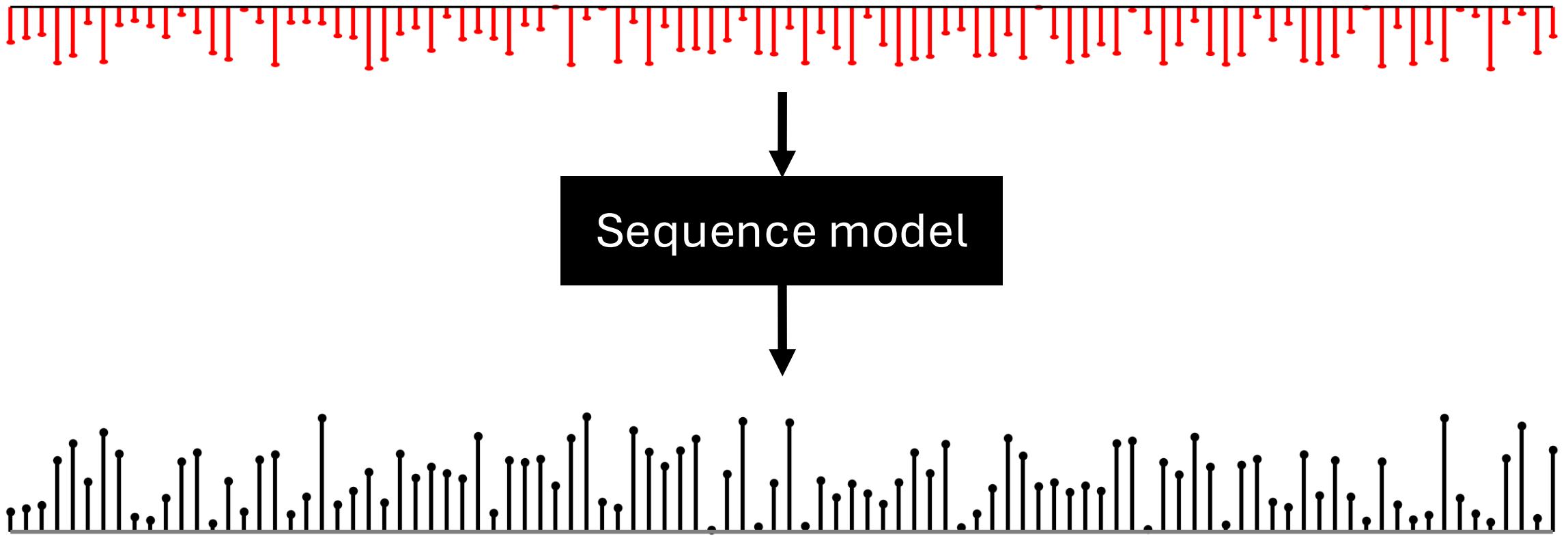




Unifying ChIP-exo DNA-Binding and RNA-Seq Coverage Predictions with a Multi-Species Fungal Language Model

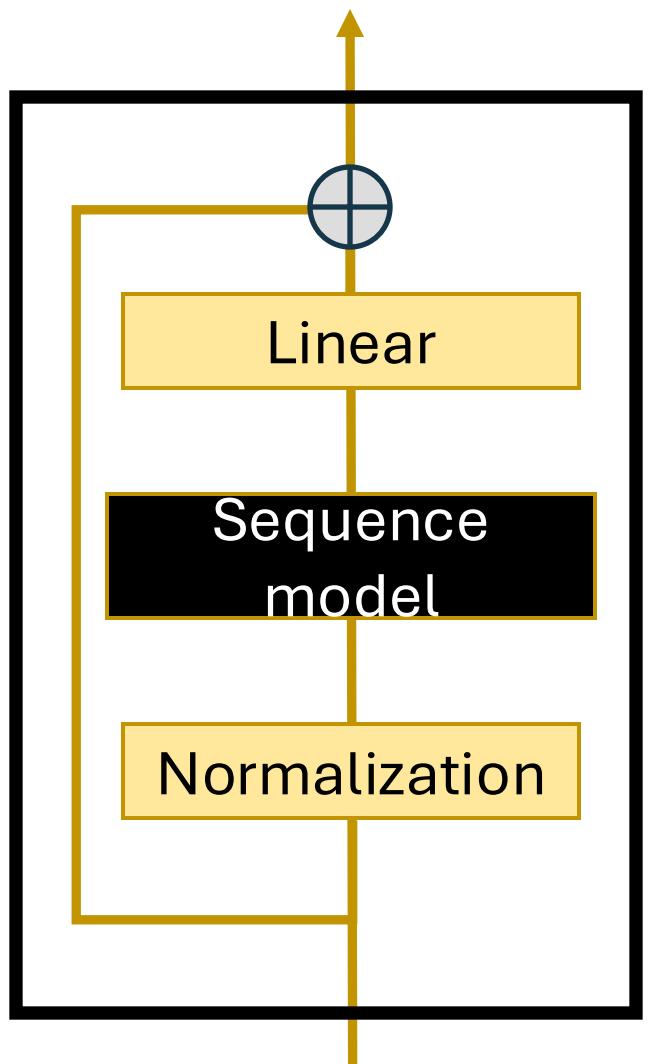
Kuan-Hao Chao

2025.01.22



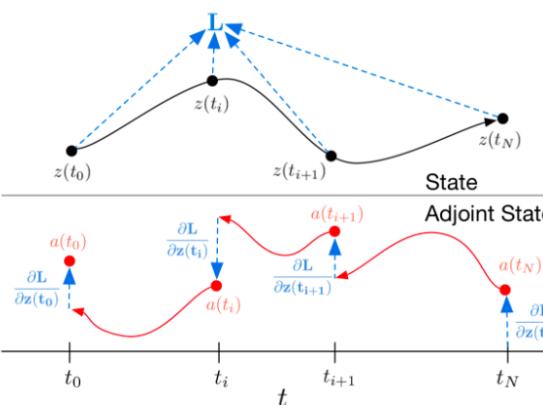
Sequence models map a sequence to a sequence

(batch, length, dim)



(batch, length, dim)

Neural ODEs



RNN

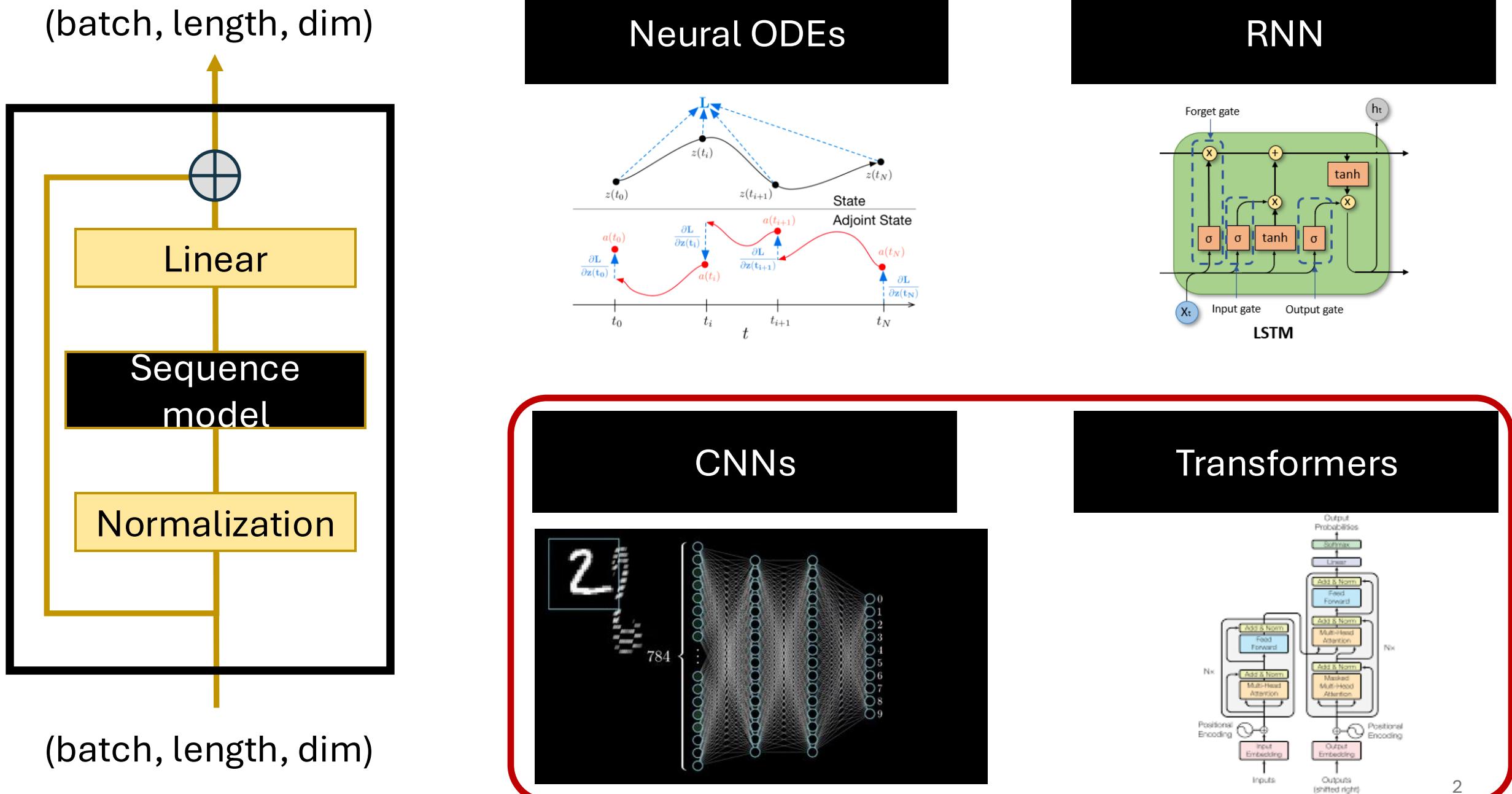
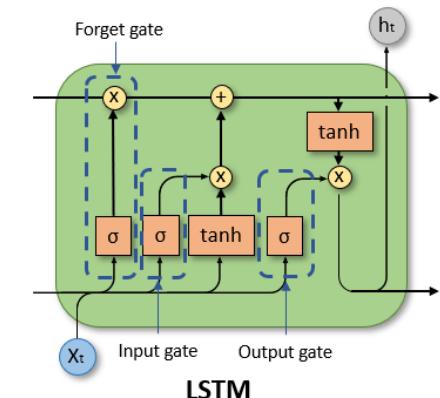
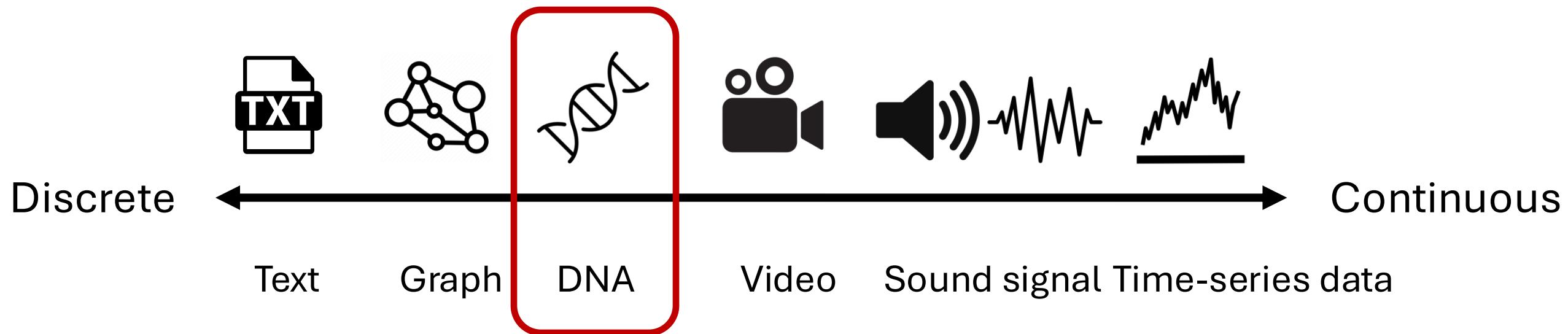


Figure 1: The Transformer - model architecture.



Spectrum of Sequential Data



Deep learning-based DNA sequence model

nature methods

[View all journals](#) | [Search](#)

Explore content

Troyanskaya Lab Princeton

DeepSEA
2015

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

Jian Zhou & Olga G Troyanskaya [✉](#)

CSH PRESS GENOME
RESEARCH

HOME | ABOUT | ARCHIVE | SUBMIT | SUBSCRIBE | ADVERTISE | AUTHOR

Institution: MILTON S EISENHOWER LIBRARY Sign In

Calico

Basset
2016

Basset: learning the regulatory
accessible genome with deep co
neural networks

David R. Kelley¹, Jasper Snoek² and John L. Rinn¹

Cell

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

Article

Predicting Splicing from Pri
with Deep Learning

Kishore Jagannathan^{1,6}, Sofia Kyriazopoulou Pangiatopoul¹,
Siavash Fazel Darbandi², David Knowles³, Yang Li¹, Jack
Wenwu Cui¹, Grace B. Schwartz², Eric D. Chow³, Efstratios
Serafim Batzoglou¹, Stephan J. Sanders², Kyle Kai-How Forh^{1,7} [✉](#) [✉](#)

Illumina

SpliceAI
2019

Agarwal 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 mar

Vikram Agarwal^{1,2} [✉](#) and David R. Kelley^{1*}

Calico

Saluki
2022

nature biotechnology

Explore content | About the journal | Publish with us

FUToronto

DeepBind
2015

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

Babak Alipanahi, Andrew Delong, Matthew T Weirauch & Brendan J Frey [✉](#)

CSH PRESS GENOME
RESEARCH

HOME | ABOUT | ARCHIVE | SUBMIT | SUBSCRIBE | ADVERTISE | AUTHOR

Institution: MILTON S EISENHOWER LIBRARY Sign In

Calico

Basenji
2018

Sequential regulatory activity
across chromosomes with co
neural networks

David R. Kelley¹, Yakir A. Reshef², Maxwell Bileschi³, Da
Cory Y. McLean³ and Jasper Snoek³

Calico

Akita
2020

nature

methods

Predicting 3D genome folding from
with Akita

Geoff Fudenberg [✉](#), David R. Kelley [✉](#) and Katherine S. Pollard [✉](#)

nature methods

Explore content | About the journal | Publish with us

Calico

scBasset
2022

Bioinformatics

Gifford Lab MIT

DNA-TF binding
2016

Article 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

ExPecto
2018

Article | Published: 16 July 2018
Deep learning sequence-based
variant effects on expression an

Jian Zhou, Chandra L. Theesfeld, Kevin Yao, Kathleen M. C

ARTICLES

<https://doi.org/10.1038/ng.41592-021-01252-x>

OPEN

Effective gene express
sequence by integratin

Ziga Avsec [✉](#), Vikram Agarwal^{2,4}, Daniel Vis
Agnieszka Grabska-Barwińska¹, Kyle R. Taylor
and David R. Kelley [✉](#) [✉](#)

DeepMind + Calico

Enformer
2021

nature genetics

Article

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

Johannes Linder [✉](#), Divyanshi Srivastava, Han Yuan, Vikram Agarwal [✉](#)

Calico

Borzoi
2025

Introduction

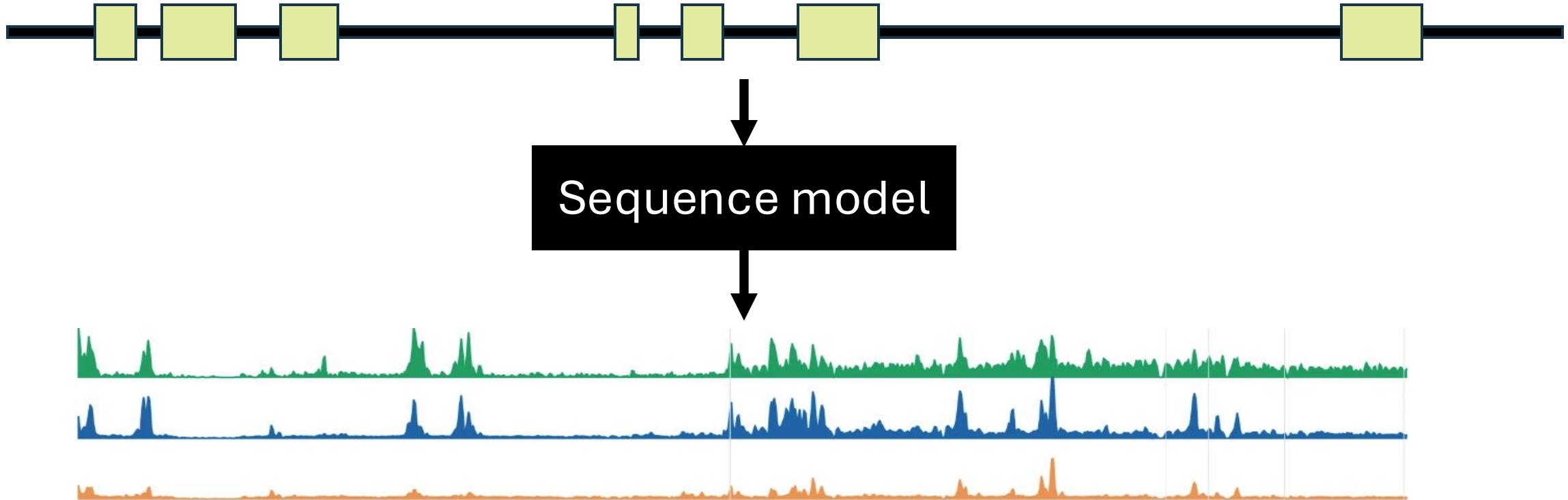
Pre-training LM

Fine-tuning LM

Applications & Conclusions

Input: DNA sequences
 $(N_{batch_size} * L_{input} * 4)$

DeepMind + Calico	Calico
Enformer 2021	Borzoi 2025



Output: Genomics tracks
 $(N_{batch_size} * L_{output} * T_{track_number})$

Supervised learning

Input: DNA sequences

Output: Genomics tracks



Stage 1 Self-supervised learning

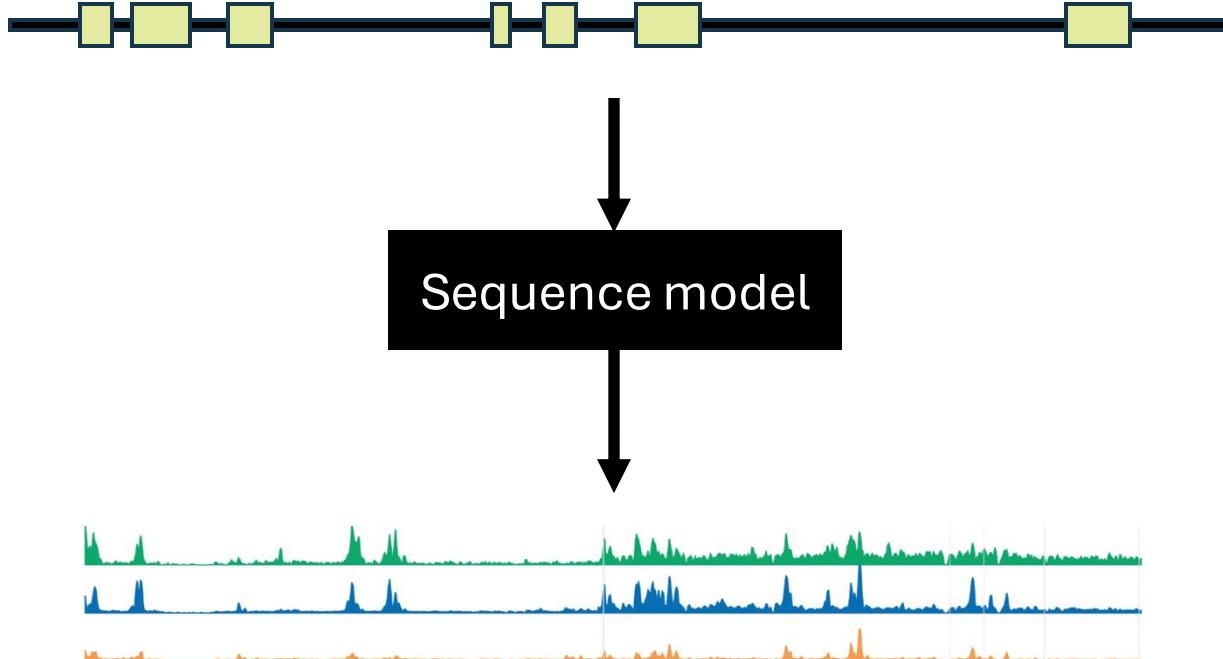


Language model (LM) /
Foundation model

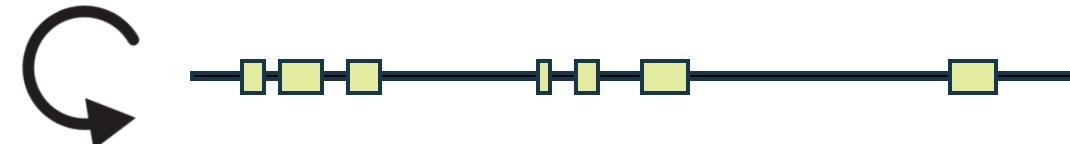
Stage 2 Fine-tuning LM



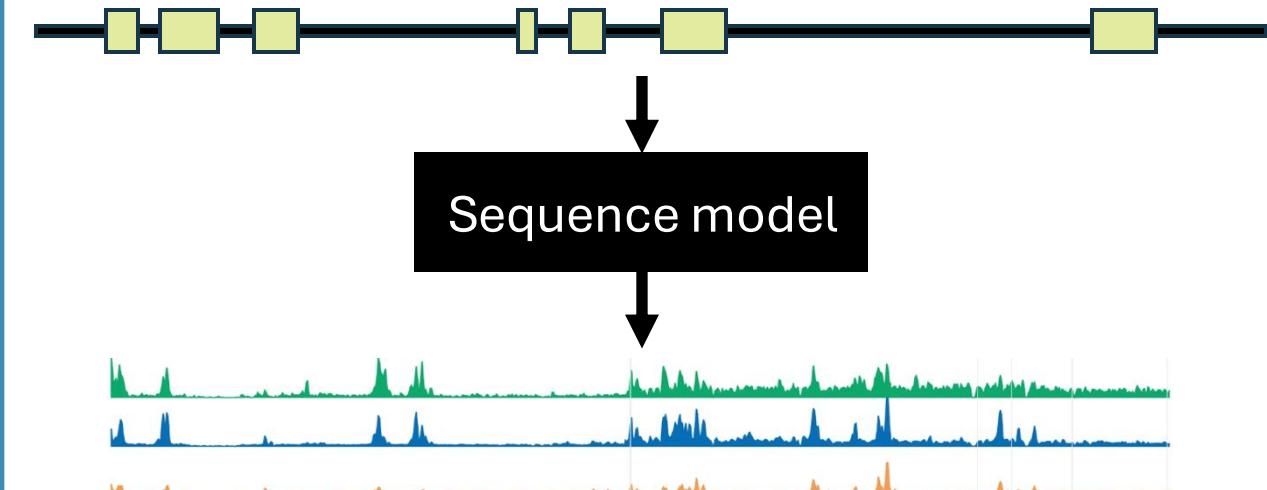
Supervised learning



Stage 1 Self-supervised learning



Stage 2 Fine-tuning LM



Protein Language model

RESEARCH ARTICLE | BIOLOGICAL SCIENCES | 8
Biological structure and function from scaling up to one million protein structures
Facebook
First LM attempt
PNAS 2020
Alexander Rives, Joshua Meier, Tom Sercu, +7 others, and Rob Fergus
Authors info & Affiliations

Article | Open access | Published: 27 July 2022
ProtGPT2 is a deep unsupervised learning model for protein design
Höcker Lab
ProtGPT2
Nat Commun 2022
Noelia Ferruz, Steffen Schmid
Nature Communications 13, Article number: 1022 (2022)
Article | Published: 03 October 2022
Single-sequence protein structure prediction using a language model
Harvard + Columbia
AminoBERT
Nat Biotechnol 2022
Ratul Chowdhury, Nazim Belghazi, Anant Kharkar, Koushik Roy, Zhenhua Zhang, George M. Church, et al.
Primer | Published: 15 February 2024
Designing proteins with AI
Profluent
PLM
Nat Biotechnol 2024
Jeffrey A. Ruffolo & Ali Madani
Nature Biotechnology 42, 200–202 (2024) | [View this article online](#)

arXiv:2006.15222 (cs)
[Submitted on 26 Jun 2020 (v1), last revised 28 Mar 2021]
BERTology Meets Biology: How Attention in Protein Language Models Works
Salesforce + UIUC
BERTology
ICLR 2021
Jesse Vig, Ali Madani, Lav R. Varshney, Caiming Xiong, Richard Socher, James Zou, et al.

