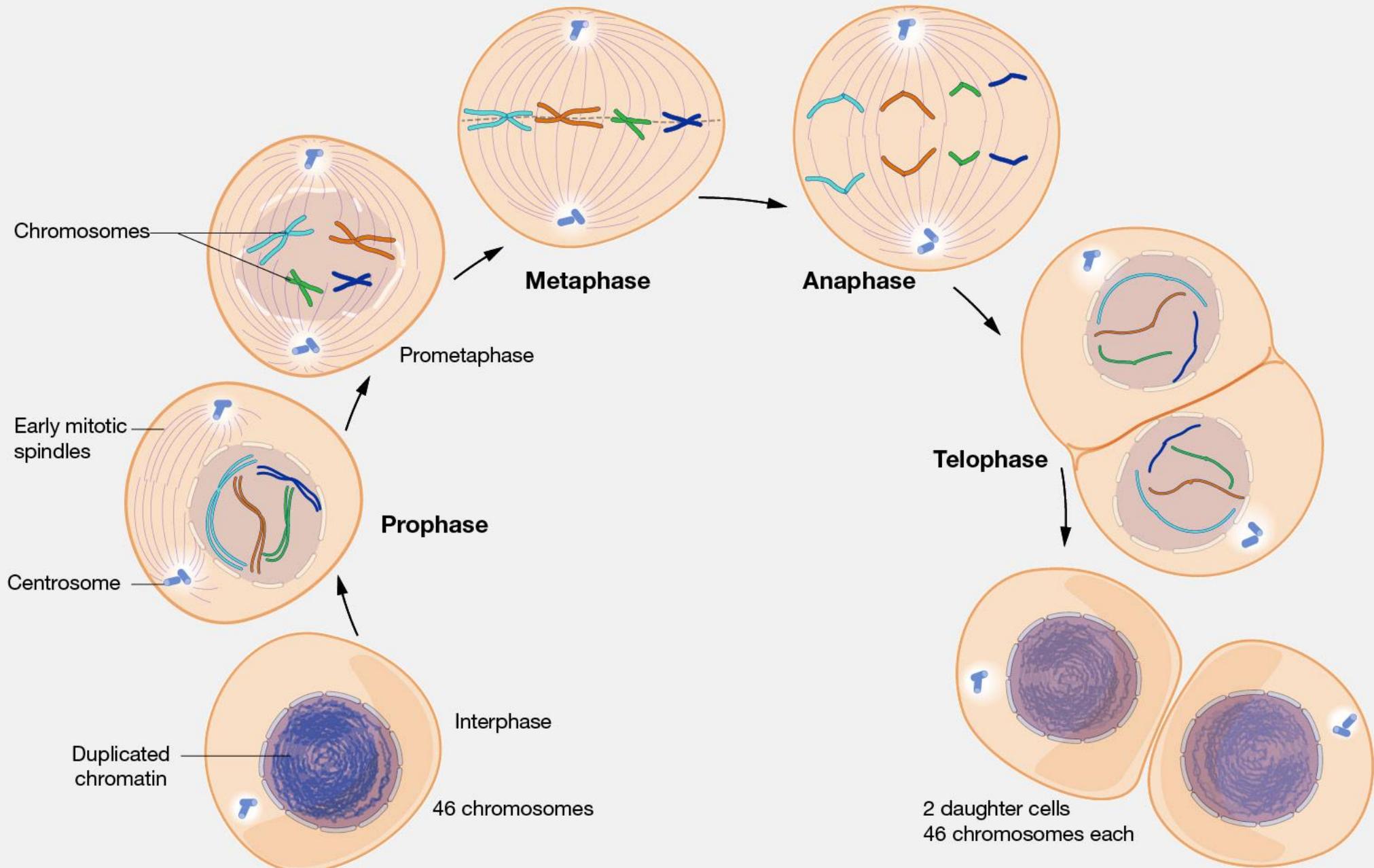
Decoding the Language of Genomes: Bridging Sequences and Function through Deep Learning

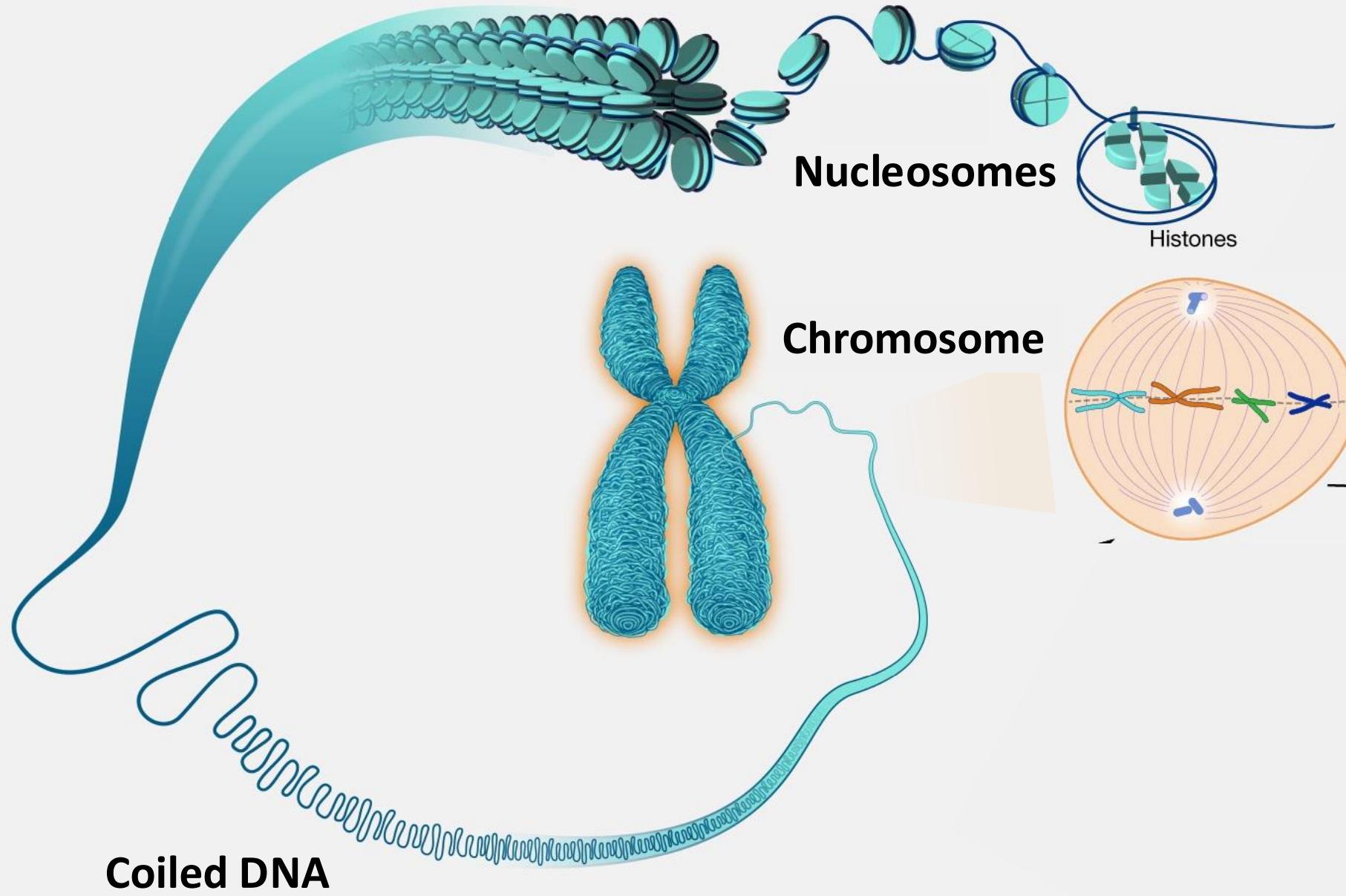
Kuan-Hao Chao

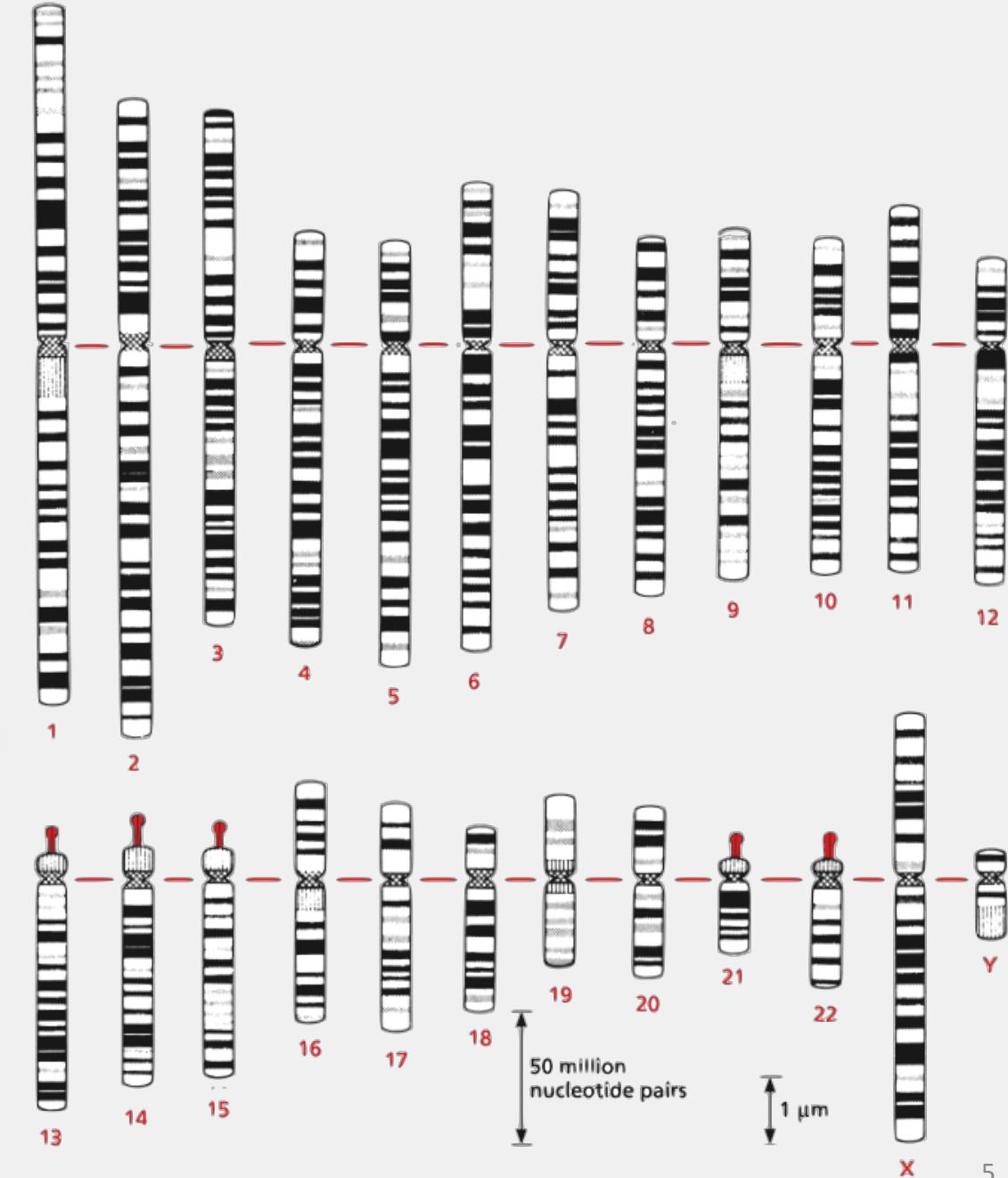
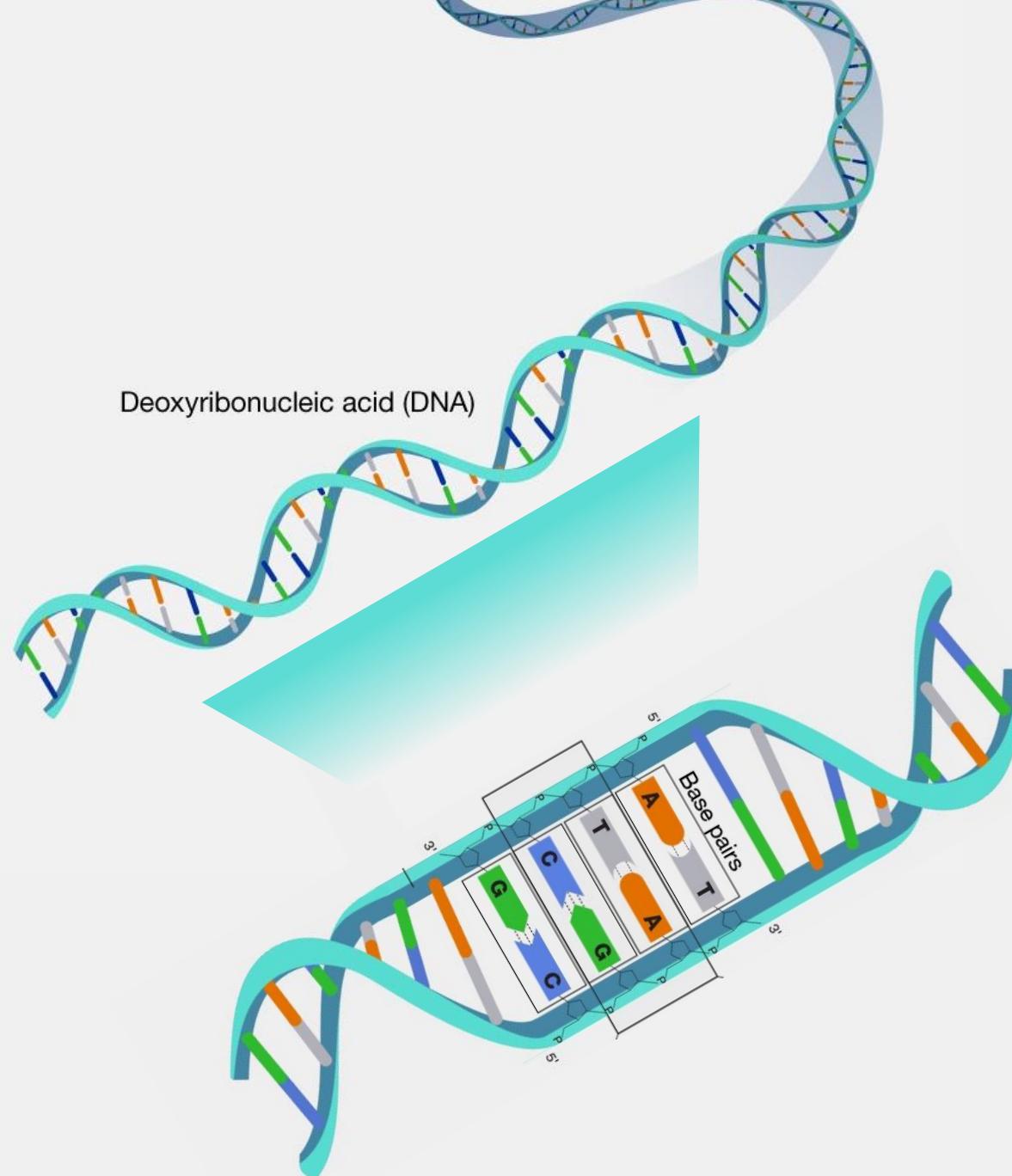
2025.08.25

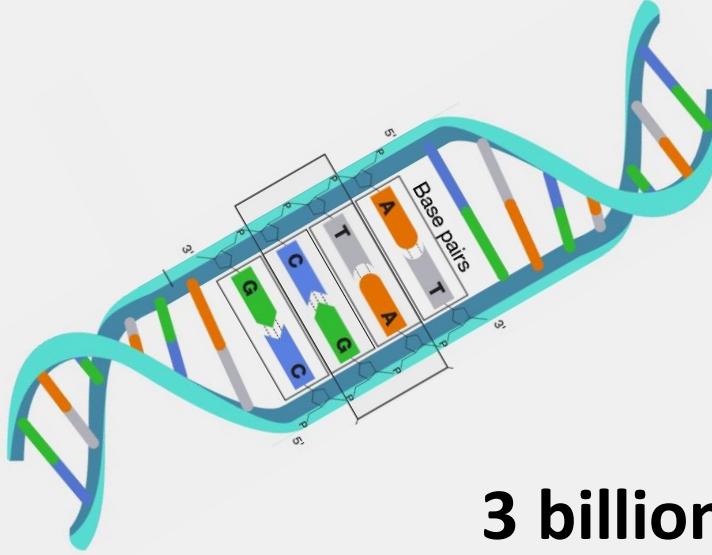




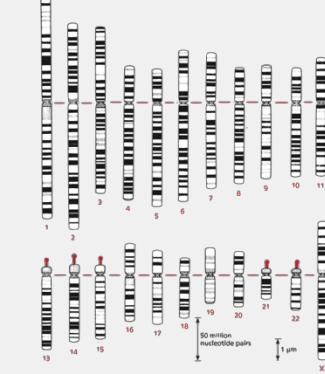




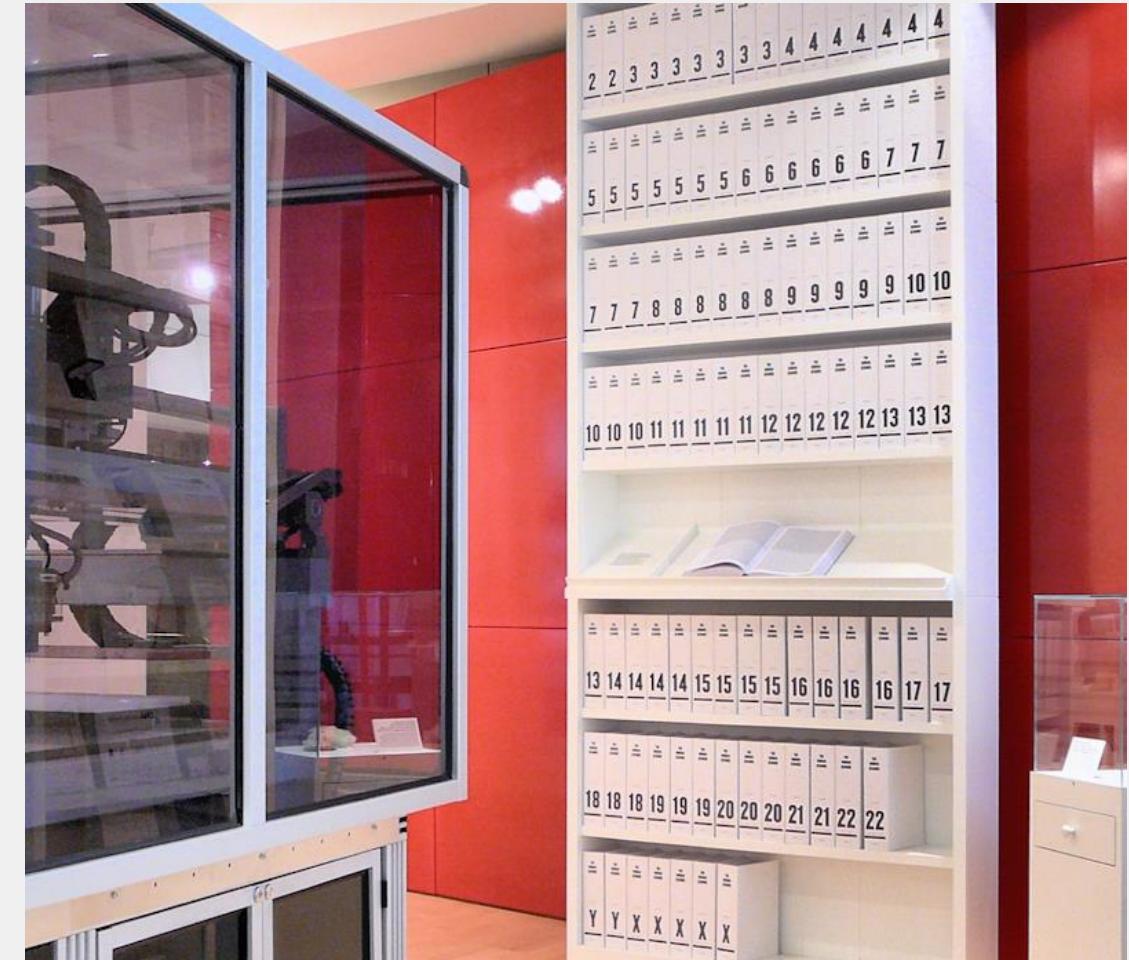




3 billions of ACGT



Wellcome genome bookcase

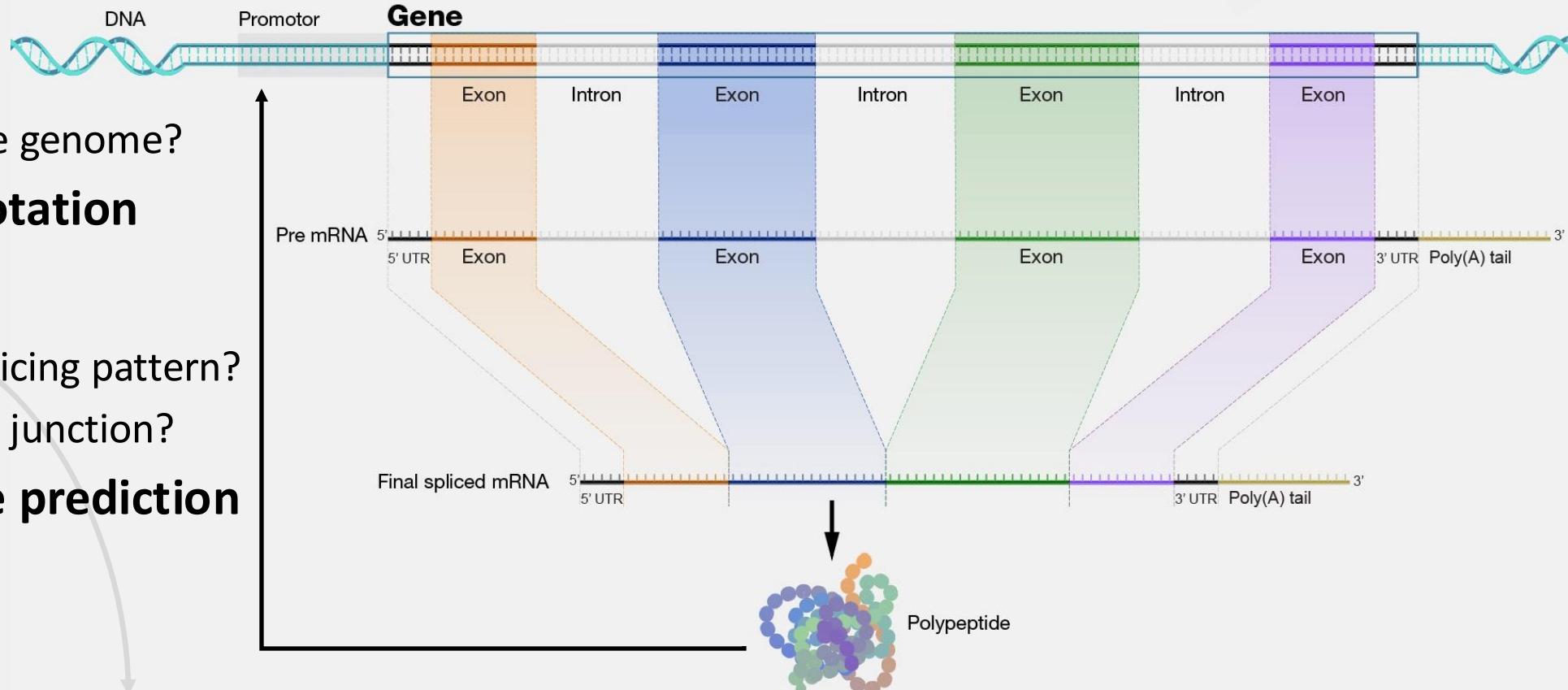


How do we assemble the complete 3-billion-nucleotide genome?

How can we efficiently represent multiple genomes for fast pattern matching?



Part I & II: Genome Assembly & Indexing



Where are the genes in the genome?

Part III: Genome Annotation

What are the canonical splicing patterns?

Alternative splicing? Splice junction?

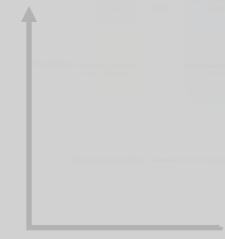
Part IV & V: Splice site prediction

Can we predict gene expression by learning the regulatory grammar in the genome?

Part VI: Shorkie. RNA-Seq coverage prediction



I Part I & II



Steven Salzberg



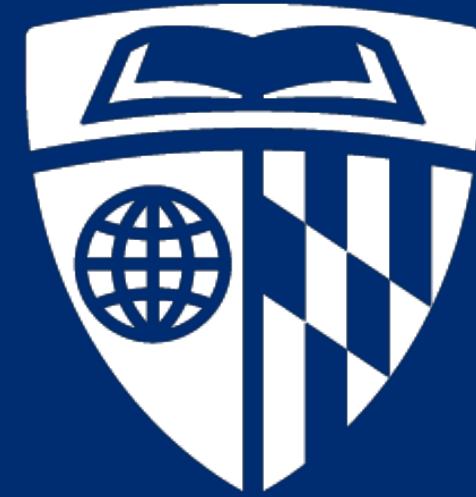
Mihaela Pertea



Ben Langmead

Genome Assembly & Indexing

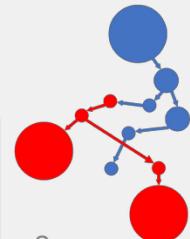
- Han1: first gapless Southern Han Chinese genome
- WGT: Wheeler graph recognition algorithm



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.

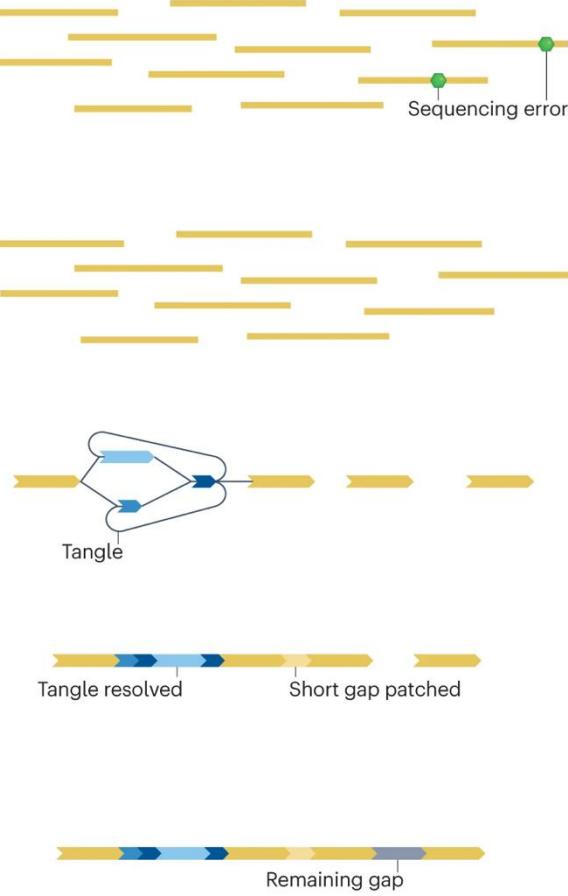


Chao, K. H., Chen, P. W., Seshia, S. A., & Langmead, B. (2023). WGT: Tools and algorithms for recognizing, visualizing, and generating Wheeler graphs. *iScience*, 26(8).

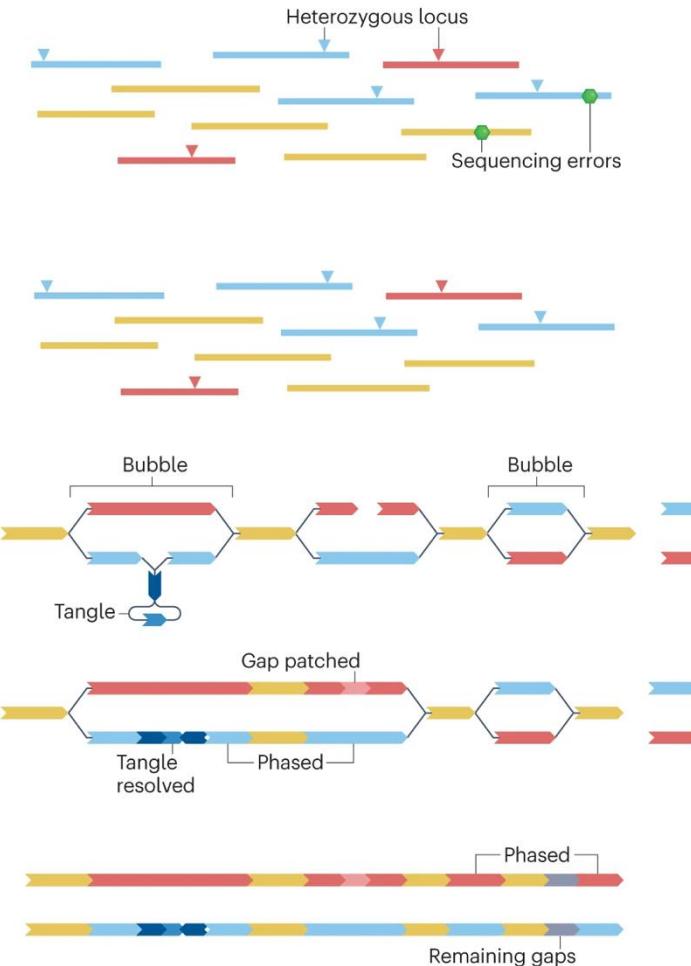


Genome Assembly

a Homozygous genome



b Heterozygous genome



1972 Bacteriophage MS2 coat gene

1976 Bacteriophage MS2 ~3.5 kb

1977 Sanger sequencing

1990-2000: BAC-by-BAC shot gun strategy

2003 near-complete human genome

2000s second-generation sequencing
short-read sequencing
(Illumina)

2010s third-generation sequencing
long-read sequencing
(PacBio; ONT)

2019 Telomere-to-Telomere (T2T) era

2022 First complete human genome

Li, H., Durbin, R. Genome assembly in the telomere-to-telomere era. Nat Rev Genet 25, 658–670 (2024).

Genome Assembly: Han1

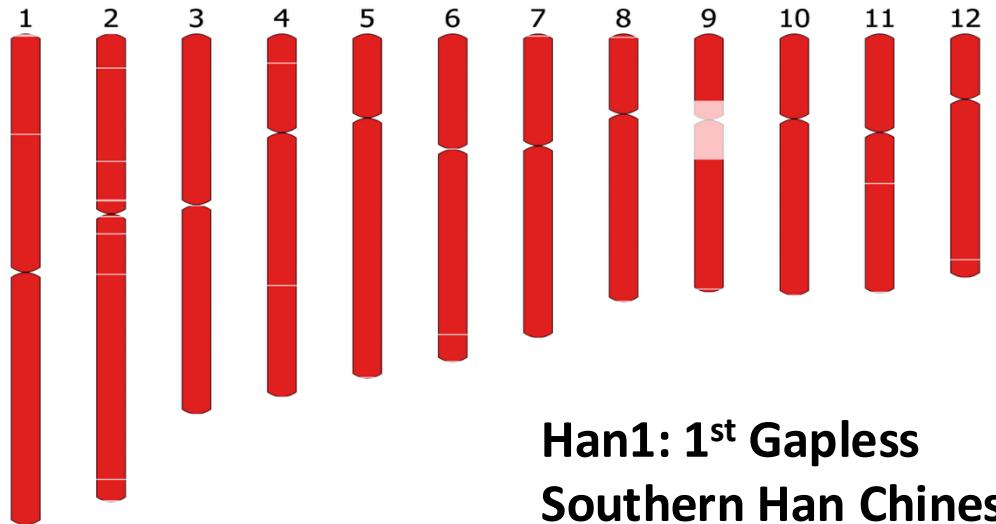


Table 3.

Statistics for the preliminary Han1 assemblies.

Assembler	Sequencing data	assembled sequence (bp)	Contig N50	Number of contigs	Quality value
Hifiasm v0.16.1-r375	PacBio HiFi	3,110,501,483	95,769,069	182	57.8a
Flye v2.5	ONT Ultralong	2,974,205,132	40,850,737	1,658	25.6a

Table 1.

A comparison among Han1, GRCh38, and the assemblies of previously released Chinese genomes.