Article | Published: 26 January 2023
Large language models generate functional protein sequences across diverse families
Salesforce
ProGen
Nat Biotechnol 2023
Ali Madani, Ben Krause, Daniel Mohr, James M. Holton, Jason Sun, Richard Socher, James Zou, et al.

SYNTHESIS · Volume 12, Issue 6, P654-669.E3, June 16, 2020 | **MIT**
Open Access
Learning the protein structure, and function
Cell Systems 2021
Tristan Bepler, Bojan Filipovic, +12 others

ARTICLE · Volume 14, Issue 11, November 2023 | **Salesforce + Profluent**
ProGen2: Exploring the limits of protein language models
Erik Nijkamp, Jeffrey A. Ruffolo, Ali Madani, et al.

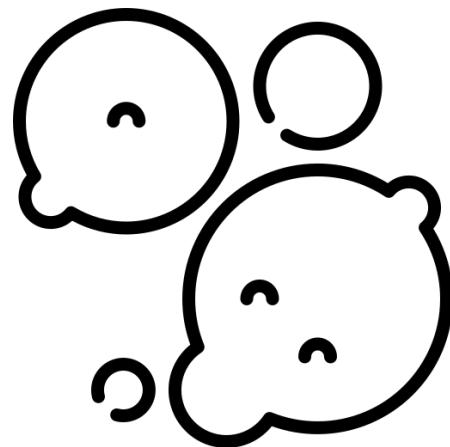
~ 6000 citations

DNA Language model

<p>DNABERT: pre-trained Bidirectional Encoder Representations from Transformers model DNA-language in genomic context</p> <p>Yanrong Ji, Zihhan Zhou, Han Liu</p>	<p>SBU</p> <p>DNABERT Bioinformatics 2021</p>	<p>DNABERT-2: Efficient Foundation Model Benchmark For Multi-Species</p> <p>Zhihan Zhou, Yanrong Ji, Weijian Li, Pratik Dandekar</p>	<p>SBU + NU</p> <p>DNABERT-2 ICLR 2024</p>	<p>Sequence-based molecular modeling</p> <p>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 Ré</p> <p>doi: https://doi.org/10.1101/2024.02.27.582234</p>	<p>Stanford + Arc Inst + TogetherAI</p> <p>Evo bioRxiv 2024</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 Authors & Affiliations</p>	<p>UC Berkeley</p> <p>GPN PNAS 2023</p>	<p>Article Open access Published: 23 July 2024</p> <p>DNA language models predict gene context in the human genome</p> <p>Melissa Sanabria, Jonas Hirsch, Pierre Lefebvre, and others Authors & Affiliations</p>	<p>TUD</p> <p>GROVER Nat Mach Intell 2024</p>	<p>[Submitted on 5 Mar 2024]</p> <p>Caduceus: Bi-Directional Range DNA Sequence Model</p> <p>Yair Schiff, Chia-Hsiang Kao, Aaron Gokaslan, Tripti Agarwal, and others Authors & Affiliations</p>	<p>Cornell + Princeton + CMU</p> <p>Caduceus ICML 2024</p>
<p>The Nucleotide Transformer: Foundation Models for</p> <p>Hugo Dalla Favera, Adam Henne, Bernardo P. Ribeiro, and Marie Lopez</p>	<p>InstaDeep + Nvidia + TUM</p> <p>Nucleotide Transformer bioRxiv 2023</p>	<p>Article Open access Published: 03 April 2024</p> <p>Genomic language model predicts gene regulation and function</p> <p>Yunha Hwang, Andre L. Cornman, Elizabeth J. C. Gaglio, and others Authors & Affiliations</p>	<p>Harvard + MIT</p> <p>Genomic LM Nat Commun 2024</p>	<p>Cross-specific prediction of nucleotide dependency</p> <p>Jingjing Zhai, Aaron Gokaslan, Yair Schiff, and others Authors & Affiliations</p> <p>doi: https://doi.org/10.1101/2024.06.06.536712</p>	<p>Cornell + USDA-ARS + Simons</p> <p>PlantCaduceus bioRxiv 2024</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</p> <p>Eric Nguyen, Michael Poli, Marjan Faizi, Armin Thomas, Aman Patel, Clayton Rabideau, Stefano Massaroli, Yoshua Bengio, and others Authors & Affiliations</p>	<p>Stanford</p> <p>HyenaDNA NeurIPS 2023</p>	<p>Research Open access Published: 02 April 2024</p> <p>Species-aware DNA language models capture regulatory elements across species</p> <p>Alexander Karolchik, Julien Gagneur, and others Authors & Affiliations</p>	<p>TUM</p> <p>Species-aware DNA LM Genom Biol 2024</p>	<p>Nucleotide dependency analysis of DNA language models</p> <p>TUM</p>	<p>Nucleotide dependency bioRxiv 2024</p>

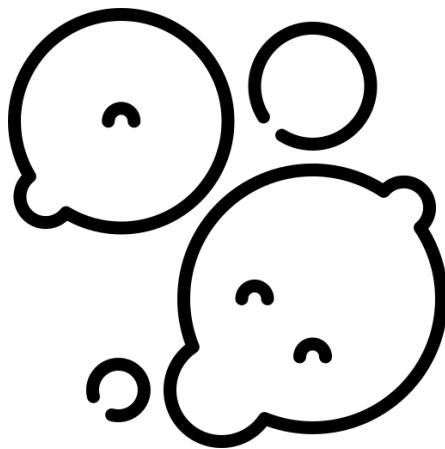
Study Goals

- Building a SOTA **yeast** gene expression model.
- **Part I:** Exploring DNA **Language Model (LM)**.
- **Part II:** Fine-tuning DNA LM to predict gene expression (RNA-Seq tracks).
 - Is ***self-supervised learning with fine-tuning*** better than ***training from scratch***?
- Understanding what models learn at each stage: (LM, fine-tuning, and scratch-training)
- **Part III:** Applications:
 - Predicting the influence of context and distal regulatory elements on gene expression.
 - Assessing the variant effects on eQTLs and negatively selected eQTLs



Why yeast?

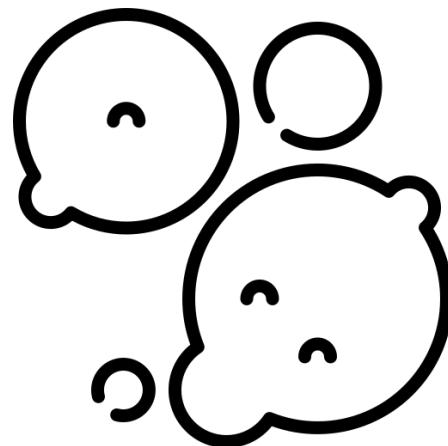
- Computation is expensive for human (human: 3B nt; yeast: 12 M nt)



Human + Mouse genomes = **5.7 B nt**

	1 sequence (All tracks)	All sequences (All tracks)	Practical runtime
Enformer	(# 128-resolution bin) × (# tracks) 896 × 6,956 6,232,576 float32 (4 bytes) ≈ 25 MB	$5.7B \div 196,608 \approx 28992$ $28992 * 2.5MB \approx \textcolor{red}{724.8 GB}$	1 model: 64 TPU v3 cores (16GB GPU memory each) (~ 3 days)
Borzoi	(# 32-resolution bin) × (# tracks) 16,384 × 10,219 167,428,096 float32 (4 bytes) ≈ 670 MB	$5.7B \div 524,288 \approx 10872$ $10872 * 670 MB \approx \textcolor{red}{7.28 TB}$	1 model: 2 A100 GPUs (40 GB GPU memory each) (~ 25 days)

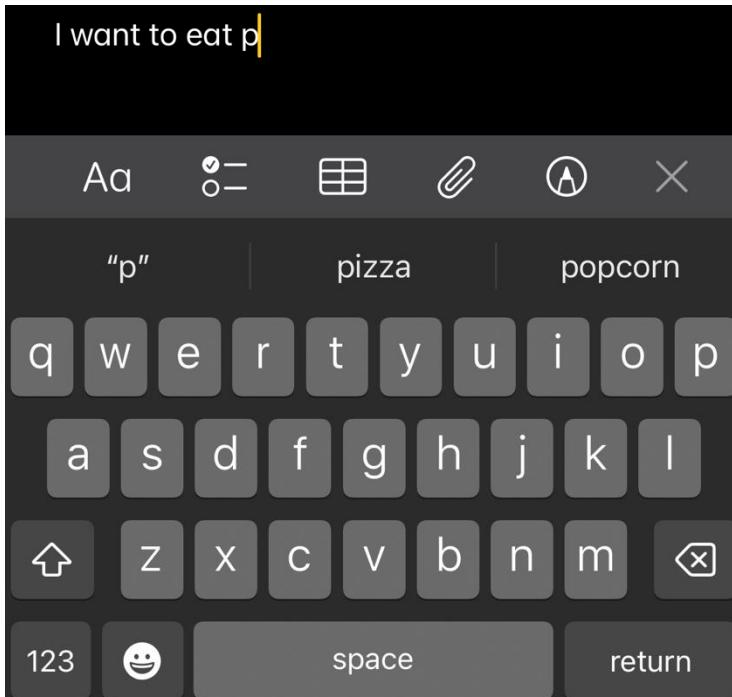
Why yeast ?



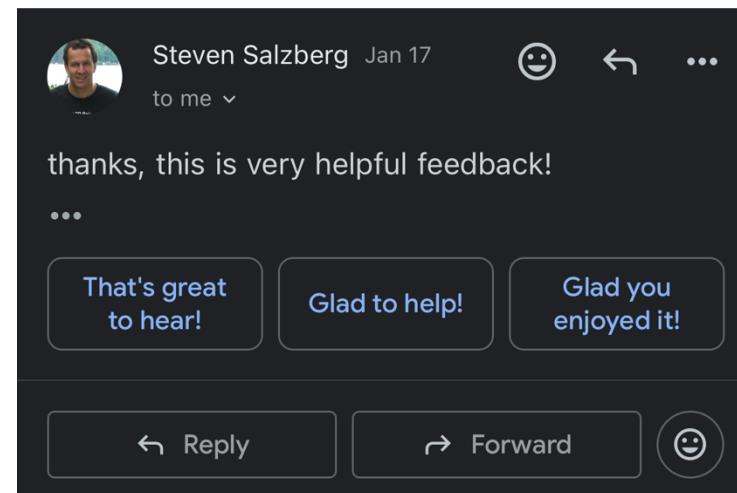
- Computation is expensive for human (human: 3B nt; yeast: 12 M nt)
- In human, we can't do large scale TF perturbation study
- Yeast is a great model organism to generate data & training models
 - Simple Eukaryotic Model: cost-effectiveness and scalability
 - Rapid Growth and Easy Culturing and Quick Lift Cycle
 - Genetic Manipulability
 - Well-Characterized Genome
 - Conserved Regulatory Mechanisms
 - Great species to study aging! (bud scar / cell size)

Language Model (LM) in our daily life

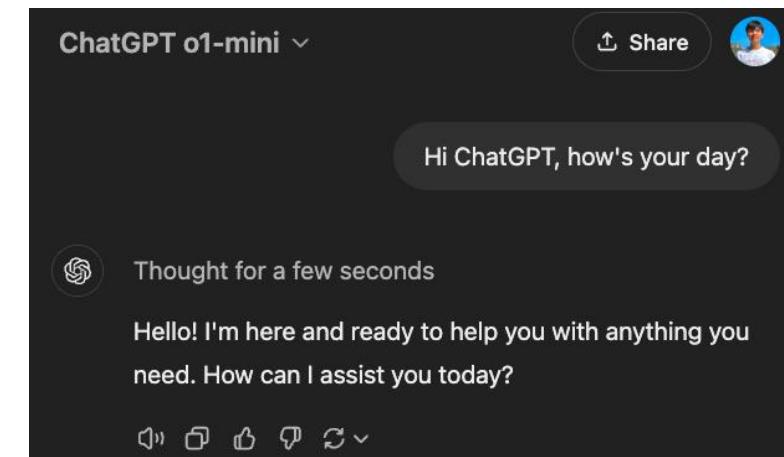
Next word prediction



Email suggested reply



ChatGPT



What is a **language model (LM)**?

$$P(X_t \mid X_1, X_2, \dots, X_{t-1})$$


Next word

Context or prefix

What is a **language model (LM)**?

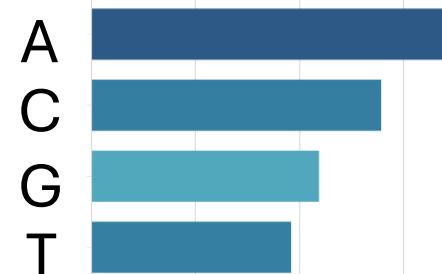
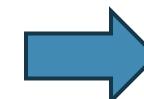
- Directly we train models on “**marginals**”
- We are implicitly learning the **full/joint distribution** of language.

$$P(\underbrace{X_t}_{\text{Next word}} | \underbrace{X_1, X_2, \dots, X_{t-1}}_{\text{Context or prefix}})$$

The chain rule:

$$P(X_1, \dots, X_t) = P(X_1) \prod_{i=1}^t P(X_i | X_1, X_2, \dots, X_{i-1})$$

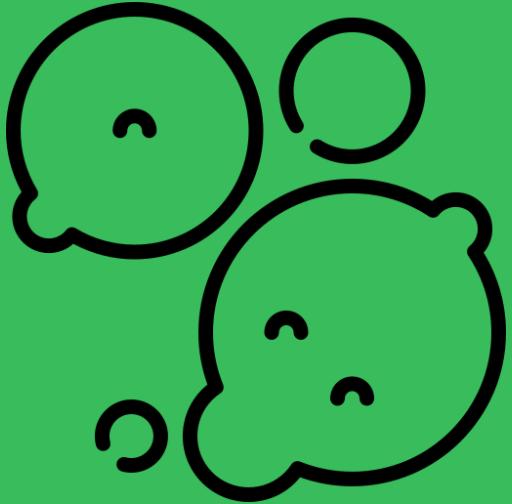
“A C T T A C T A G A [MASK]”



Language Modeling \triangleq learning prob distribution over language sequence.

Part I

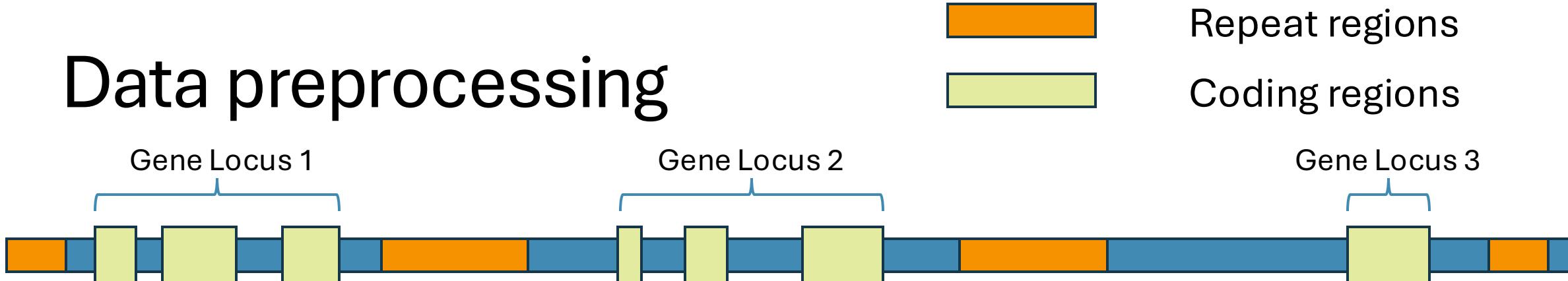
Fungal Language Model



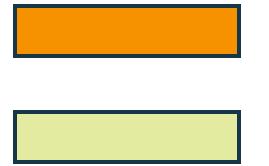
Summary

We need to understand the genome & carefully preprocess genome to correctly train a Fungi language model.

Data preprocessing



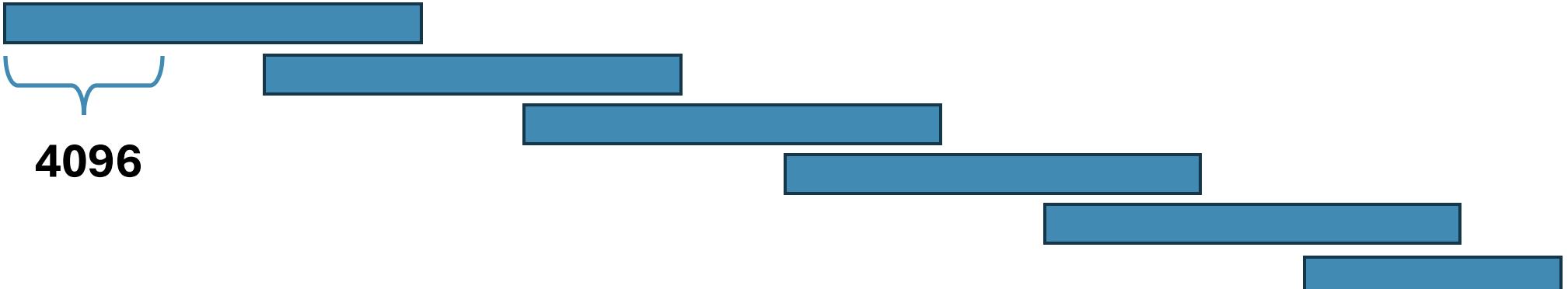
Data preprocessing



16384



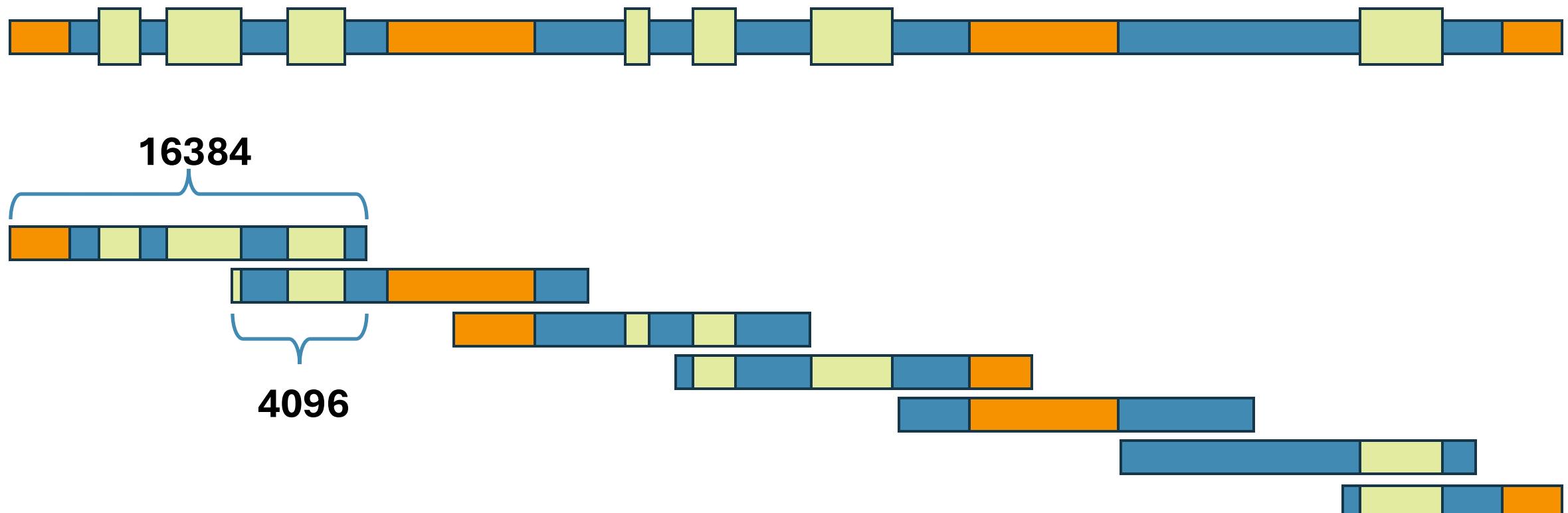
4096



~ 8 genes per window

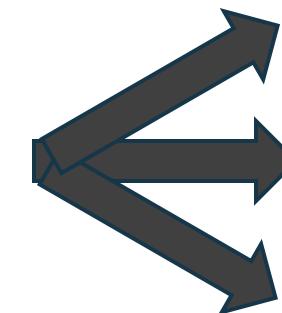
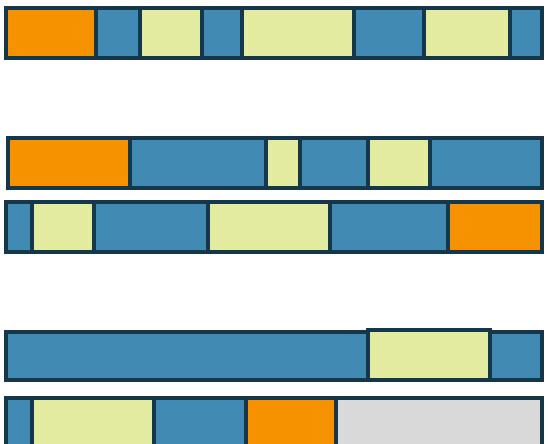
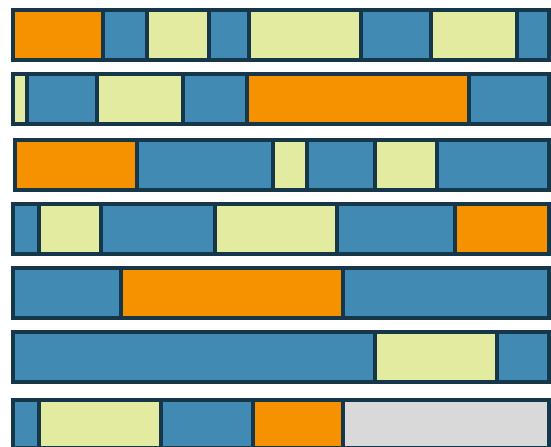
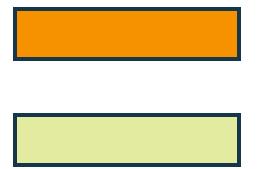
Data preprocessing

Repeat regions
Coding regions



~ 8 genes per window

Data preprocessing

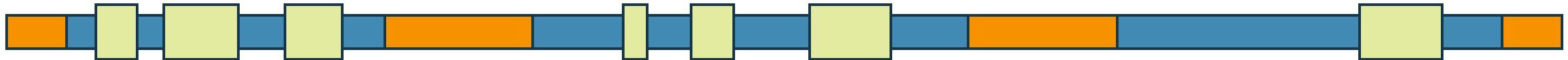


Testing
(chrXII, chrXIV, chrXVI)

Why building a Fungal Language Model (LM)?

- Yeast genome is small. 12Mbps.
- Thousands of fungal genomes with high quality. No supervised measurements
- Language model pre-training on all available genomes followed by transfer learning to the smaller yeast genome.



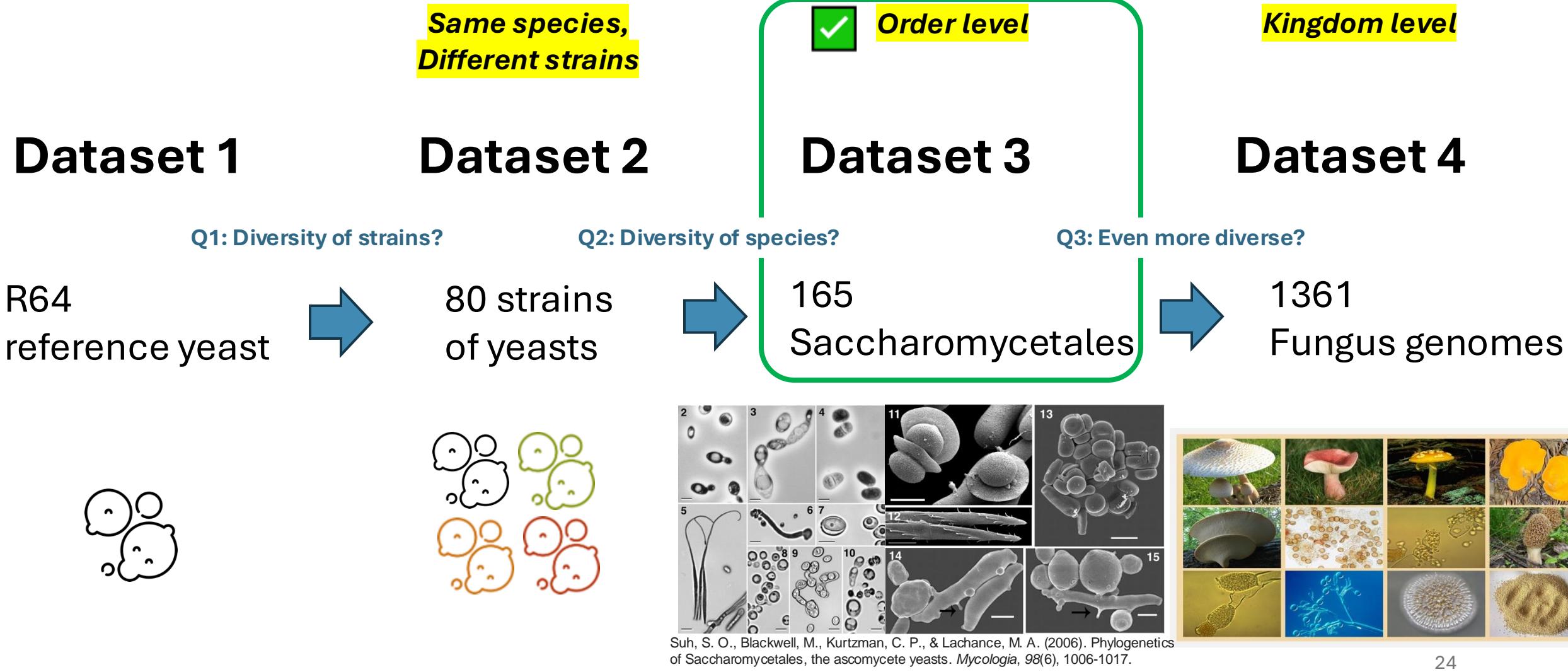


Q1: To what evolutionary distance should we include in our LM?



Selected Genomes for LM

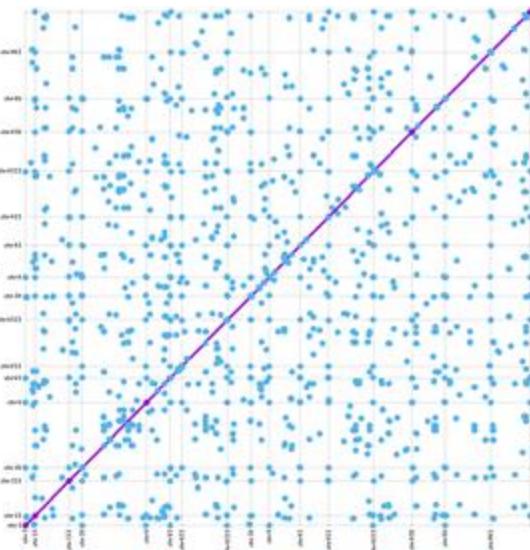
Fungi diverged from other life around 1.5 billion years ago



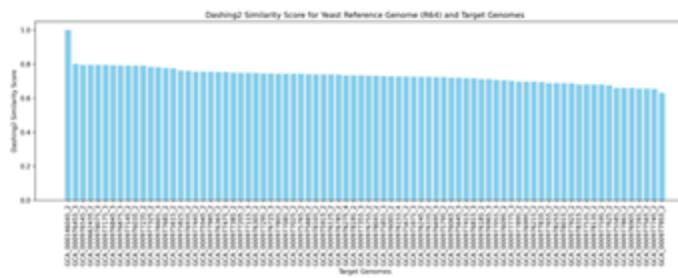
Genome distance evaluation

R64 Reference Yeast

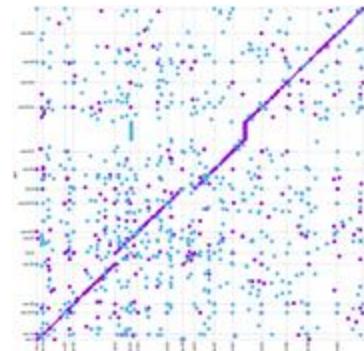
Mummerplot



Dashing2
similarity

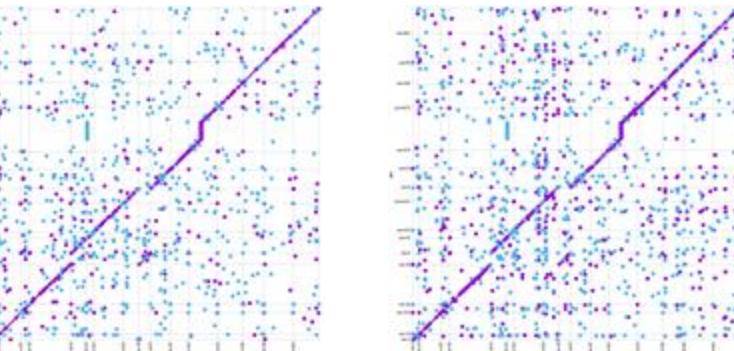
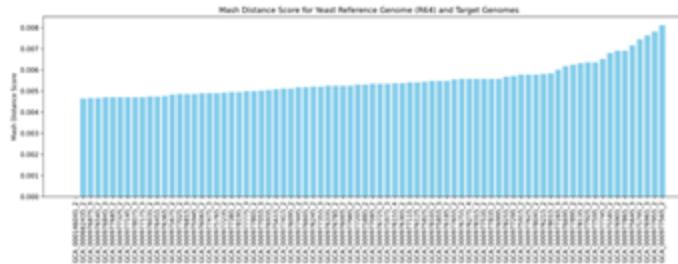


Mummerplot



80 strains of yeasts

Mash
distance

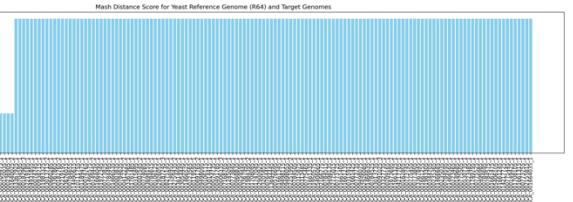


✓ 165 Saccharomycetales

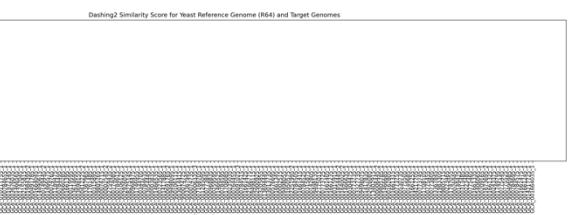
Mummerplot

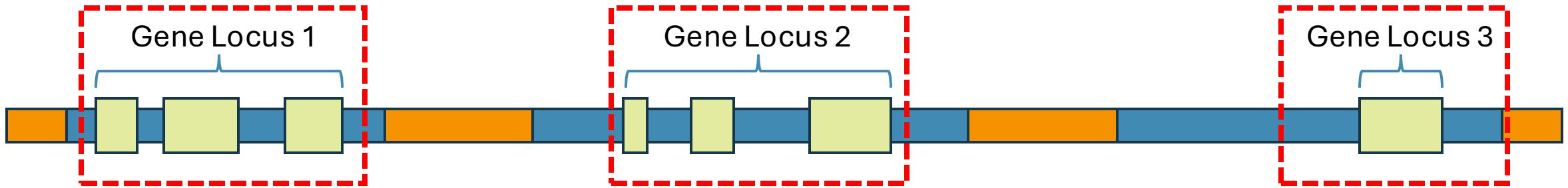


Mash
distance



Dashing2
similarity





Q2: How many genes per window?

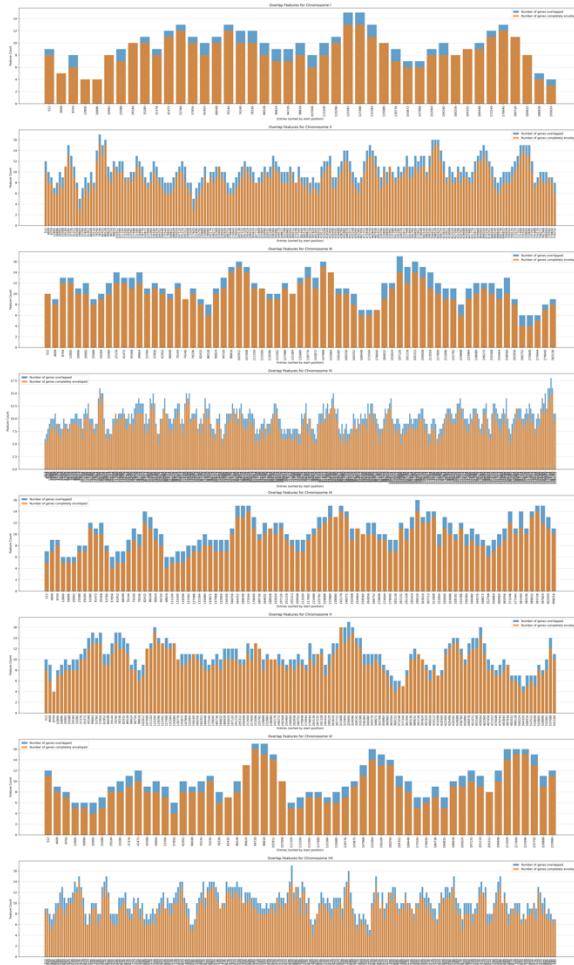
Q3: What is the quality of the annotation?

Q4: What is the protein-coding region ratio?



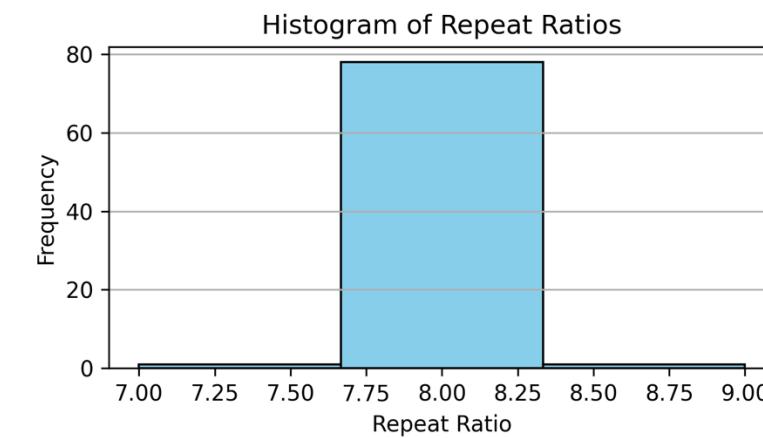
Genome evaluation – # genes per 16K window

R64 Reference Yeast

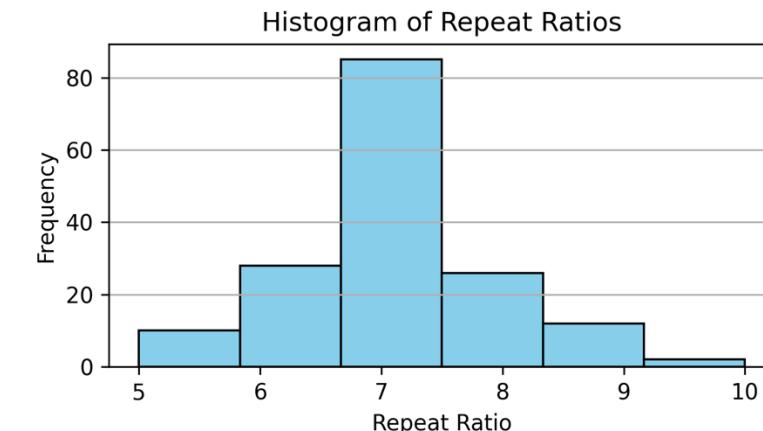


Median: 9.0; Mean: 8.98

80 strains of yeasts

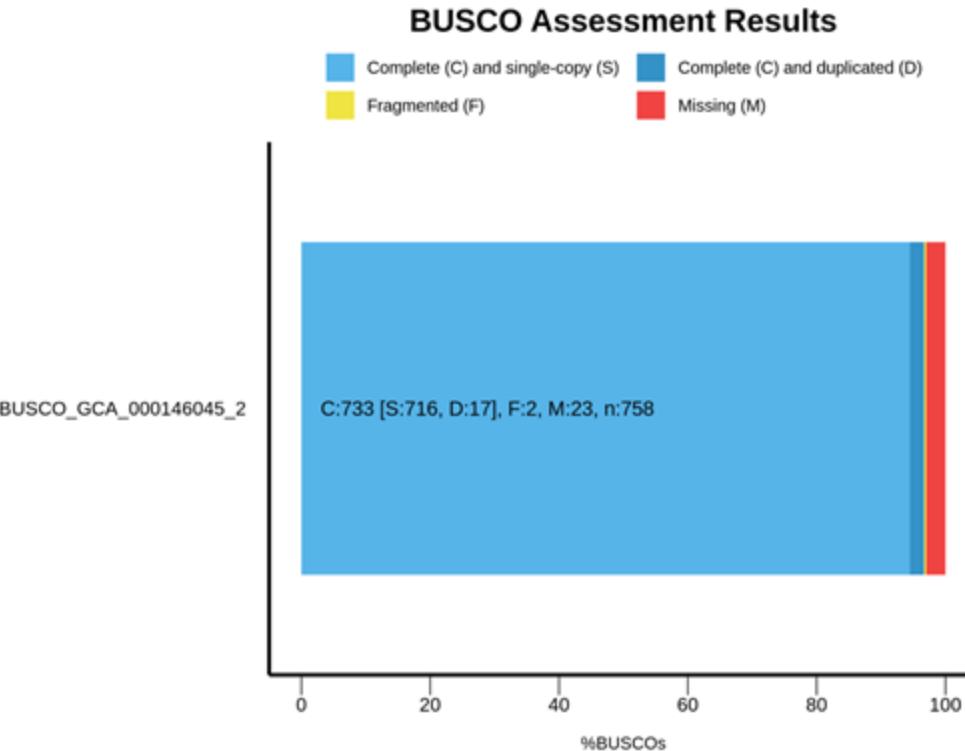


165 Sachramonycetales



Genome annotation completeness evaluation

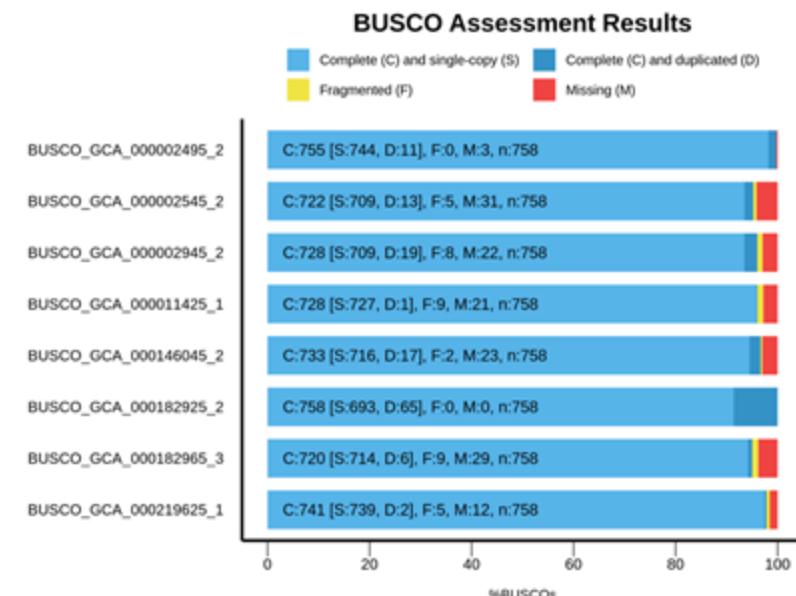
R64 Reference Yeast



80 strains of yeasts



165 Sachramonycetales



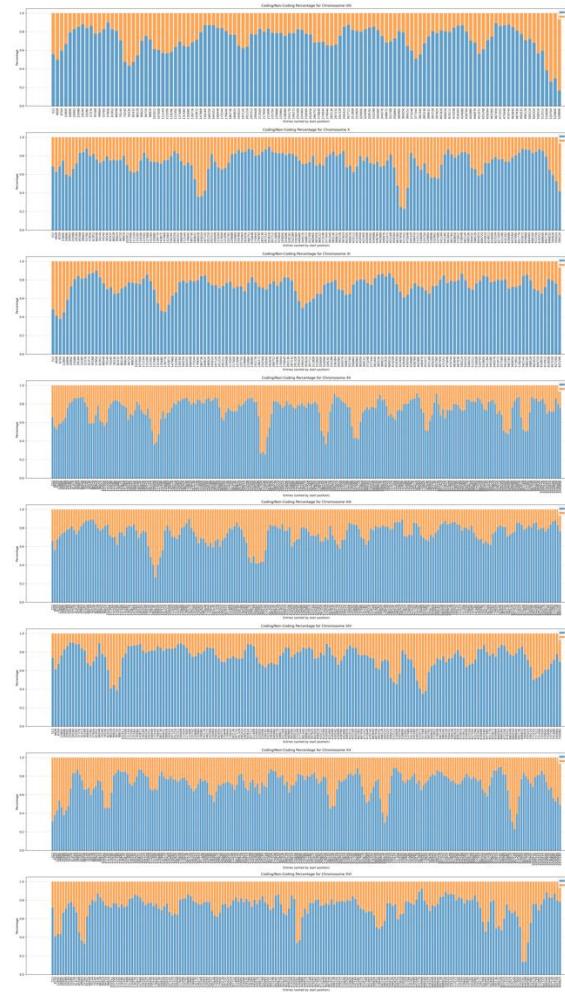
Conclusion:
BUSCO ~95% completeness

Genome evaluation – coding / noncoding regions

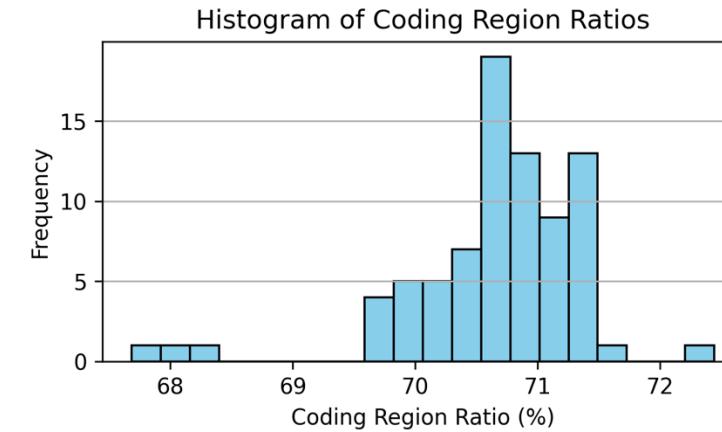
R64 Reference Yeast



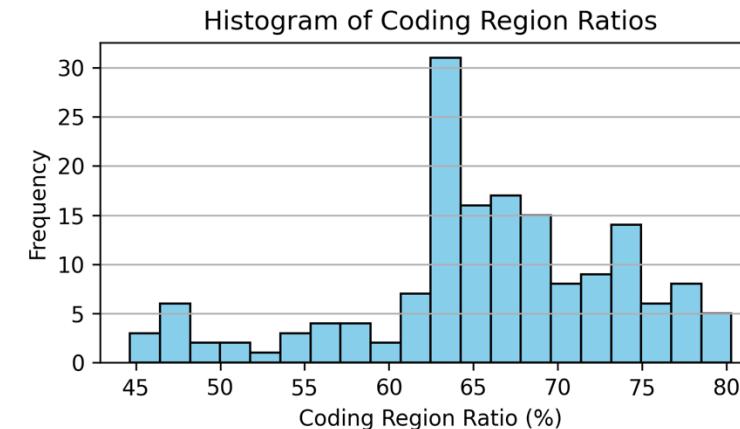
72.46% coding regions

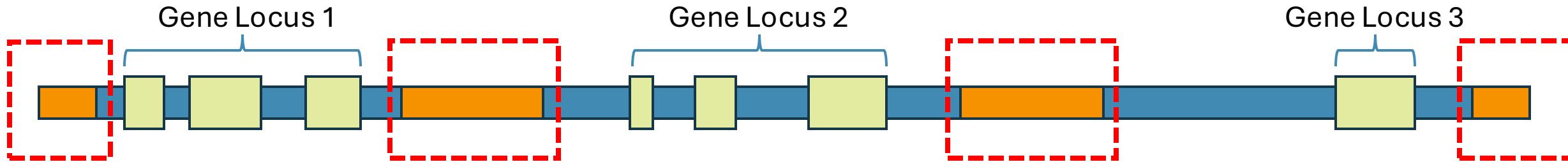


80 strains of yeasts

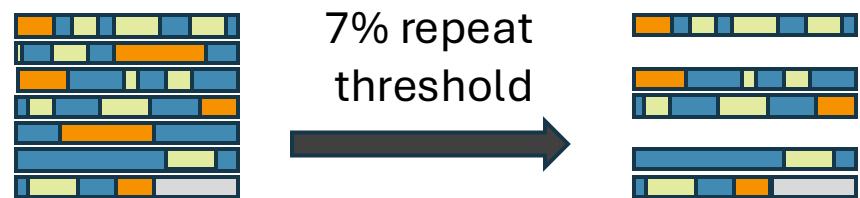


165 Sachamonycetales



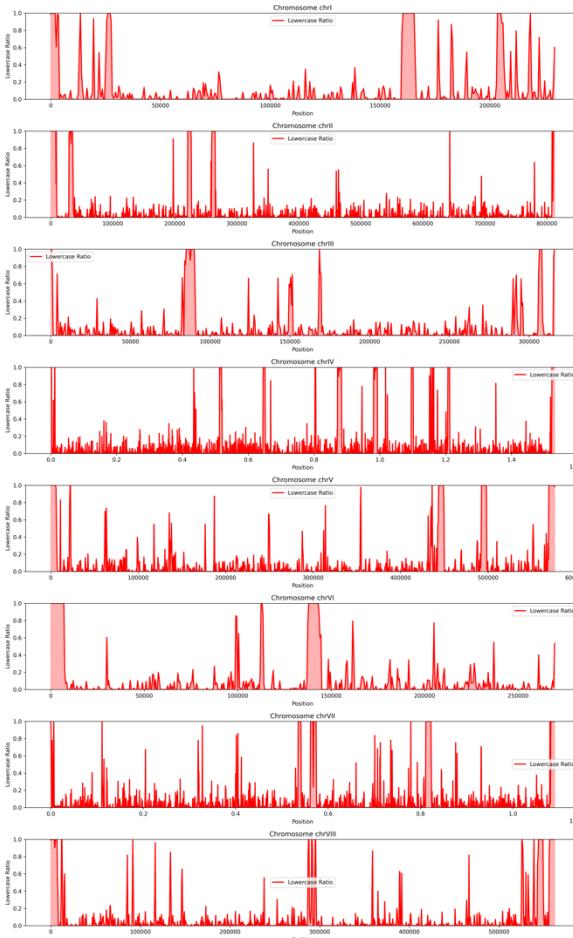


Q5: How repetitive are the genomes?