Genome	Ethnicity	Contig N50 (Mb)	Number of contigs	Number of gaps	Assembly size (Gb)
Han1a	SH Chinese	148.02	25	0	3.10
HG00621 (hifiasm)b	SH Chinese	95.77	182	157	3.11
T2T-CHM13v2.0c	Northern European	150.62	25	0	3.12
HJ-H1d	NH Chinese	28.15	1,330	427	3.07
HJ-H2d	NH Chinese	25.90	896	390	2.91
NH1d	NH Chinese	3.60	11,019	8,484	2.89
HX1d	SH Chinese	8.33	5,843	4,025	2.93
YH2.0e	SH Chinese	0.02	361,157	235,514	2.91
TJ1.p0f	Tujia	13.67	1,430	907	2.87
TJ1.p1f	Tujia	13.70	1,426	873	2.87
ZF1g	Tibetan	23.62	1,384	1360	2.85
GRCh38.p14h	Mixed	57.88	994	804	3.10

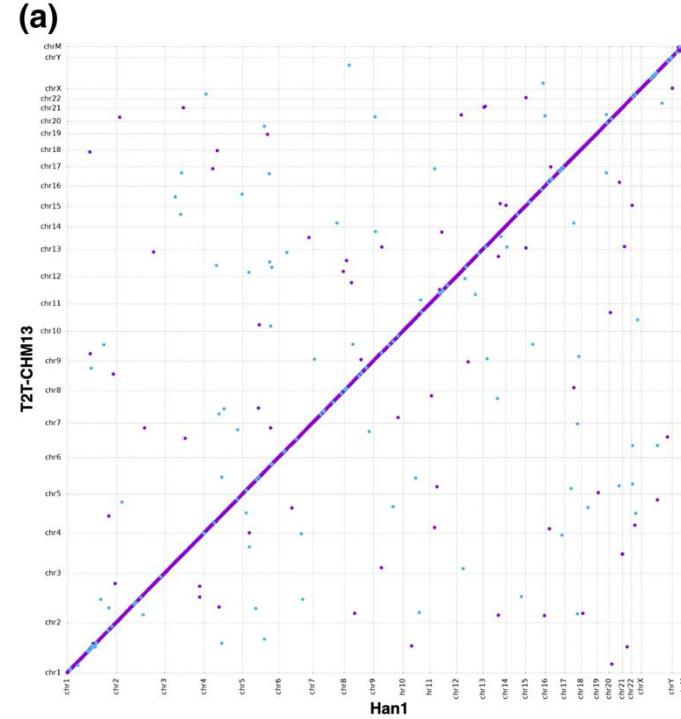


GenBank assembly:
GCA_024586135.1

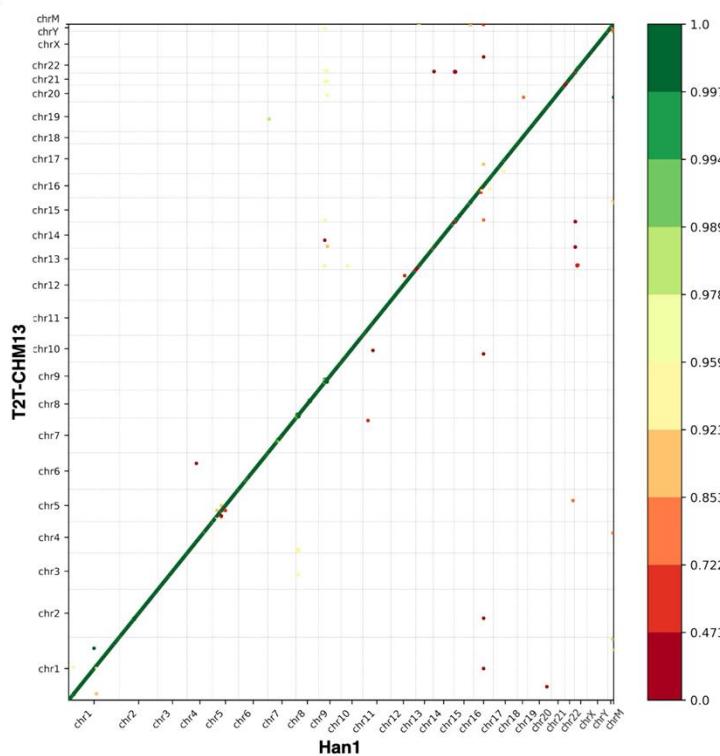


GRCh38, T2T-CHM13, and other assemblies of Chinese genome

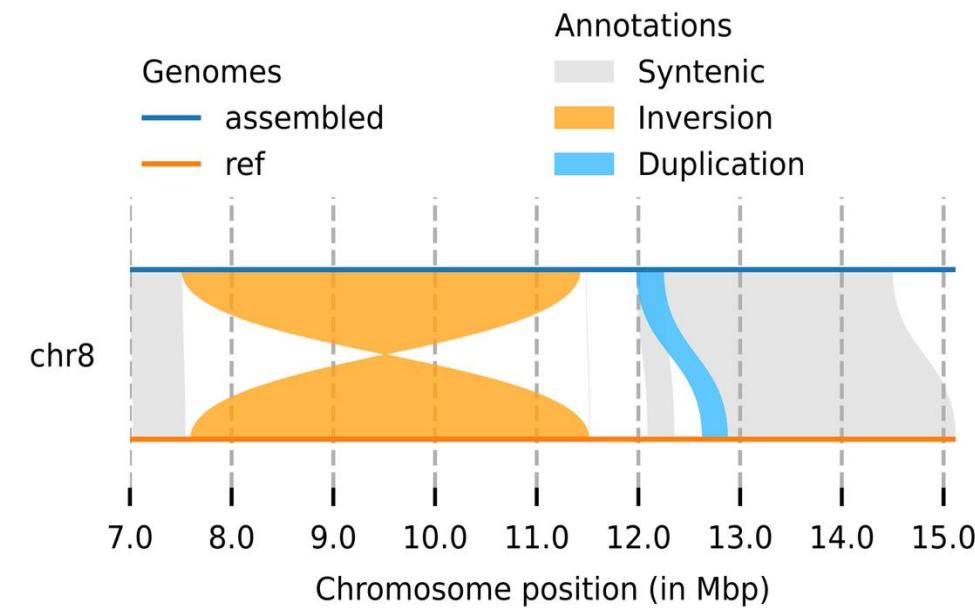
Mummer dot plot



(b)

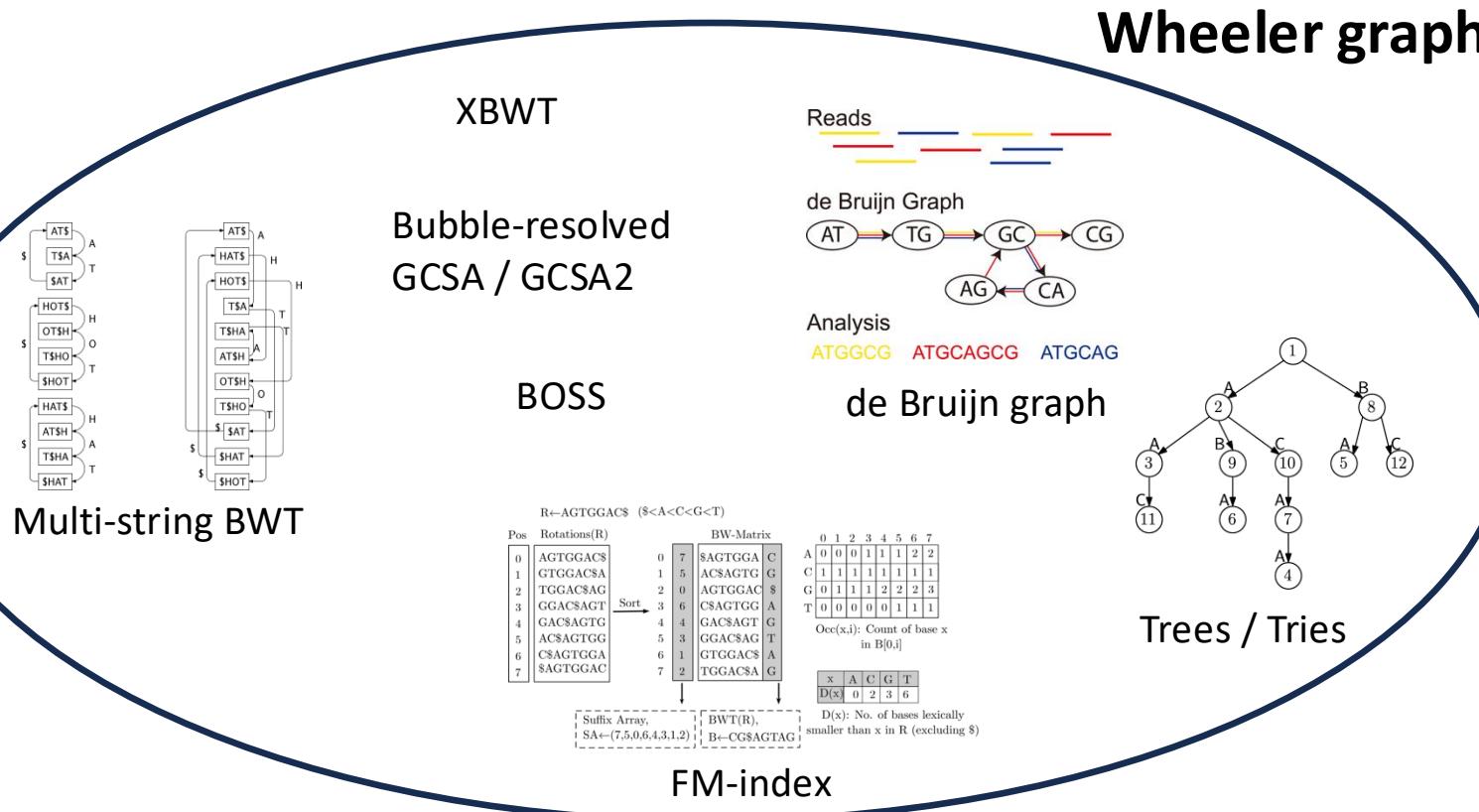


Inversion on β defensin gene cluster



Wheeler Graph Recognition

- Wheeler graphs are the basis of many pangenome and other sequence analysis tools



- 1994 Burrows-Wheeler Transform (BWT)
- 2000 FM-index
- 2009 Bowtie
- 2012 colored de Bruijn graphs
- 2012 BOSS data structure
- 2014 GCSA
- 2017 GCSA2
- 2017 VARI
- 2017 Wheeler graph**
- 2018 Vg toolkit
- 2020 Graph Burrows-Wheeler Transform (GBWT)
- 2021 vg Giraffe

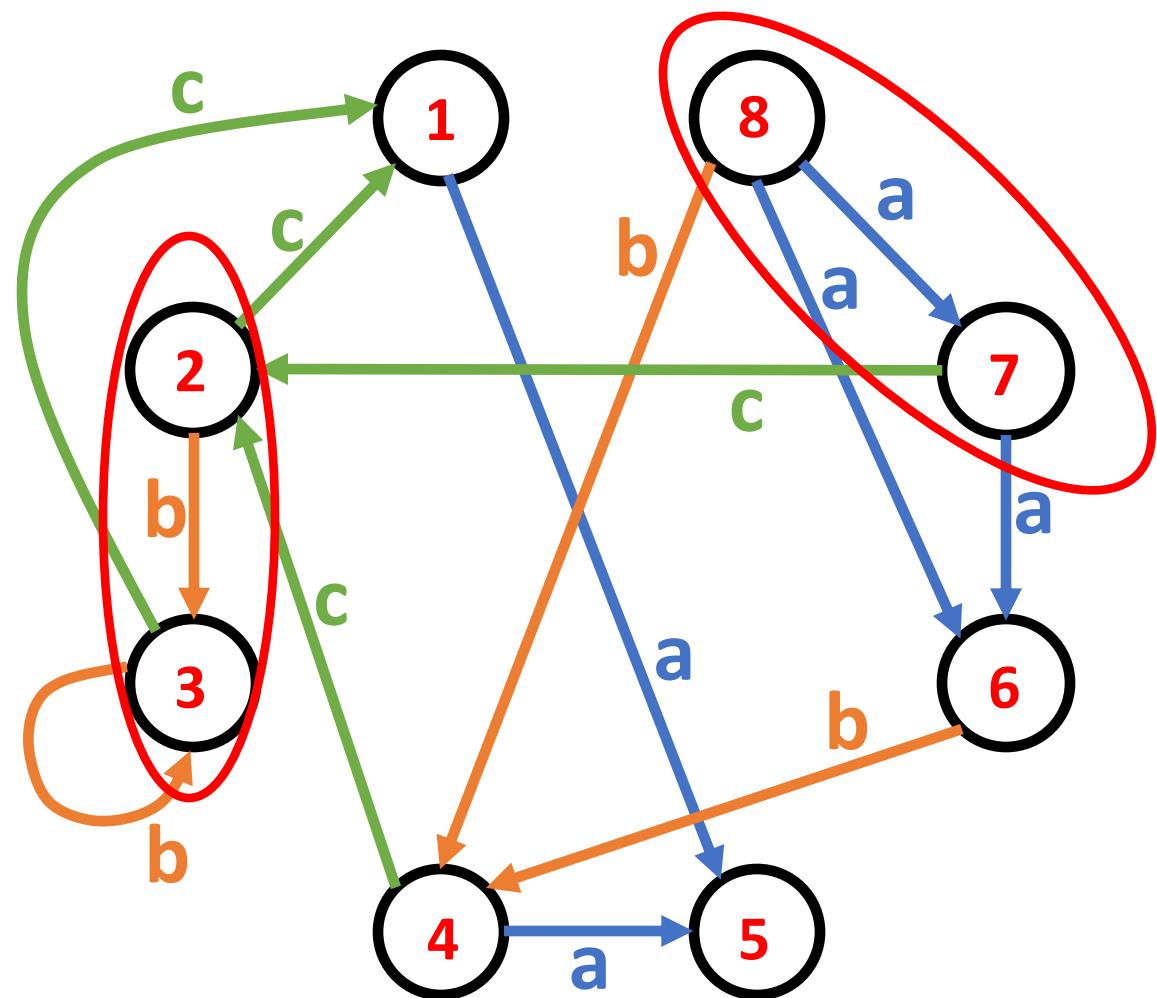
Wheeler Graph Recognition



RECOMB-SEQ

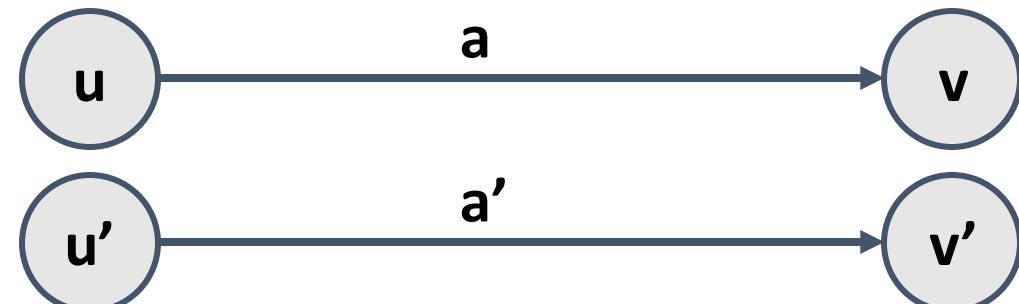


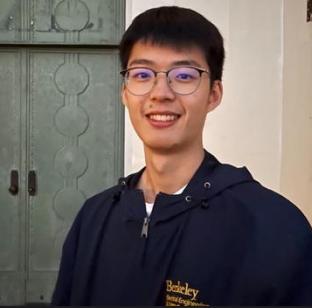
iScience



Wheeler graph recognition problem

1. Node with indegree 0 comes before every other nodes.
2. $a \prec a' \implies v < v'$
3. $(a = a') \wedge (u < u') \implies v \leq v'$





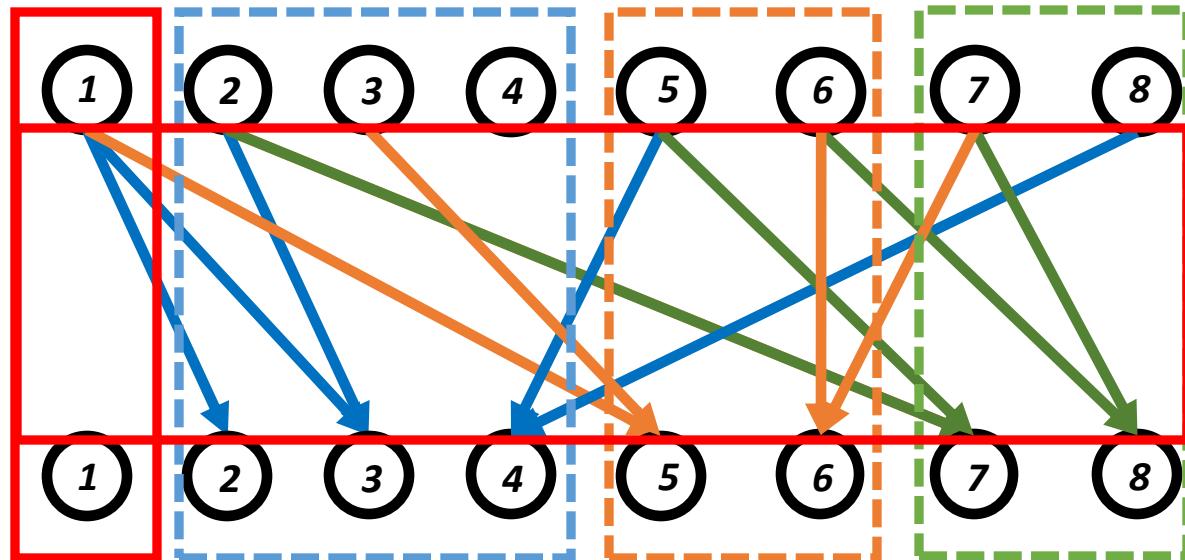
Satisfiability modulo theories (SMT)

- SMT solver = SAT solver + IDL theory solver

1. Node with indegree 0 comes before every other nodes.

$$2. \quad a \prec a' \implies v < v'$$

$$3. \quad (a = a') \wedge (u < u') \implies v \leq v'$$



```
#####
## Adding Wheeler graph constraints
#####
# SMT solver initialization
s = SolverFor('QF_IDL')

# 1st WG rule: nodes without incoming edge should be in front
for u in no_incoming_edge:
    s.add(u <= len(no_incoming_edge))

# 2nd WG rule: a < a' => v < v'
for us, (lb, ub) in group_info:
    for u in us:
        s.add(And(u > lb, u <= ub))
    s.add(Distinct(*us))

# 3rd WG rule: (a = a') ^ (u < u') => v <= v'
for i in range(len(edges)):
    # for j in range(len(edges)):
    for j in range(i+1, len(edges)):
        if i != j:
            ui, vi, wi = edges[i]
            uj, vj, wj = edges[j]
            if wi == wj:
                s.add(Implies(ui < uj, vi <= vj))
                s.add(Implies(ui > uj, vi >= vj))
```

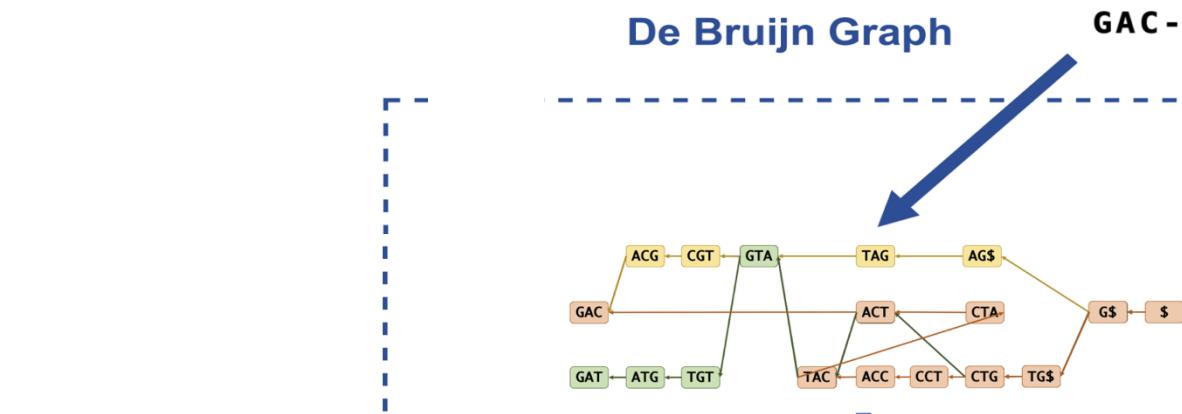
Wheeler Graph Recognition

WGT
Wheeler Graph Toolkit



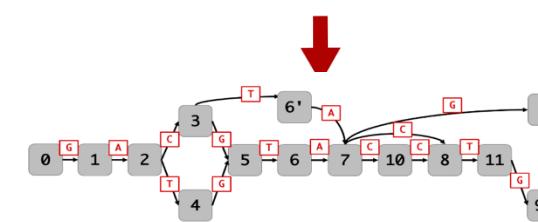
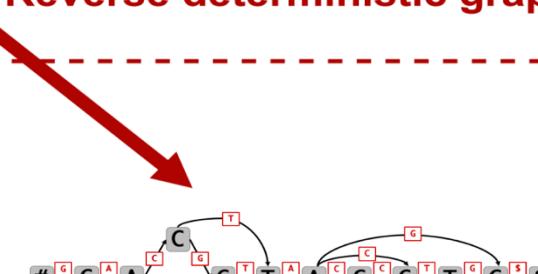
GCSA (Sirén et al., 2014)

Reverse deterministic graph



18 nodes, 21 edges
1 “distant gluing”

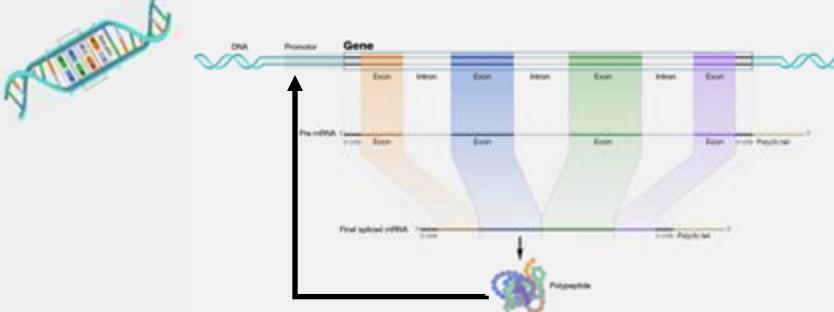
Original label	Wheeler order	Original label	Wheeler order
0	1	4	10
11	2	1	11
15	3	10	12
6	4	17	13
3	5	8	14
16	6	9	15
7	7	5	16
12	8	2	17
14	9	13	18



14 nodes, 16 edges
No “distant gluing”

Original label	Wheeler order	Original label	Wheeler order
0	1	9	8
2	2	5	9
7	3	9'	10
3	4	4	11
10	5	6	12
8	6	11	13
1	7	6	14

Part III



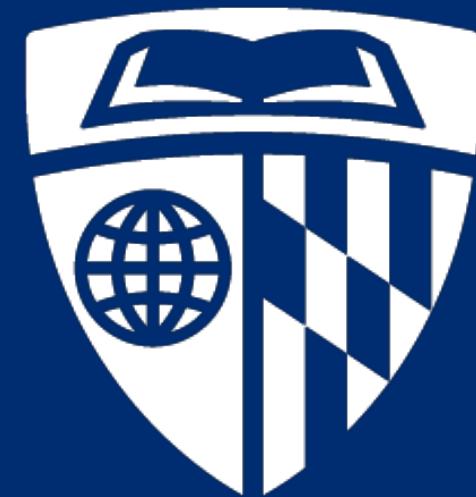
Steven Salzberg



Mihaela Pertea

Genome Annotation

- LiftOn: genome annotation lift-over
- Application: CHESS (GRCh38 -> CHM13; GRCh38 -> Han1)



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



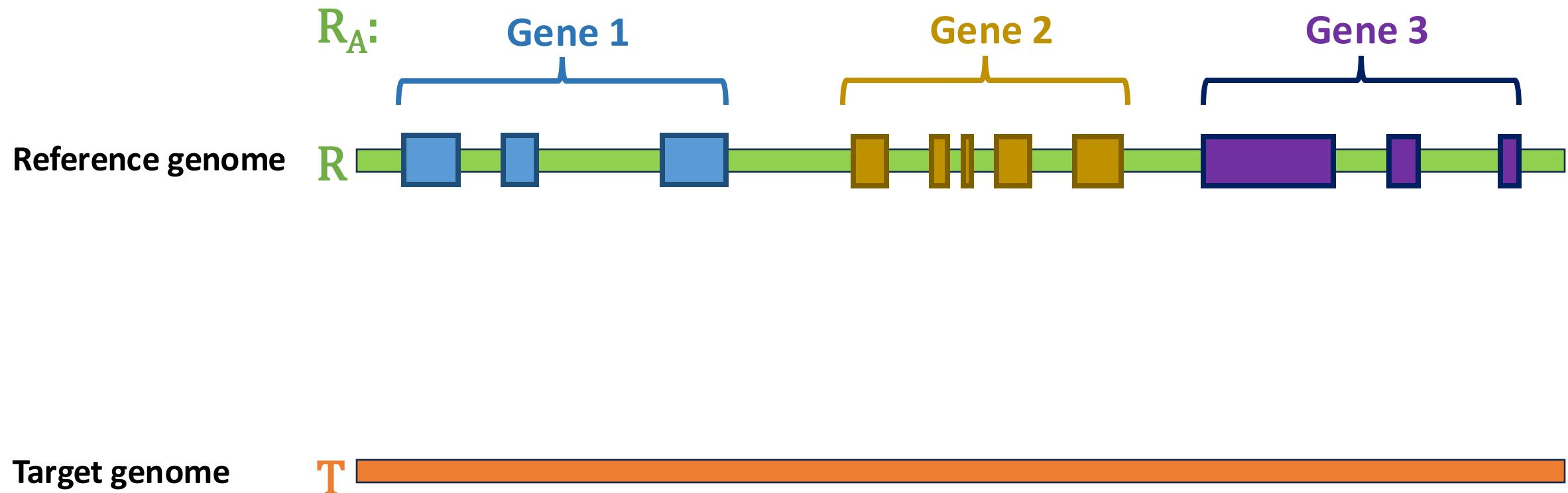
Genome (FASTA)

```
CAGCCCCCGGAGACTtaaatacaggaagaaaaaggCAGGACAGAATTACAAGGTGCTGGCCAGGGCGGGCAGCGGCCCT
GCCTCCTACCCTTGCCTCATGACCAGCTTGAAGAGATCCGACATCAAGTGCCACCTGGCTCGTGGCTCTCACT
GCAACGGAAAGCCACAGACTGGGGTGAAGAGAGTTCAAGTCACATGCGACCGGTgactccctgtccccaccccatgACACT
CCCCAGCCCTCCAAGGCCACTGTGTTCCAGTTAGCTCAGAGCCTCAGTCGATCCCTGACCCAGCACCAGGGCACTGATG
AGACAGCGGCTGTTGAGGagccacctcccagccacctcggggcccagggccagggtgtGCAGCACCACTGTACAATGGGG
AAACTGGCCCAGAGAGGTGAGGCAGCTTGCCTGGGTACAGAGCAAGGCAAAAGCAGCGCTGGGTACAAGCTAAAACC
ATAGTGCCCAGGGCACTGCCGCTGCAGGGCGAGGCATCGCATCACACCAGTGTCTGCCTCACAGCAGGCATCATCAGTA
```

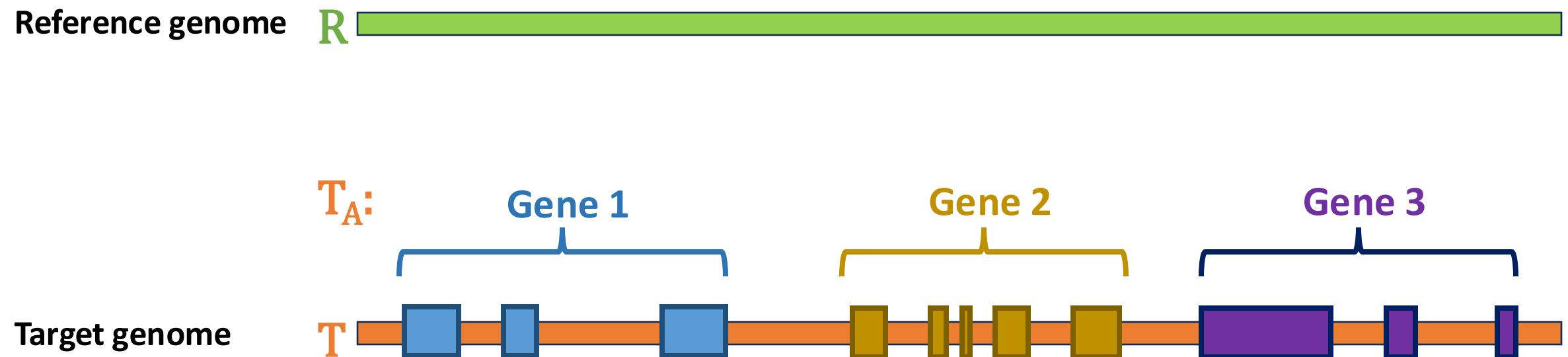
Annotation (GFF / GTF)

```
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;
```

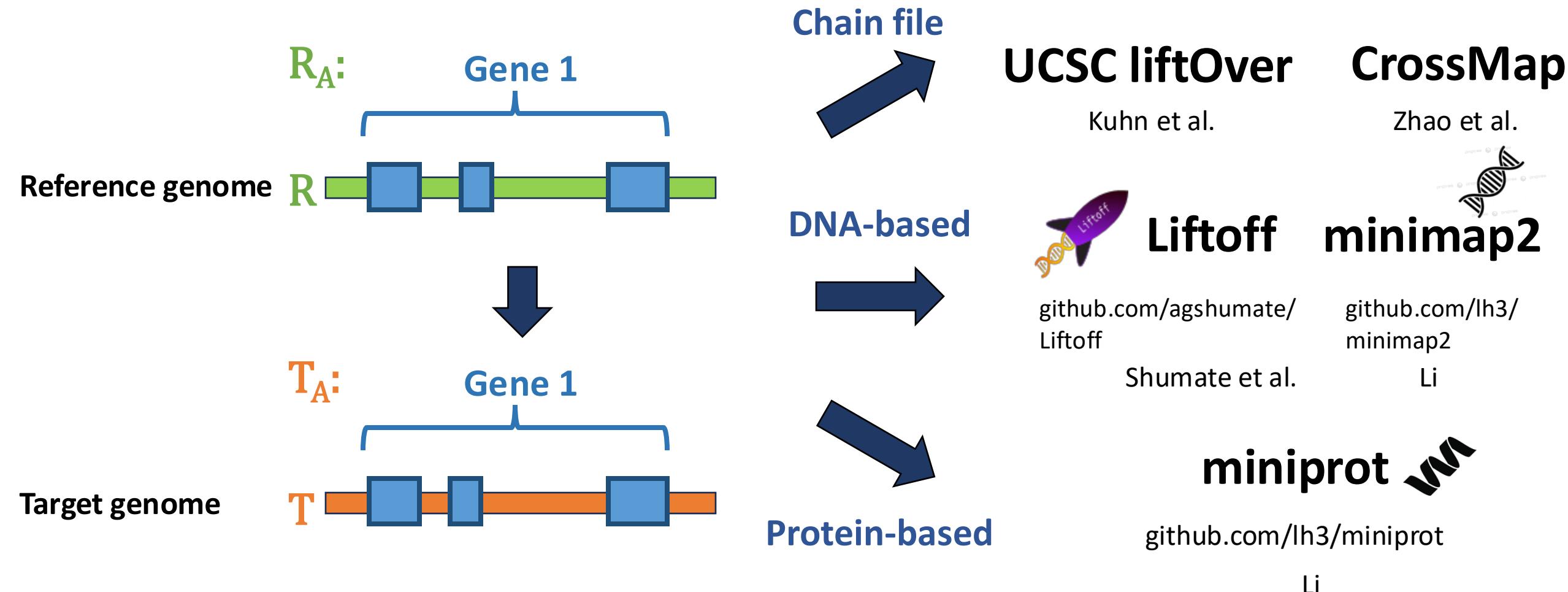
Lift-over Problem Definition:



Lift-over Problem Definition:



Lift-over problem, what methods are available?



Application: GRCh38 to T2T-CHM13 lift-over

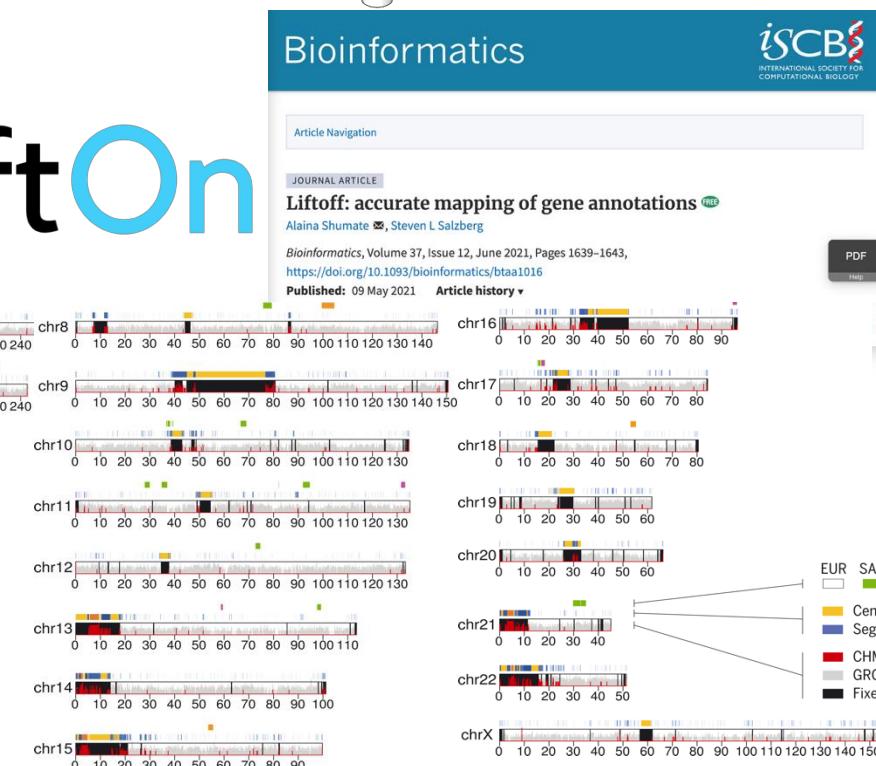
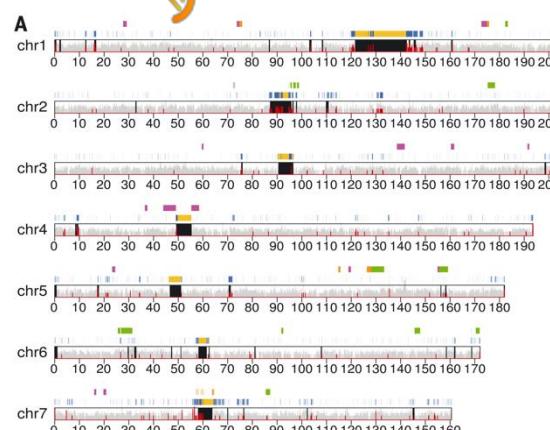
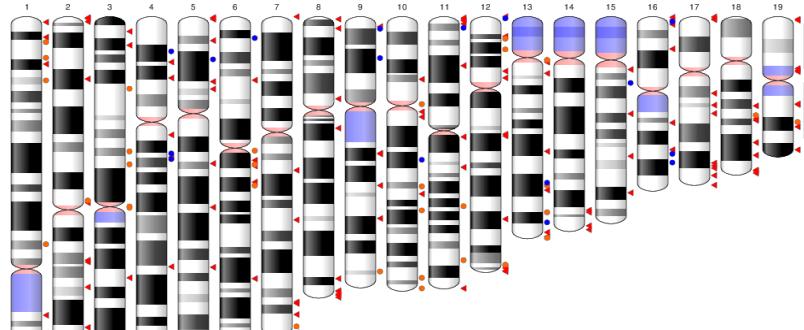
NCBI Build 34
NCBI Build 35
NCBI Build 36.1
GRCh37



GRCh38
(Dec 2013)



T2T-CHM13
(Jan 2022)



nature

Explore content ▾ About the journal ▾ Publish with us ▾

nature > article

Original Article | Published: 01 February 2001

Initial sequencing and analysis of the human genome

International Human Genome Sequencing Consortium

Science

HOME > SCIENCE > VOL. 291, NO. 5507 > THE SEQUENCE OF THE HUMAN GENOME

Special Reviews

The Sequence of the Human Genome

J. CRAIG VENTER, MARK D. ADAMS, EUGENE W. MYERS, PETER W. LI, [...] AND XIACHONG ZHU +269 authors Authors Info & Affiliations



Bioinformatics

Article Navigation

JOURNAL ARTICLE

Liftoff: accurate mapping of gene annotations

Alaina Shumate □, Steven L Salzberg

Bioinformatics, Volume 37, Issue 12, June 2021, Pages 1639–1643,

<https://doi.org/10.1093/bioinformatics/btaa1016>

Published: 09 May 2021 Article history ▾



PDF

Help



ISCB

INTERNATIONAL SOCIETY FOR COMPUTATIONAL BIOLOGY

bioinformatics.oxfordjournals.org

ISSN 1367-4813

© 2021 The Author(s)

Journal compilation © 2021 Association for Computing Machinery. All rights reserved.

DOI: 10.1093/bioinformatics/btaa1016

Published online first 09 May 2021

Received 12 January 2021; revised 12 April 2021; accepted 12 April 2021

Editorial handling: Michael Schatz

Reviewing editor: Michael Schatz

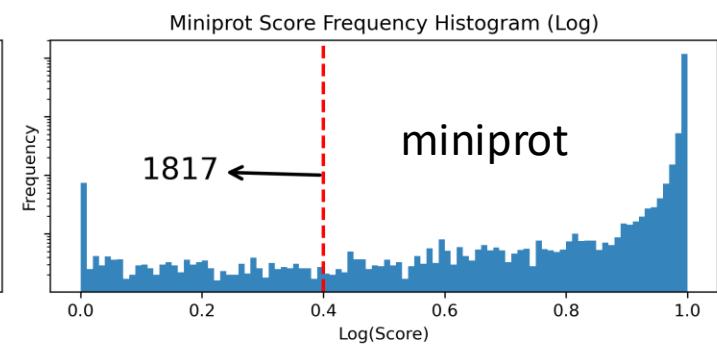
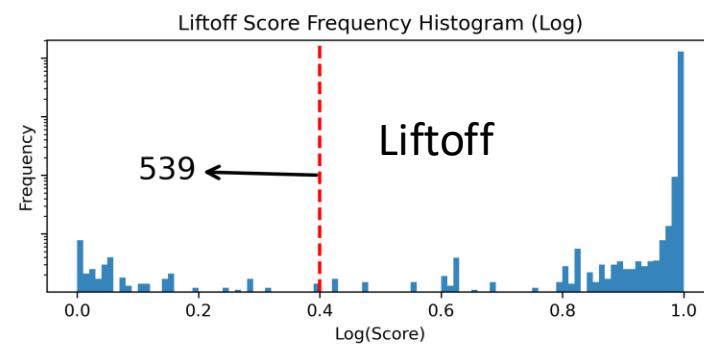
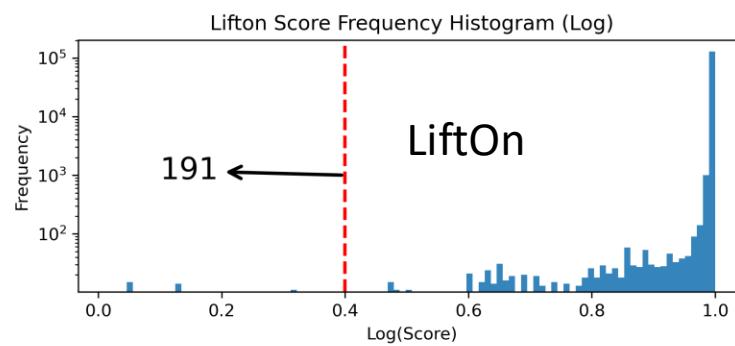
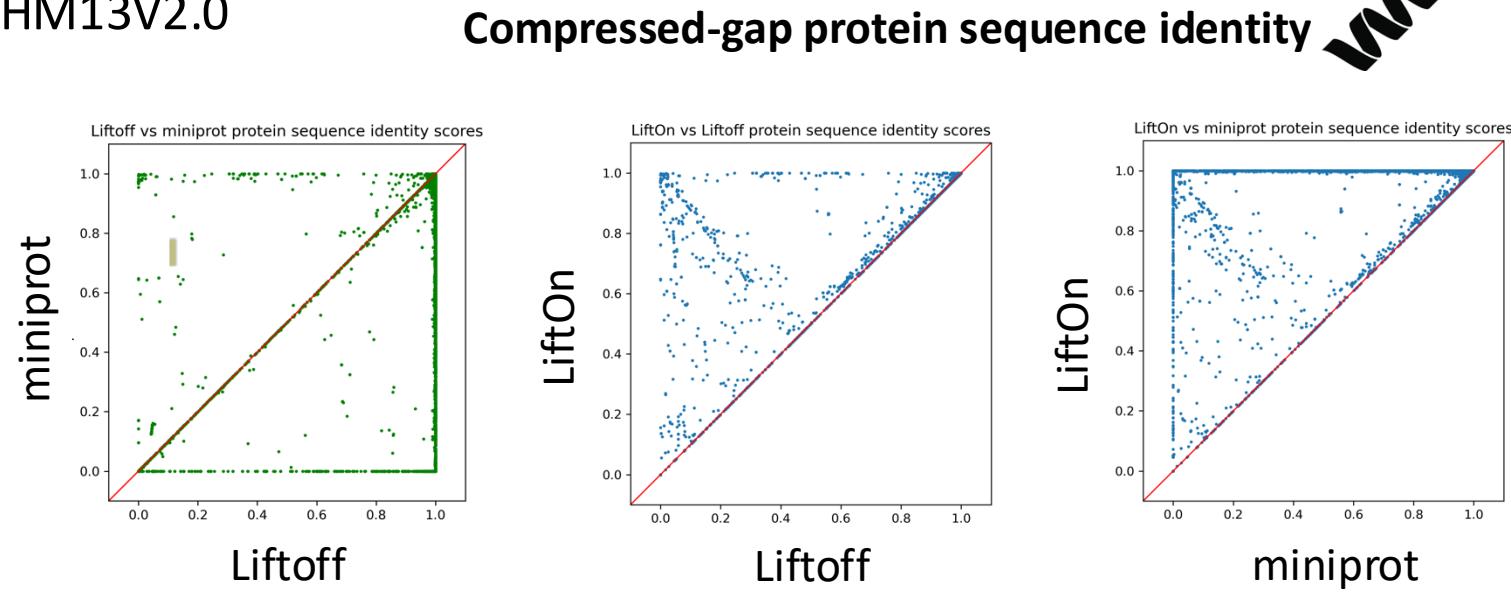
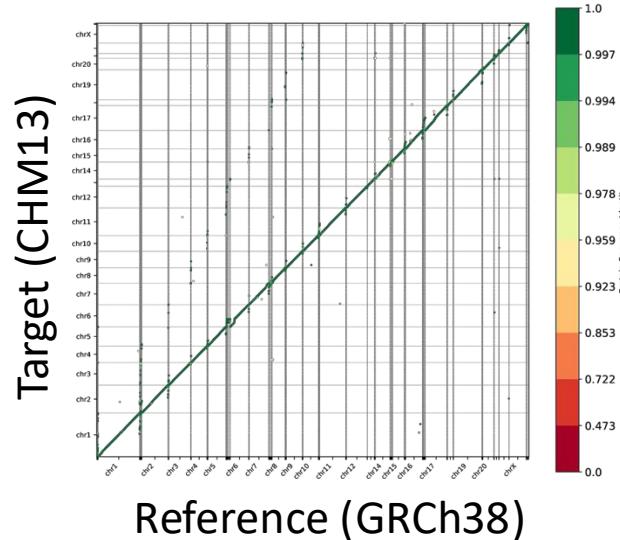
Associate editor: Michael Schatz

Editorial office: Michael Schatz

Result 1: improves DNA & protein-based lift-over



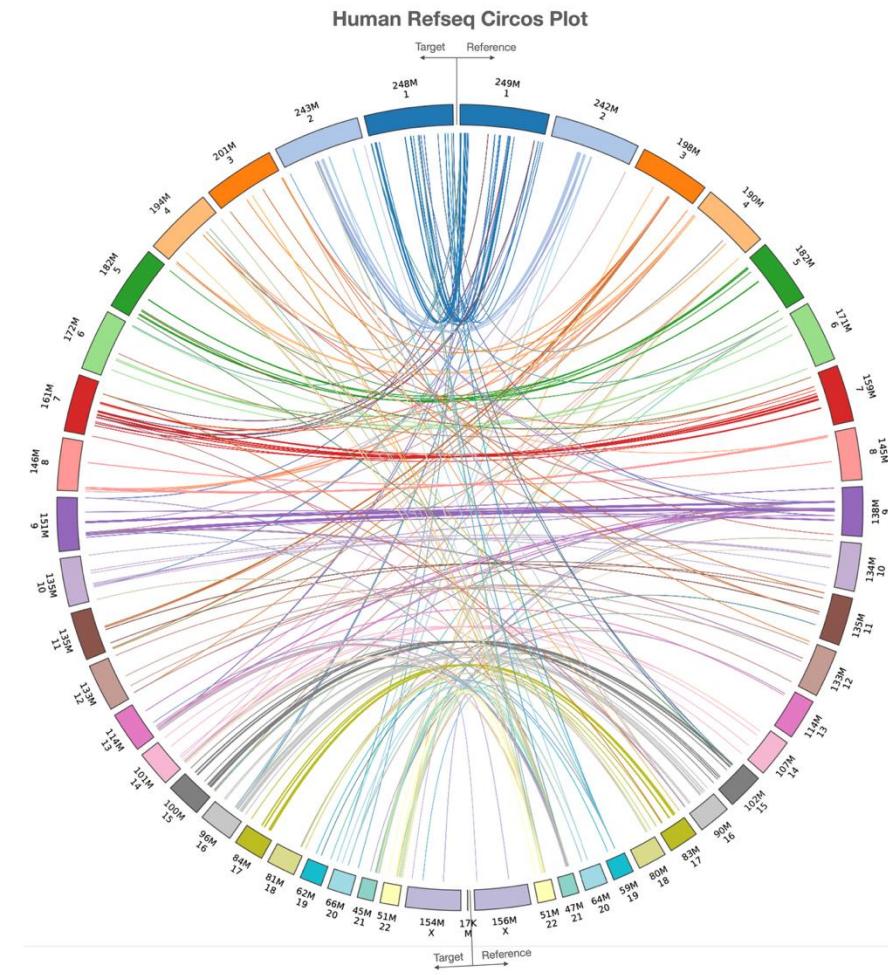
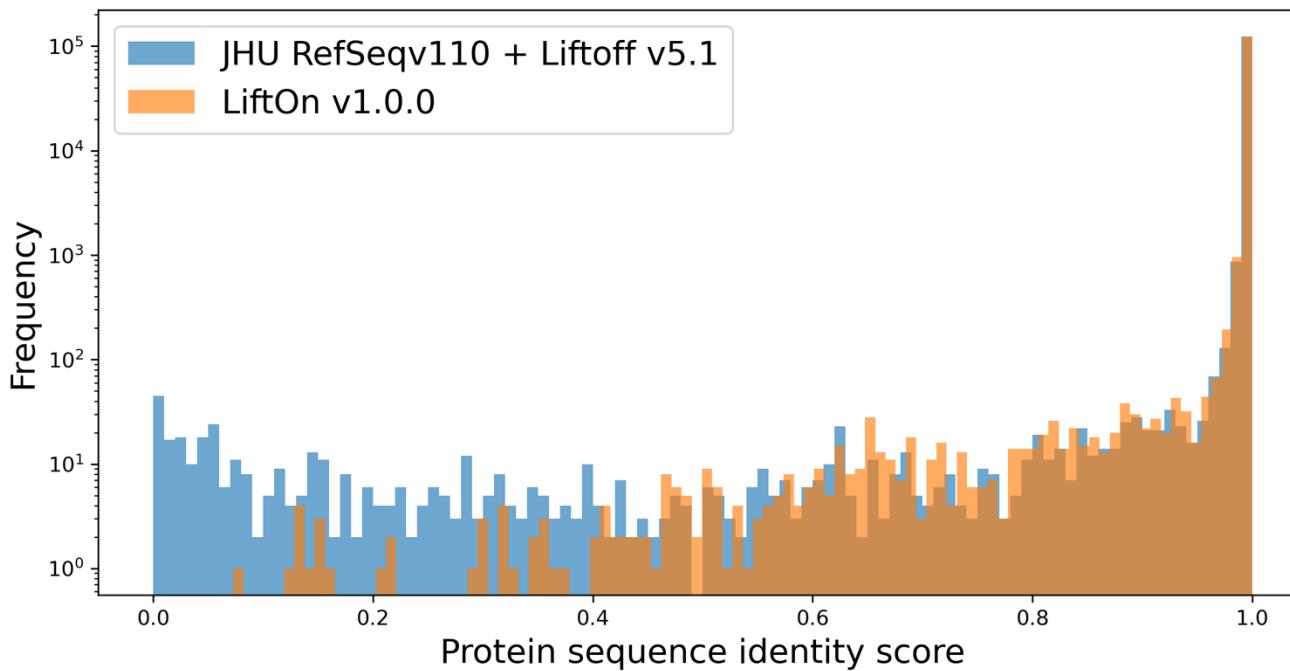
Map RefSeq v220 from GRCh38 -> CHM13V2.0



Result 2: improve CHM13 protein annotations

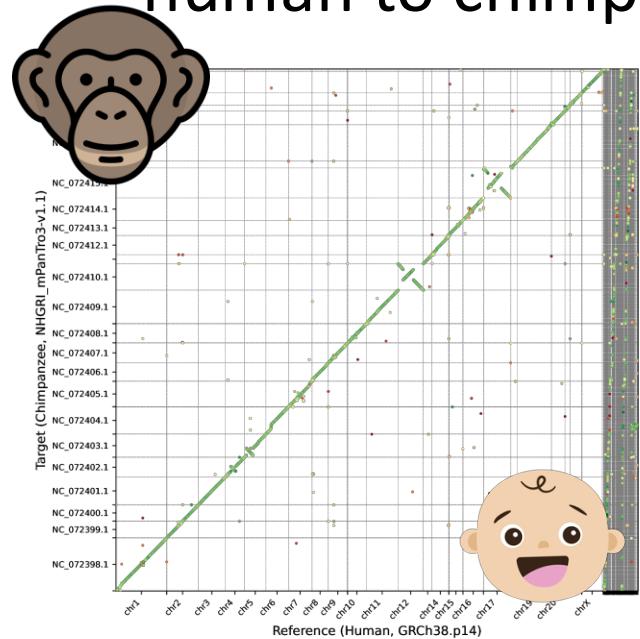


Protein sequence identity score frequency histogram

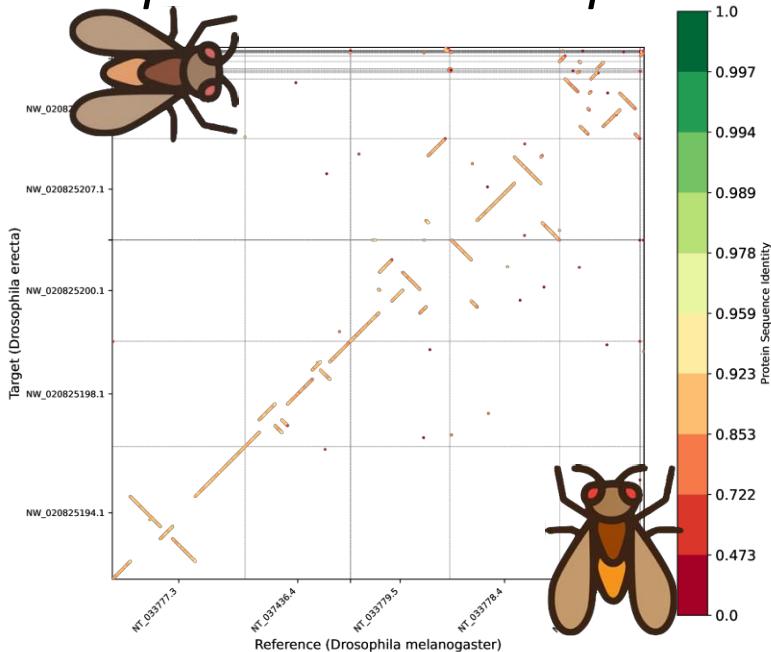


Result 3: improve distant species lift-over

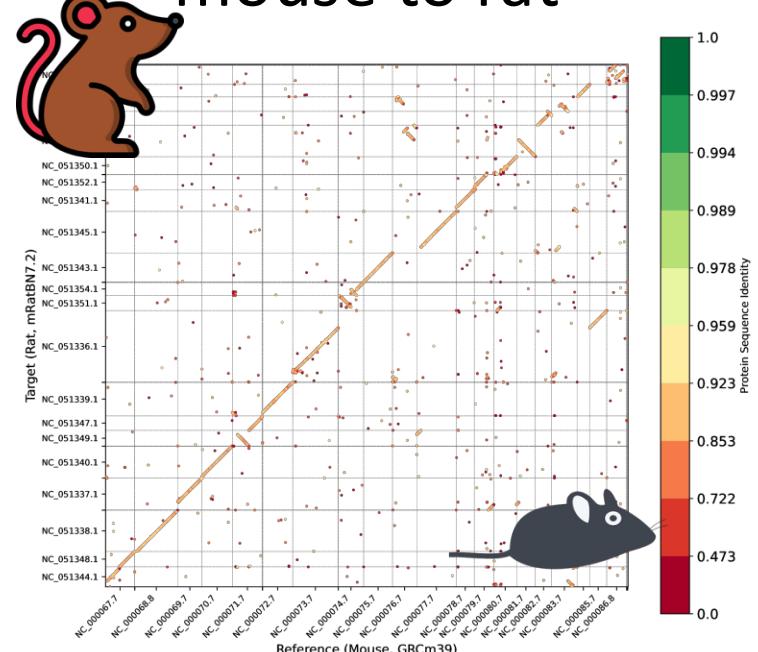
human to chimp



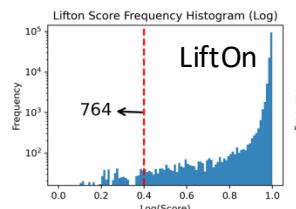
Drosophila m. to *Drosophila e.*



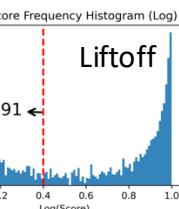
mouse to rat



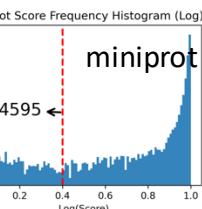
DNA + Protein



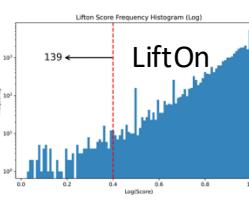
DNA-only



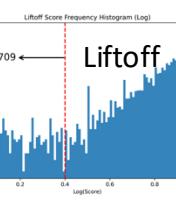
Protein-only



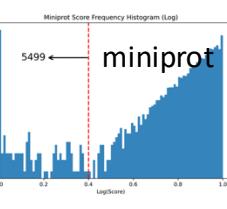
DNA + Protein



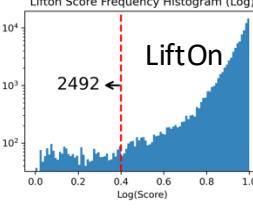
DNA-only



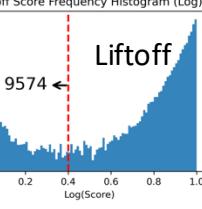
Protein-only



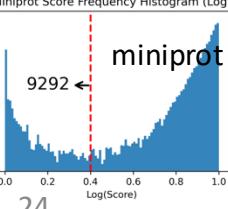
DNA + Protein



DNA-only



Protein-only



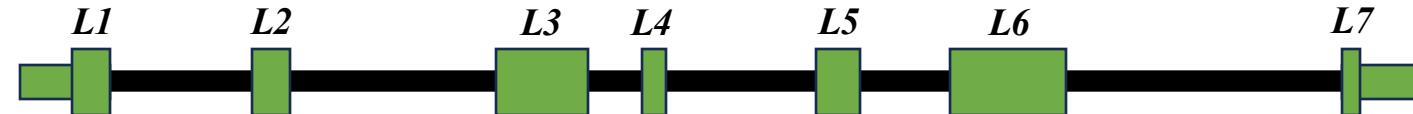
LiftOn: Protein-maximization algorithm

A

Target genome +
Expected annotation



1. Liftoff annotation

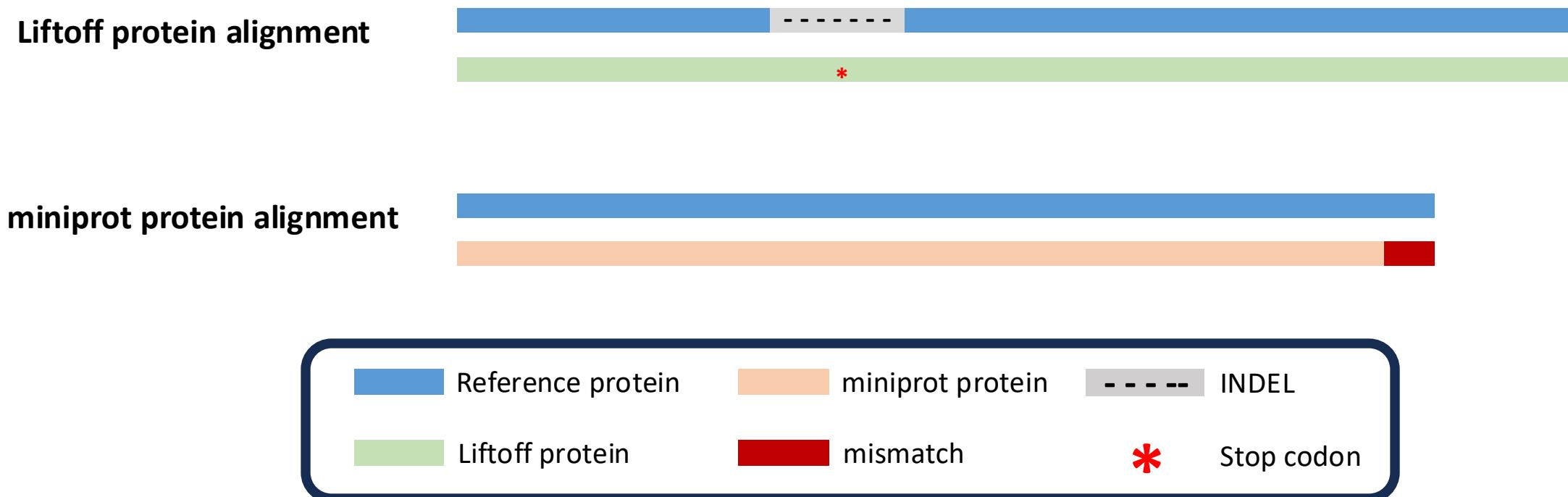


2. miniprot annotation



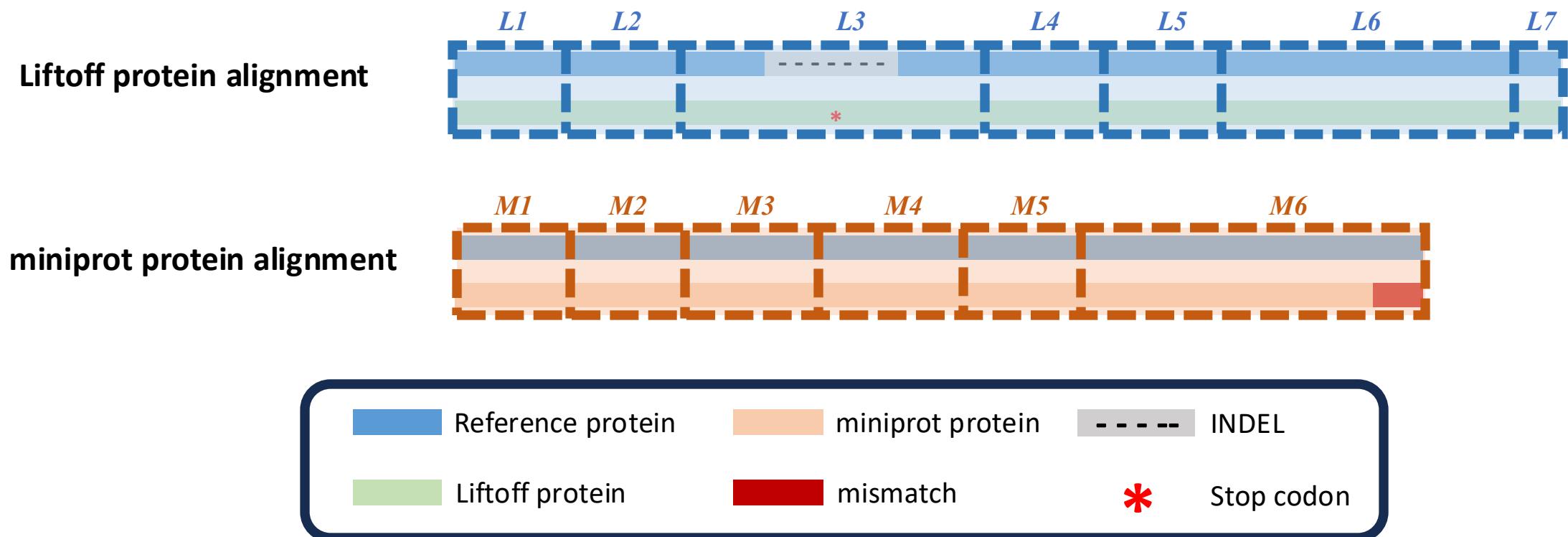
LiftOn: Protein-maximization algorithm

B Step 1: Align Liftoff & miniprot proteins to reference protein



LiftOn: Protein-maximization algorithm

C Step 2: Mapped CDS boundaries onto Liftoff & miniprot protein alignments

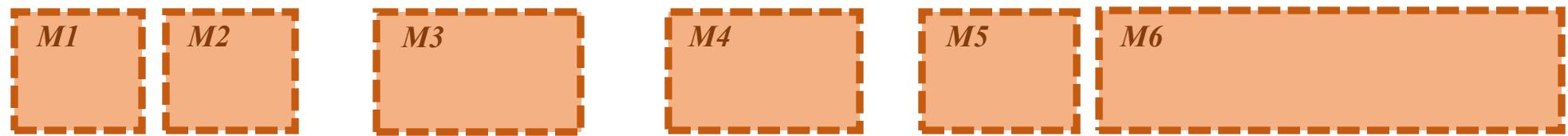


LiftOn: Protein-maximization algorithm

D Step 3: group CDSs by “accumulated AA in the reference protein”

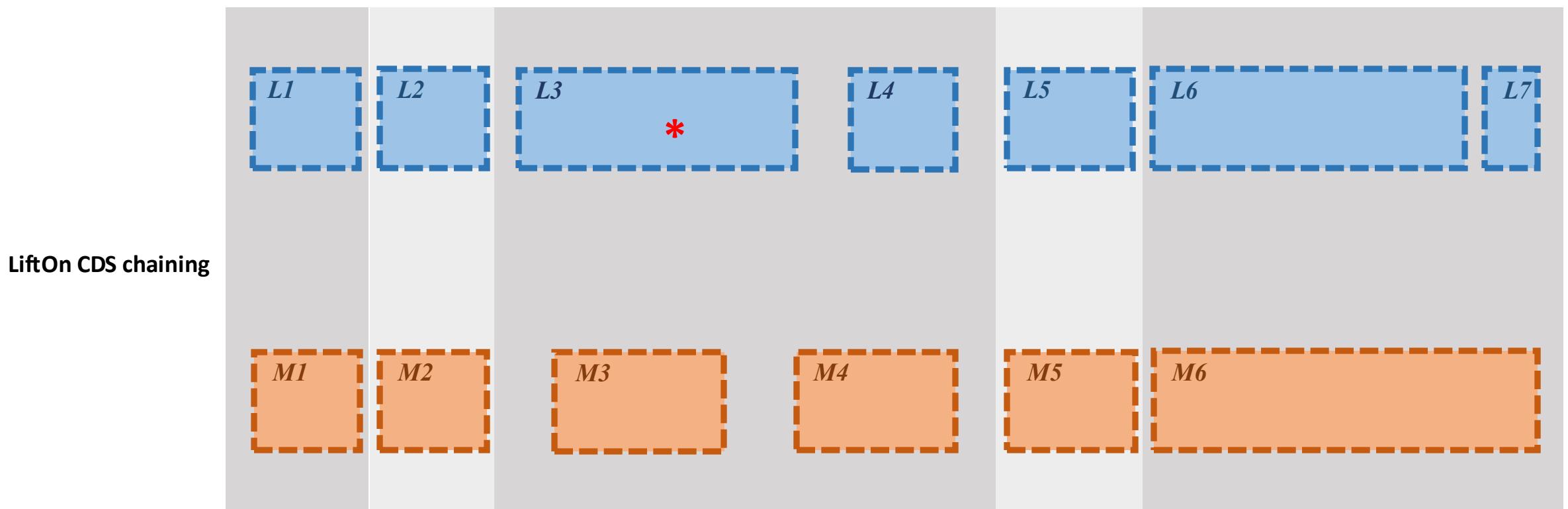


LiftOn CDS chaining



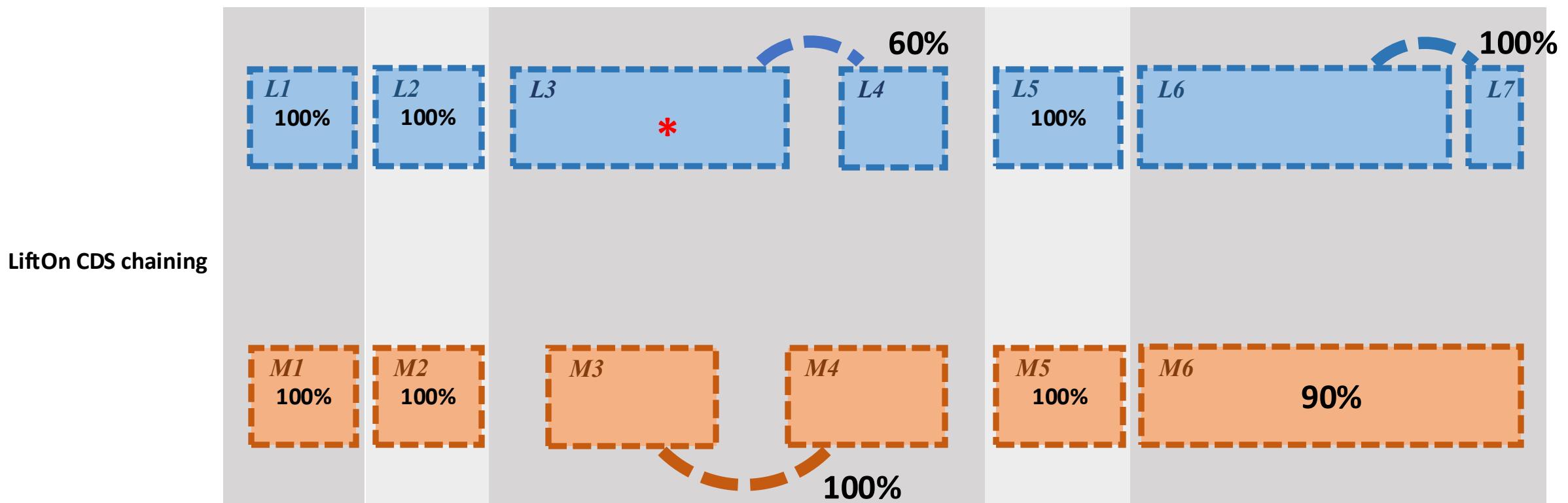
LiftOn: Protein-maximization algorithm

D Step 3: group CDSs by “accumulated AA in the reference protein”