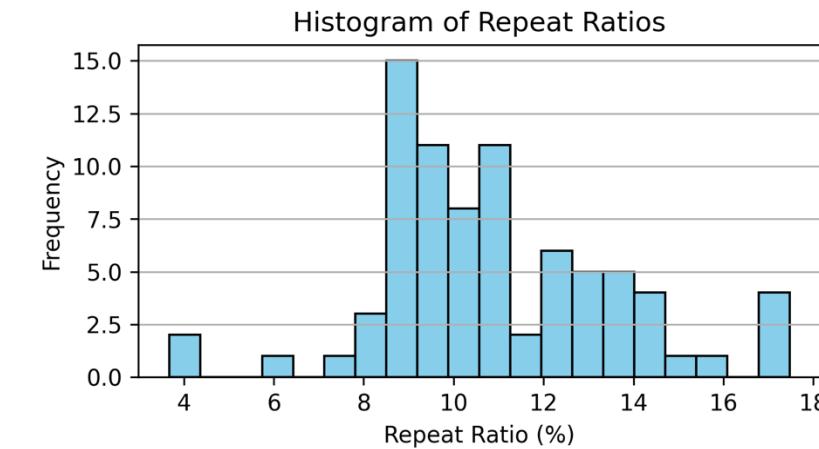
Genome evaluation – repeat regions

R64 Reference Yeast



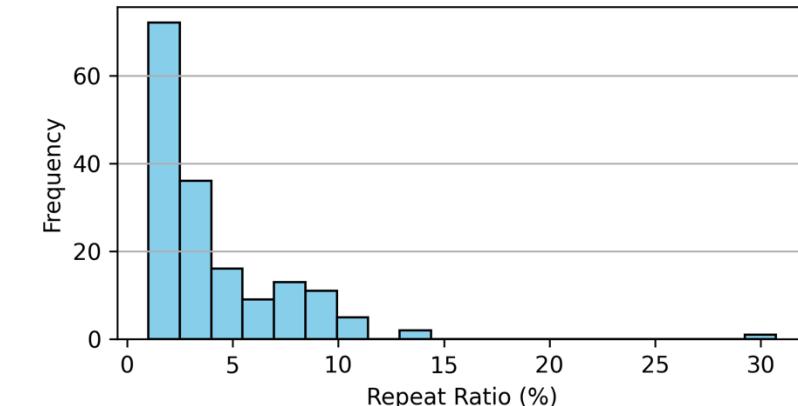
7.39% repeat regions

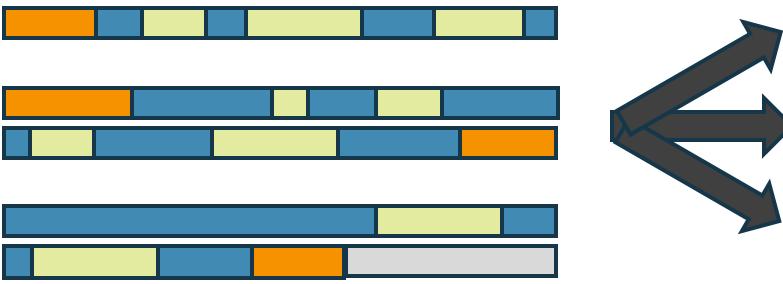
80 strains of yeasts



165 Sachamonomycetales

Histogram of Repeat Ratios





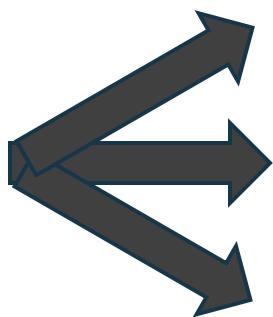
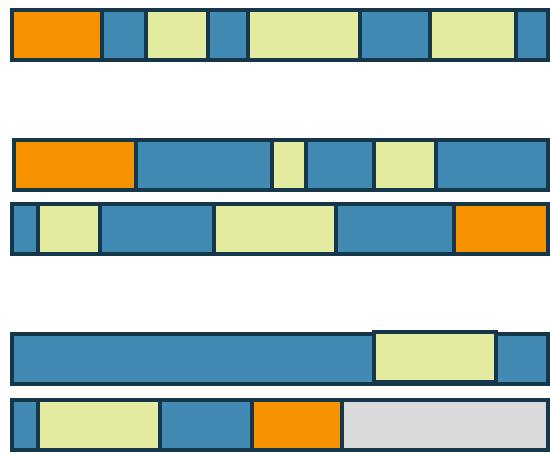
Training

Validation (chrXI, chrXIII, chrXV)

Testing (chrXII, chrXIV, chrXVI)

Q6: How many homologous sequences are there between training and testing?

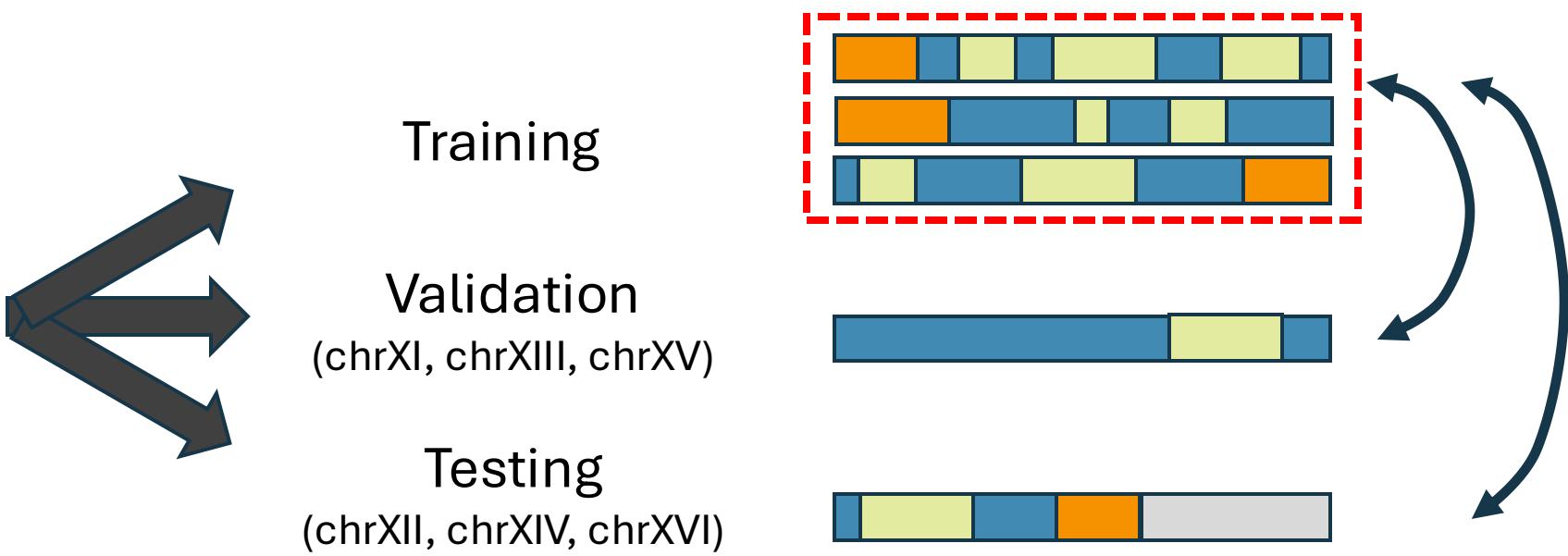




Training

Validation
(chrXI, chrXIII, chrXV)

Testing
(chrXII, chrXIV, chrXVI)



Detect homologous sequence using **Minimap2** (5% sequence divergence)

Final sequence for training / testing / validation

Before cleaning

r64

Train : 1440
Test : 608
Validation : 576

-377 (-14.4%)

80 strains

Train : 108960
Test : 608
Validation : 576

-6783 (-6.16%)

165 *Saccharomycetales*

Train : 404608
Test : 608
Validation : 576

-19195 (-4.73%)

After cleaning

Train : 1201
Test : 528
Validation : 518

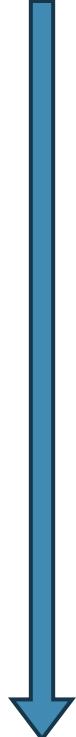
Train : 102315
Test : 528
Validation : 518

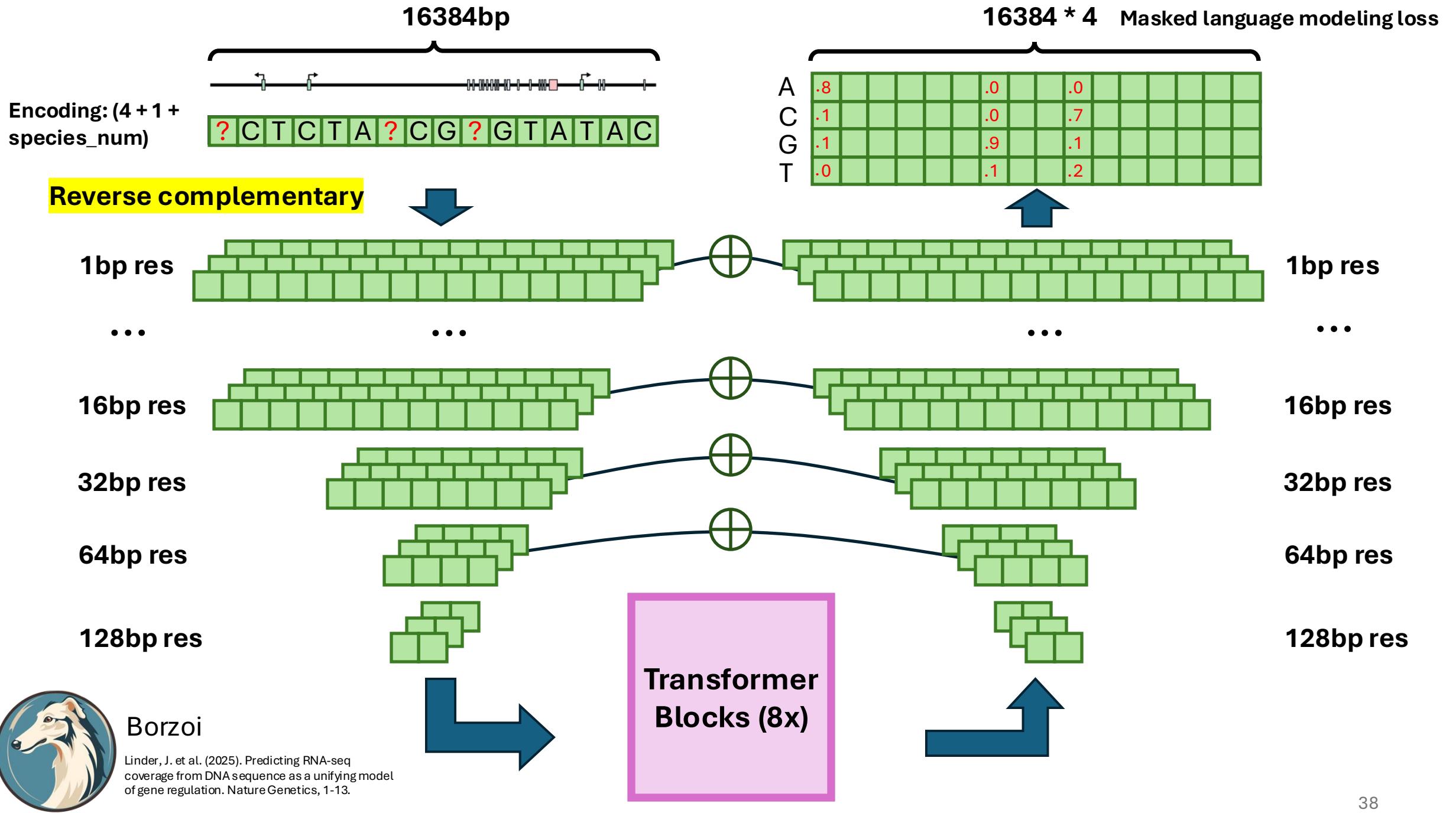
Train : 385551
Test : 528
Validation : 518

Fungal Language Model Architecture



Different model architecture we've tried

- Dilated convolutional neural network (small) Total params: 320,708 (1.22 MB)
 - Dilated convolutional neural network (large) Total params: 3,642,116 (13.89 MB)
 -  Transformer-based unet (small) Total params: 13,665,828 (52.13 MB)
 - Transformer-based unet (large) Total params: 71,790,564 (273.86 MB)
- 
- Model gets bigger**



Fungal Language Model

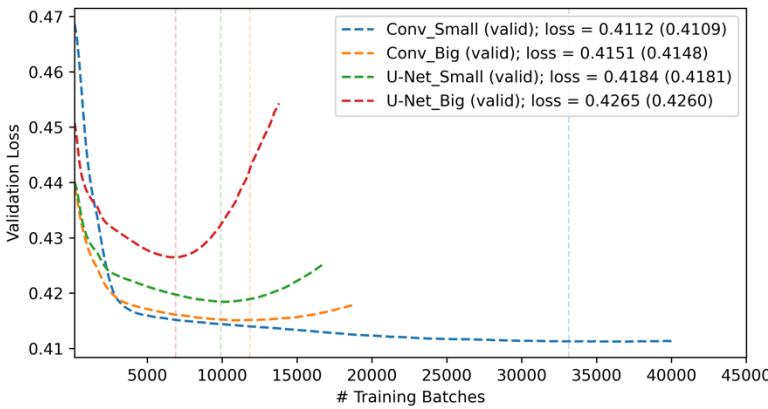
Self-supervised training Results



Model comparison

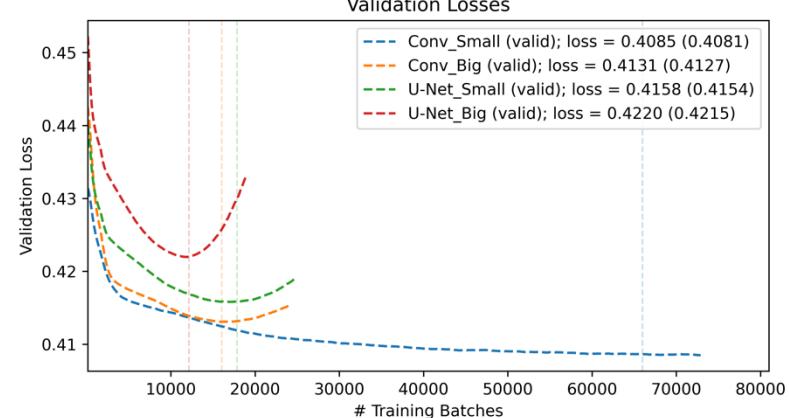
r64

Validation Losses



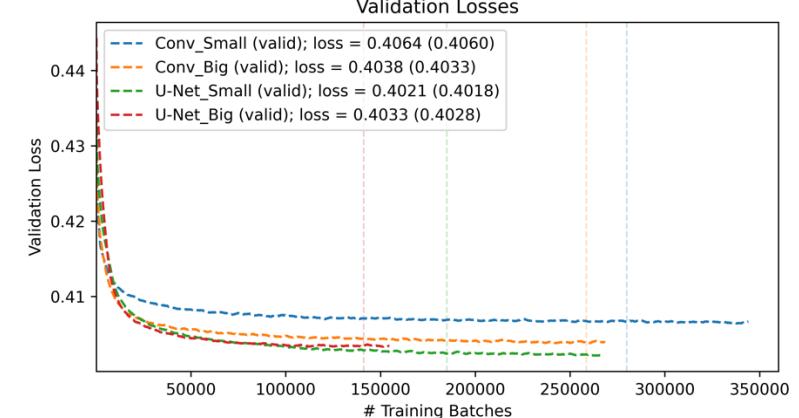
80 strains

Validation Losses

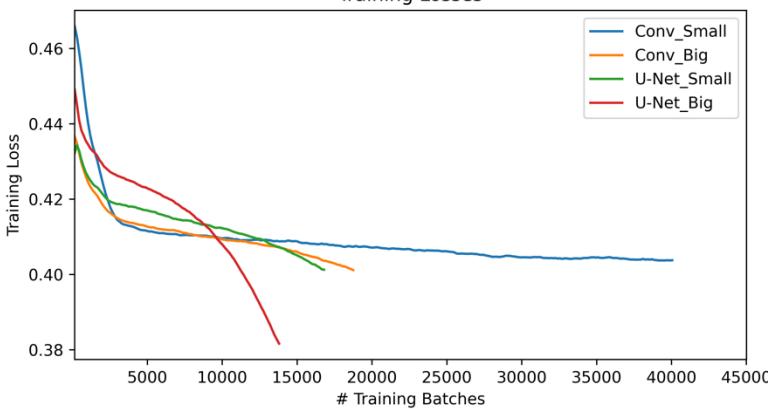


165 *Saccharomycetales*

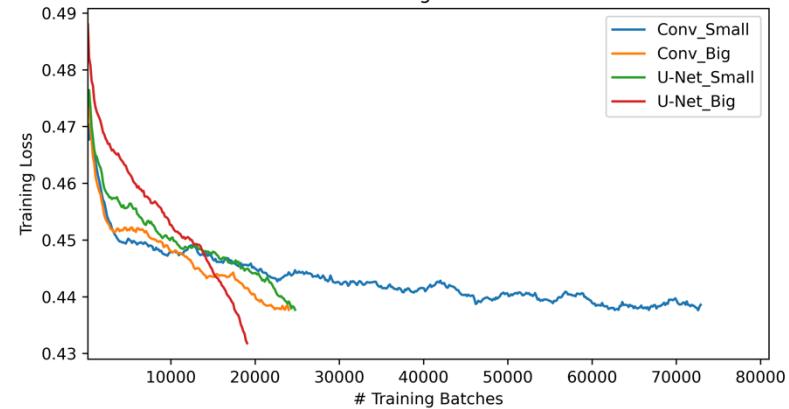
Validation Losses



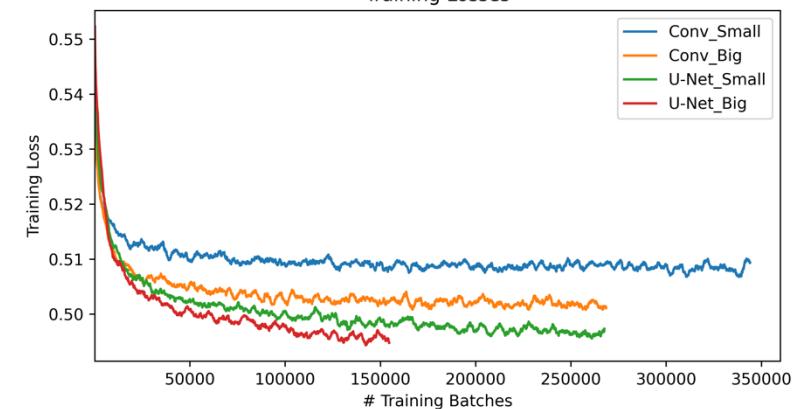
Training Losses



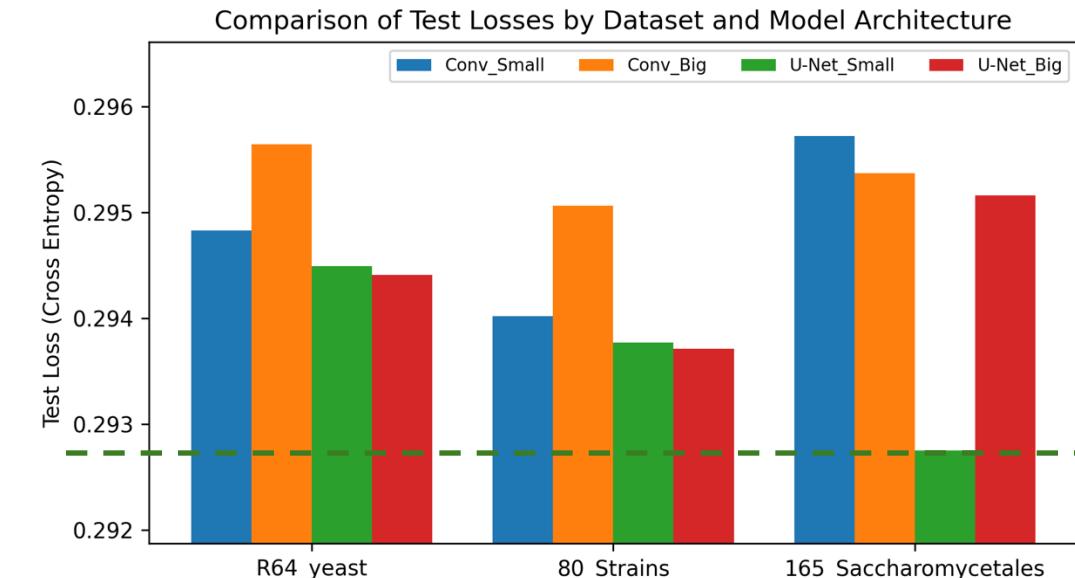
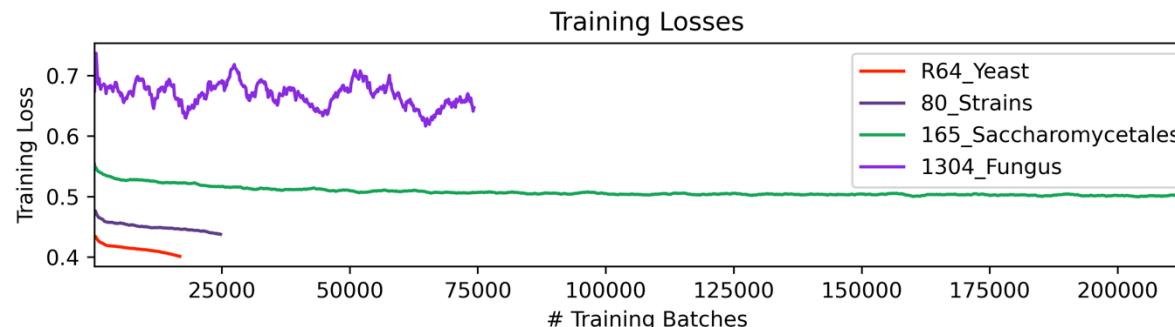
Training Losses



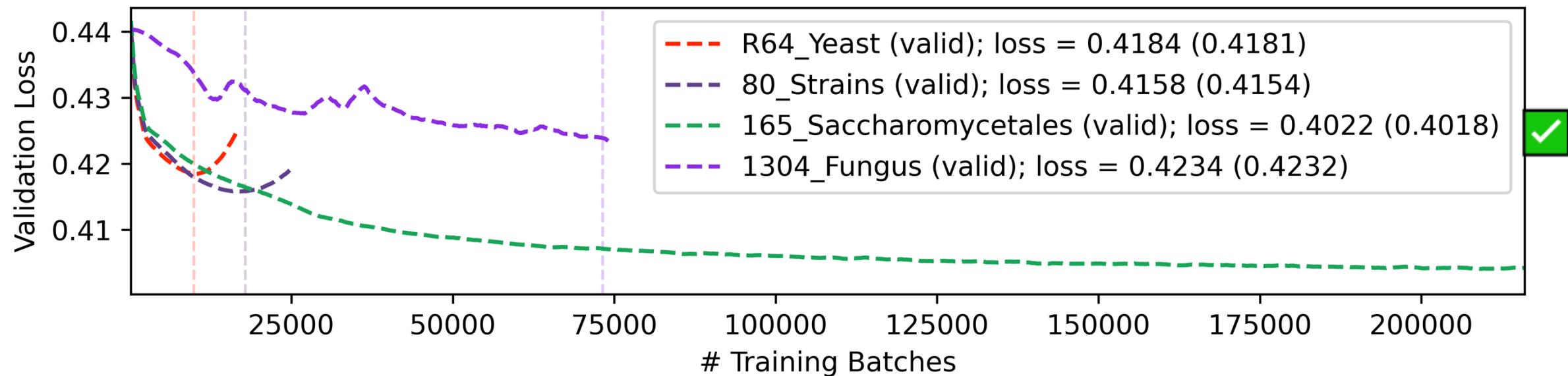
Training Losses



Dataset comparison



Validation Losses



Fungal Language Model Interpretability

1. Motif inference
2. Attention map visualization

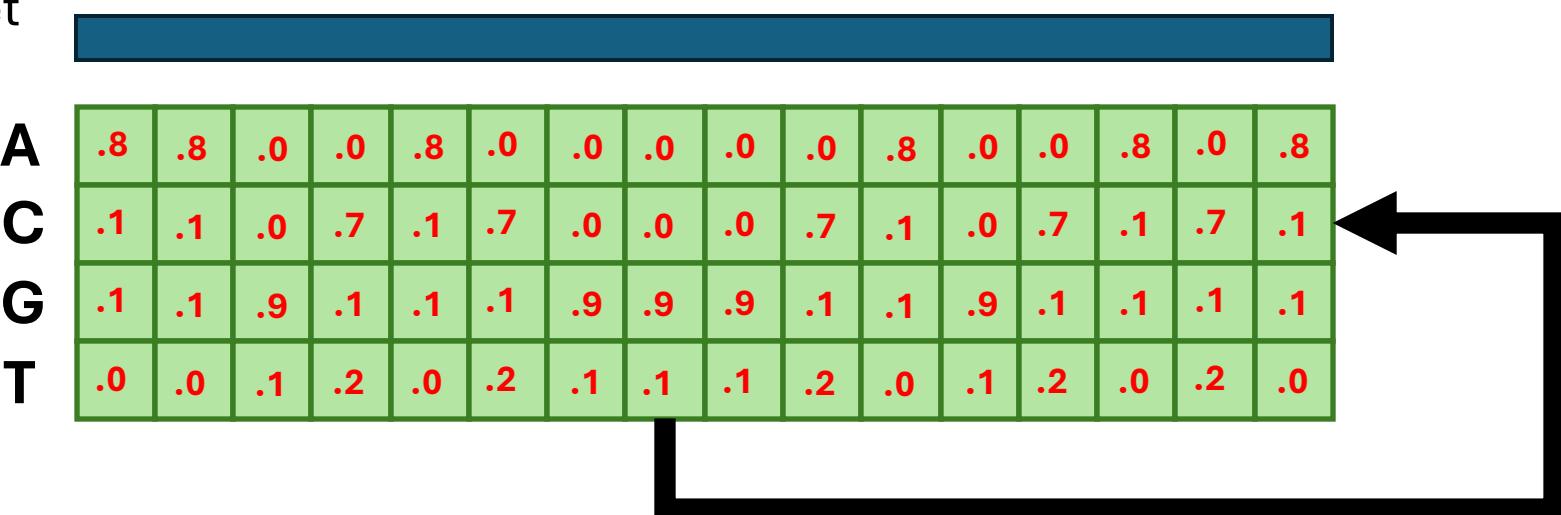


Fungal LM learns sequence conservation

Step 1: Construct PWM from test set

Input sequences in test set
(chrXII, chrXIV, chrXVI)

Predicting **15 %** masked
regions for each iteration



Step 2: PWM Normalization

1. Pseudocount Addition

$$p_{i,j} = p_{i,j} + \epsilon$$

$p_{i,j}$ is the read count for
nucleotide i at position j

2. Row Normalization

$$p_{j,i}^{norm} = \frac{p_{j,i}}{\sum_{k=1}^4 p_{j,k}}$$

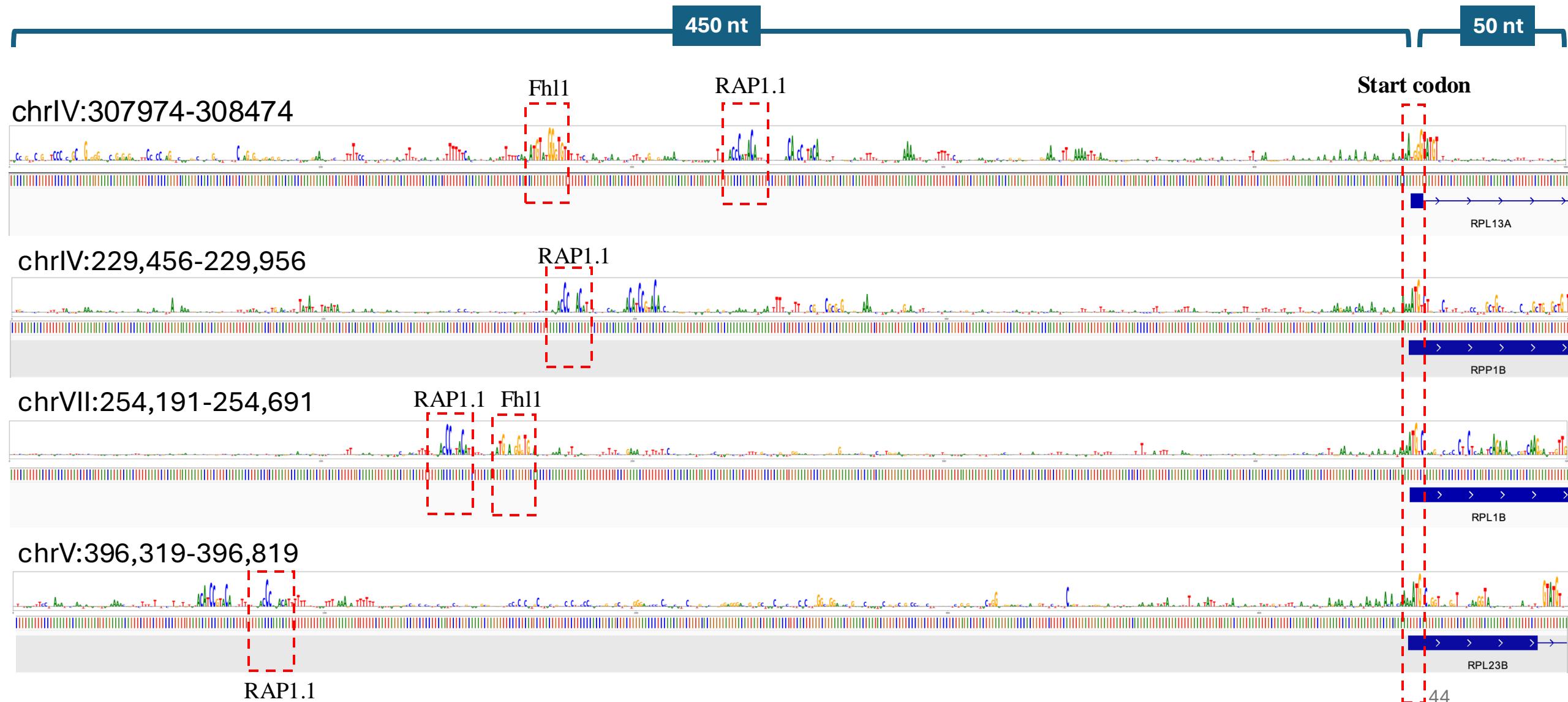
3. Entropy Calculation

$$H_j = - \sum_{i=1}^4 p_{j,i}^{norm} \cdot \log_2(p_{j,i}^{norm})$$

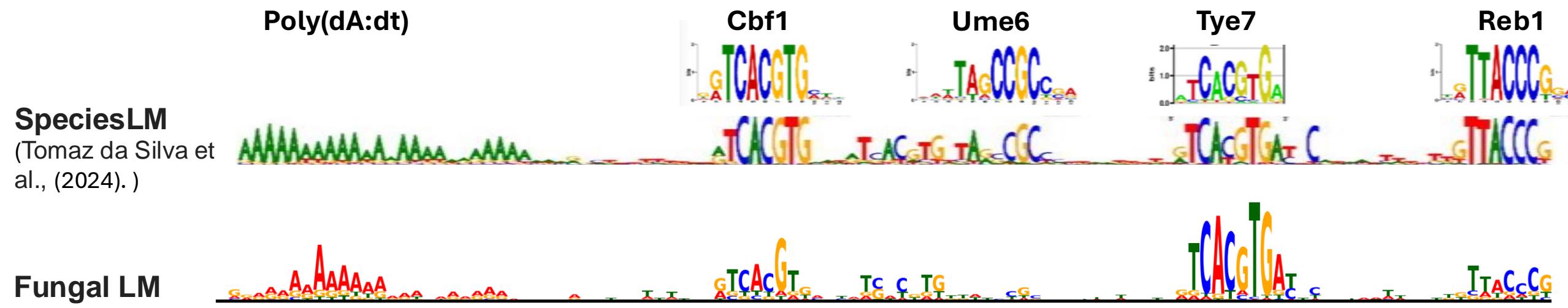
4. Conservation Calculation

$$C_j = 2 - H_j$$

Ribosomal Protein Upstream Promoter regions



Upstream Promoter region of SMT3 gene

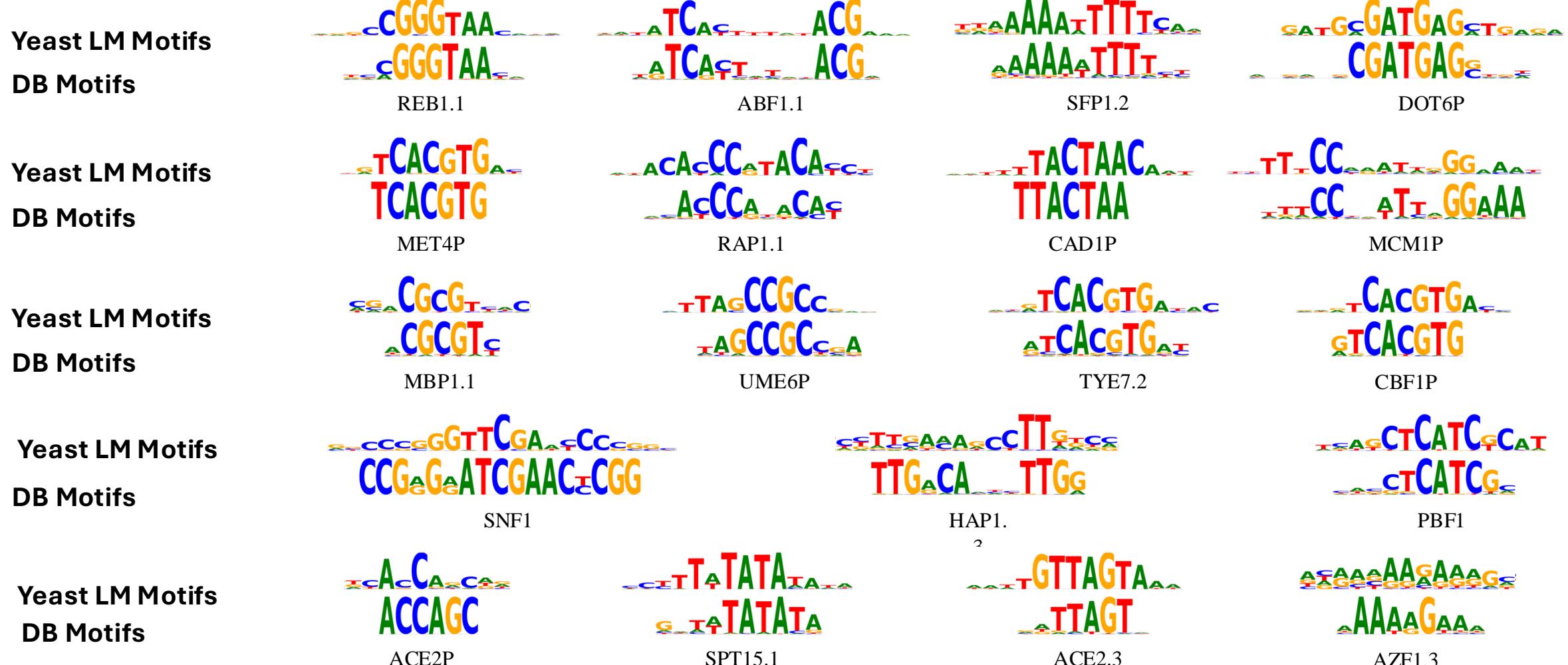


SpeciesLM: Trained with 5' and 3' regions only.

Fungal LM: Trained with 165 full Saccharomycetales genomes (**harder approach**)

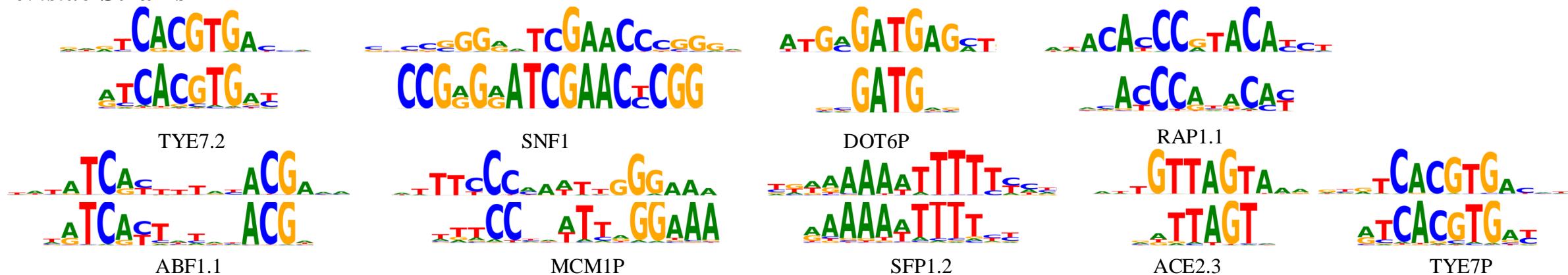
Tomaz da Silva et al., (2024). Nucleotide dependency analysis of DNA language models reveals genomic functional elements. bioRxiv

Constructing motifs in *S. cerevisiae* genome

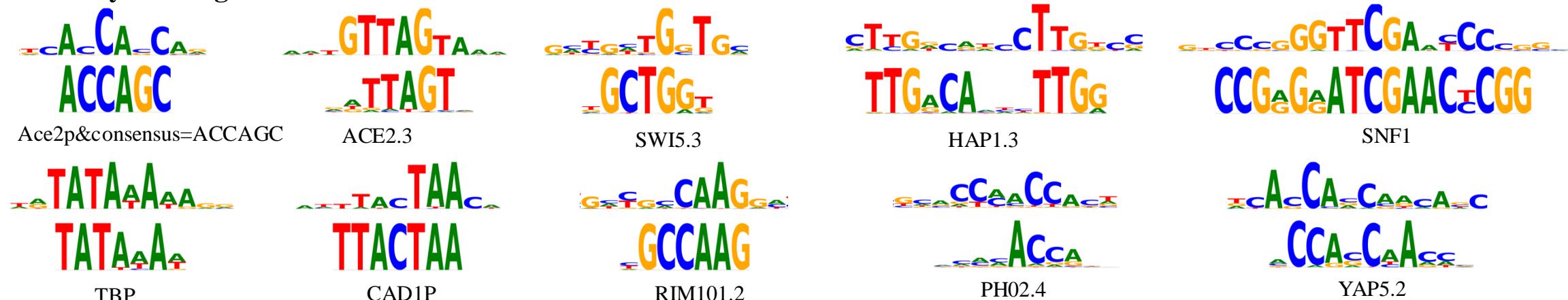


Constructing motifs in unseen genomes

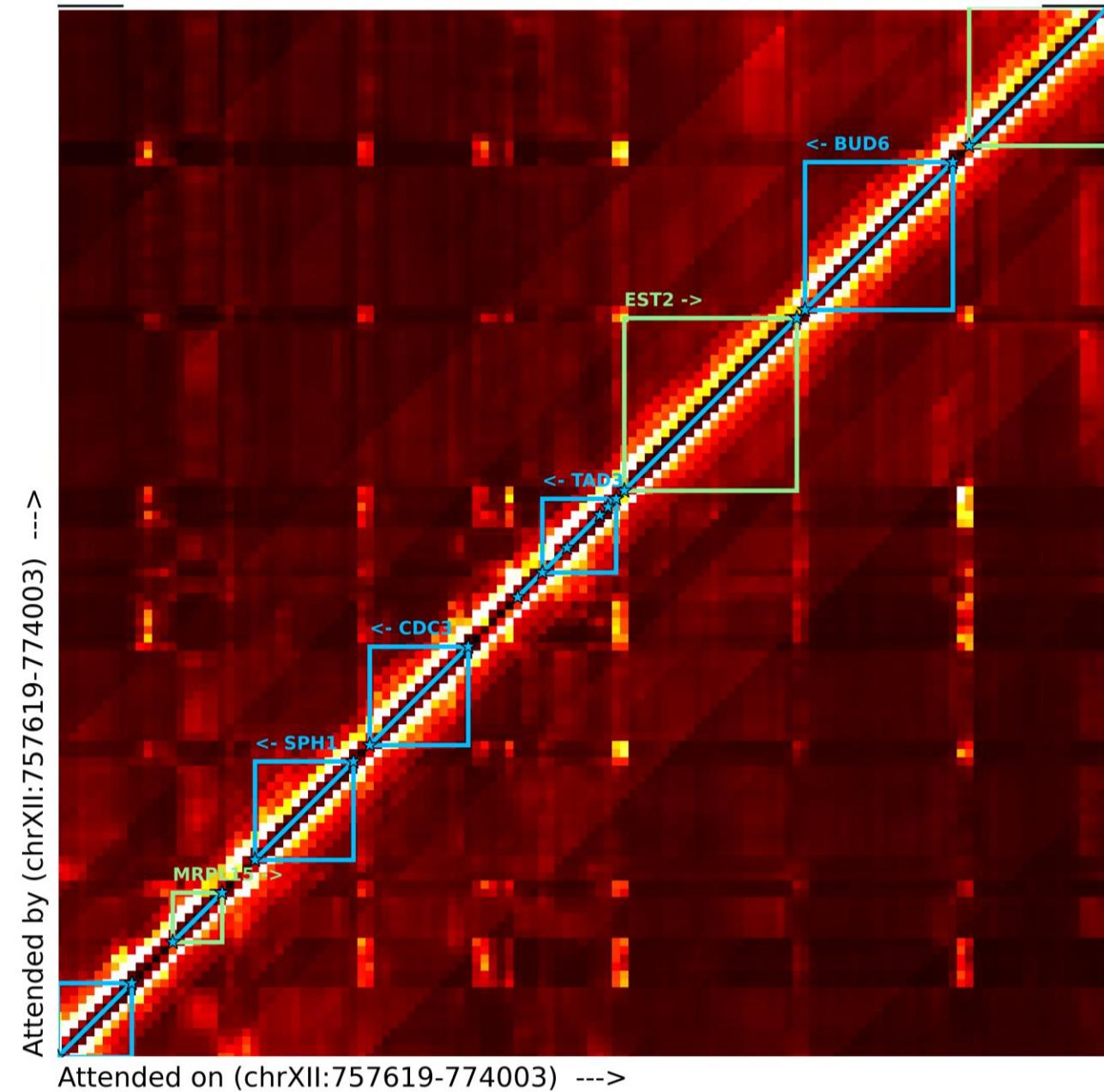
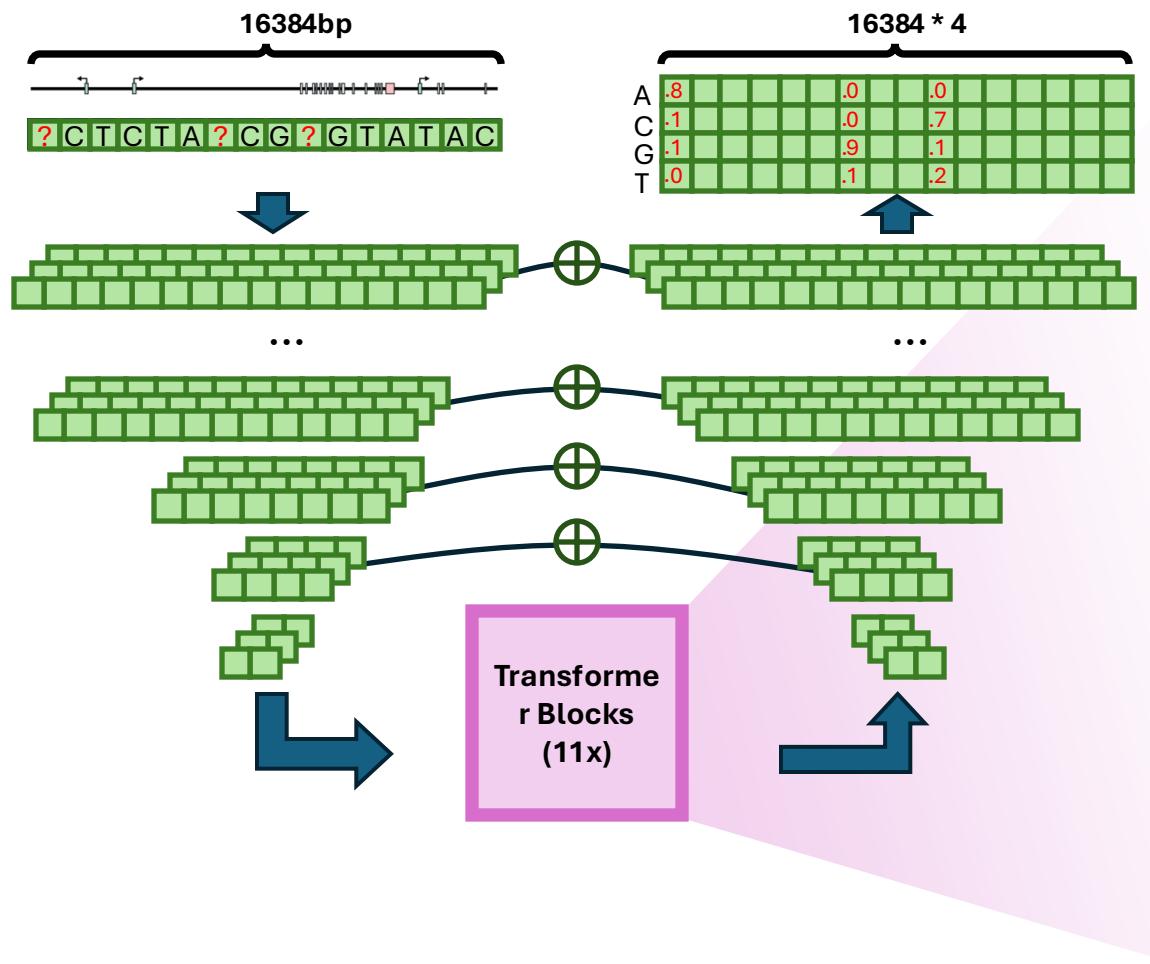
8 *S. cerevisiae* Strains



8 Schizosaccharomycetales genomes



Self-attention maps



Fungal LM: Summary

1. The **Saccharomycetales order** is a good evolutionary distance, offering good species diversity.
2. Thoroughly investigate genomes (protein-coding / repetitive / # gene per window)
3. Homologous sequence removal between train-test/validation is crucial
4. Transformer-based U-Net architecture is the best
5. Model interpretability
 1. LM can capture cis-regulatory motifs
 2. Attention maps highlights potential regulatory elements

Part II

Fine-tuning Fungal Language Model

ChIP-exo, histone marks, RNA-Seq prediction

Q: Does fine-tuning a pretrained LM outperform training a new model from scratch under the exact model architecture?