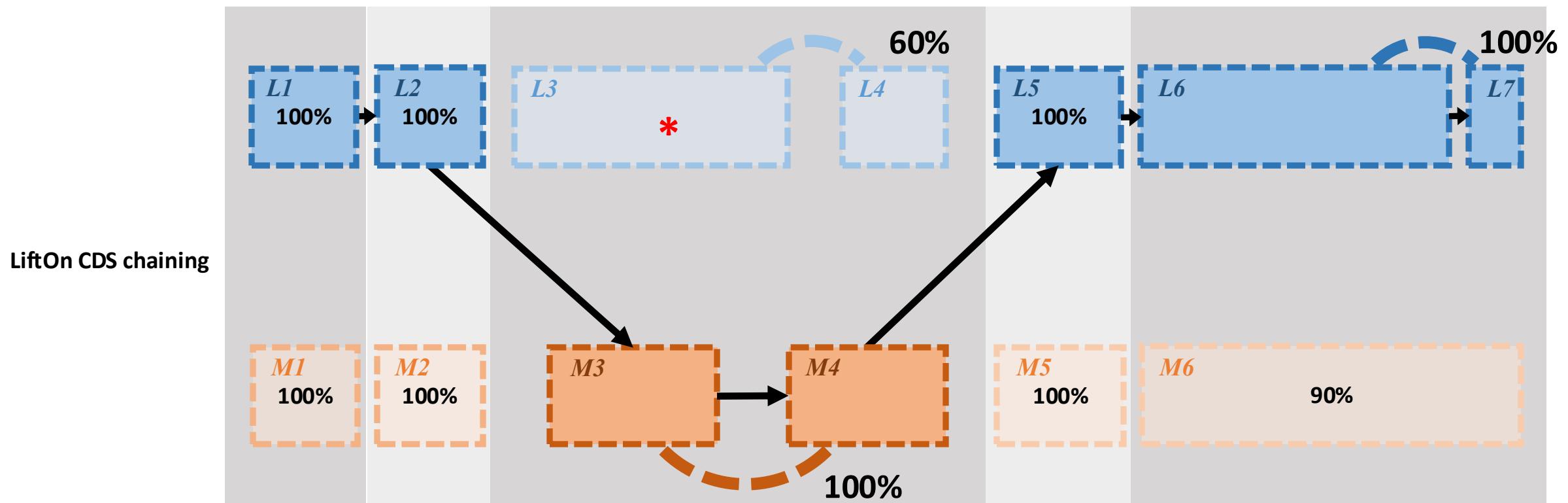
LiftOn: Protein-maximization algorithm

D Step 3: group CDSs by “accumulated AA in the reference protein”

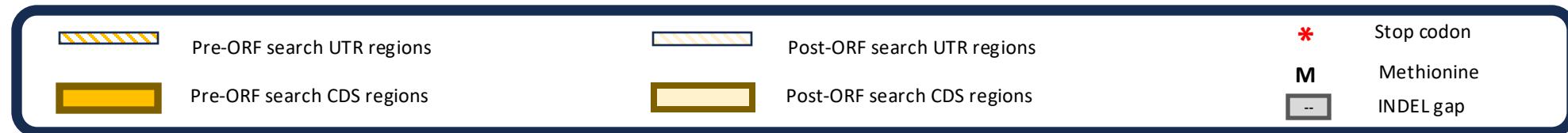


LiftOn: Protein-maximization algorithm

D Step 3: group CDSs by “accumulated AA in the reference protein”

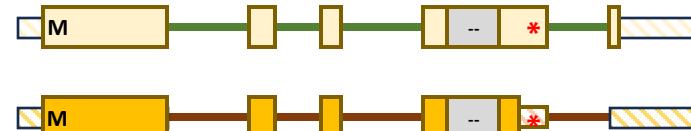


LiftOn: Protein-maximization algorithm



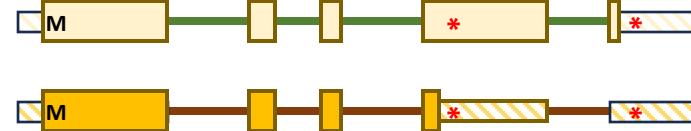
A

Frameshift



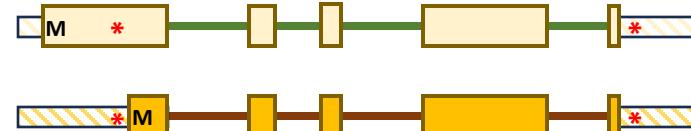
B

Stop codon gain:
Early
translation stop



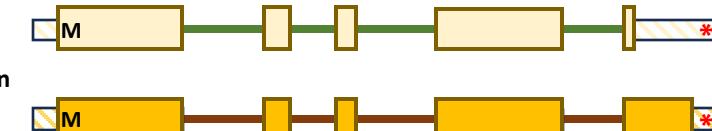
C

Stop codon gain:
Switching
translation start



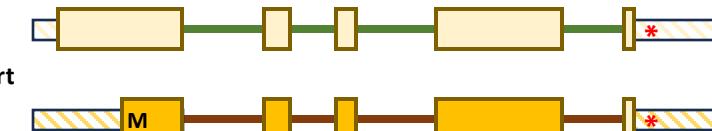
D

Stop codon lost:
Protein extension



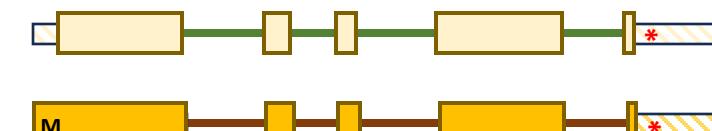
E

Start codon lost:
Downstream start



F

Start codon lost:
Upstream start



LiftOn: more lift-over applications



House mouse
(*Mus musculus*)



Yeast
(*Saccharomyces cerevisiae*)



Thale cress
(*Arabidopsis thaliana*)



Honey bee
(*Apis mellifera*)



fruit fly
(*Drosophila melanogaster*)

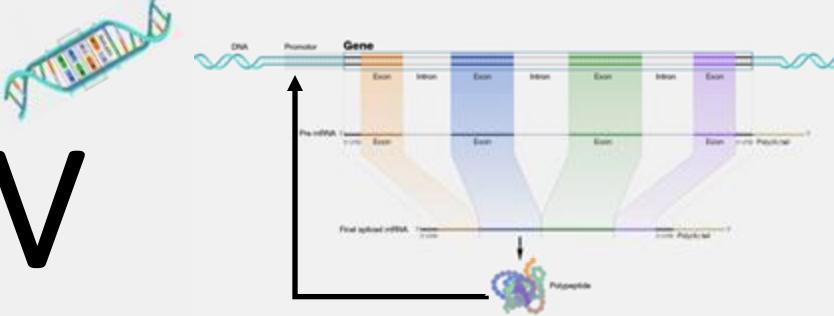


Rice
(*Oryza sativa*)

GENOME
RESEARCH

LiftOn

I Part IV & V



Mihaela Pertea



Steven Salzberg



Anqi Liu

Splice Site Prediction

- OpenSpliceAI
- Splam

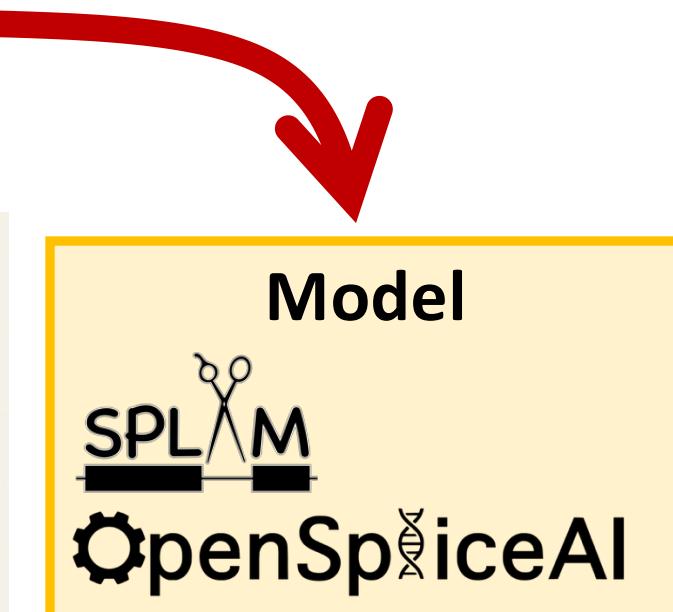
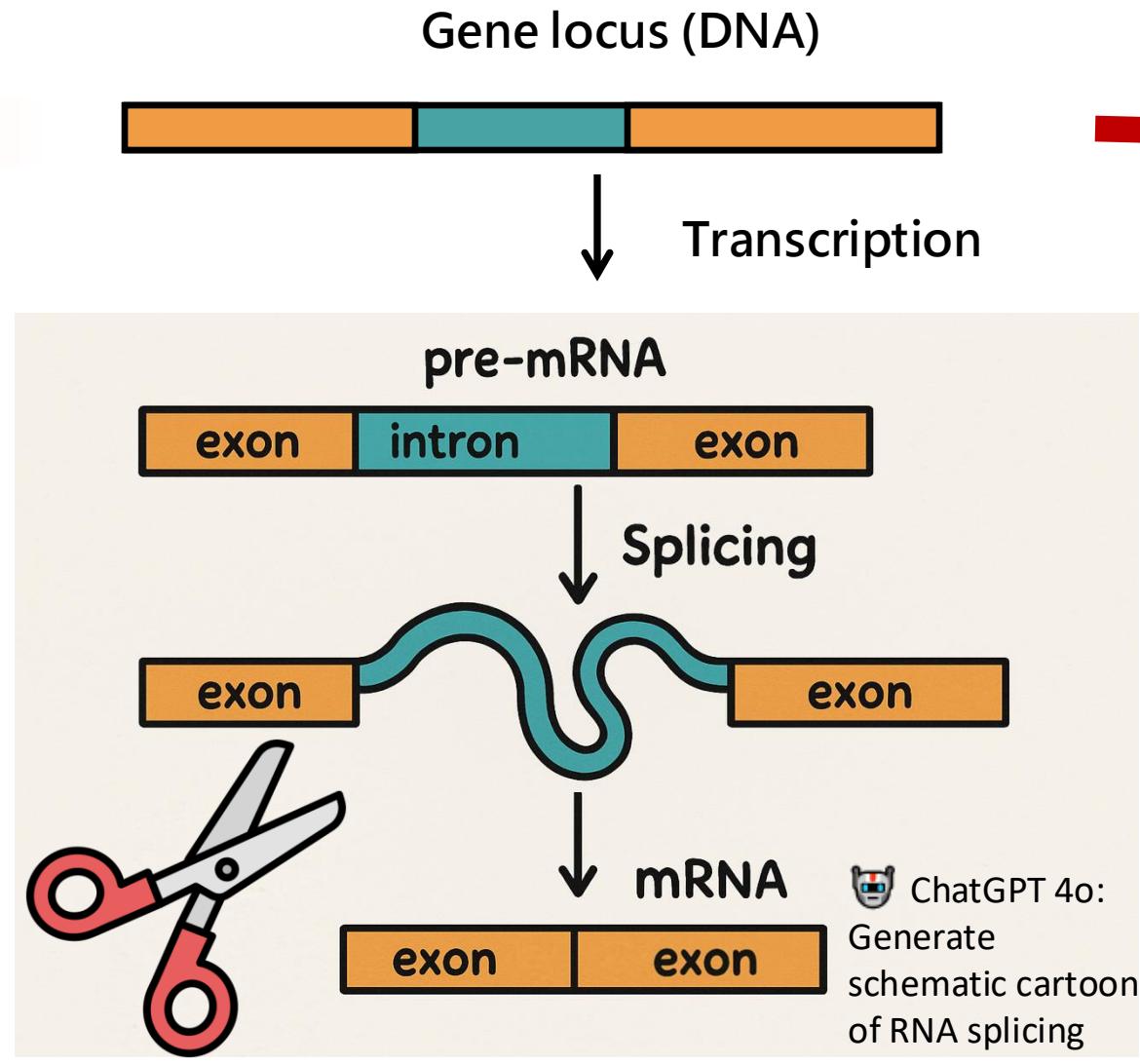
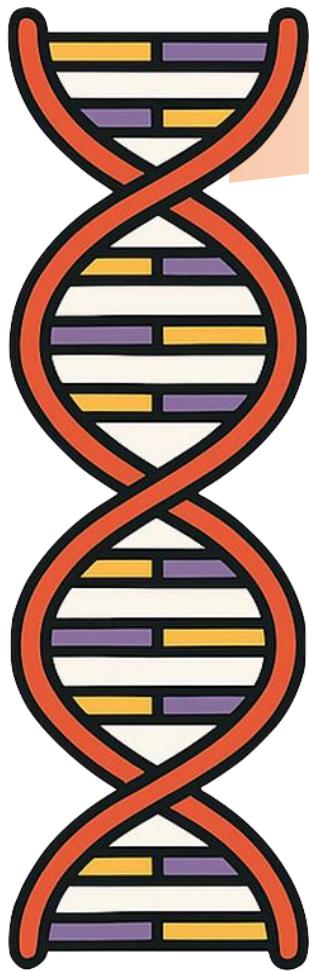


Chao, K. H., Mao, A., Salzberg, S. L., & Pertea, M. (2024). Splam: a deep-learning-based splice site predictor that improves spliced alignments. *Genome biology*, 25(1), 243.



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. *eLife*, 2025-03.

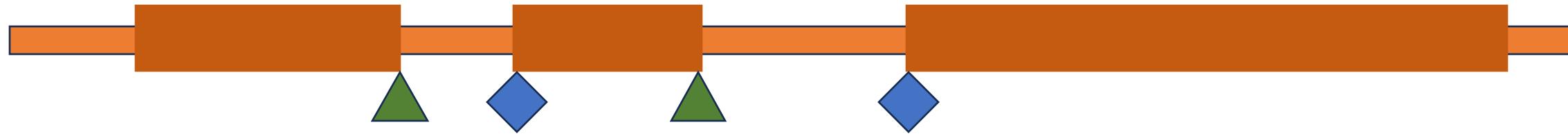
Splice Site Prediction



Where are the splice sites?

OpenSpliceAI: Splice site prediction

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. *eLife*



X

AGACTCAGCCCCCGGAGACTTAGTTAGAGGAAGAAAAGGTAGGACAGAAGAAAAAGGCAGGACATACAAGGTGCTGGCCCAGGGCGG

Y



Donor: 2

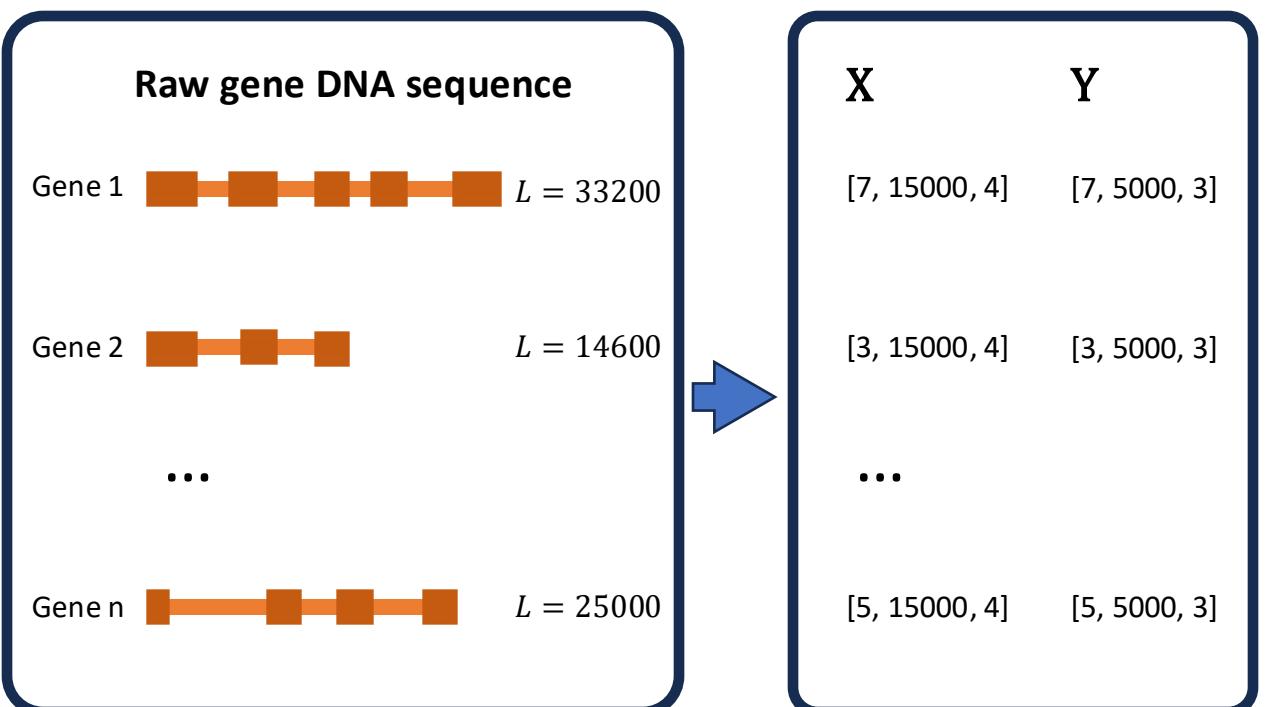


Acceptor: 1

Neither: 0

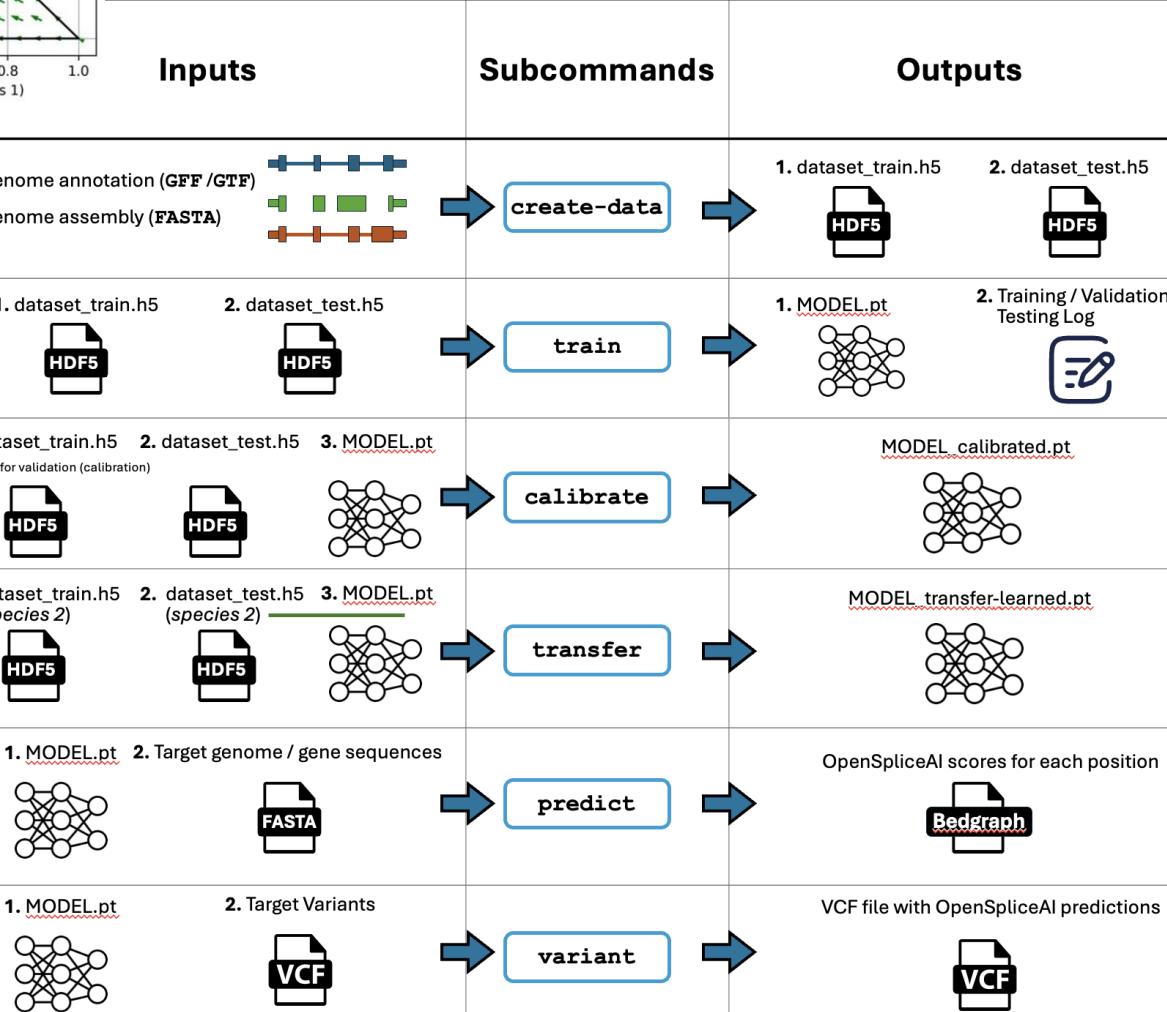
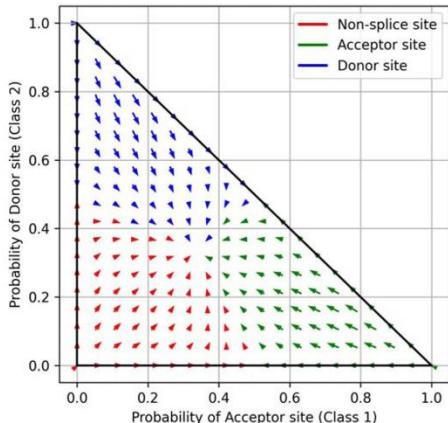
openSpliceAI

~20k protein-coding genes



$$W = 5000$$

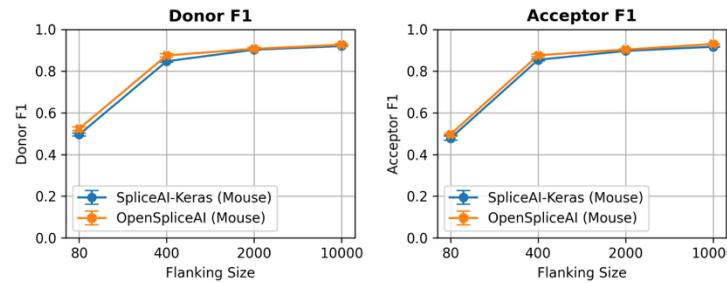
$$F = 10,000$$



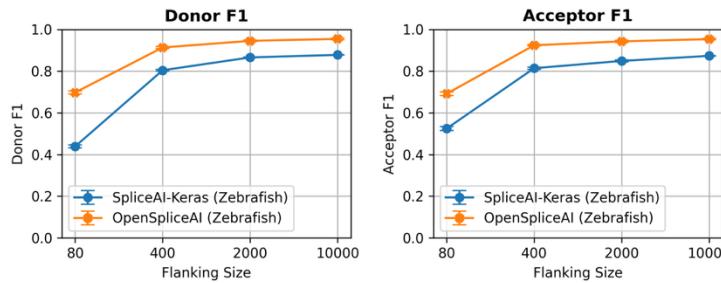
openSpliceAI: retrain on different species



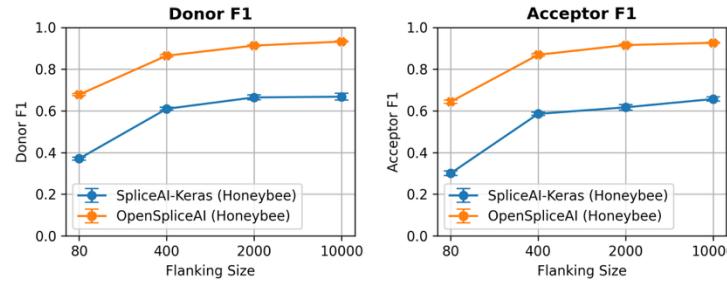
Splice site prediction metrics for Mouse



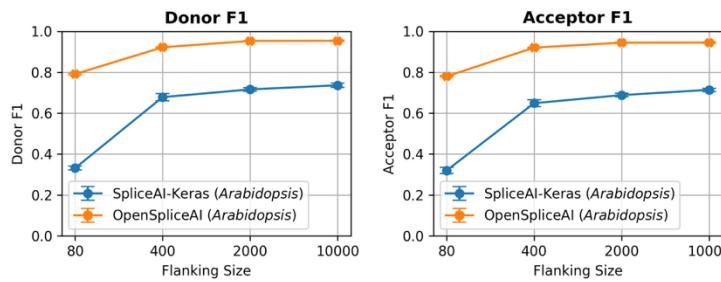
Splice site prediction metrics for Zebrafish



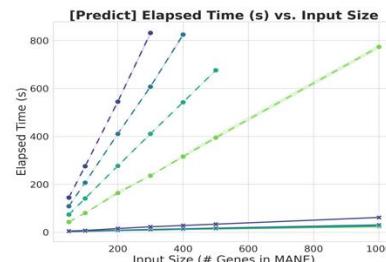
Splice site prediction metrics for Honeybee



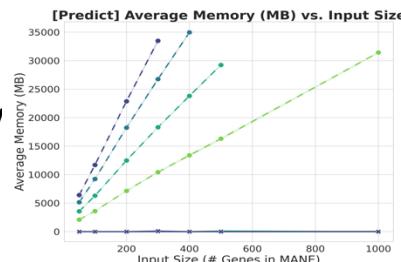
Splice site prediction metrics for Arabidopsis



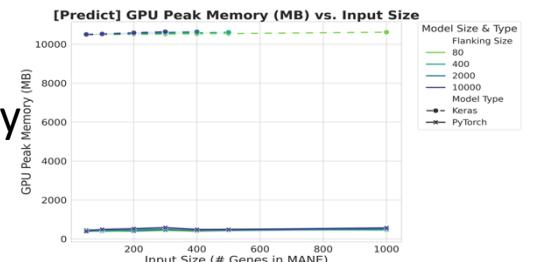
System time



Memory

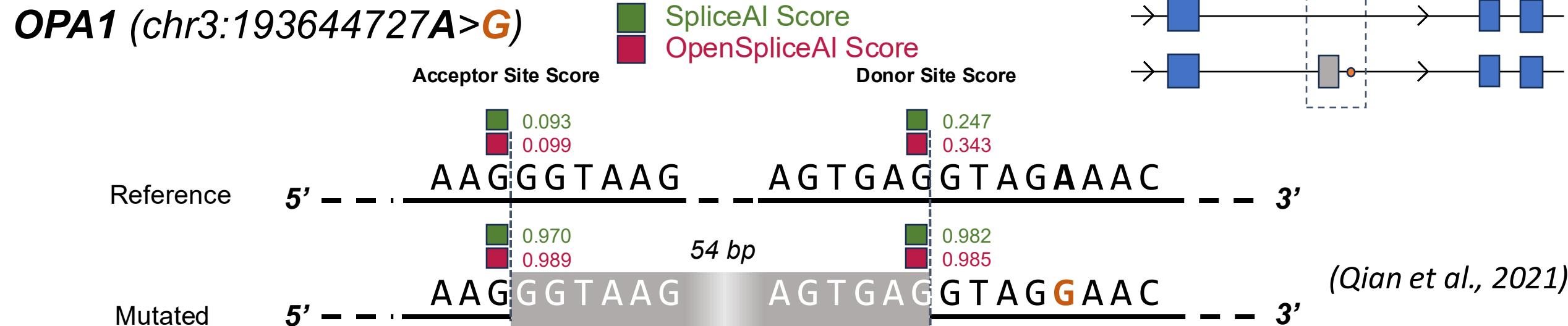
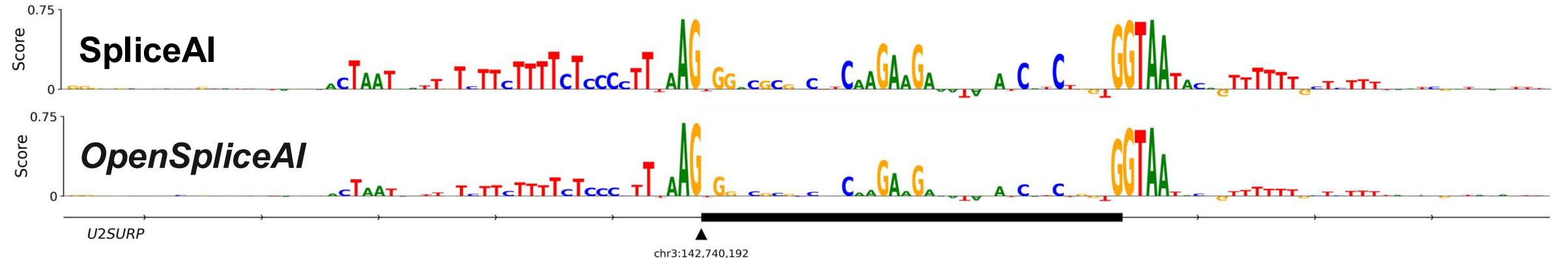


GPU Avg memory





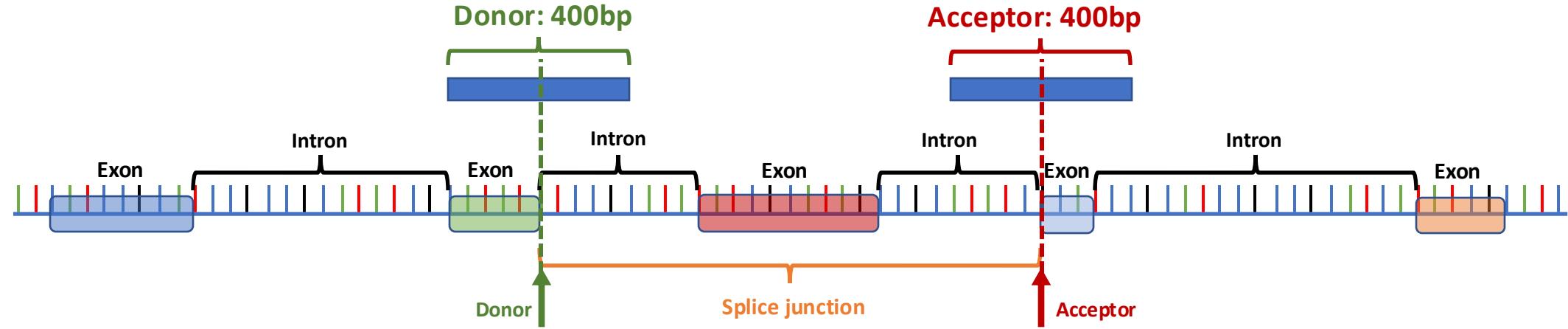
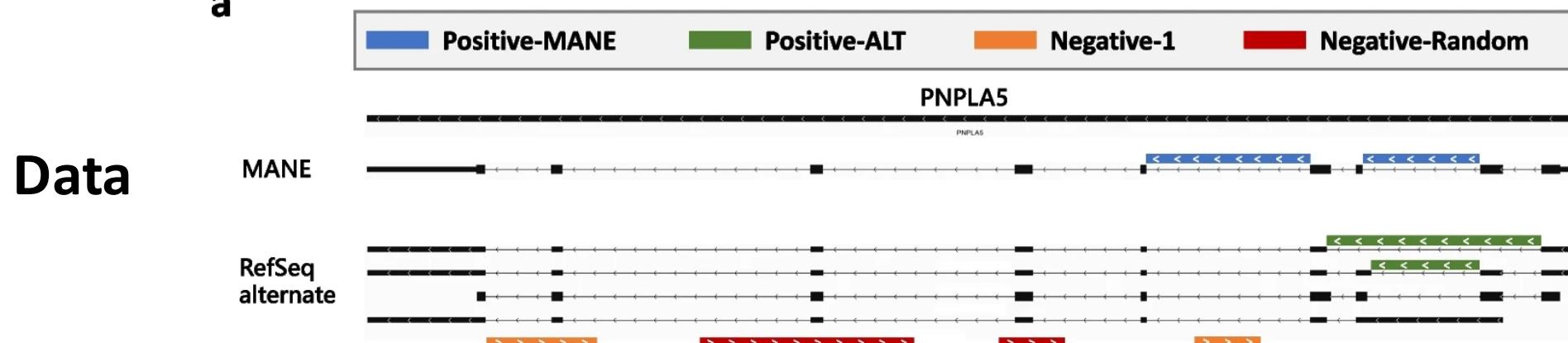
openSpliceAI: ISM & variant prediction



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.