Genomic tracks intro & Data preprocessing



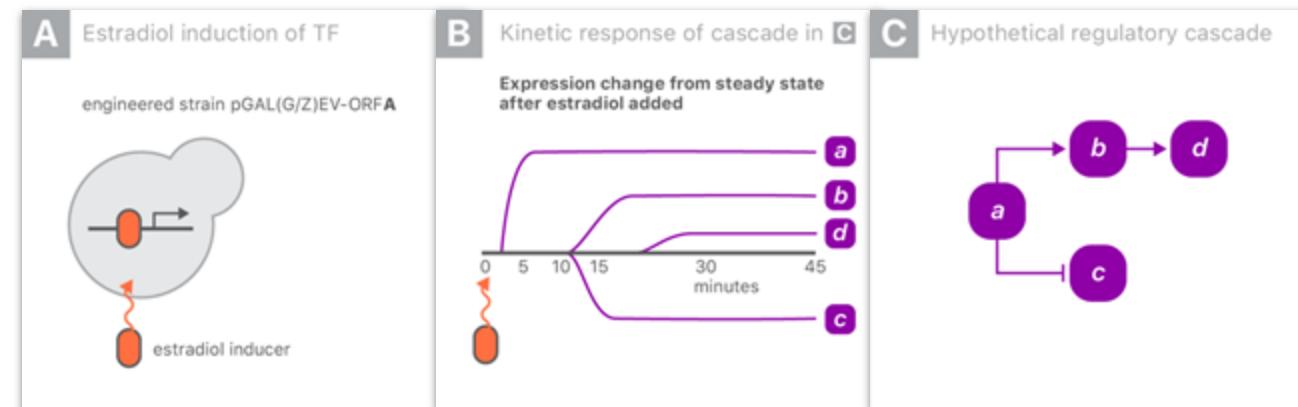
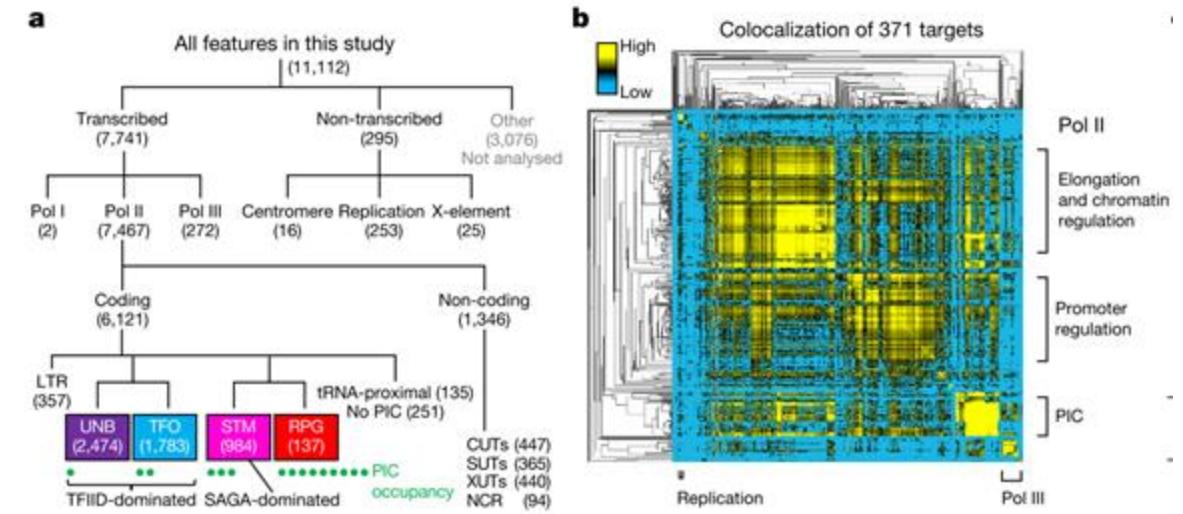
ChiP-exo + Histone Marks + RNA-Seq

- ChiP-exo provides high res view of protein-DNA binding across the yeast genome.
- Dataset includes 1128 ChiP-exo experiments
- Histone Mods MNase-ChIP-seq

Rossi, M. J., Kuntala, P. K., Lai, W. K., Yamada, N., Badjatia, N., Mittal, C., ... & Pugh, B. F. (2021). A high-resolution protein architecture of the budding yeast genome. *Nature*, 592(7853), 309-314.

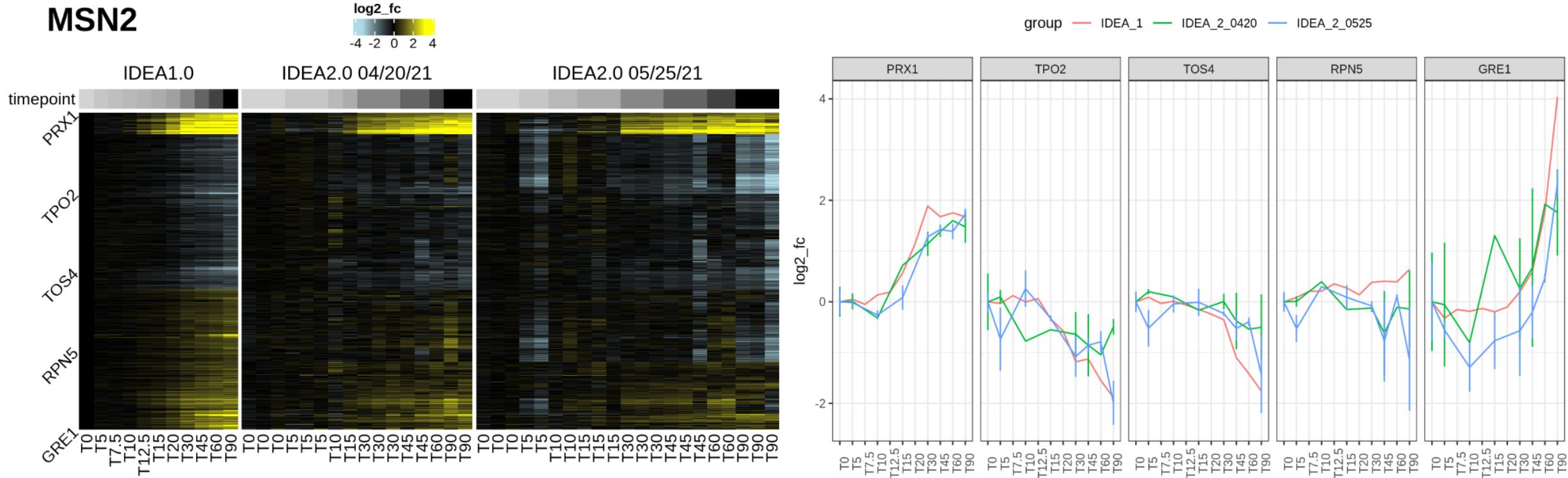
- Genome-scale **perturbation dynamics** propagate signals across regulatory networks (1340 experiments)
- Aggregating dynamics across many **time-courses** enables disambiguation of cause > effect relationships

Hackett, S. R., Baltz, E. A., Coram, M., Wranik, B. J., Kim, G., Baker, A., ... & McIsaac, R. S. (2020). Learning causal networks using inducible transcription factors and transcriptome-wide time series. *Molecular systems biology*, 16(3), e9174.



RNA-Seq

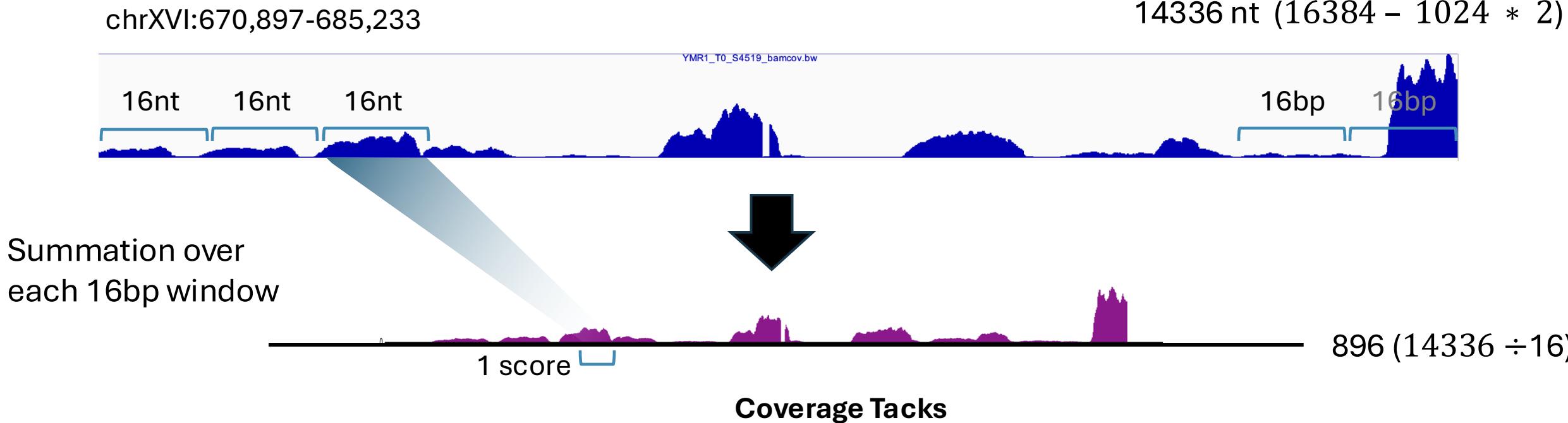
- IDEA (the Induction Dynamics gene Expression Atlas)



Hackett, S. R., Baltz, E. A., Coram, M., Wranik, B. J., Kim, G., Baker, A., ... & McIsaac, R. S. (2020). Learning causal networks using inducible transcription factors and transcriptome-wide time series. *Molecular systems biology*, 16(3), e9174.

Borzoi: squashed scale $y_{j,t}^{(\text{squashed})} = \begin{cases} y_{j,t}^{(3/4)} & \text{if } y_{j,t}^{(3/4)} \leq 384, \text{ otherwise } 384 + \sqrt{y_{j,t}^{(3/4)} - 384} \end{cases}$

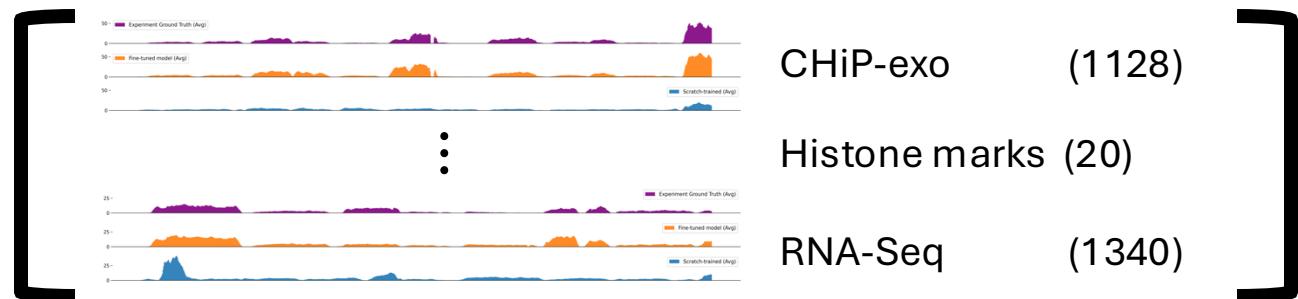
Track transformation



No normalization across tracks

Loss function: Poisson Loss.

The depth of the track somehow reflects the quality of the data.



Scratch-trained model vs.

Fine-tuned Fungal LM:

Initialization & Training





Borzoi

16384bp

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

16384 * 4

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





Borzoi

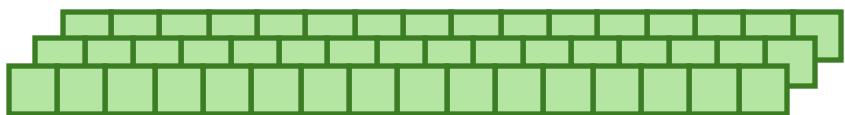
16384bp



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



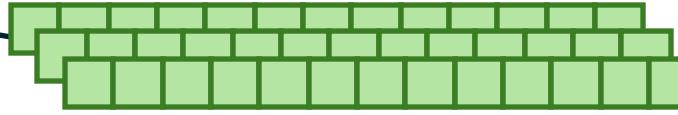
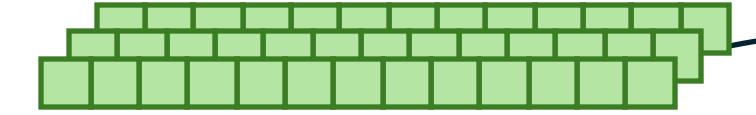
1bp res



...

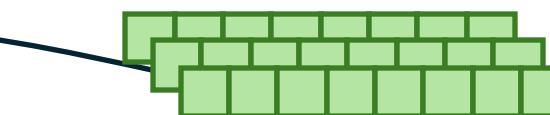
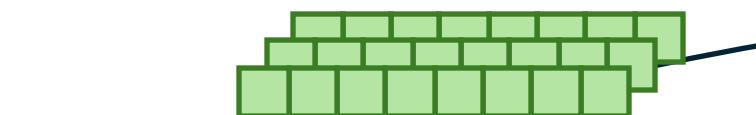
...

16bp res



16bp res

32bp res



32bp res

64bp res



64bp res

128bp res



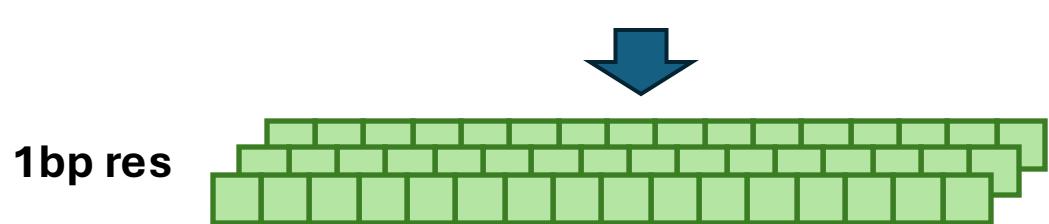
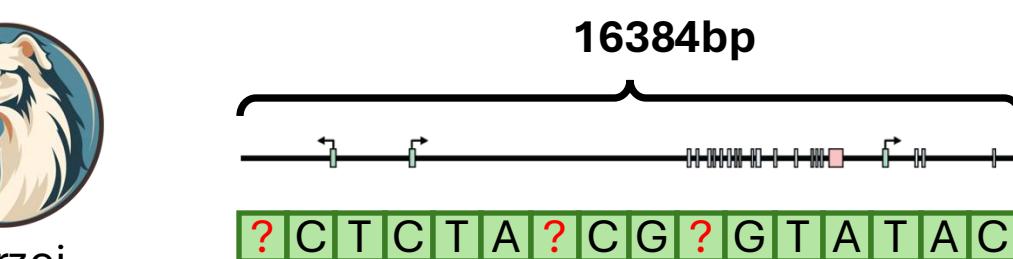
Transformer
Blocks (8x)



128bp res



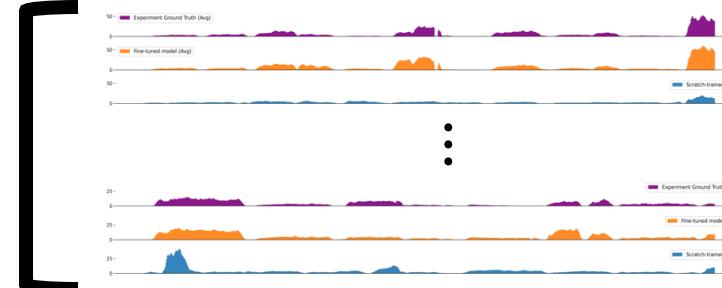
Borzoi



...



Coverage Tacks (896 * 2488)



ChIP-exo (1128)

Histone marks (20)

RNA-Seq (1340)

32bp res



64bp res



128bp res

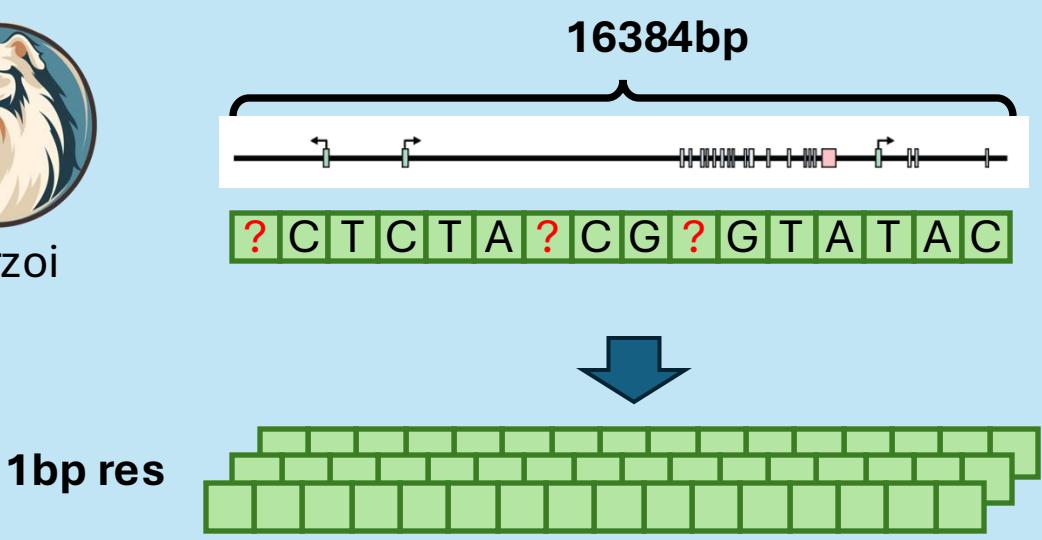


Transformer
Blocks (8x)

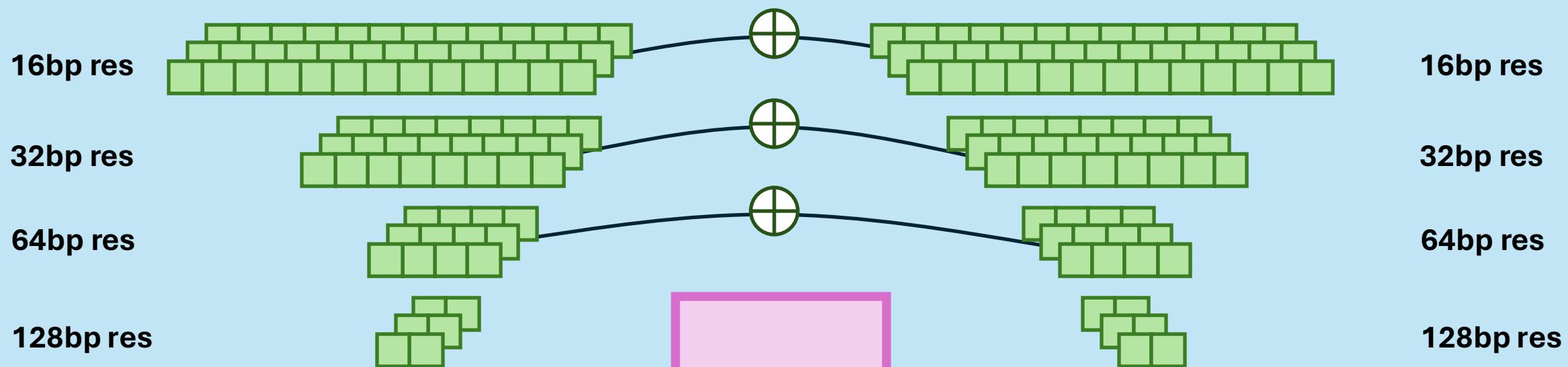
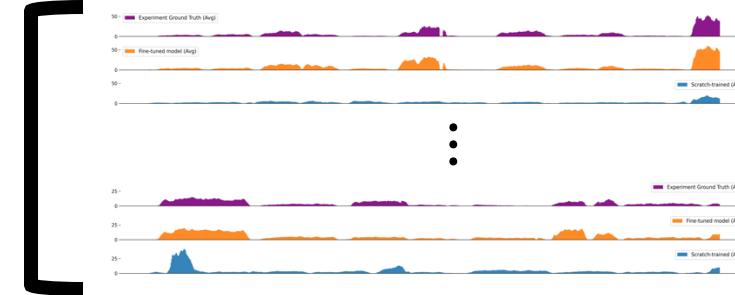




Borzoi



Coverage Tacks (896 * 2488)



1. Initialized with LM weights (Fine-tuned)
2. Random initialization (Scratch-trained)

Model training: 8-fold cross validation

- Divide genome into 8 folds.
- Train 8 models with distinct validation and test folds.



Fold0: 743 seq, 1406020 nt (0.1244)

chrXIV: 0-628758
chrX: 0-436307
chrXI: 440246-666816
chrIII: 0-114385

Fold1: 736 seq, 1433427 nt (0.1268)

chrXI: 0-440129
chrV: 0-151987
chrV: 152104-576874
chrXIII: 0-268031
chrVI: 0-148510

Fold2: 806 seq, 1521492 nt (0.1346)

chrII: 238323-813184
chrVII: 0-496920
chrIV: 0-449711

Fold3: 755 seq, 1408276 nt (0.1246)

chrXVI: 0-555957
chrIV: 449821-990877
chrVI: 48627-270161
chrVIII: 0-105586
chrIX: 355745-439888

Fold4: 732 seq, 1444997 nt (0.1278)

chrIV: 990877-1531933
chrXII: 614562-1078177
chrII: 0-238207
chrIII: 114501-316620

Fold5: 742 seq, 1284157 nt (0.1136)

chrVII: 497038-1090940
chrX: 436425-745751
chrI: 0-151465
chrI: 151582-230218
chrXII: 0-150828

Fold6: 785 seq, 1446481 nt (0.1280)

chrXIII: 268149-924431
chrXII: 150947-614562
chrXV: 0-326584

Fold7: 733 seq, 1360020 nt (0.1203)

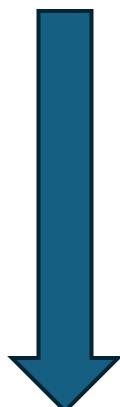
chrVIII: 105703-562643
chrXVI: 556073-948066
chrIX: 0-355629
chrXIV: 628875-784333

Scratch-trained

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

Input

16,384 * 4



Model

Fine-tuned

ACTCTACCCGGGTATAC

$$16,384 * (4 + 1 + 165)$$

A 5x20 grid representing a sequence alignment. The rows are labeled A, C, G, T, and ending. A vertical bracket on the left indicates the sequence starts at the second column.

Masked encoding

Species en (r64 : **109**)

Model

Scratch-trained model vs.

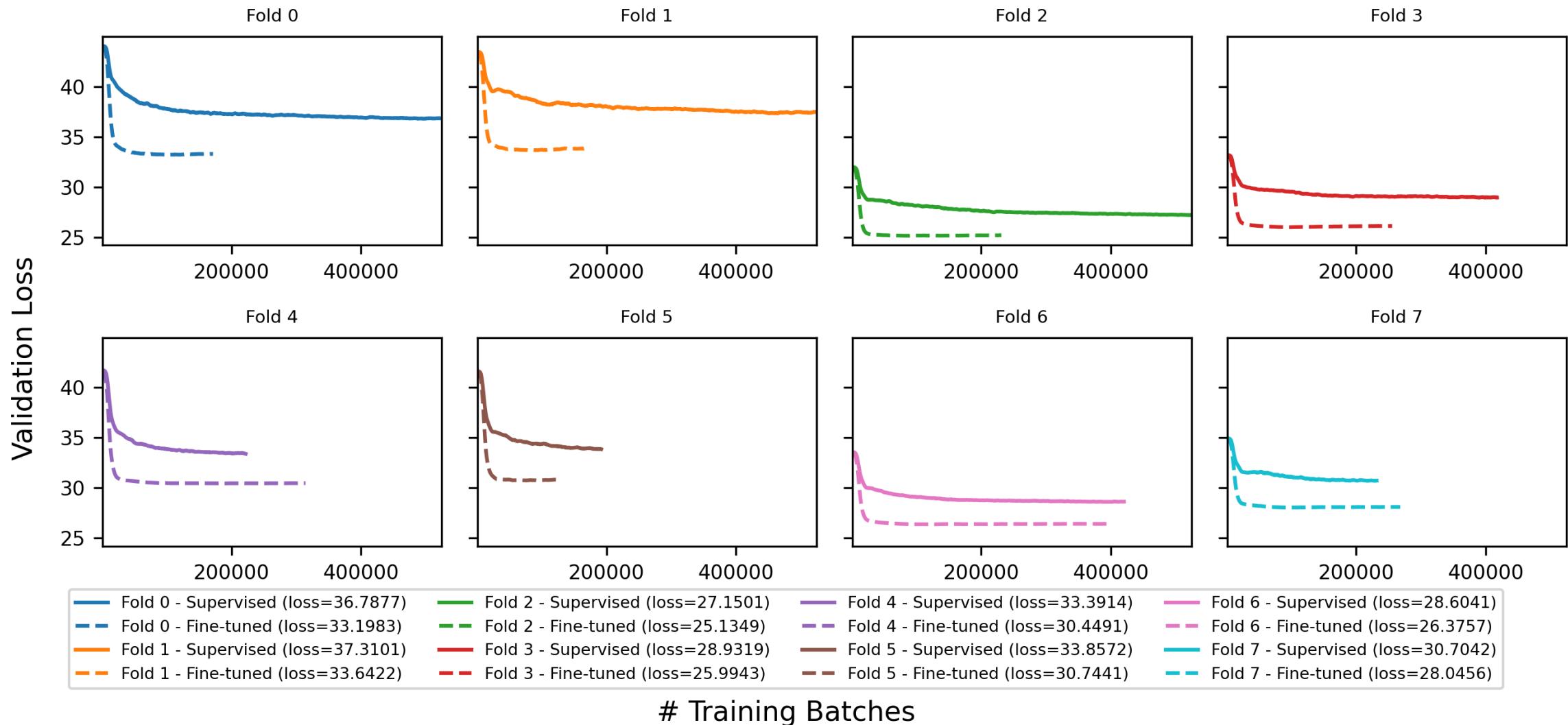
Fine-tuned Fungal LM:

Training Results



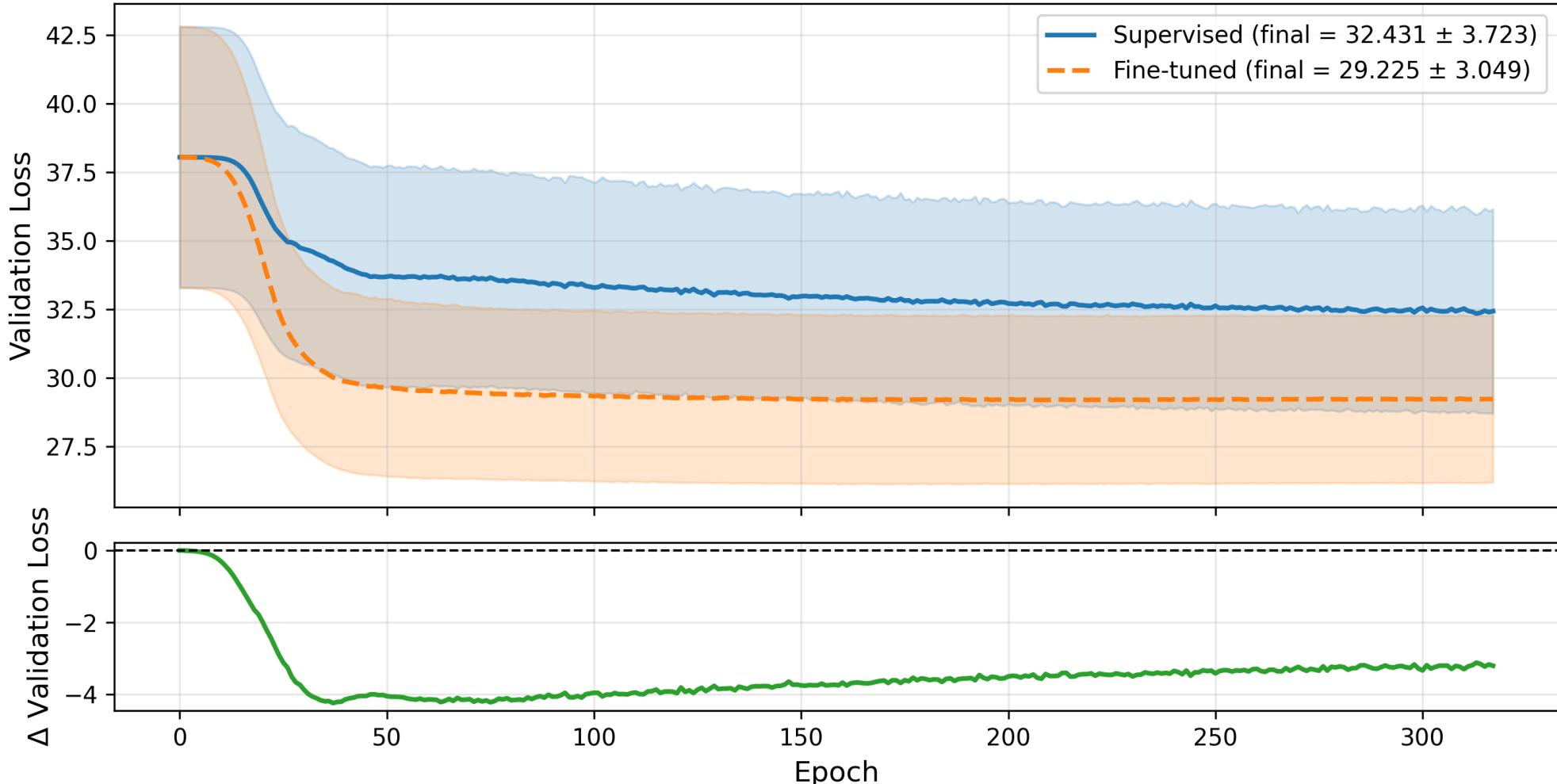
Fine-tuned vs Scratch-Trained

(16 bp resolution)



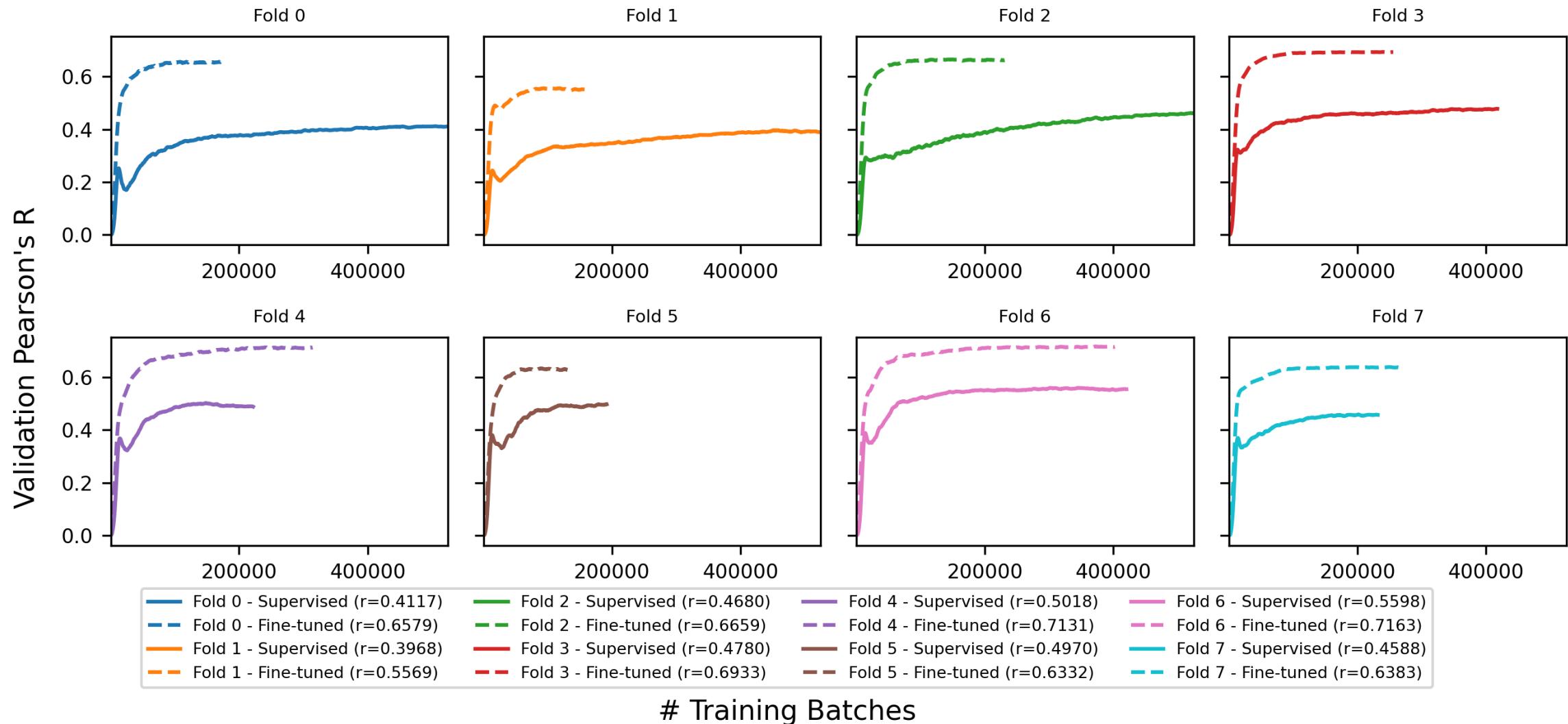
Fine-tuned vs Scratch-Trained

Supervised vs. Fine-Tuned (Validation Loss)



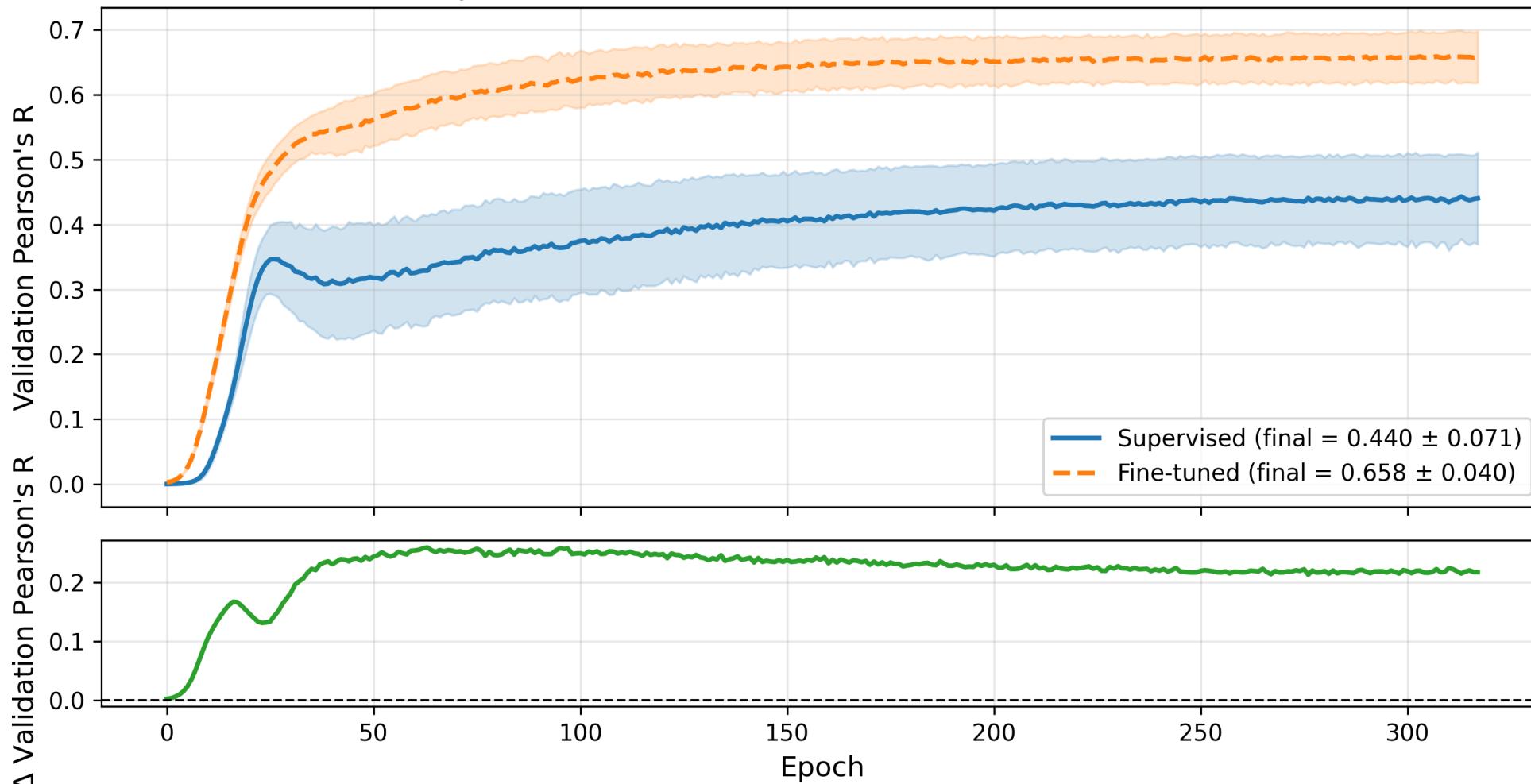
Fine-tuned vs Scratch-Trained

(16 bp resolution)



Fine-tuned vs Scratch-Trained

Supervised vs. Fine-Tuned (Validation Pearson's R)



(16 bp resolution)

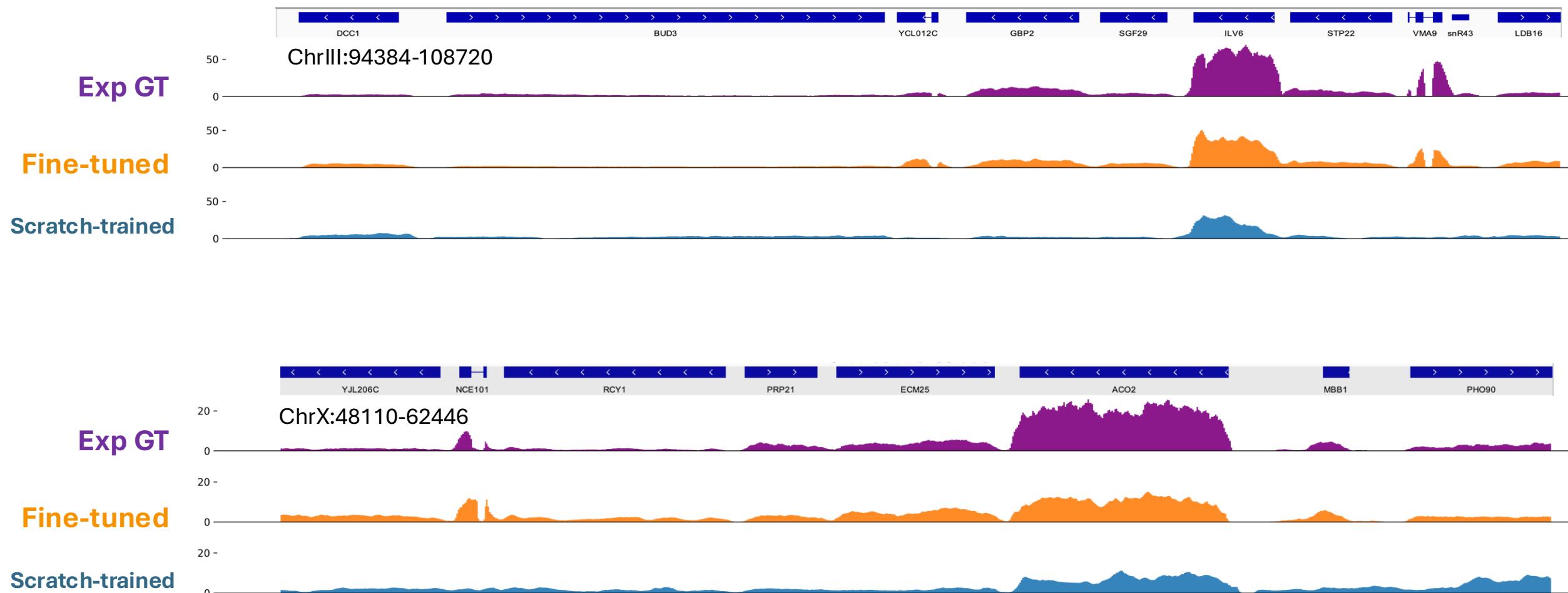
Scratch-trained model vs.

Fine-tuned Fungal LM:

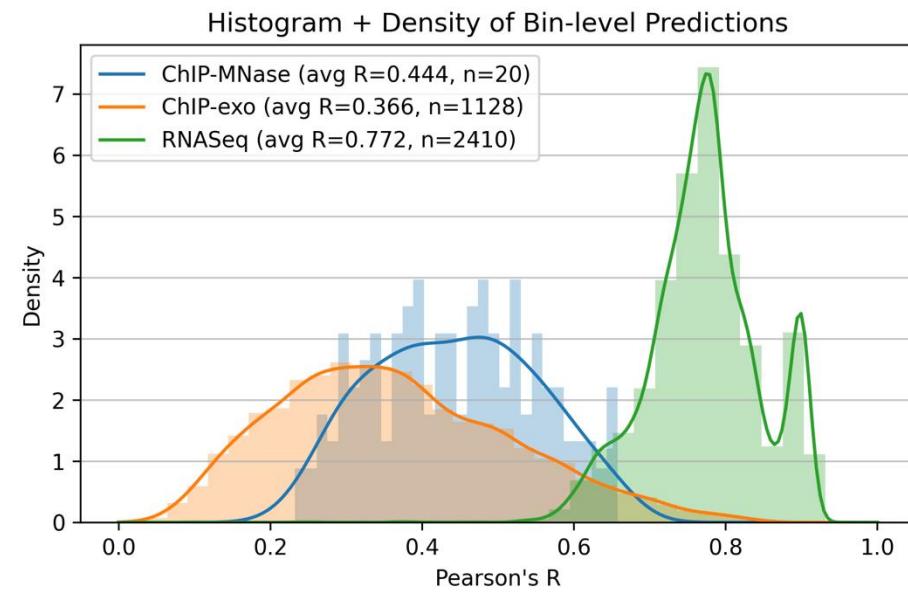
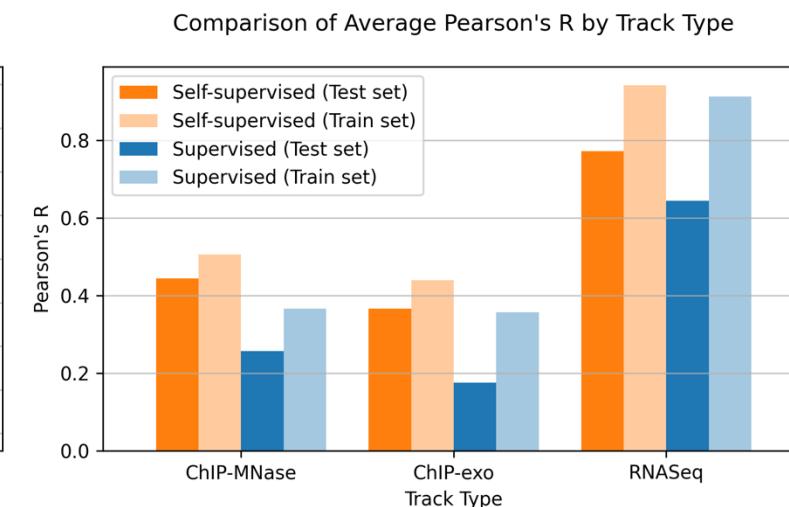
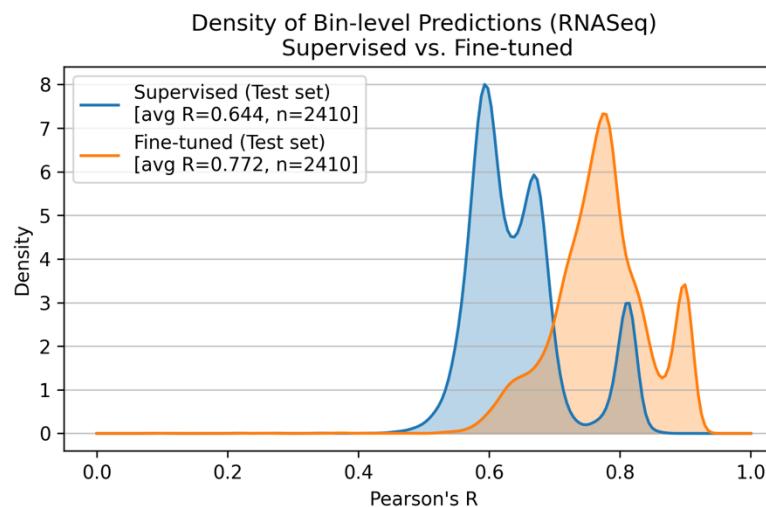
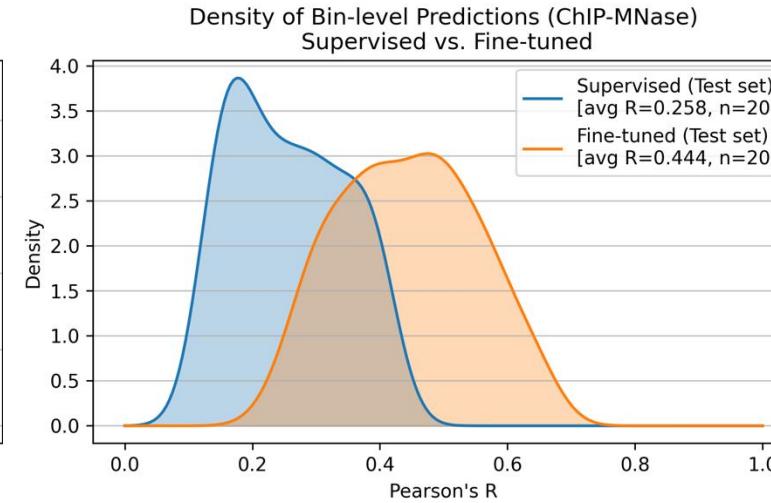
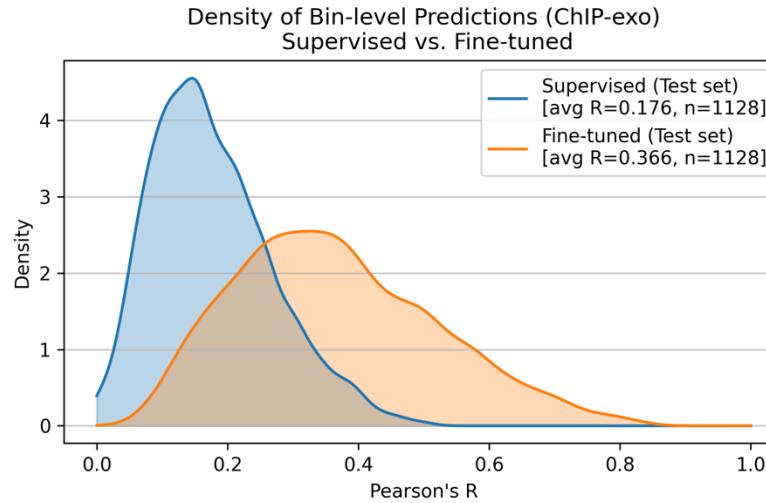
Track-level prediction evaluation



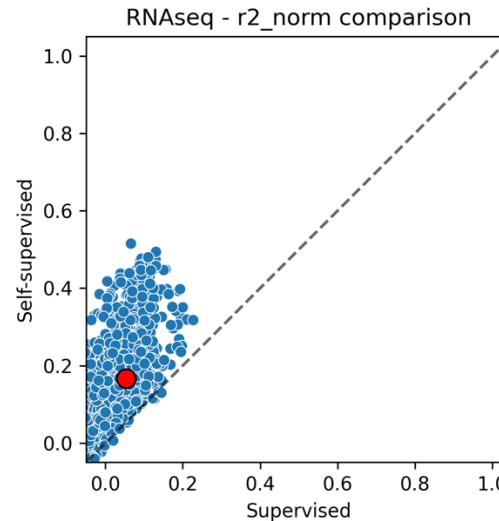
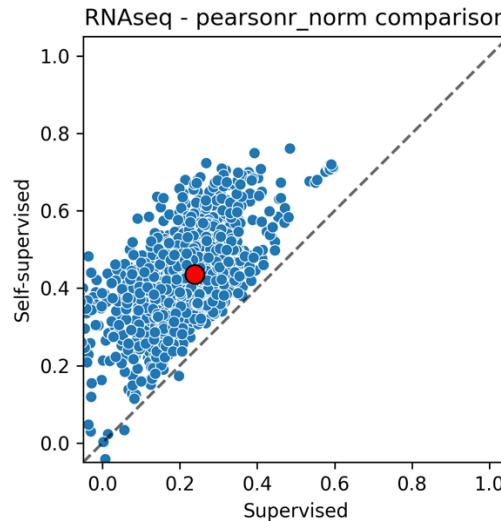
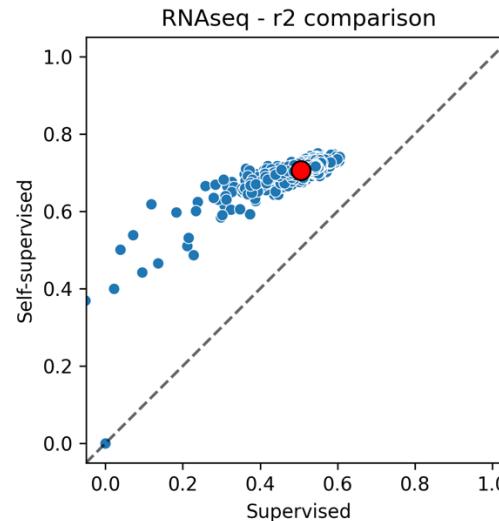
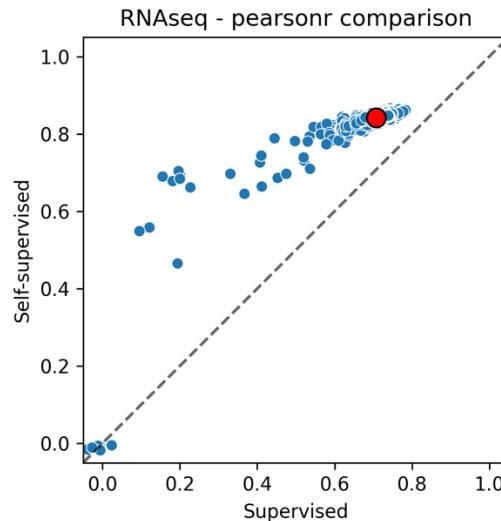
Fine-tuned vs Scratch-Trained (Test set)