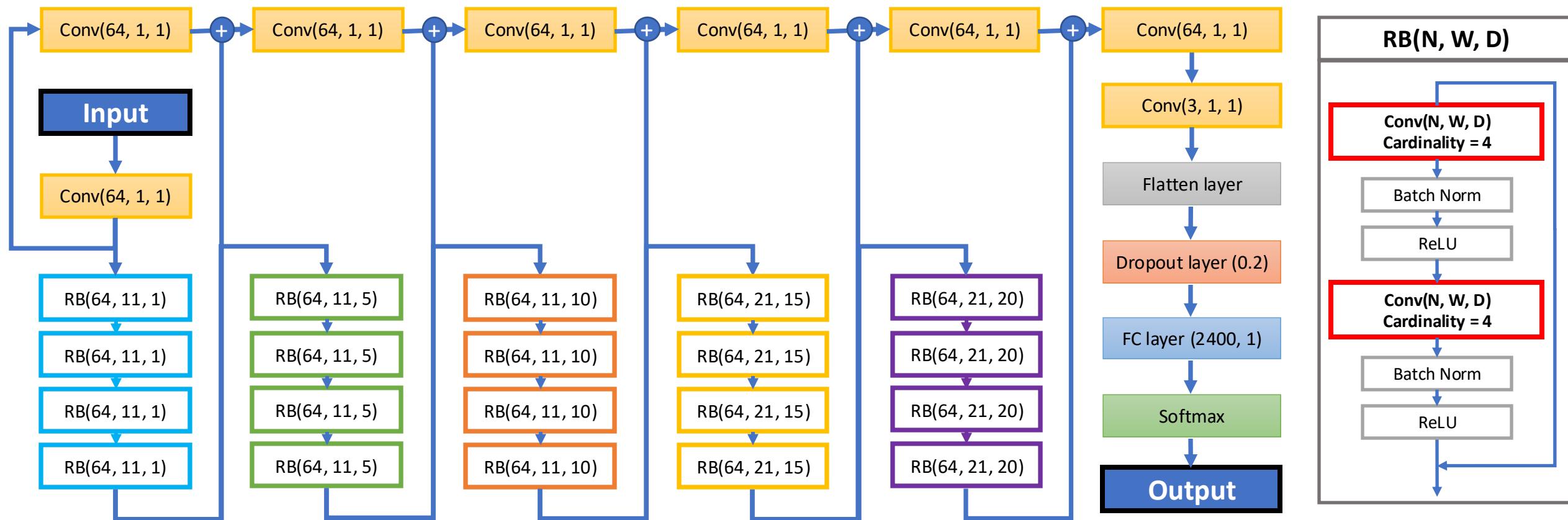


: predicting splice sites in a pair

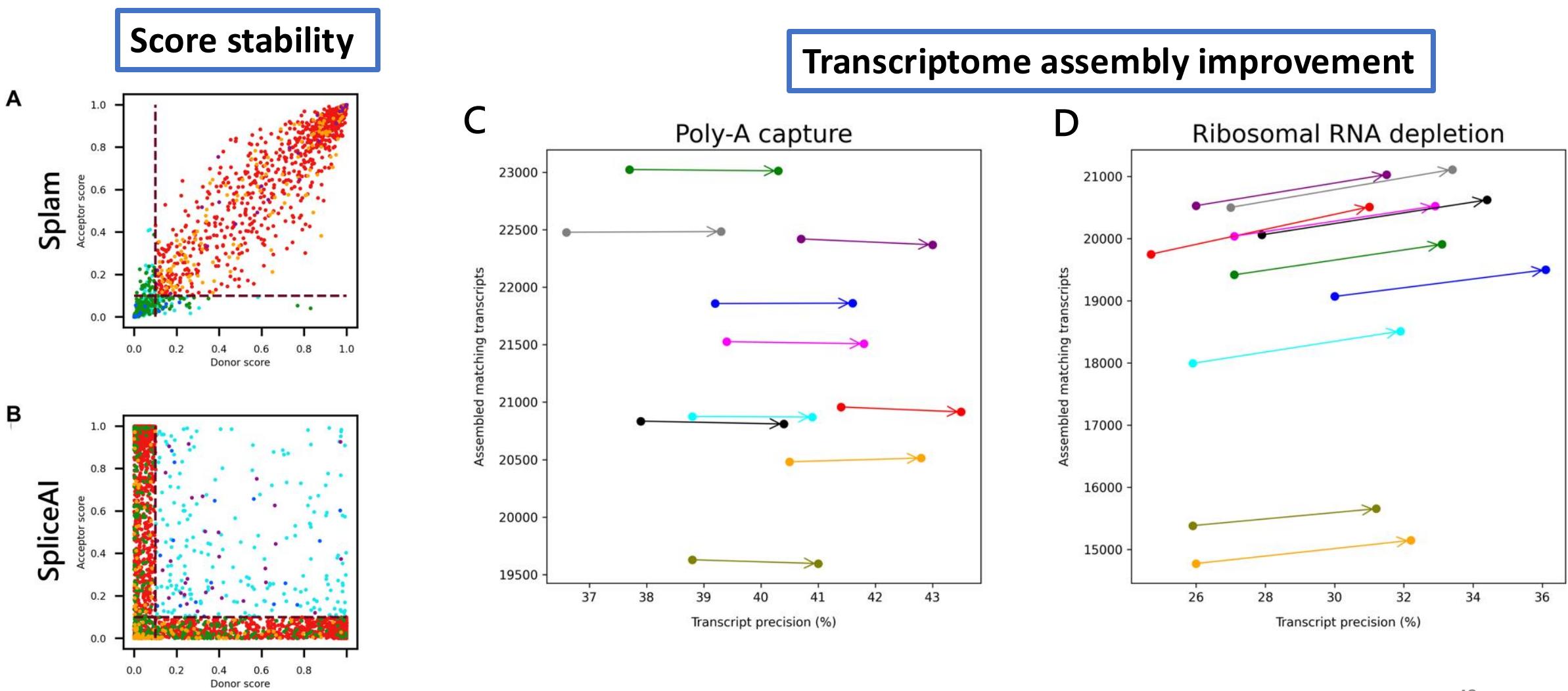
**a**



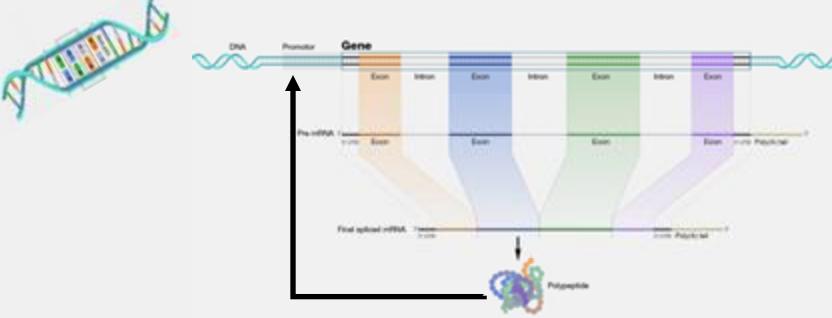
SPLAM : Splam Model Architecture



SPLAM: improving transcriptome assembly



Part VI



David Kelley



Johannes Linder



Majed Mohamed Magzoub

Yeast RNA-Seq Prediction

- Fungal language model (LM)
- **Shorkie**: models to predict yeast RNA-Seq coverages
- Model interpretability



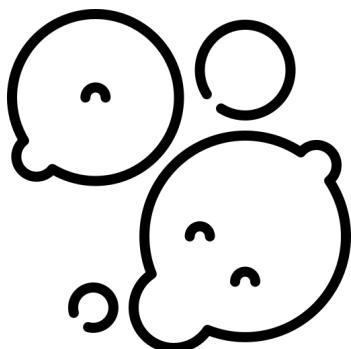
Chao, K. H., Magzoub, M., Stoops, E., Hackett, S., Linder J., * and Kelley, D. R.,
(manuscript in preparation). Predicting dynamic expression patterns in
budding yeast with a fun-gal DNA language model

Manuscript in preparation.
Expected to be on bioRxiv soon

Questions we answer in this study:

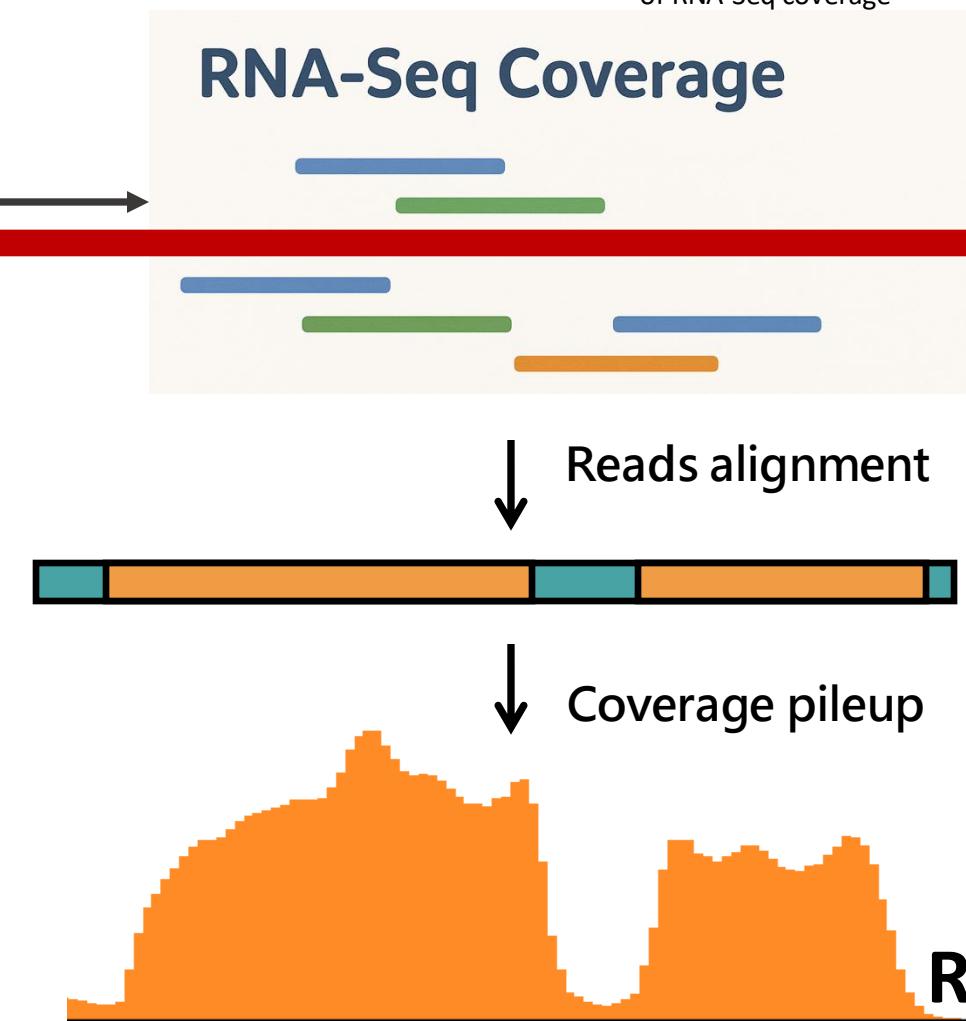
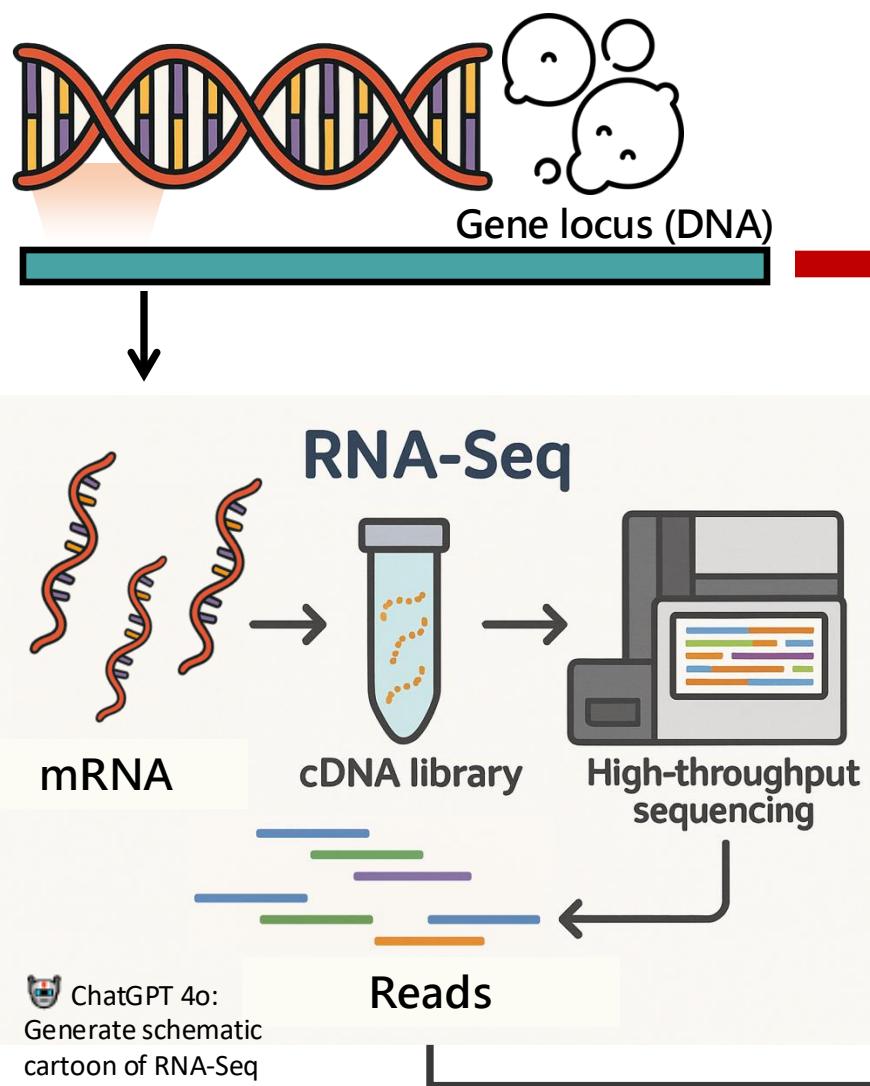
Main Goal

- Predict yeast gene expression (RNA-Seq coverage) directly from DNA sequences.

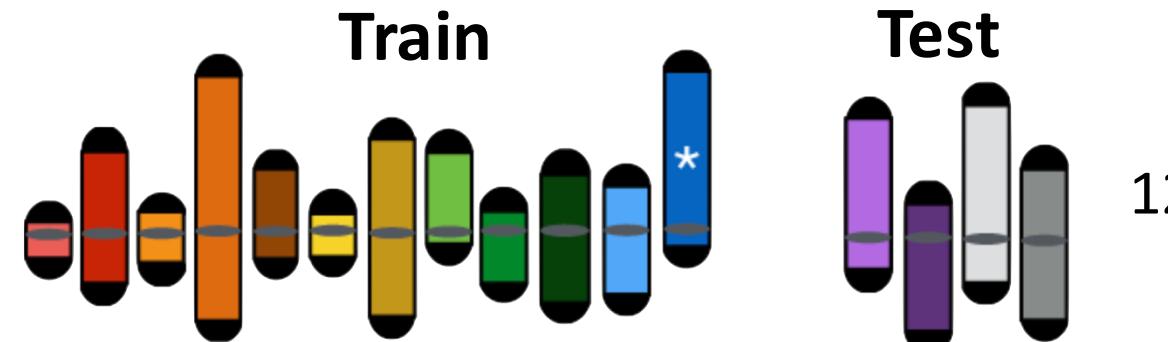
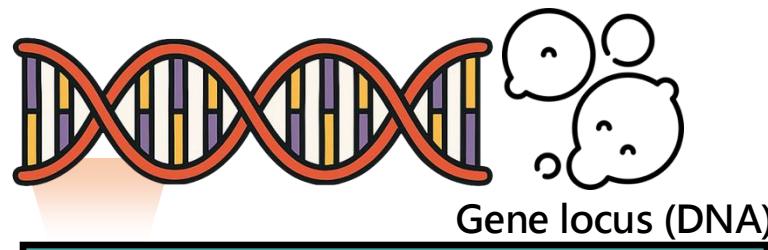


I Yeast RNA-Seq Prediction

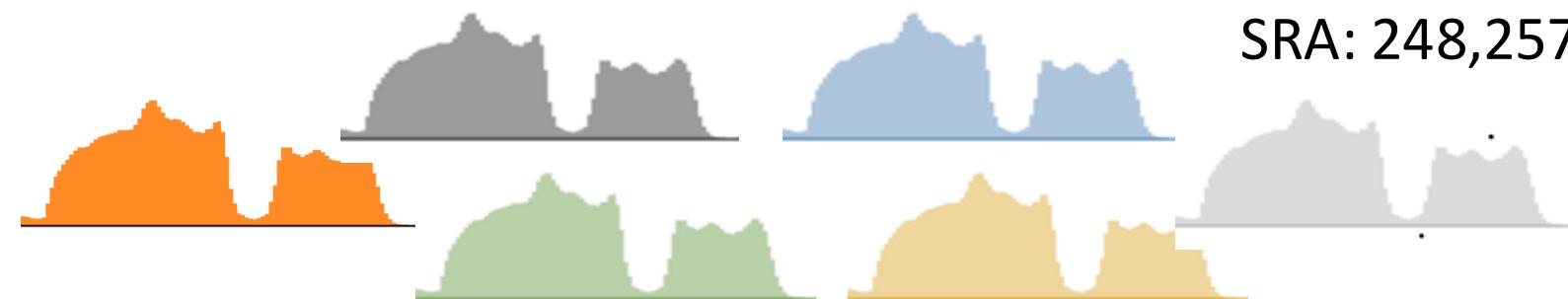
ChatGPT 4o:
Generate schematic cartoon
of RNA-Seq coverage



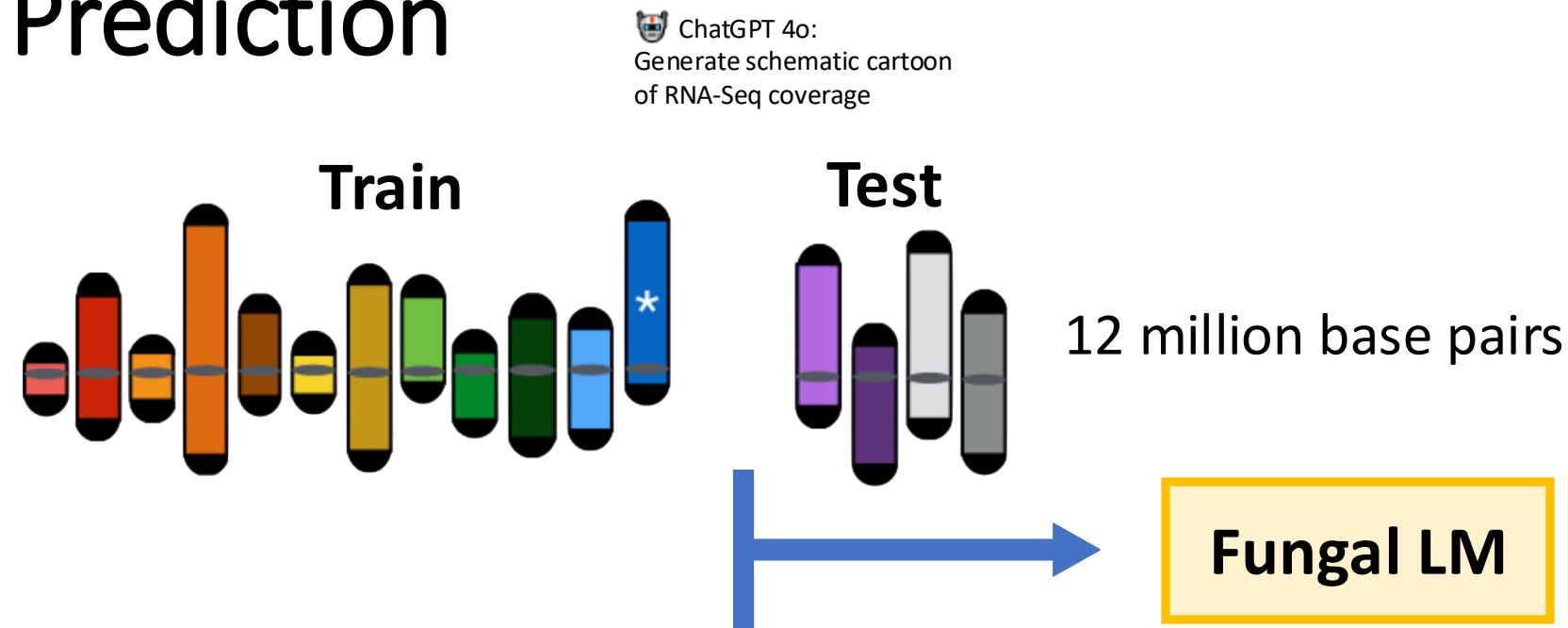
I Yeast RNA-Seq Prediction



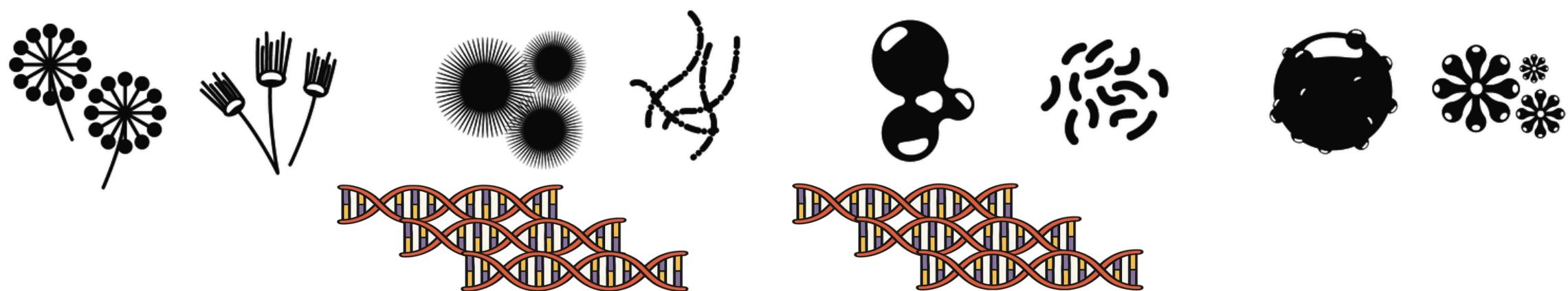
Genome size is too small, leading to overfitting.



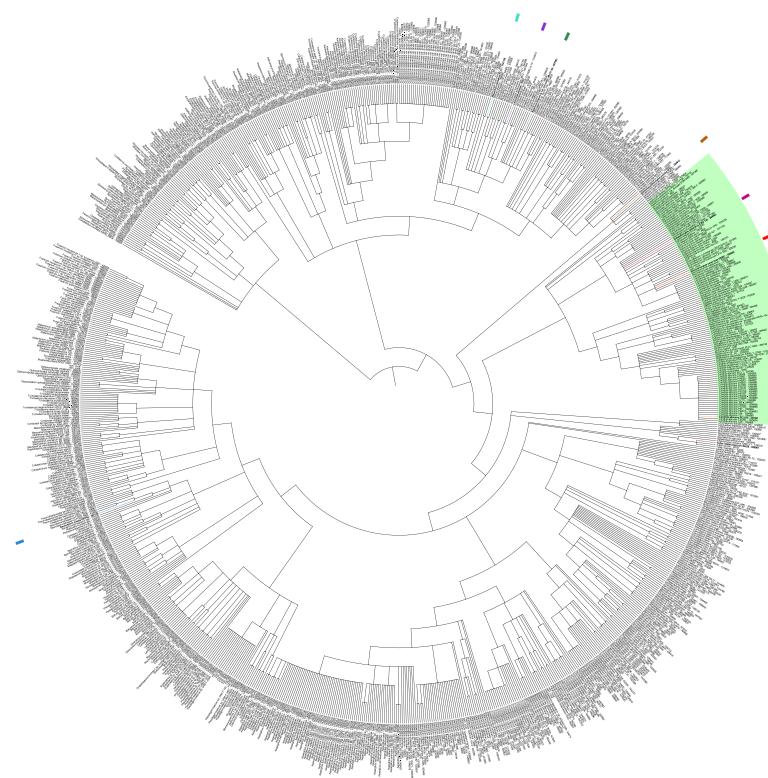
I Yeast RNA-Seq Prediction



Ensembl fungi: 1500 genomes



I Yeast RNA-Seq Prediction



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

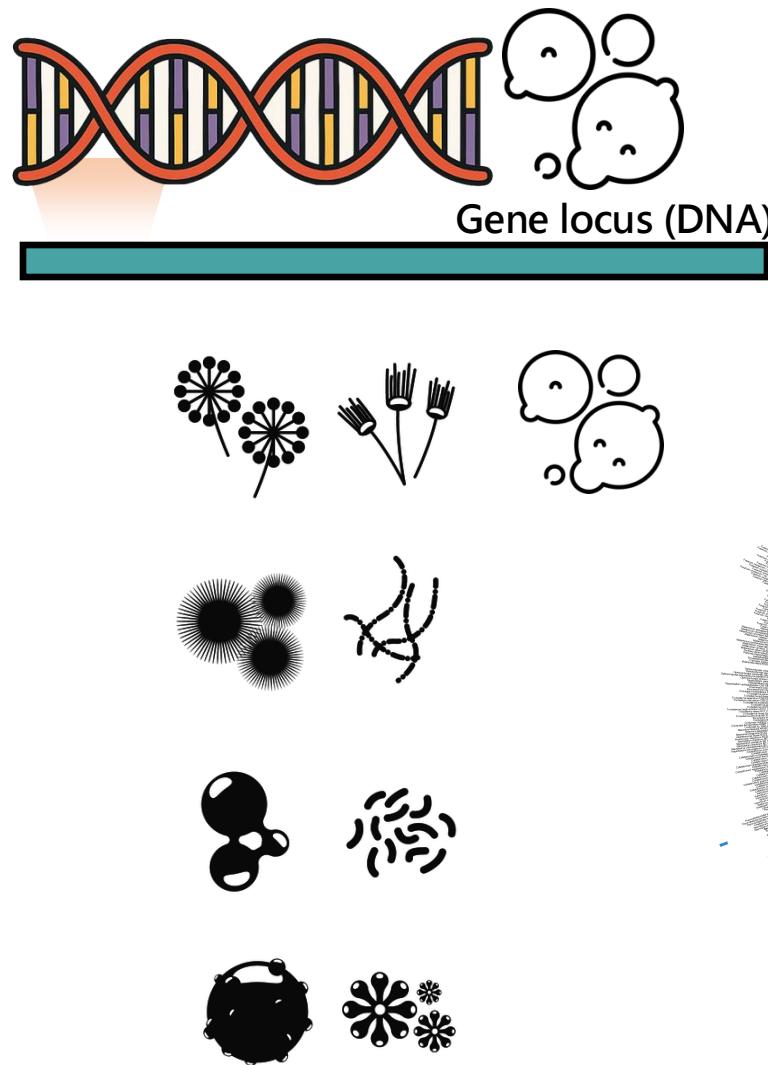
Fungal LM

Gene expression
model

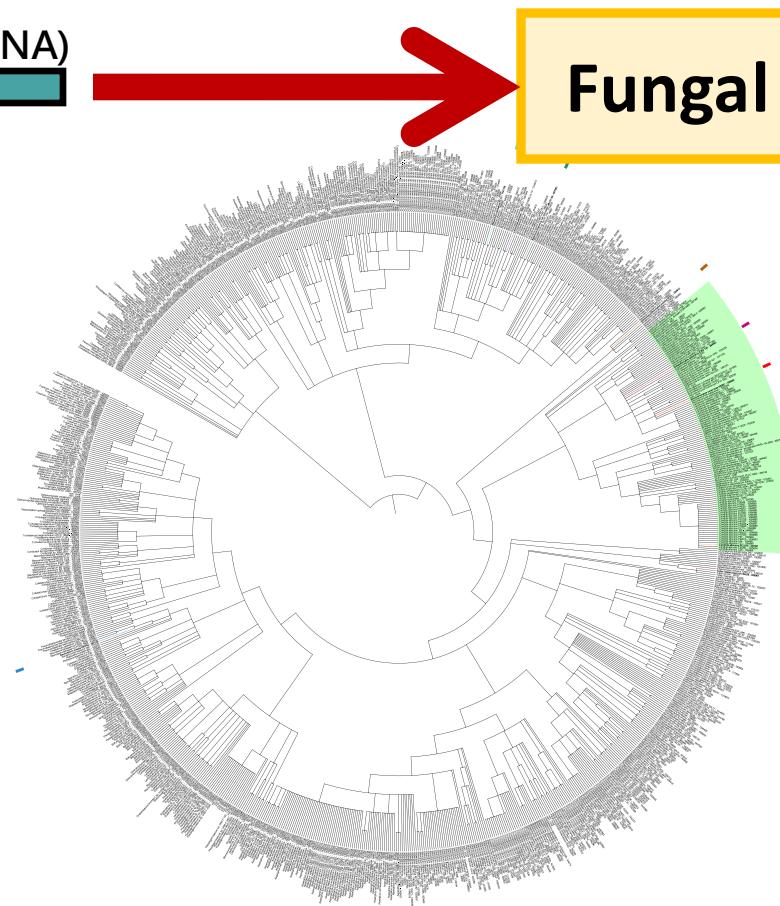


RNA-Seq expression?

I Yeast RNA-Seq Prediction



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



Fungal LM

Model
Fungal Model



RNA-Seq expression?

Deep learning-based DNA sequence model

nature methods

View all journals | Search

Explore content

Troyanskaya Lab Princeton

DeepSEA
2015Brief Communication | Published: 24 August 2015
Predicting effects of noncoding
learning-based sequence mo

Jian Zhou & Olga G Troyanskaya

GENOME
RESEARCH

HOME | ABOUT | ARCHIVE | SUBMIT | SUBSCRIBE | ADVERTISE | AUTH

Institution: MILTON S EISENHOWER LIBRARY Sign In

Basset: learning the regulatory
accessible genome with deep co
neural networksDavid R. Kelley¹, Jasper Snoek² and John L. Rinn¹

Calico

Basset
2016

Cell

Volume 176, Issue 3, 24 January 2019, Pages 535-548.e24

Article

Predicting Splicing from Pri
with Deep LearningKishore Jagannathan^{1,6}, Sofia Kyriazopoulou Panagiotopoul¹,
Sivash Fazel Darbandi², David Knowles³, Yang J. Li^{1,7}, Jack
Wenwu Cui¹, Grace B. Schwartz², Eric D. Chow³, Efstratios
Serofim Batzoglou¹, Stephan J. Sanders², Kyle Kai-How Farh^{1,7}

Illumina

SpliceAI
2019Agarwal and Kelley Genome Biology (2022) 23:245
https://doi.org/10.1186/s13059-022-02811-x

RESEARCH

The genetic and biochemical deter
of mRNA degradation rates in marVikram Agarwal^{1,2*} and David R. Kelley^{*}

Calico

Saluki
2022

nature biotechnology

Explore content | About the journal | Publish with us

FUToronto

DeepBind
2015Analysis | Published: 27 July 2015
Predicting the sequence sp
DNA- and RNA-binding pro
learning

Babak Alipanahi, Andrew Delong, Matthew T Weirauch & Brendan J Frey

GENOME
RESEARCH

HOME | ABOUT | ARCHIVE | SUBMIT | SUBSCRIBE | ADVERTISE | AUTH

Institution: MILTON S EISENHOWER LIBRARY Sign In

Sequential regulatory activity
across chromosomes with co
neural networksDavid R. Kelley¹, Yakir A. Reshef², Maxwell Bileschi³, Da
Cory Y. McLean³ and Jasper Snoek³

Calico

Basenji
2018

nature methods

Explore content | About the journal | Publish with us

Calico

Akita

2020

nature methods

Predicting 3D genome folding from
with AkitaGeoff Fudenberg^{1,15}, David R. Kelley^{2,3,5} and Katherine S. Pollard⁴

nature methods

Explore content | About the journal | Publish with us

Calico

scBasset
2022

Bioinformatics

Gifford Lab MIT

DNA-TF binding
2016Article Navigation
JOURNAL ARTICLE
Convolutional neural
protein binding

Haoyang Zeng, Matthew D. Edwards, Ge Liu, David K. Gifford

nature genetics

Explore content

nature

>

Troyanskaya Lab Princeton

Article | Published: 16 July 2018

Deep learning sequence-based
variant effects on expression an

Jian Zhou, Chandra L. Theesfeld, Kevin Yao, Kathleen M. O'Byrne

ExPecto
2018

ARTICLES

https://doi.org/10.1038/s41592-021-01252-x

OPEN

Effective gene expressio
sequence by integratinŽiga Avsec^{1,15}, Vikram Agarwal^{2,4}, Daniel Visontai¹,
Agnieszka Grabska-Barwińska¹, Kyle R. Taylor³
and David R. Kelley^{2,3,5}

DeepMind + Calico

Enformer
2021

nature genetics

Article

Predicting RNA-seq coverage
DNA sequence as a unifying n
of gene regulation

Calico

Borzoi
2025

Received: 28 August 2023

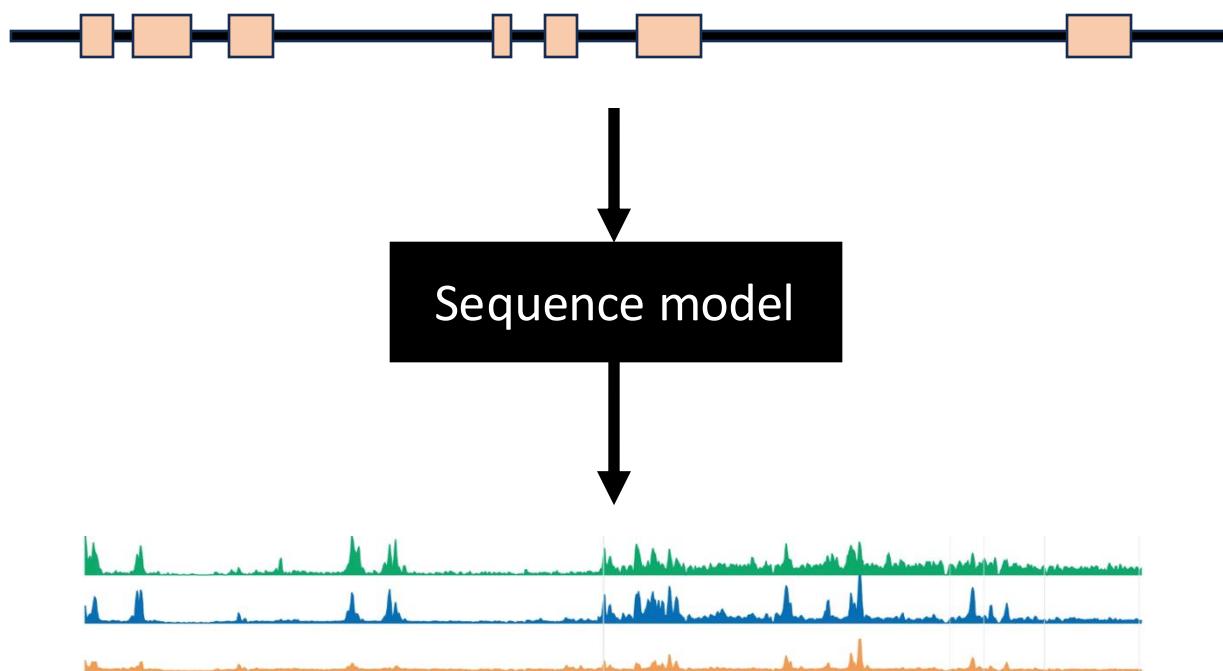
Johannes Linder^{1,15}, Divyanshi Srivastava¹, Han Yuan¹, Vikram Agarwal^{2,4}

Accepted: 4 December 2024

Supervised learning

Input: DNA sequences

Output: Genomics tracks



Protein Language model (PLM)

RESEARCH ARTICLE | BIOLOGICAL SCIENCES | Open access | Published: 27 July 2022

Facebook

Biological structure and function from scaling up to one million protein sequences

First LM attempt

PNAS 2020

Alexander Rives, Joshua Meier, Tom Sercu, and Rob Fergus [Authors info & affiliations](#)

Höcker Lab

ProtGPT2 is a deep unsupervised language model for protein design

Noelia Ferruz, Steffen Schmid, and Michael Höcker [Authors info & affiliations](#)

Nature Communications 13, Article number: 1322 (2022) | [DOI: 10.1038/s41467-022-12500-0](#) | [PMID: 35832320](#)

ProtGPT2

Nat Commun 2022

Article | Published: 03 October 2022

Transformer protein language unsupervised learning

Facebook

Transformer protein LM

ICLR 2021

Roshan Rao, Alexander Rives, and Michael Höcker [Authors info & affiliations](#)

arXiv:2006.15222 (cs)

[Submitted on 26 Jun 2020 (v1), last revised 28 Mar 2021]

Saleforce + UIUC

BERTology Meets Biology: Learning Attention in Protein Language Models

ICLR 2021

Jesse Vig, Ali Madani, Lav R. Varshney, Caiming Xiong, Richard Socher, James M. Holton, Ben Krause, and James Zou [Authors info & affiliations](#)

Large language models generate functional protein sequences across diverse families

Saleforce

ProGen

Nat Biotechnol 2023

Article | Published: 26 January 2023

SYNTHESIS · Volume 12, Issue 6, P654-669.E3, June 16, 2020 | [Open Access](#)

MIT

Learning the protein structure, and function from scaling up to one million protein sequences

Review

Cell Systems 2021

Tristan Bepler, Daniel K. Sercu, and Rob Fergus [Authors info & affiliations](#)

ARTICLE · Volume 14, Issue 11, Article number: 1100 (2023) | [DOI: 10.1038/s41551-023-01500-w](#) | [PMID: 37144000](#)

Saleforce + Profluent

ProGen2: Exploring the limits of large language models

ProGen2

Cell System 2023

Erik Nijkamp, Jeffrey A. Ruffolo, Ali Madani, and Michael Höcker [Authors info & affiliations](#)

f x Meta

RESEARCH ARTICLE | STRUCTURE PREDICTION

ZEMING LIN, HALIL AKIN, [...] AND ALEXANDER RIVES [+12 authors](#)

Evolutionary-scale prediction of protein structure with a language model

ESMFold

Science 2023

Primer | Published: 15 February 2024

Designing proteins with AI

Profluent

PLM

Nat Biotechnol 2024

Nature Biotechnology 42, 200–202 (2024) | [DOI: 10.1038/s41551-023-01500-w](#) | [PMID: 37144000](#)

arXiv:2402.16445 (cs)

[Submitted on 26 Feb 2024 (v1), last revised 16 Jul 2024 (this version: v1)]

Peking Uni

ProLLaMA: A Protein Language Model for Multi-Task Protein Language Modeling

ProLLaMA

arxiv, 2024

Lizhenghao Lv, Zongying Lin, Hao Li, Yuyang Liu, Jiaxin Wang, and Yonghong Tian [Authors info & affiliations](#)

~ 6000 citations

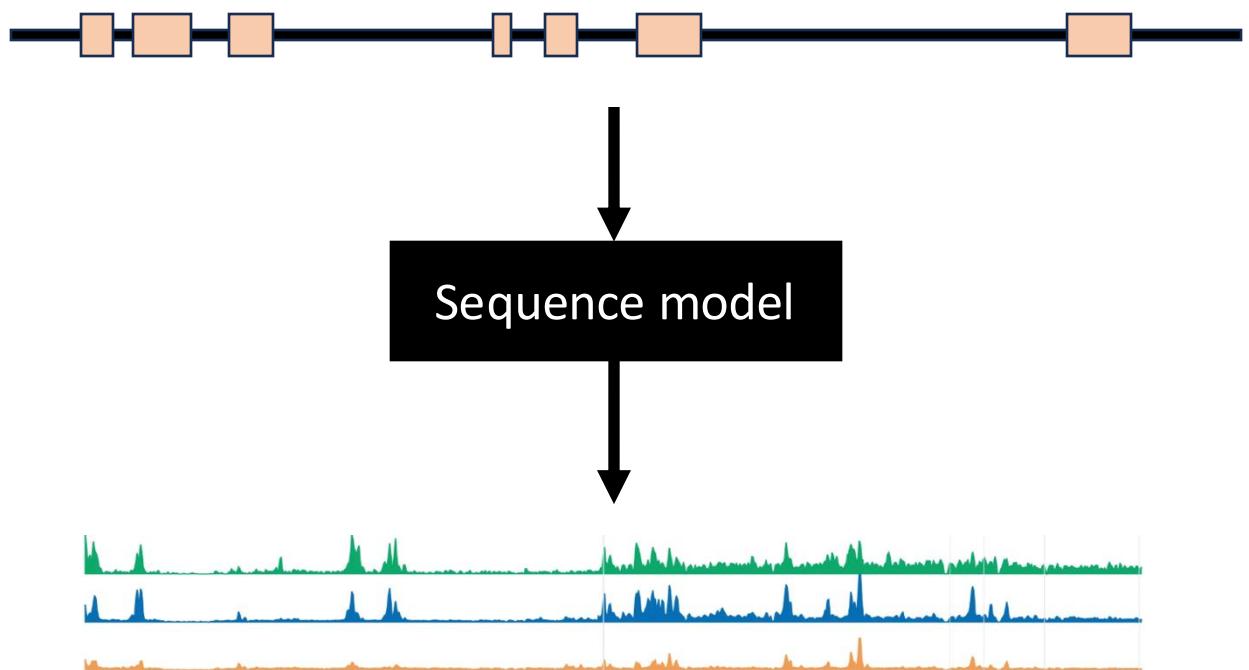
DNA Language model

<p>DNABERT Bioinformatics 2021</p> <p>UC Berkeley</p> <p>GPN PNAS 2023</p> <p>Nucleotide Transformer bioRxiv 2023</p> <p>HyenaDNA NeurIPS 2023</p>	<p>DNABERT-2 ICLR 2024</p> <p>TUD</p> <p>GROVER Nat Mach Intell 2024</p> <p>Harvard + MIT</p> <p>Genomic LM Nat Commun 2024</p> <p>TUM</p> <p>Species-aware DNA LM Genom Biol 2024</p>	<p>Stanford + Arc Inst + TogetherAI</p> <p>Evo bioRxiv 2024</p> <p>Cornell + Princeton + CMU</p> <p>Caduceus ICML 2024</p> <p>Cornell + USDA-ARS + Simons</p> <p>PlantCaduceus bioRxiv 2024</p> <p>TUM</p> <p>Nucleotide dependency analysis of DNA language models bioRxiv 2024</p>
<p>DNABERT: pre-trained Bidirectional Encoder Representations from Transformers model for DNA-language in genomics Yanrong Ji, Zihuan Zhou, Han Liu et al.</p> <p>RESEARCH ARTICLE BIOPHYSICS AND COMPUTATIONAL BIOLOGY 8</p> <p>DNA language models are powerful predictors of genome-wide variation</p> <p>Gonzalo Benegas, Sanjit Singh Batra, and Yun S. Song et al.</p> <p>The Nucleotide Transformer: Foundation Models for</p> <p>Hugo Dalla Favera, Adam Henne, Bernardo P. Ribeiro, and Marie Lopez et al.</p> <p>arXiv:2306.15794 (cs)</p> <p>[Submitted on 27 Jun 2023 (v1), last revised 14 Nov 2023 (this version, v2)]</p> <p>HyenaDNA: Long-Range Context Modeling at Single Nucleotide Resolution Eric Nguyen, Michael Poli, Marjan Faizi, Armin Thomas, Aman Patel, Clayton Rabideau, Stefano Massaroli, Yoshua Bengio, and Chris Ré et al.</p>	<p>DNABERT-2: Efficient Foundation Model Benchmark For Multi-Species Zhihan Zhou, Yanrong Ji, Weijian Li, Pratik Garg, and Han Liu et al.</p> <p>Article Open access Published: 23 July 2024</p> <p>DNA language model context in the human genome Melissa Sanabria, Jonas Hirsch, Pierre Lefebvre, and Pierre-Antoine Vial et al.</p> <p>Article Open access Published: 03 April 2024</p> <p>Genomic language model predicts gene regulation and function Yunha Hwang, Andre L. Cornman, Elizabeth J. Cade, and R. Girguis et al.</p> <p>Research Open access Published: 02 April 2024</p> <p>Species-aware DNA language models capture regulatory elements across species Alexander Karolchik, Julien Gagneur, and Daniel H. Schmid et al.</p> <p><i>Genome Biology</i> 25, Article number: 83 (2024) Cite this article</p>	<p>Sequencing molecule Eric Nguyen, Michael Poli, Matthew G. Durrant, Armin Thomas, Jeremy Sullivan, Madelena Y. Ng, Ashley Lewis, Aman Patel, Stephen A. Baccus, Tina Hernandez-Boussard, Christopher J. Burtt, and Daniel H. Schmid et al.</p> <p>doi: https://doi.org/10.1101/2024.02.27.582234</p> <p>[Submitted on 5 Mar 2024]</p> <p>Caduceus: Bi-Directional Range DNA Sequence Model Yair Schiff, Chia-Hsiang Kao, Aaron Gokaslan, Tripti Agarwal, and Jingjing Zhai et al.</p> <p>Cross-species nucleotide dependency analysis Jingjing Zhai, Aaron Gokaslan, Yair Schiff, Michelle C. Stitzer, M. Cinta Romo, and Jingjing Zhai et al.</p> <p>doi: https://doi.org/10.1101/2024.06.04.537000</p>
		<p>53</p>

Supervised learning

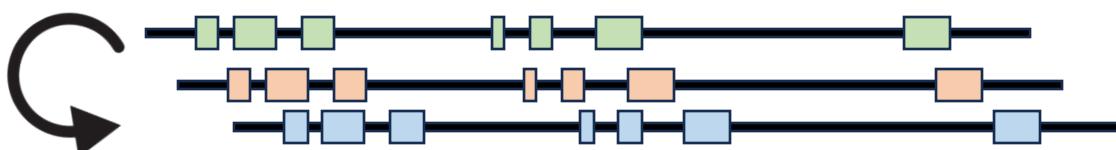
Input: DNA sequences

Output: Genomics tracks



Stage 1

Self-supervised learning



Language model (LM) /
Foundation model

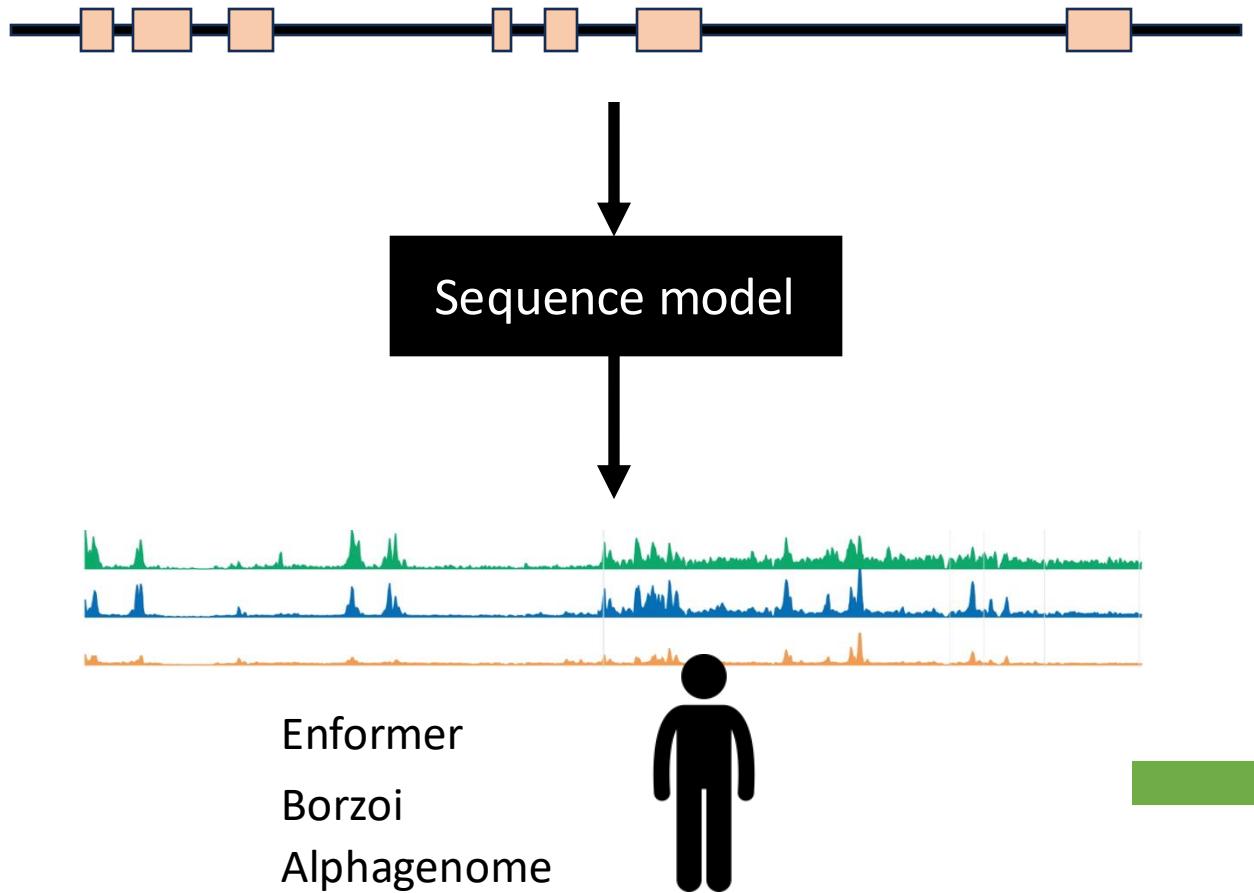
Stage 2

Fine-tuning LM



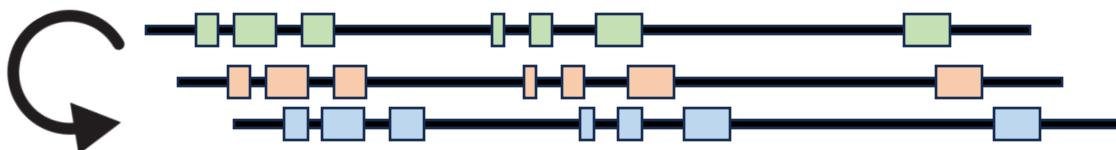
Enformer
Borzoi

Supervised learning



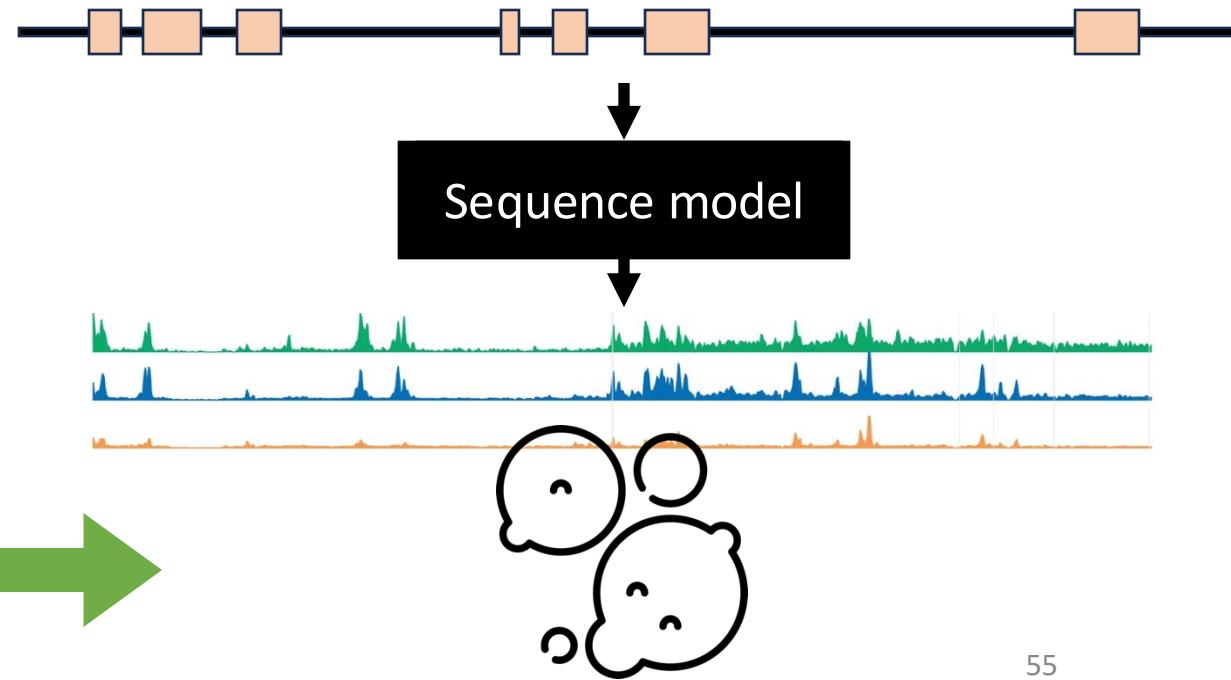
Stage 1

Self-supervised learning



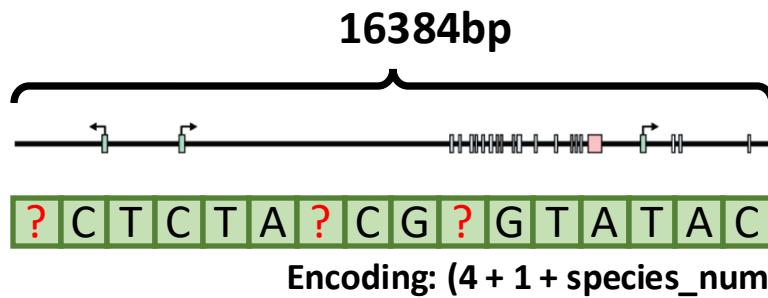
Stage 2

Fine-tuning LM



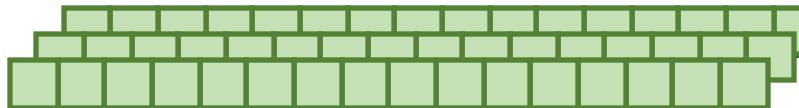


Shorkie



Reverse complementary

1bp res



16384 * 4 Masked language modeling loss

A
C
G
T

.8	.0	.0
.1	.0	.7
.1	.9	.1
.0	.1	.2



1bp res

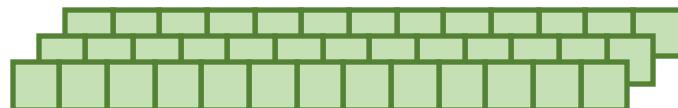
...

...

...

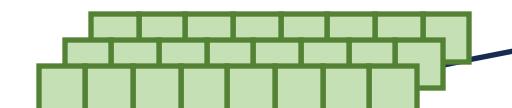
...

16bp res



16bp res

32bp res



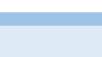
32bp res

64bp res



64bp res

128bp res



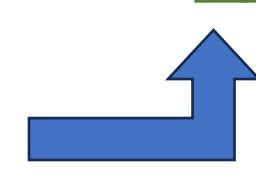
128bp res



Borzo

Linder, J. et al. (2025). Predicting RNA-seq coverage from DNA sequence as a unifying model of gene regulation. Nature Genetics, 1-13.

Transformer
Blocks (8x)



RAP1



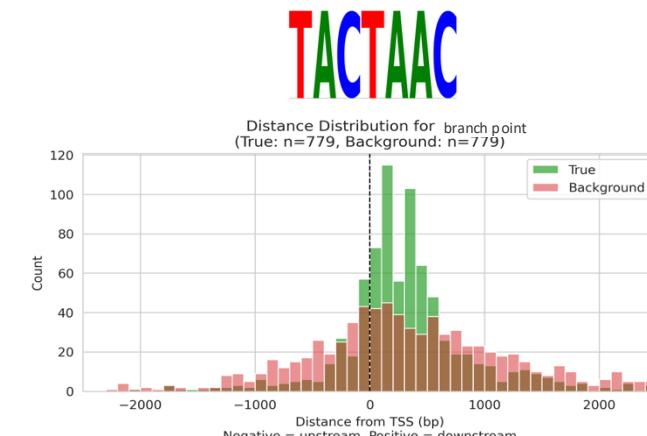
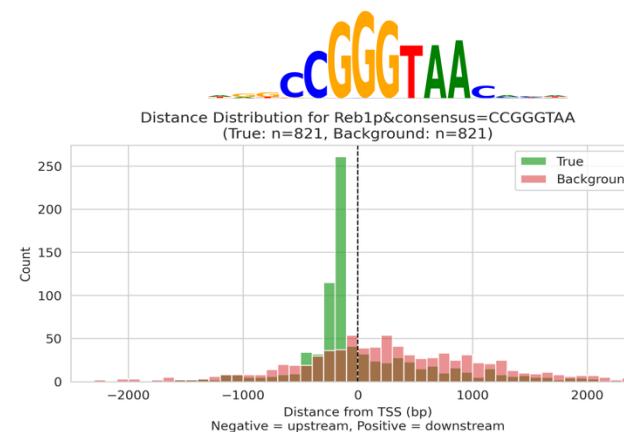
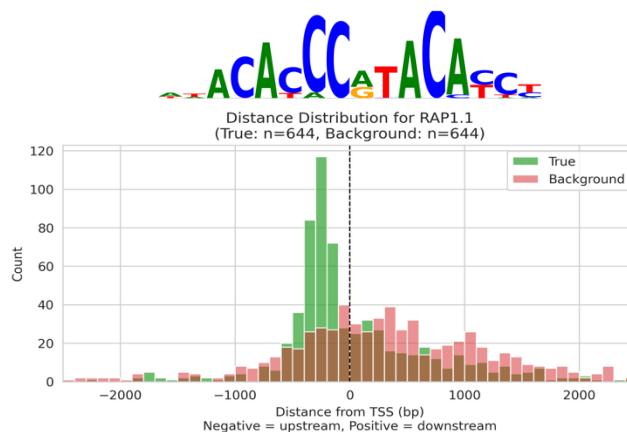
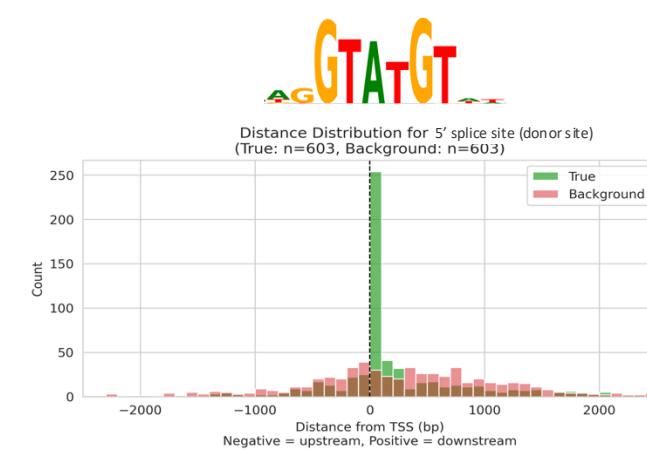
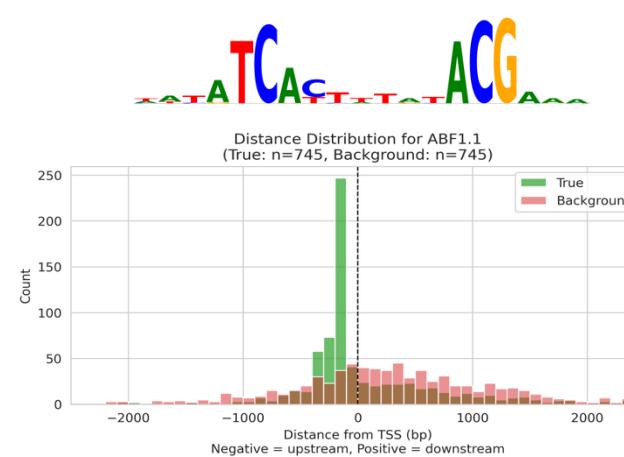
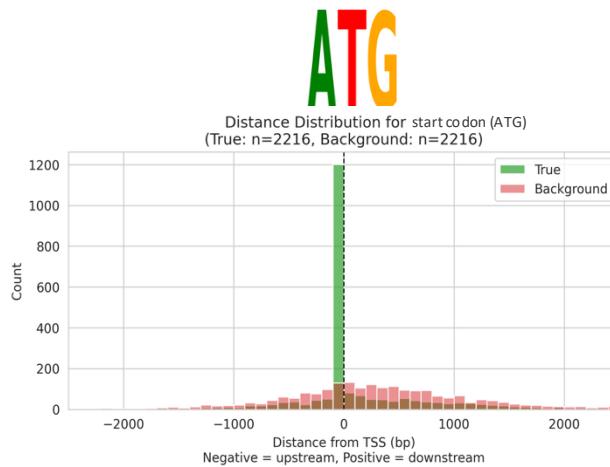
FHL1



LM predicts motifs in RP gene Promoter regions



These motifs are enriched in promoter regions





Shorkie

16384bp

? C T C T A ? C G ? G T A T A C

16384 * 4

A
C
G
T

.8	.0	.0
.1	.0	.7
.1	.9	.1
.0	.1	.2

1bp res



1bp res

...

...

...

...

16bp res



16bp res

32bp res



32bp res

64bp res



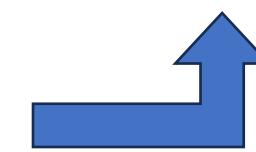
64bp res

128bp res



128bp res

Transformer
Blocks (8x)





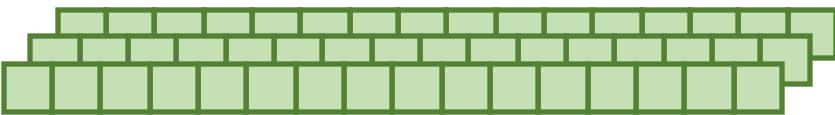
Shorkie

16384bp

? C T C T A ? C G ? G T A T A C



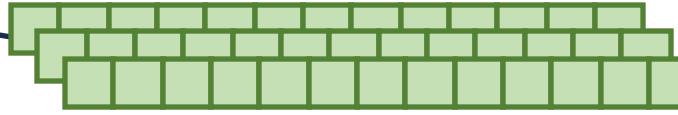
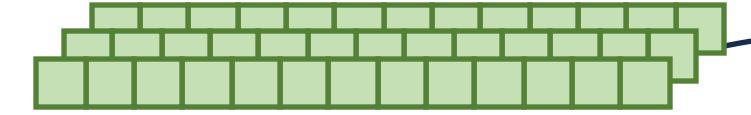
1bp res



...

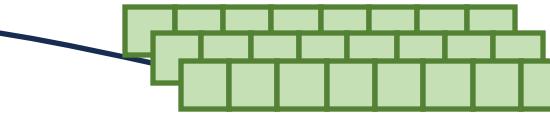
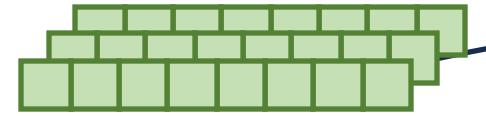
...

16bp res



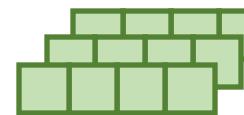
16bp res

32bp res



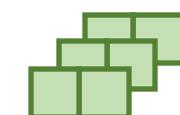
32bp res

64bp res

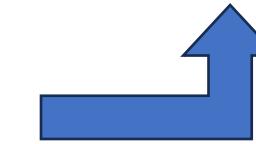


64bp res

128bp res

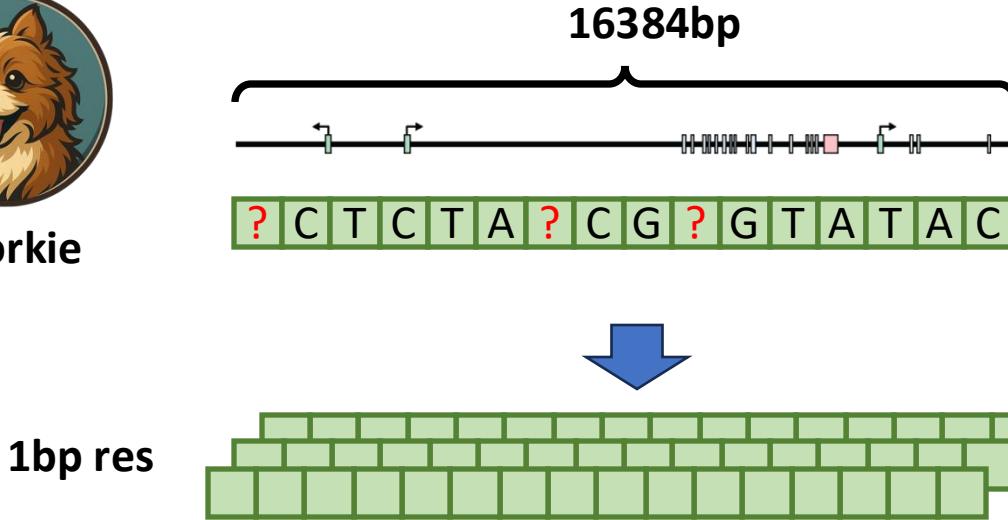


Transformer
Blocks (8x)

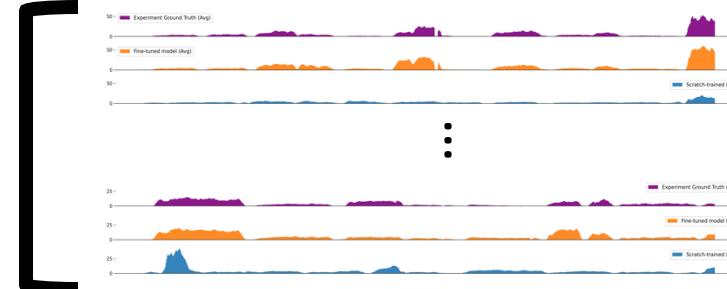




Shorkie



Coverage Tacks (896 * 2488)



ChIP-exo	(1128)
Histone marks	(20)
RNA-Seq	(3054)
1000-strains	
RNA-Seq	(1014)

16bp res



16bp res

32bp res



32bp res

64bp res



64bp res

128bp res

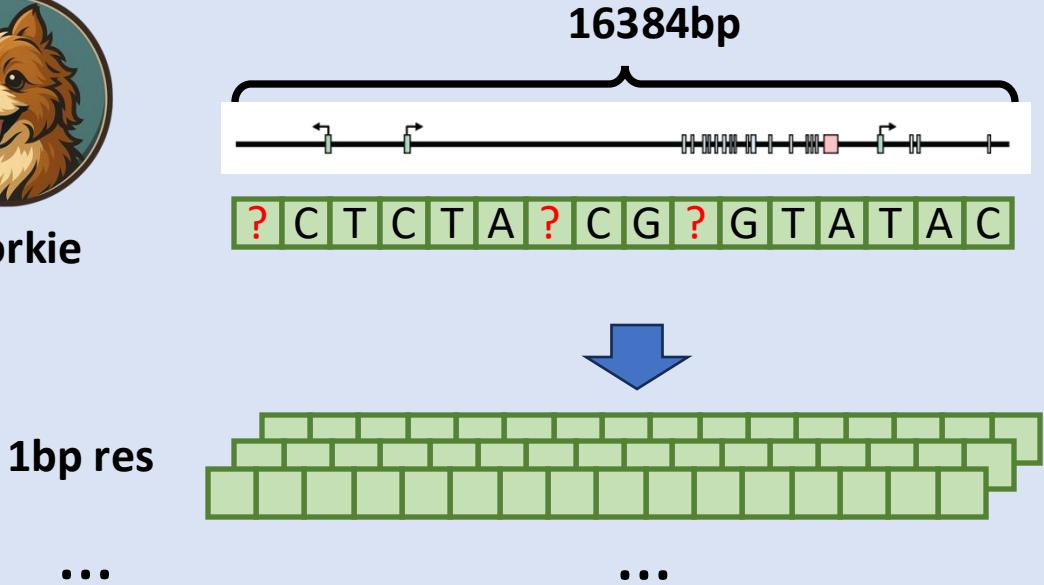


128bp res

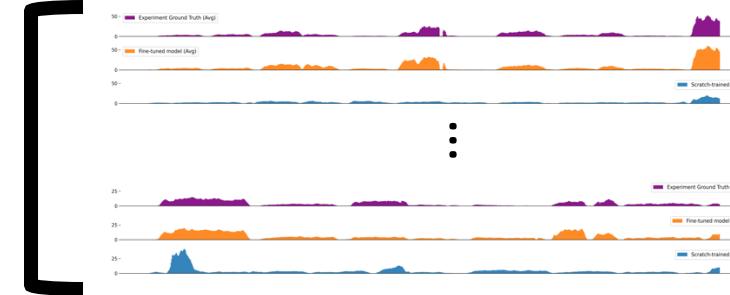
Transformer
Blocks (8x)



Shorkie



Coverage Tacks (896 * 2488)



ChIP-exo	(1128)
Histone marks	(20)
RNA-Seq	(3054)
1000-strains	
RNA-Seq	(1014)



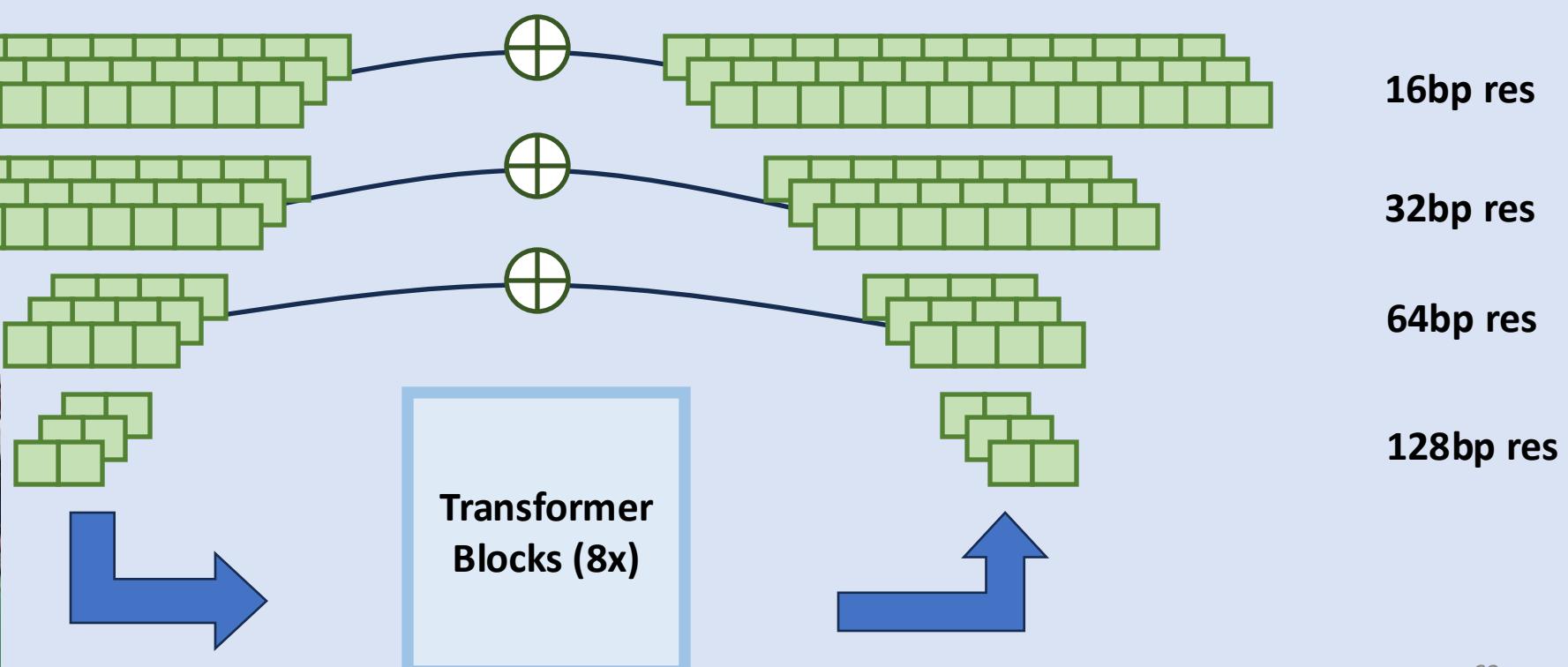
Introduction

Fungal LM

Transfer Learning

Interpretability

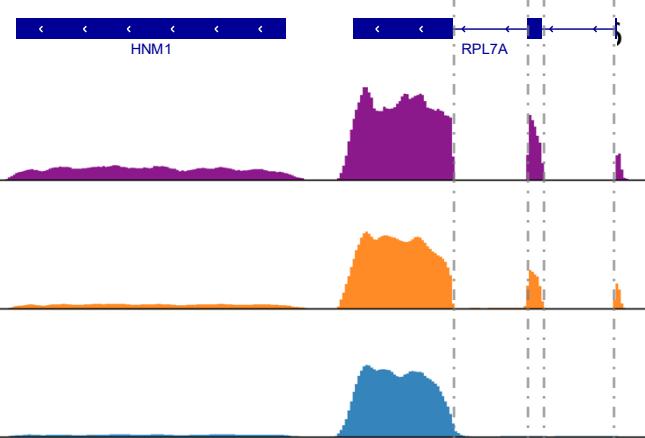
Applications



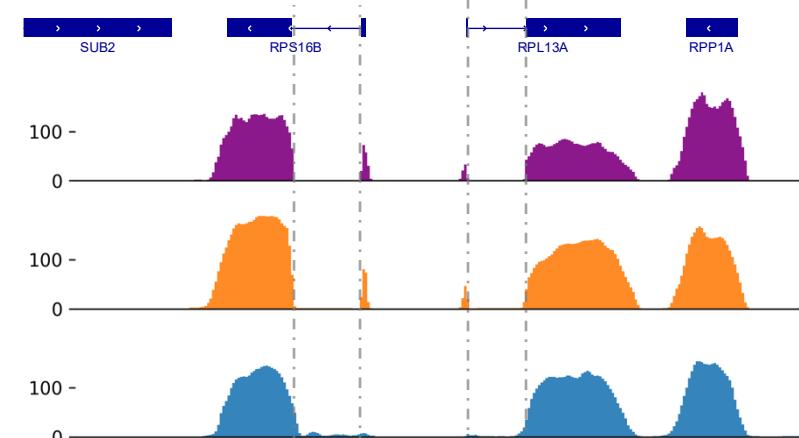
Shorkie vs Random Initialization (Test set)

RNA-Seq tracks

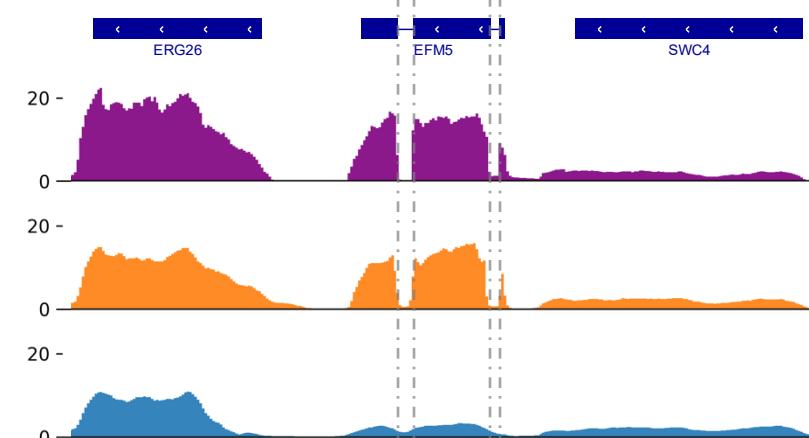
chrVII:362,180-366,023 (RNA-Seq tracks); fold 3



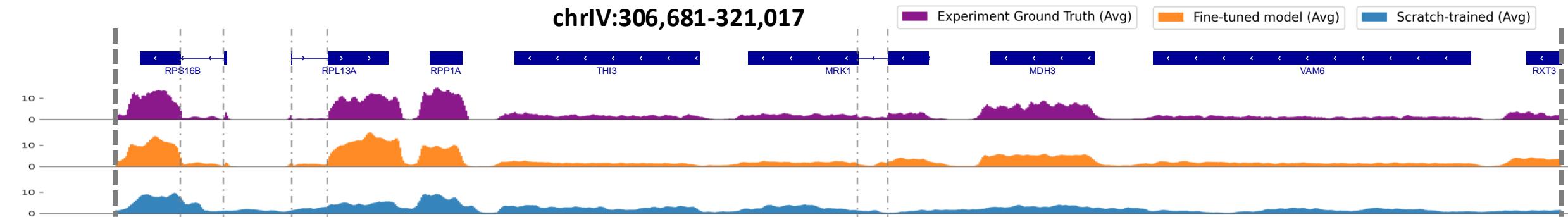
chrIV:305,657-310,505 (RNA-Seq tracks); fold 3



chrVII:495,374-499,965 (RNA-Seq tracks); fold 6



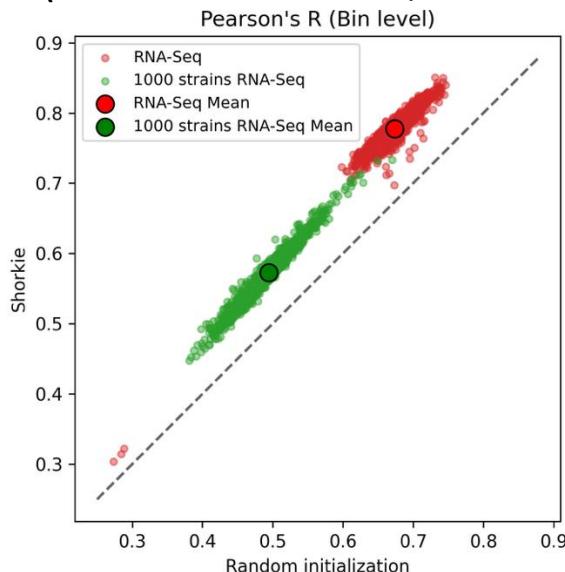
chrIV:306,681-321,017



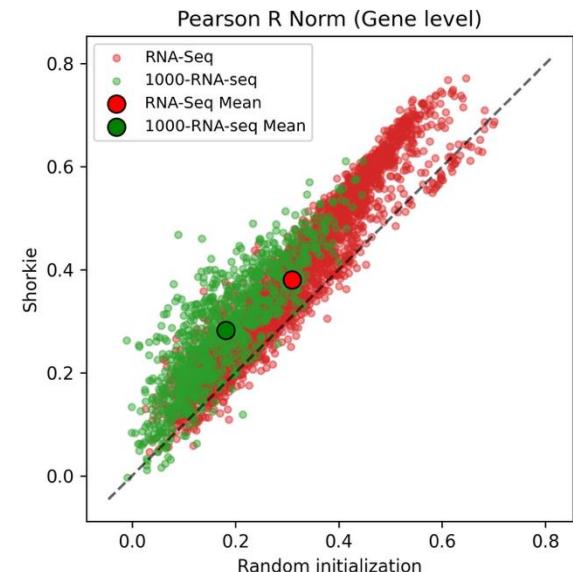
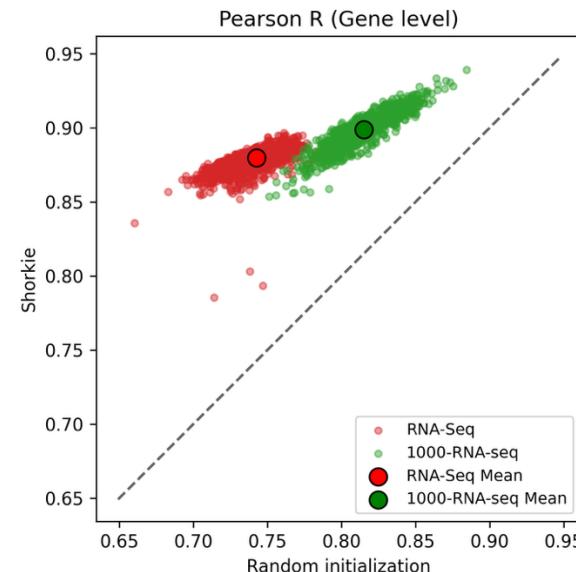
Shorkie vs Random Initialization (Test set)

Bin level

(each dot is a track; all bins)

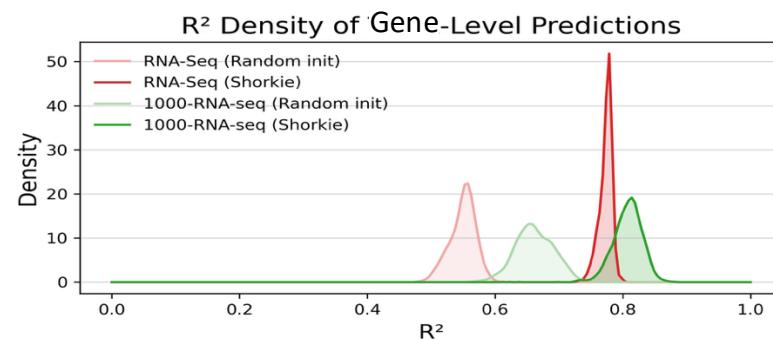
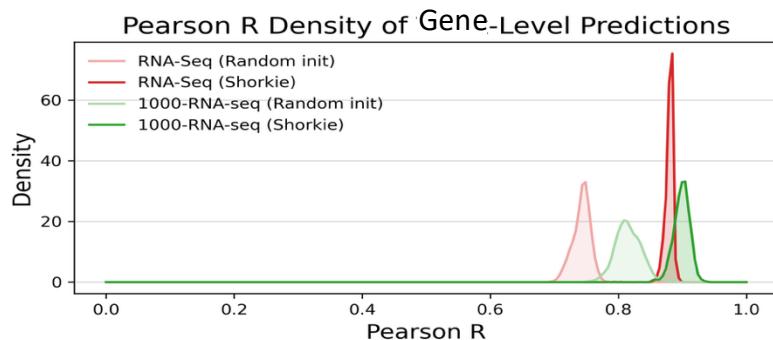
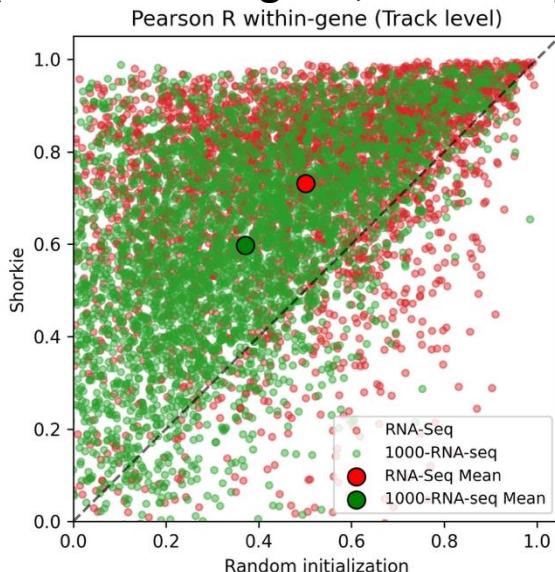


Gene level (each dot is a track; all genes)



Track level

(each dot is a gene; all tracks)



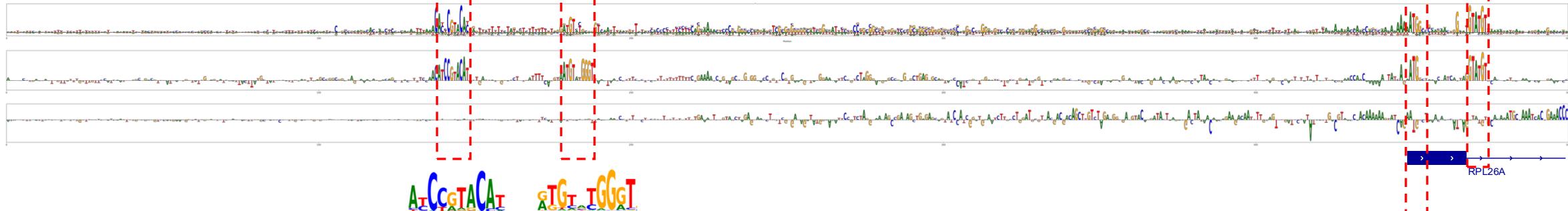
Shorkie ISM maps: Ribosomal genes & RRB genes

450 nt

50 nt

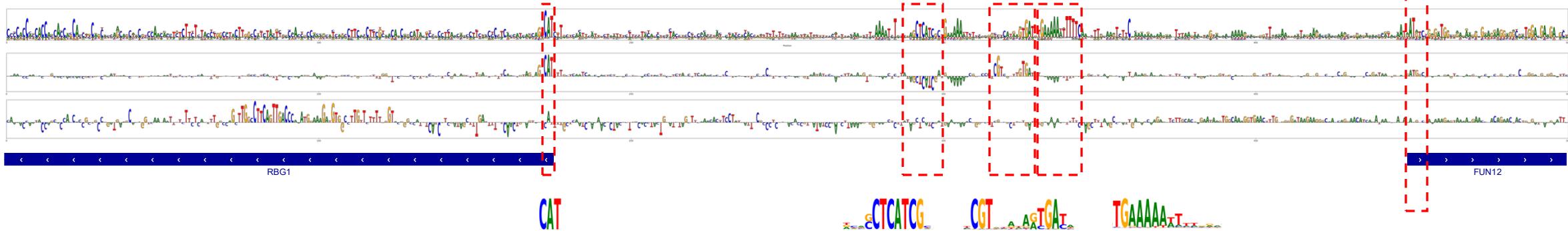
RPL26A (chrXII:818,862-819,362)

Fungal LM

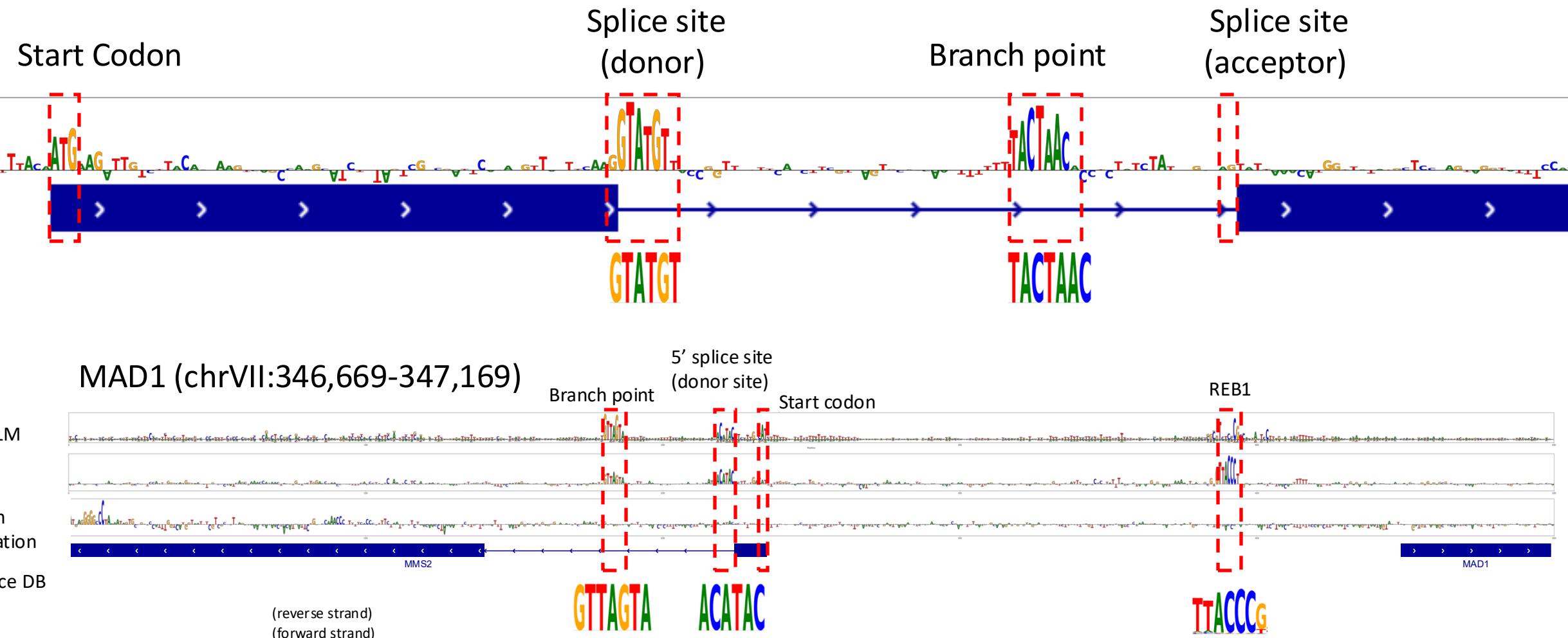


FUN12 (chrI:75,977-76,477)

Fungal LM



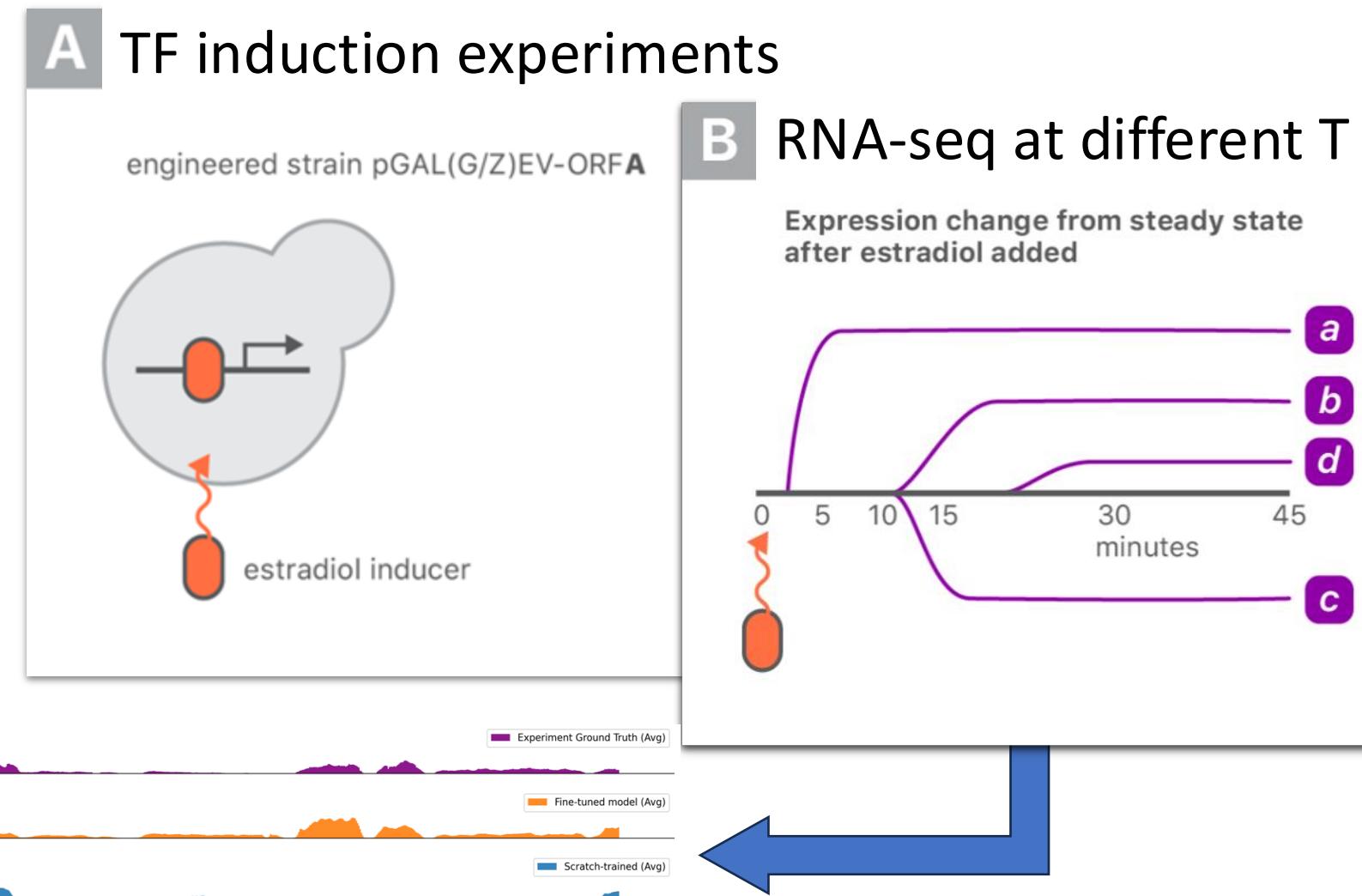
Shorkie ISM maps: Splicing motifs



Time-course TF-induced RNA-Seq

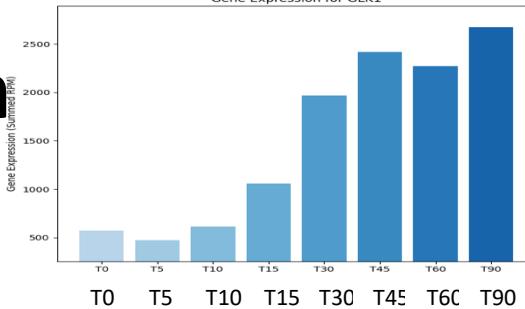
- Genome-scale transcription factor **perturbation** (1340 experiments; 3054 RNA-Seq readouts)
- Aggregating dynamics across many **time-courses**

3054 RNA-Seq readouts



Shorkie MSN2 Induction temporal RNA-Seq prediction

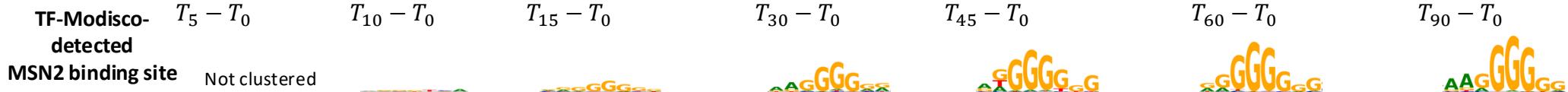
chrIII:50388-50888 (Promoter region of GLK1)



Average motif changes using TF-modisco

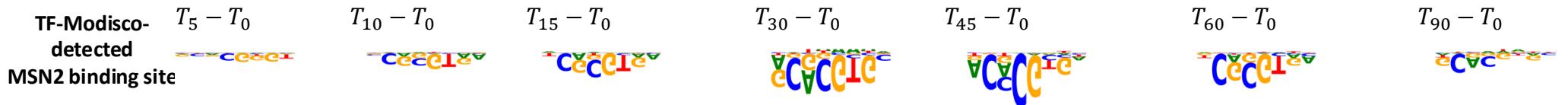
MSN2 Induction

YeTFaSCo
DB motif



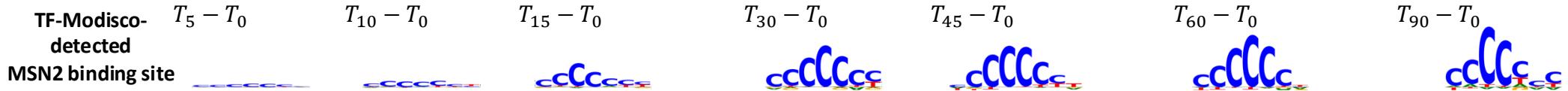
MET4 Induction

YeTFaSCo
DB motif



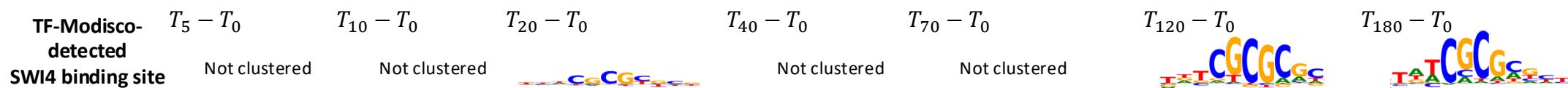
MSN4 Induction

YeTFaSCo
DB motif

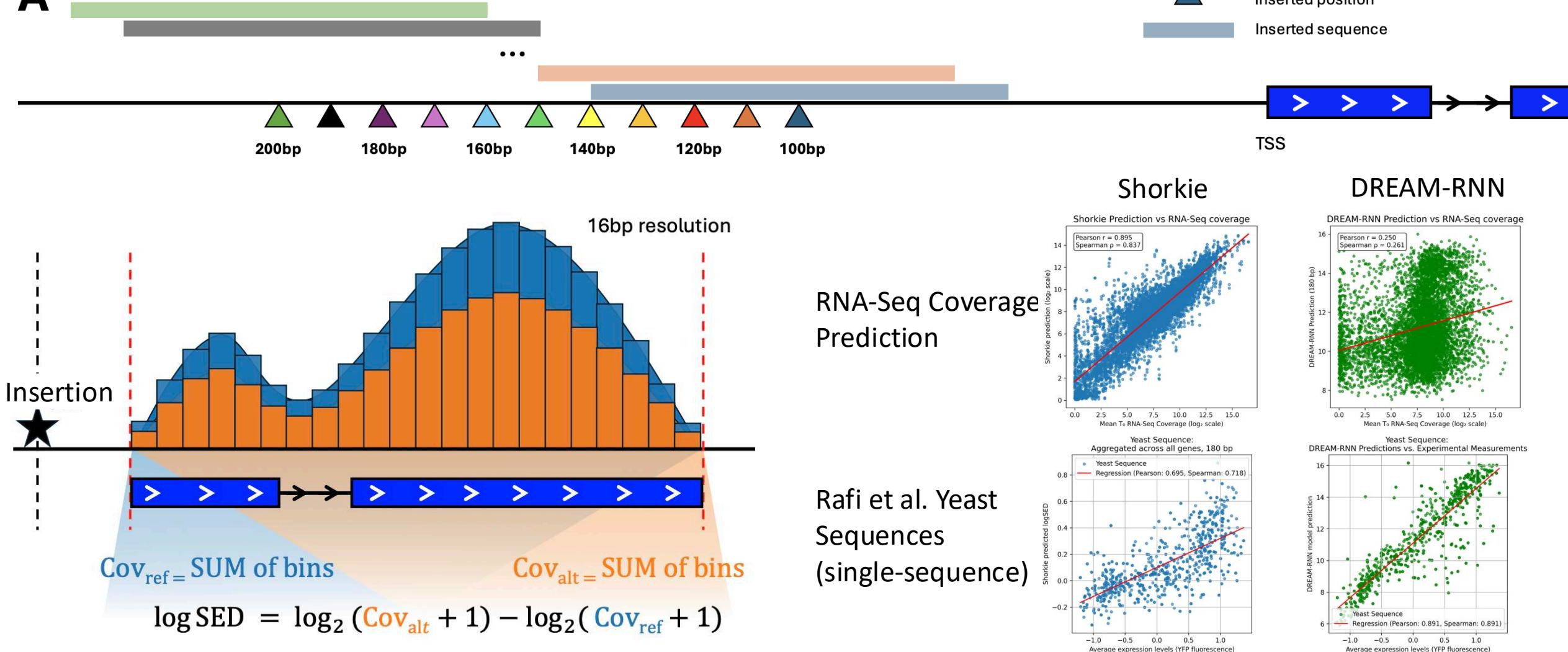


SWI4 Induction

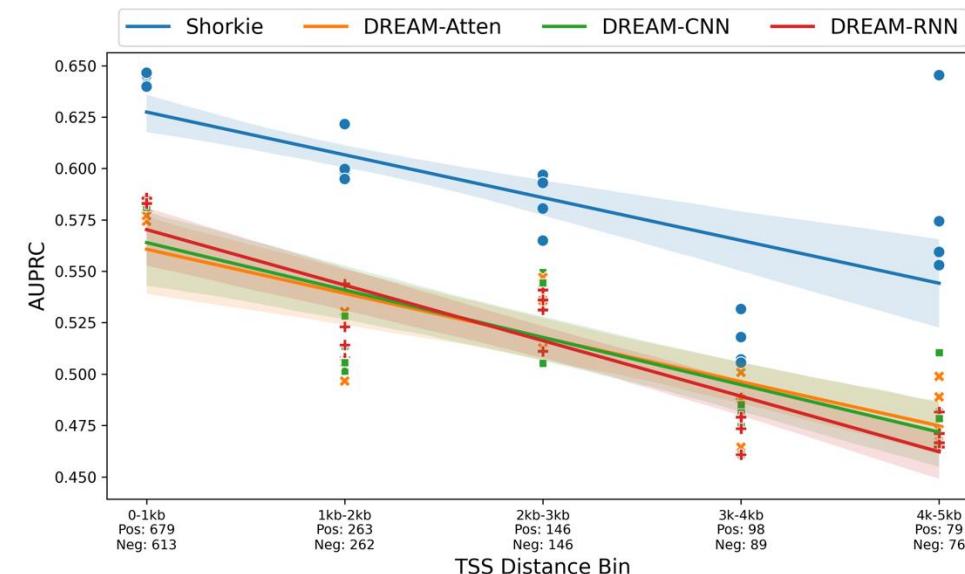
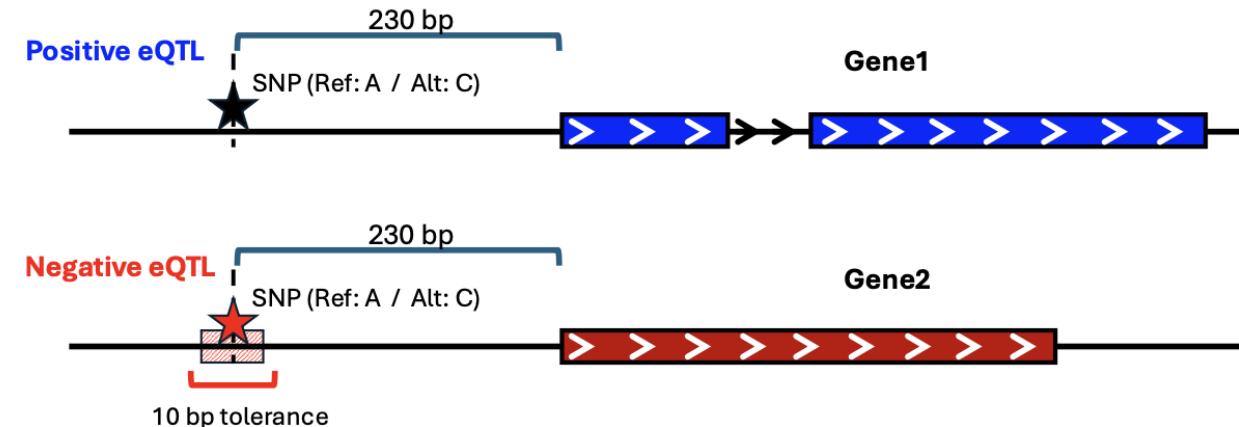
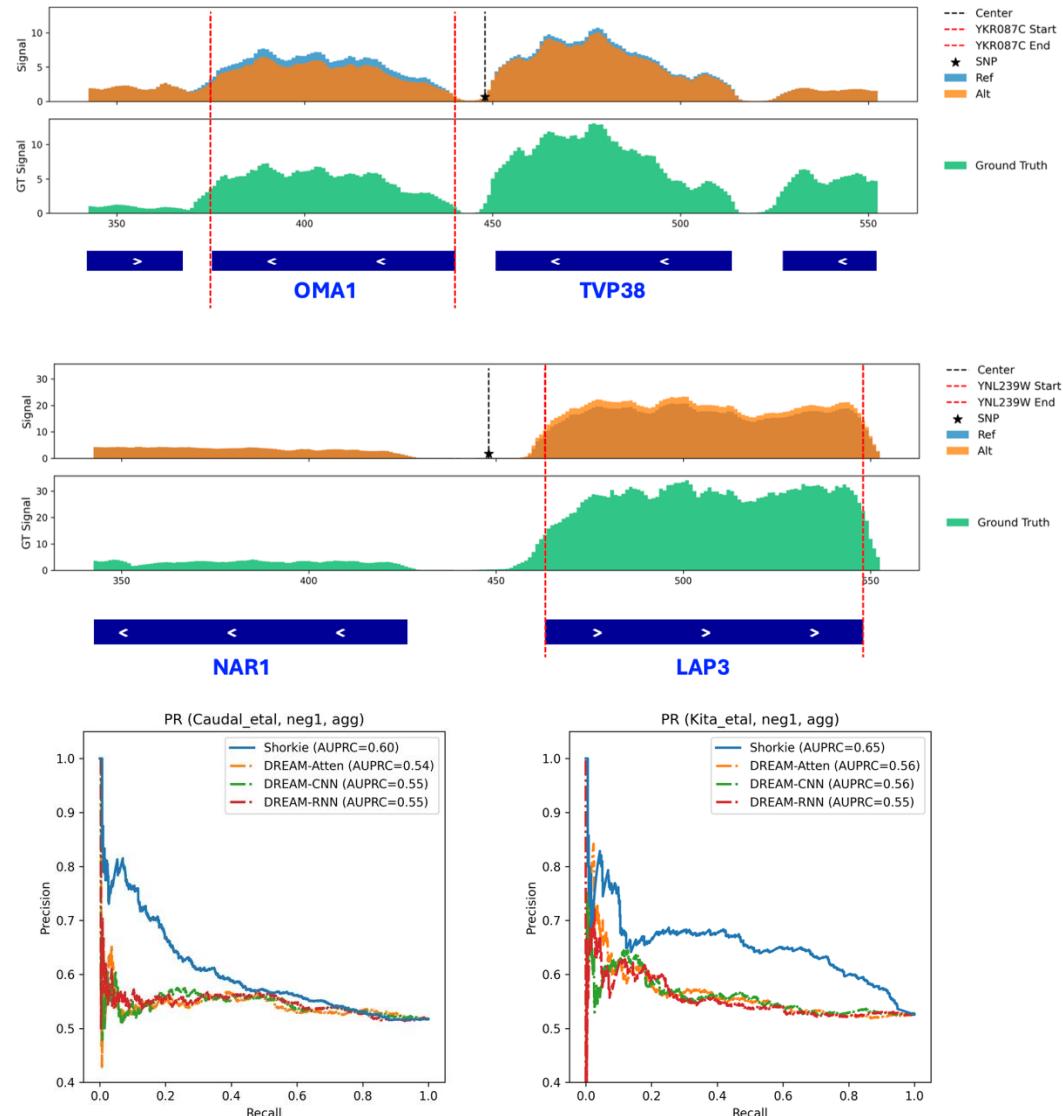
YeTFaSCo
DB motif



MPRA mutation effect prediction

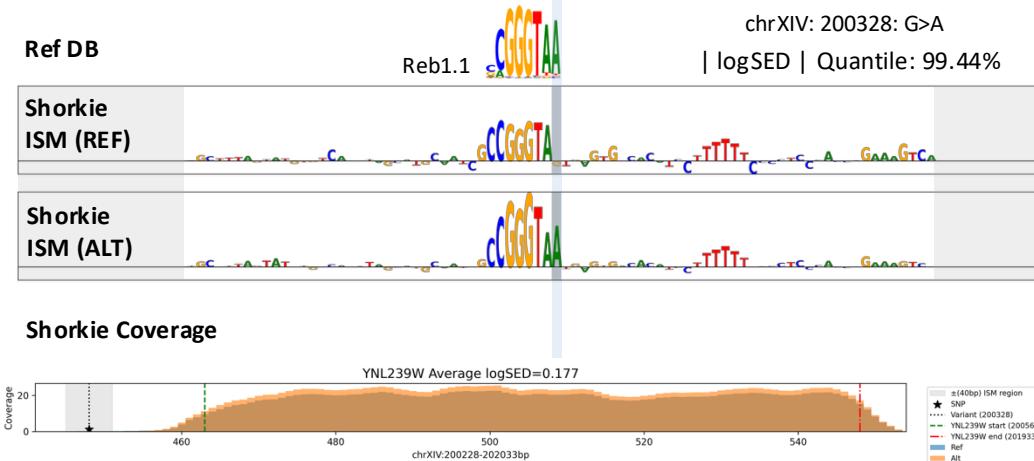
A

Variant effects on eQTLs and negatively selected eQTLs

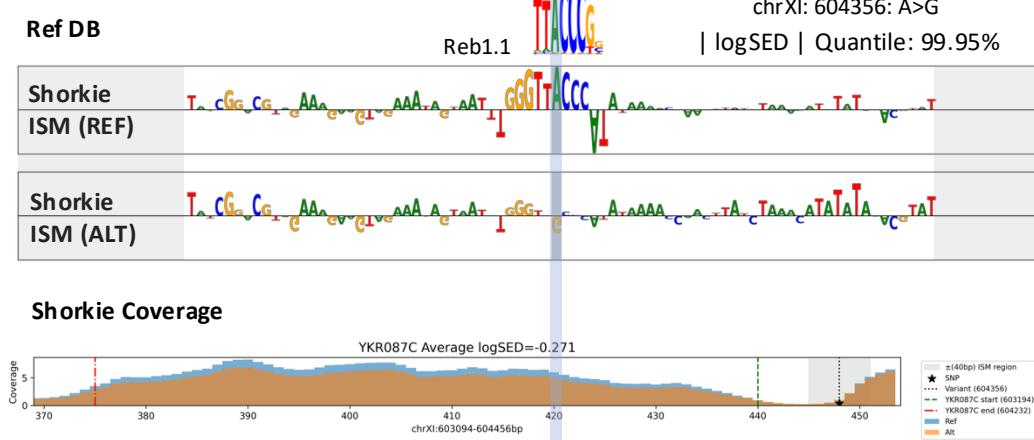


Variant effects on eQTLs and negatively selected eQTLs

chrXIV:200,288 – 200,368



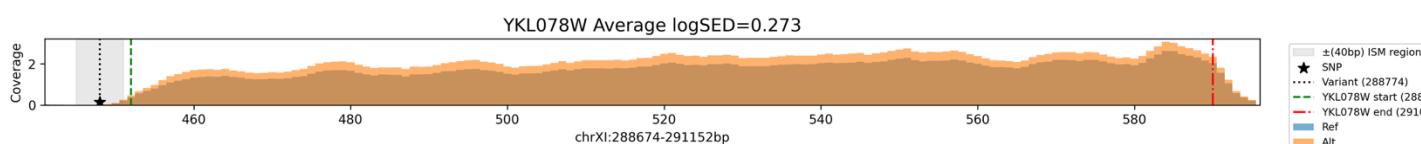
chrXI:604,316-604,396



chrXI:288,734 – 288,814



Shorkie Coverage



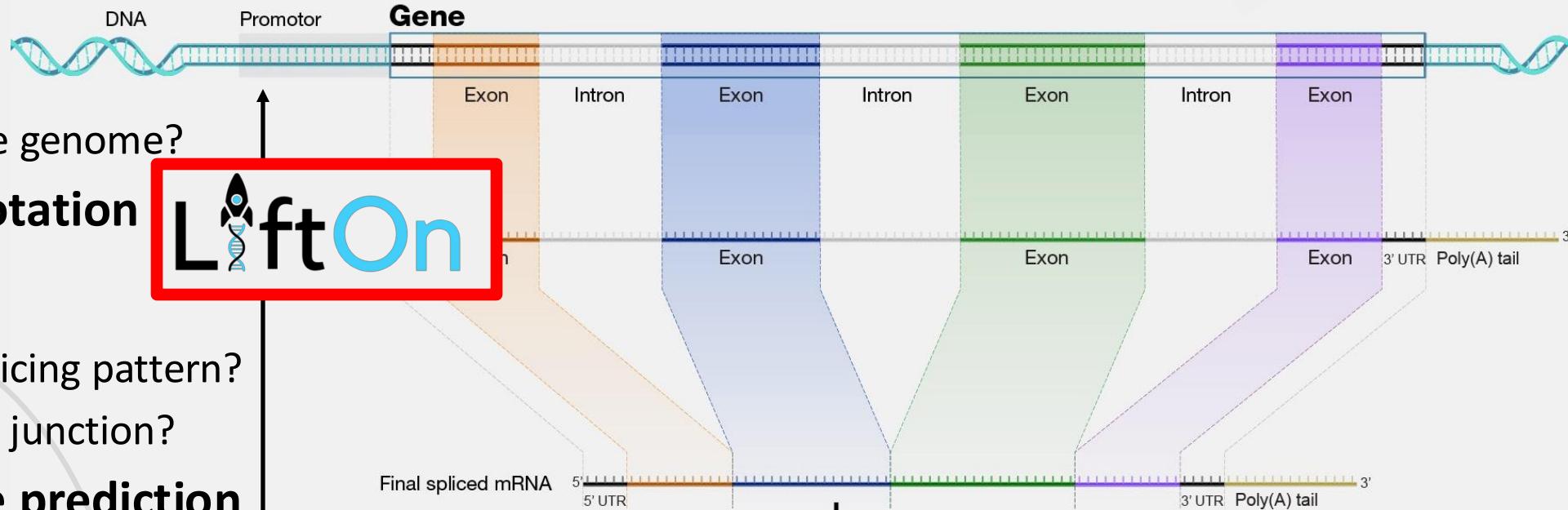
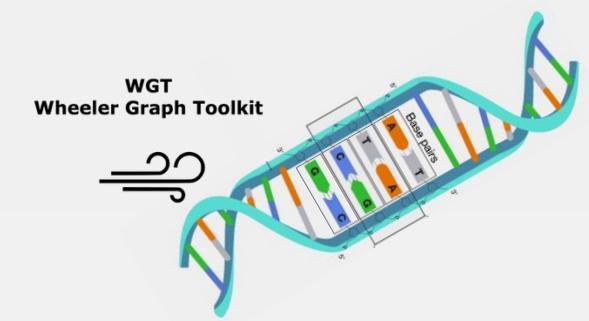
How do we assemble the complete 3-billion-nucleotide genome?

How can we efficiently represent multiple genomes for fast pattern matching?

Part I & II: Genome Assembly & Indexing

Han1

WGT



Where are the genes in the genome?

Part III: Genome Annotation

LiftOn

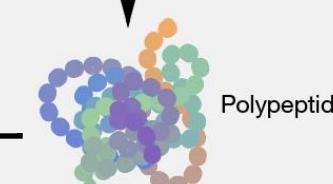
What are the canonical splicing pattern?

Alternative splicing? Splice junction?

Part IV & V: Splice site prediction

OpenSpliceAI

SPLICE



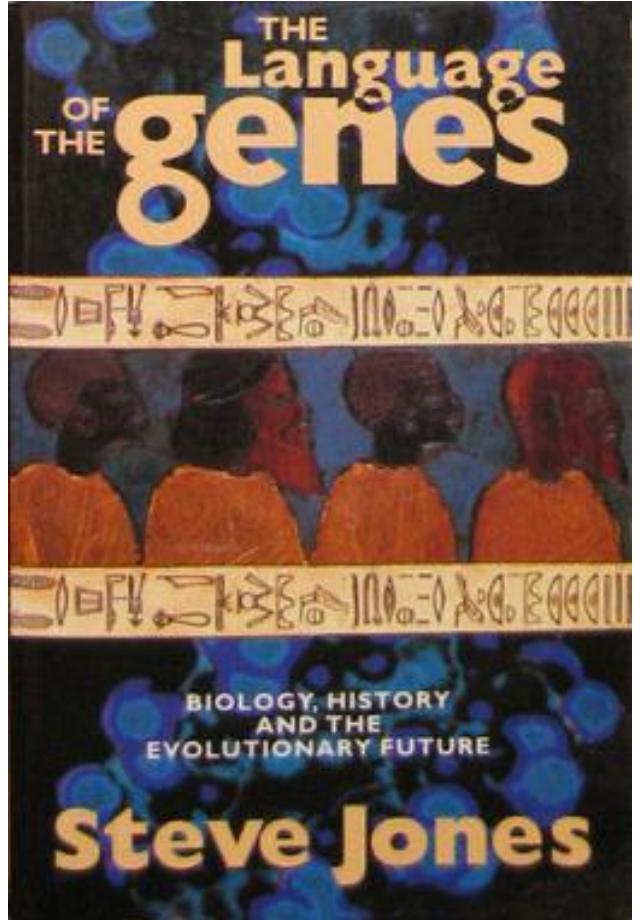
Can we predict gene expression by learning the regulatory grammar in the genome?

Part VI: Shorkie. RNA-Seq coverage prediction

Shorkie



Conclusions: Language of genomes



Steve Jones, 1993

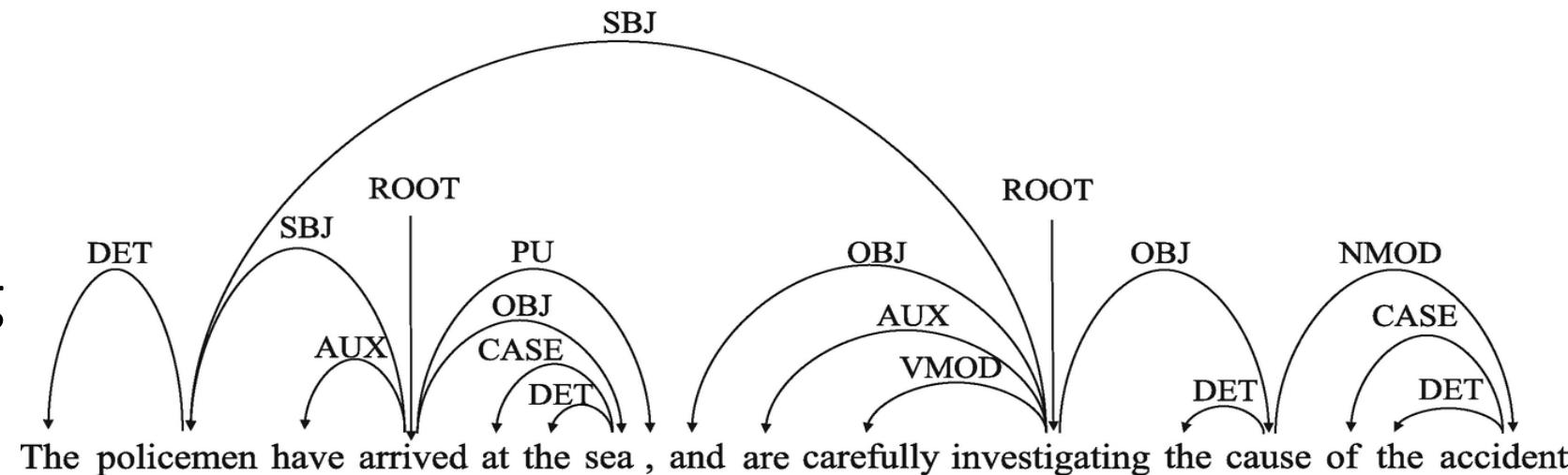
Codons as “words” and genes as “sentences”.

Zhou et al., 2023; Hwang et al., 2024; Sanabria et al., 2024; Nguyen et al., 2024; Theodoris, 2024; Dalla-Torre et al., 2025

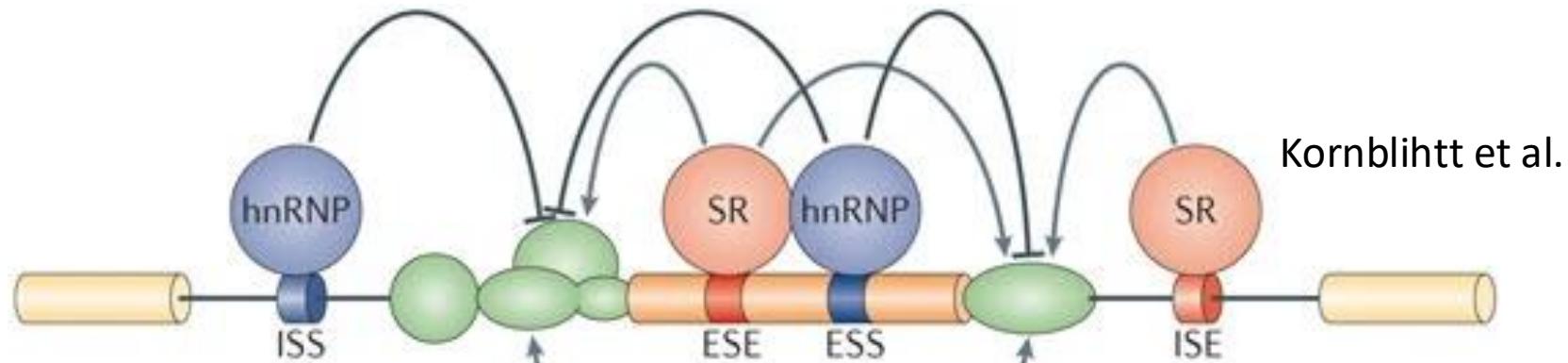
Recurring k-mer motifs as discrete tokens,
modeling long-range dependencies as grammatical rules.

Conclusions: Language of genomes

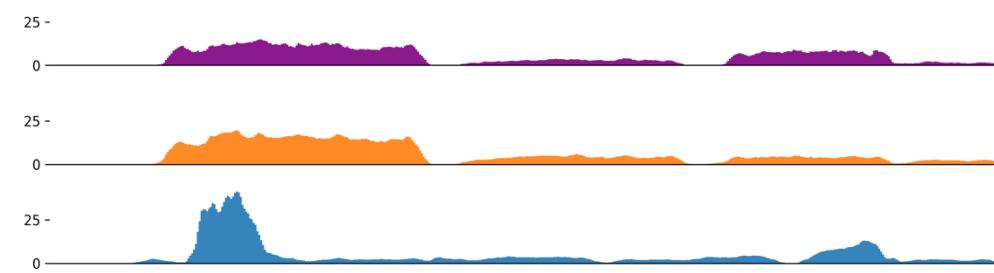
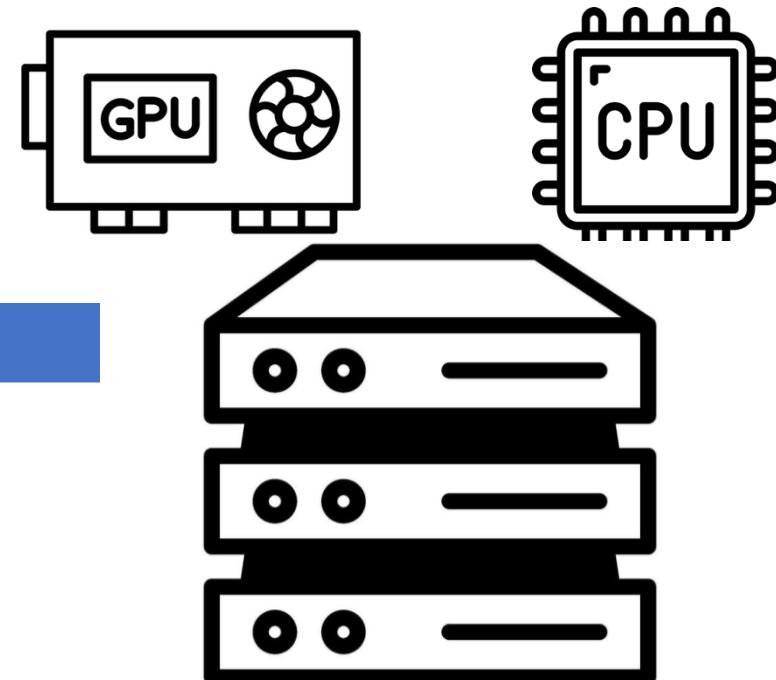
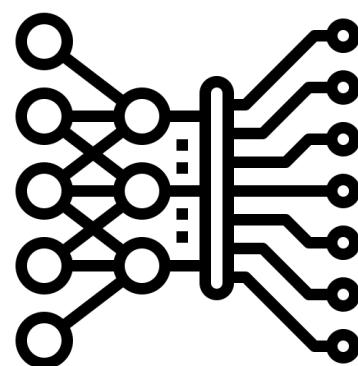
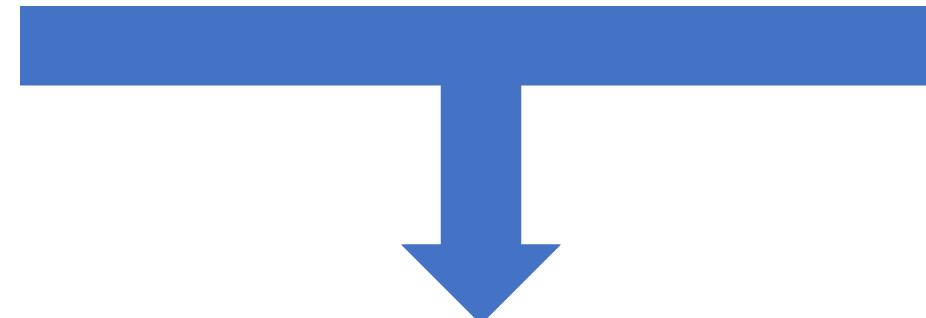
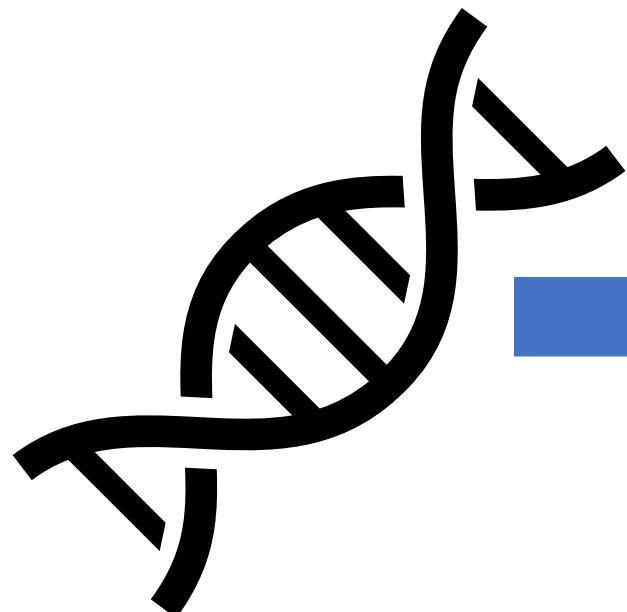
Natural Language Processing



DNA gene regulation



Conclusions: Language of genomes



Committee members



Steven Salzberg



Mihaela Pertea



Ben Langmead



David Kelley



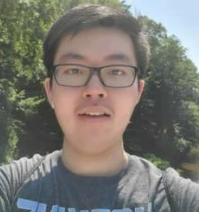
Anqi Liu

Best labmates



Salzberg Lab

Dr. Steven Salzberg

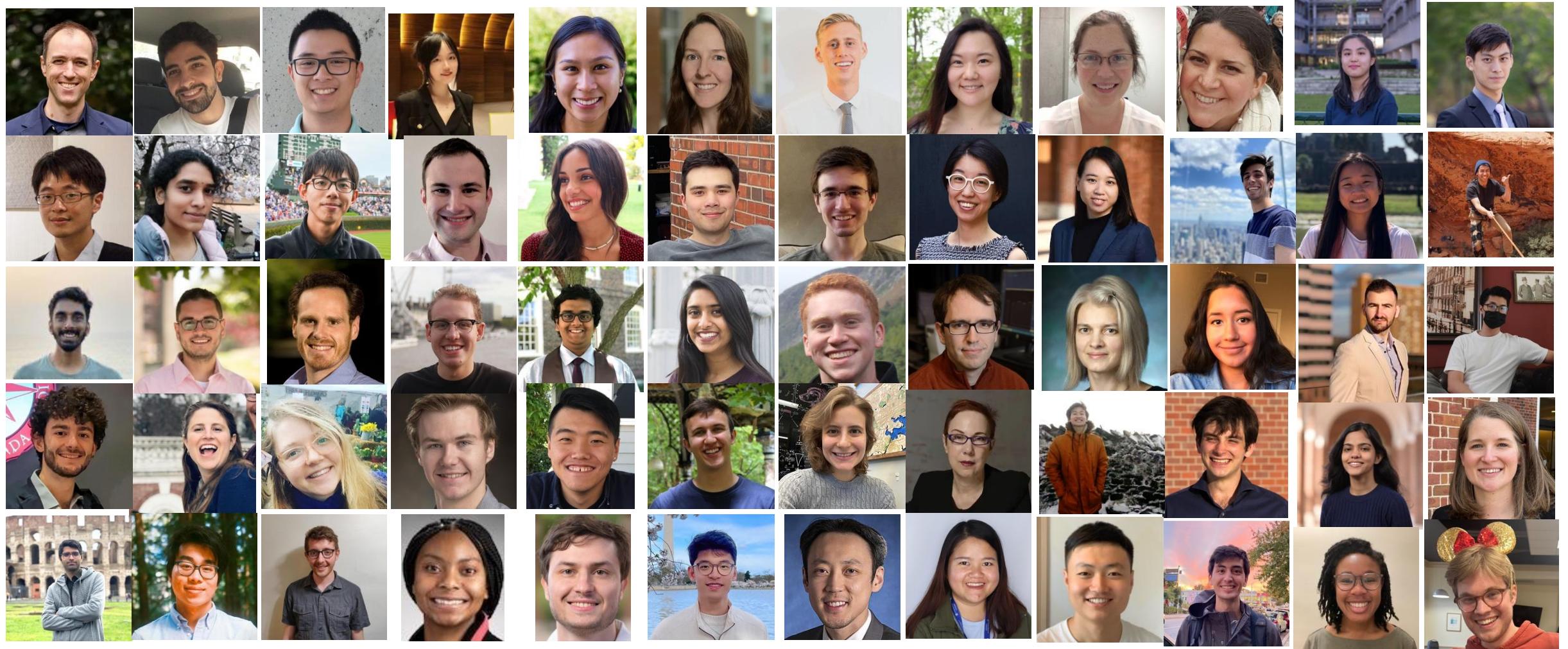


Pertea Lab

Dr. Mihaela Pertea



Hopkins Genomics + CS Friends





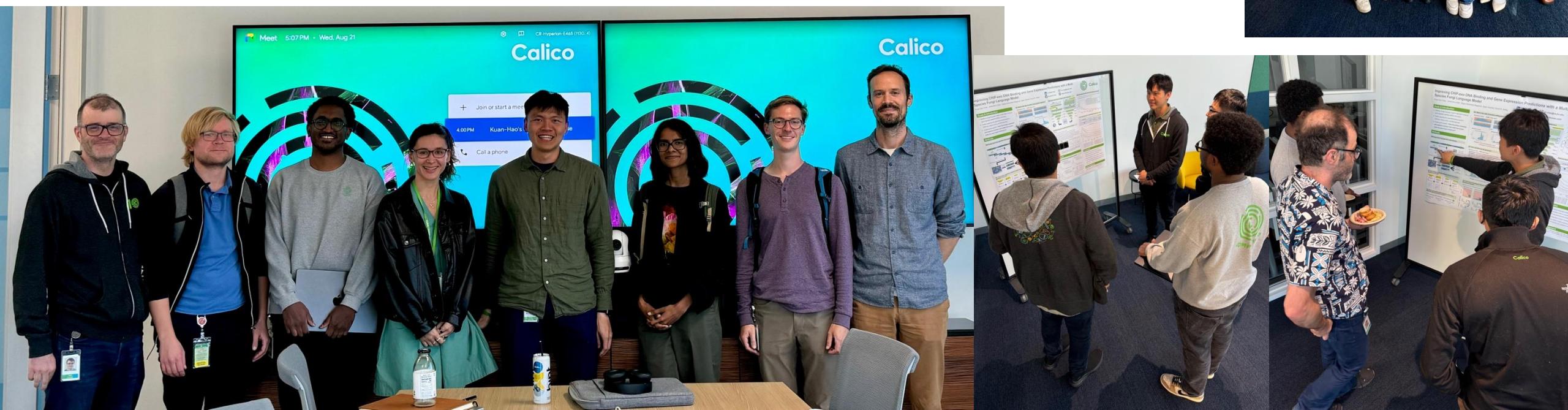
Calico Friends (2024 Summer)



Dr. David Kelley



Dr. Johannes Linder Majed Magzoub







Quarter-finals (Spring 2025)



JHU X UMD Sport Day



AC Malone (Sum 2023)



CHONVIATION

Family



Playoff (Spring 2024)



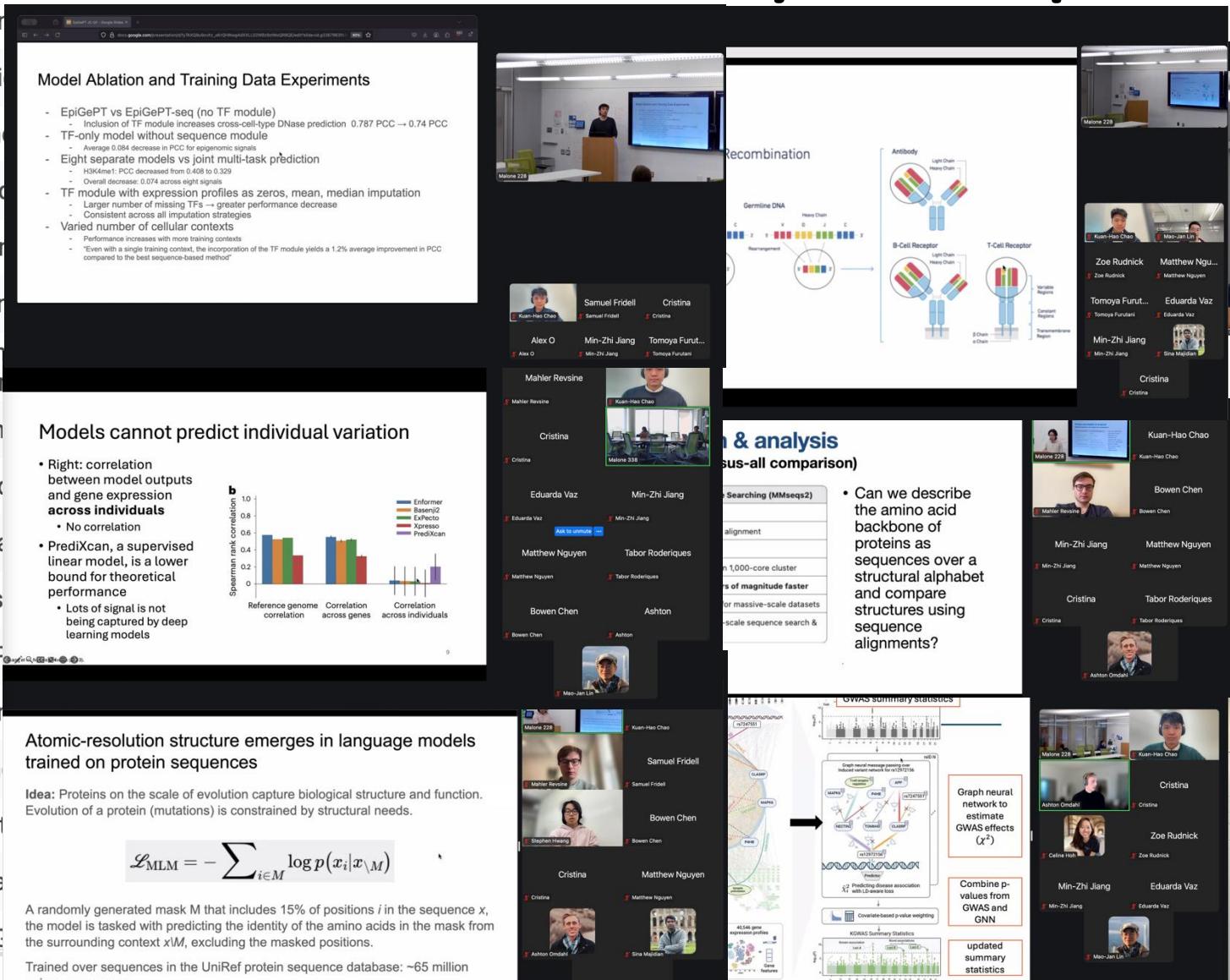
Softball (Sum 2023)



10/22	Kuan-Hao Chao	Predicting RNA- J
11/05	Mahler Revsine	HyenaDNA: Lon
11/19	Cristina Martin Linares	Machine-guided
12/03	Eduarda Vaz	Effective gene e
1/7	Gus Fridell	Applying
01/21	Stephen Hwang	Evoluti
02/11	Celine Hoh	A found
02/18	Ashton Omdahl	Small-d
03/04	Zoe Rudnick	Seque
03/25	Bowen Chen	Fast ar
04/08	Mahler Revsine	Person
04/22	Kuan-Hao Chao	Design
05/06	Cristina Martin Linares	Toward
05/27	Bohan Ni	A scal
06/03	Mao-Jan Lin	Diseas
06/17	Gus Fridell	EpiGeF
07/01	Celine Hoh	The lar
07/22	Eduarda Vaz	Param
07/29	Jake Galvin	Predict
08/12	Zoe Rudnick	Accura
08/26	Bowen Chen	AlphaC

JHU Deep Learning in Genomics Study Group

Mahler Revsine



Family



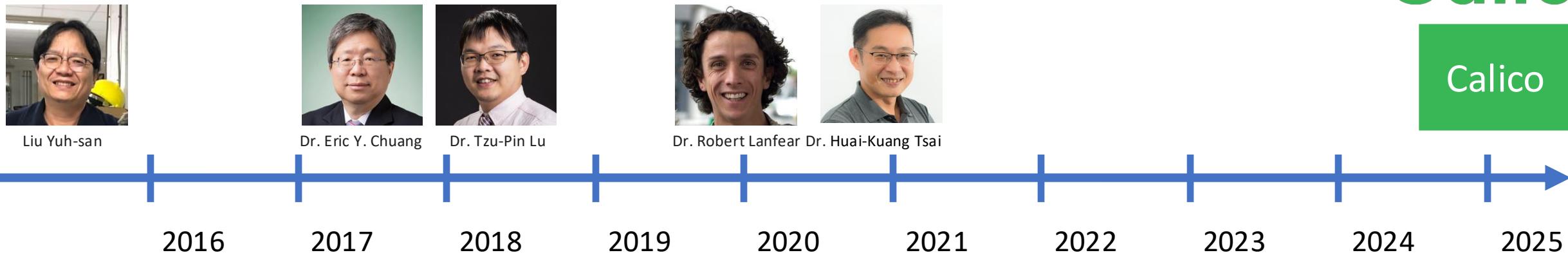
My research path

Next ...



Dr. David R Kelley Dr. Johannes Linder Majed Mohamed Magzoub

Calico



Taipei Municipal Chien Kuo High School



國立臺灣大學

National
Taiwan
University



Australian National University



Academia Sinica



Dr. Steven L Salzberg Dr. Mihaela Pertea Dr. Ben Langmead Dr. Anqi Liu



■ Next: Foster City, CA



Dr. Kyle Farh

Dr Kishore Jaganathan

Senior Deep Learning Scientist

@ Illumina AI Lab

illumina®