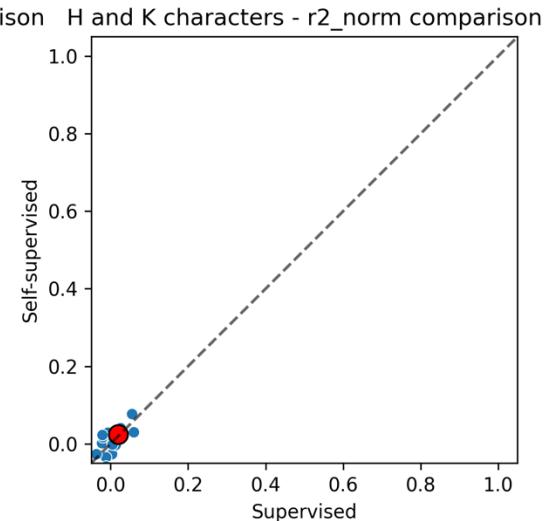
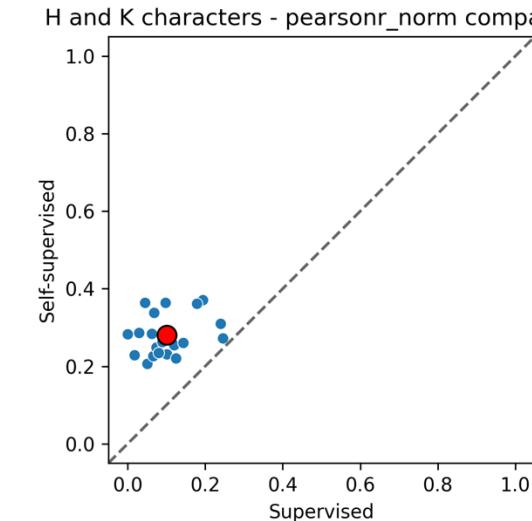
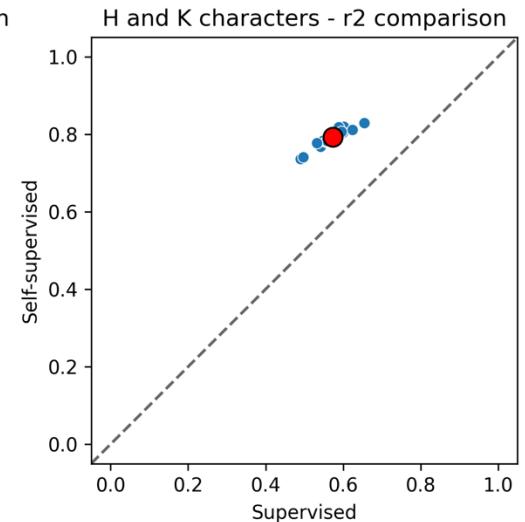
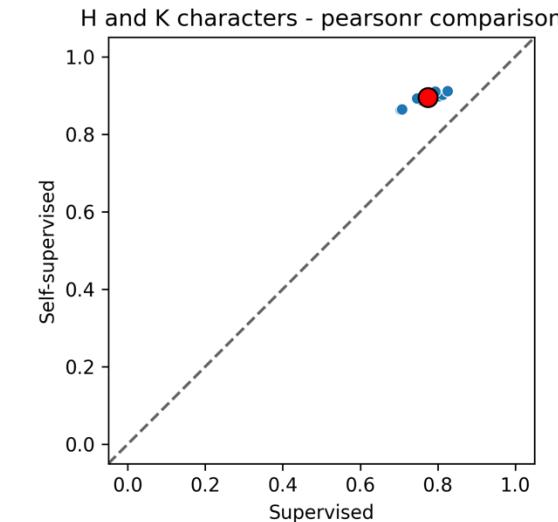
Fine-tune vs Scratch-Trained (Test set)



RNA-Seq

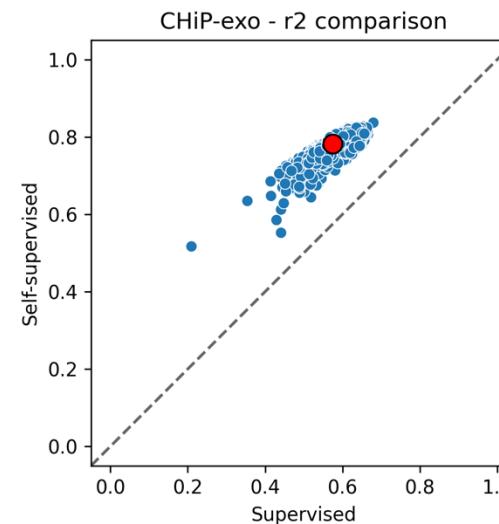
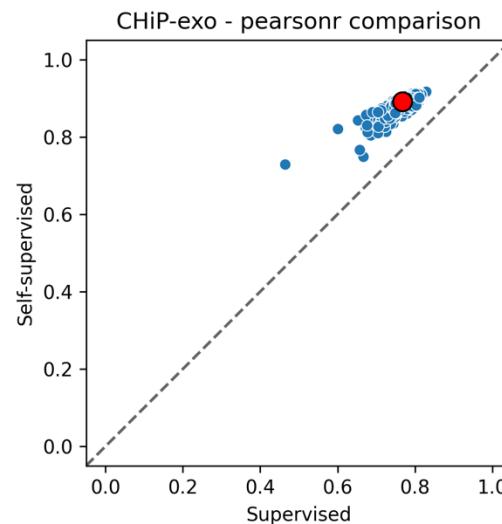


Histone Marks

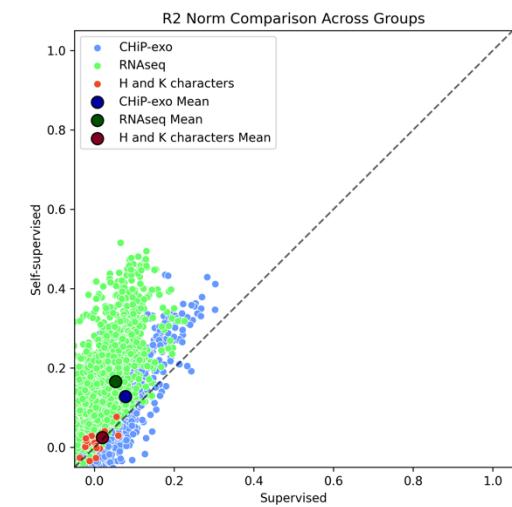
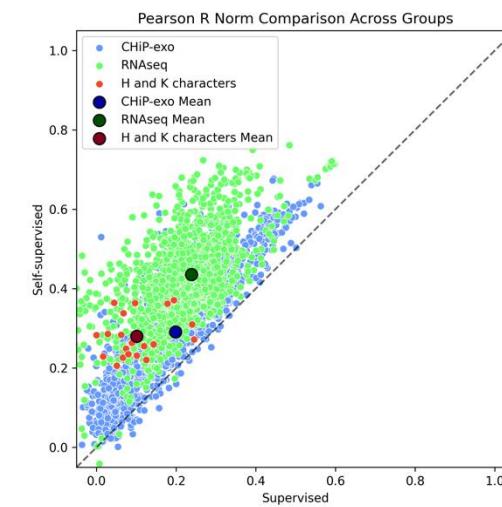
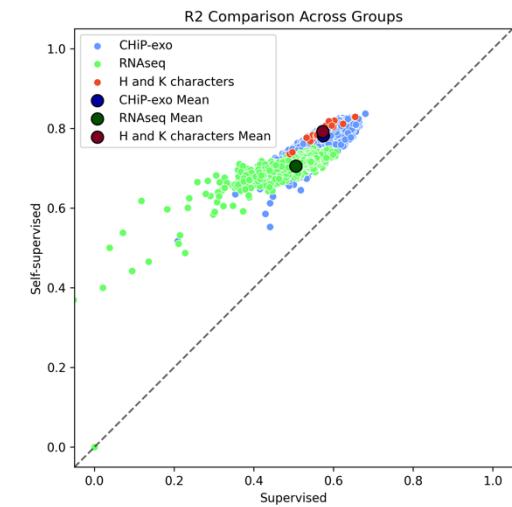
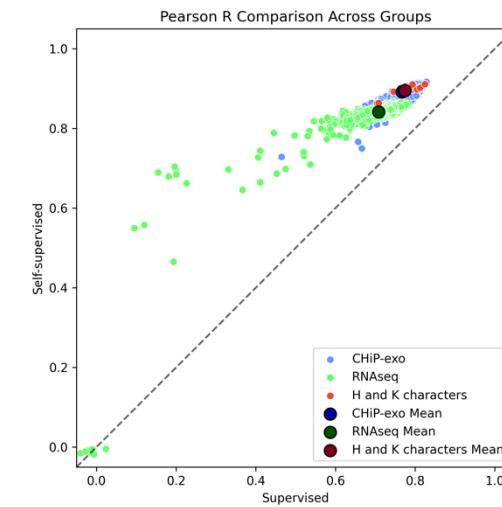


Average results across 8 folds. Each dot is a track.

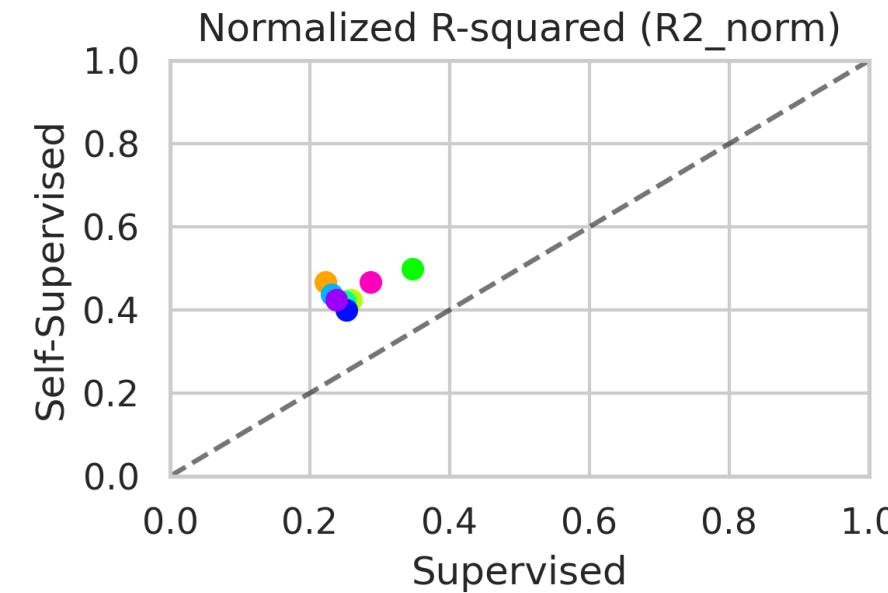
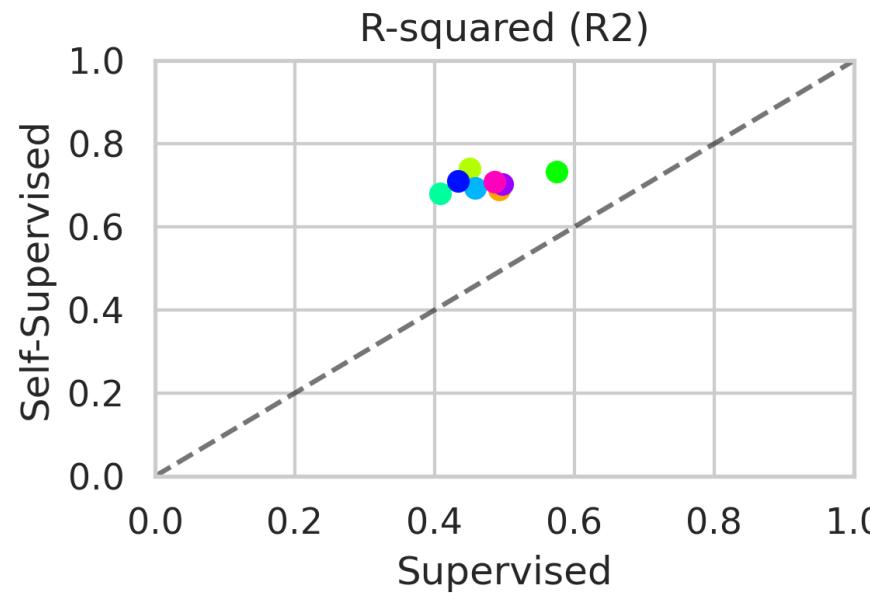
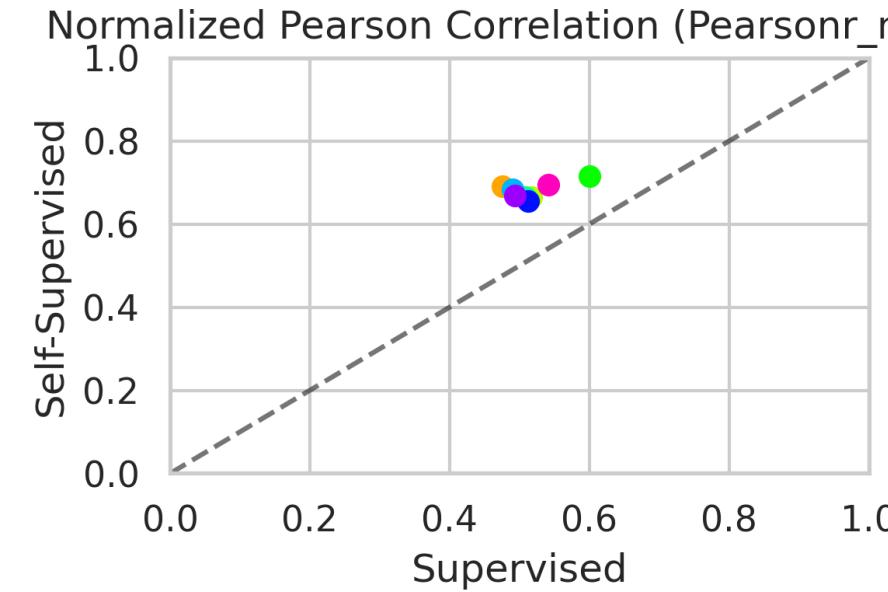
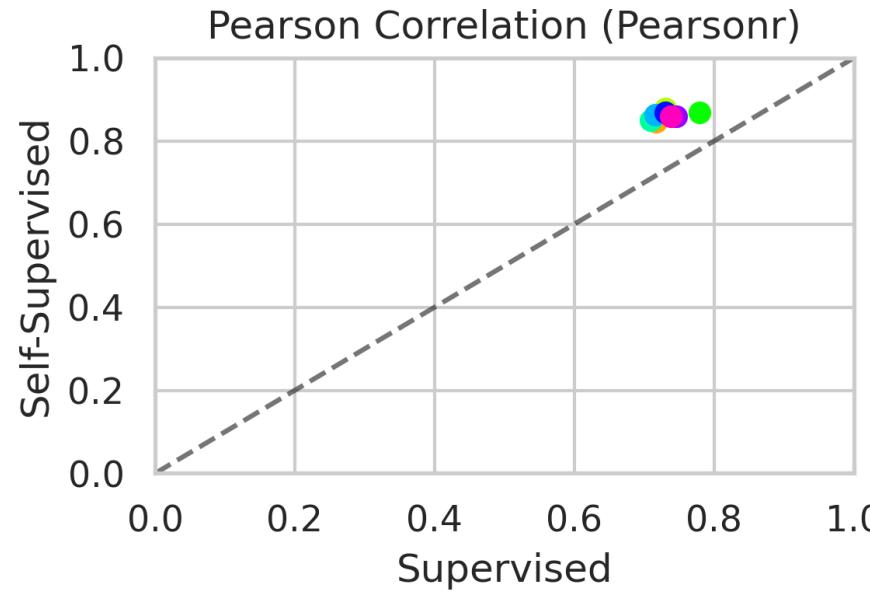
CHiP-exo



All (RNA-Seq + Histone Marks + CHiP-exo)



Average results across 8 folds. Each dot is a track.



Average results
across tracks.
Each dot is a fold



Scratch-trained model vs.

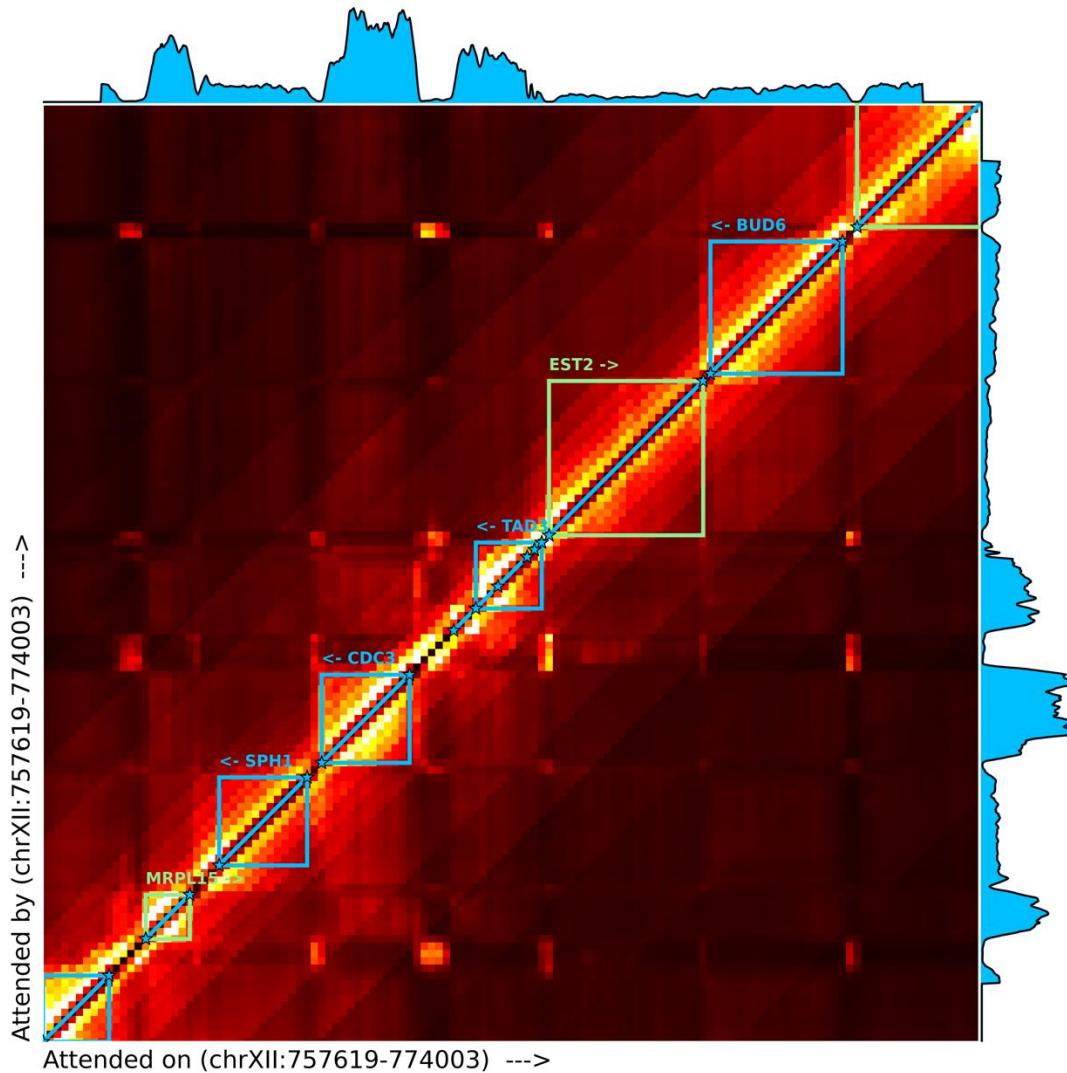
Fine-tuned Fungal LM:

Attention maps

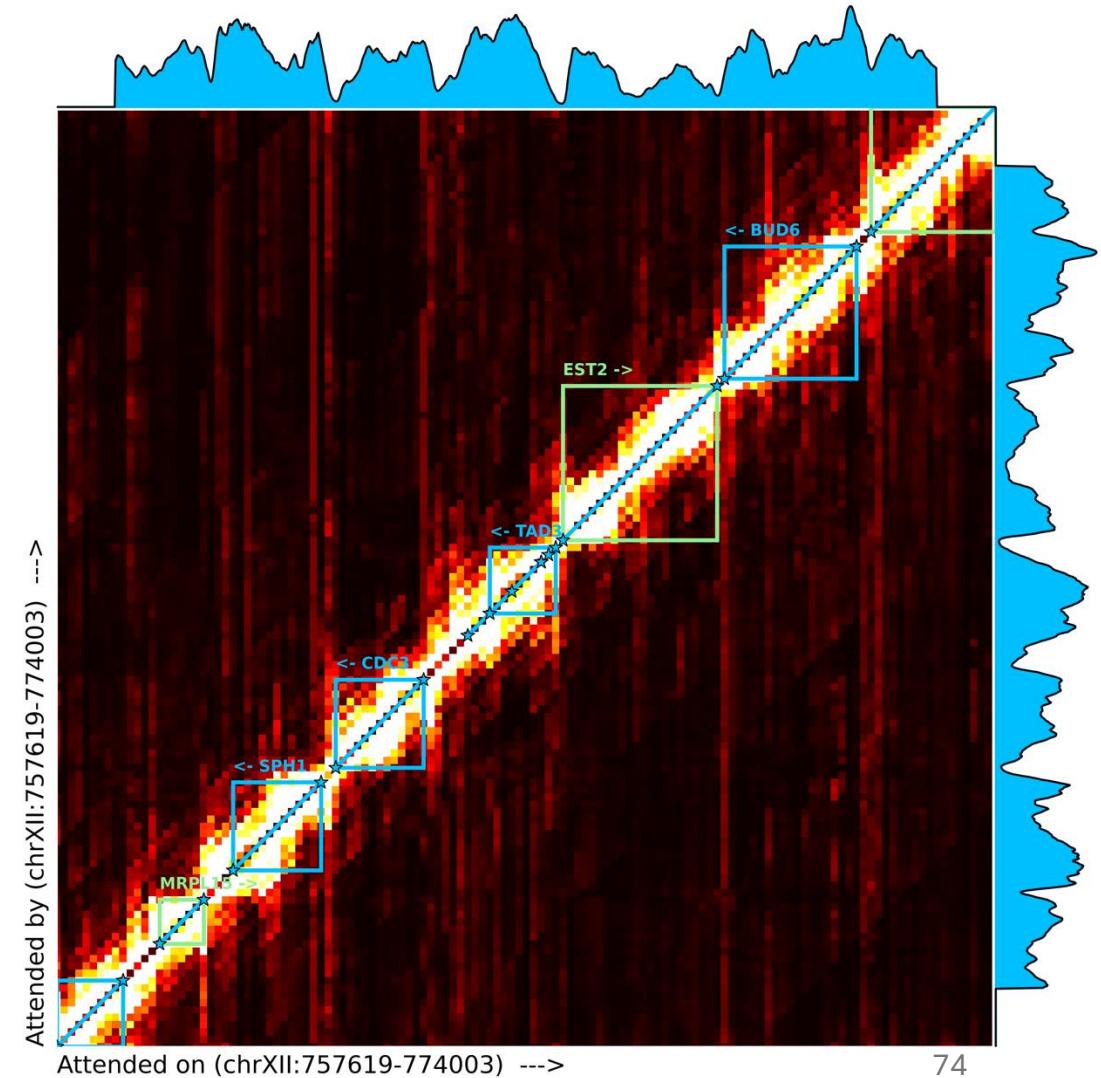


Fine-tuned vs Scratch-Trained

Fine-tuned



Scratch-Trained



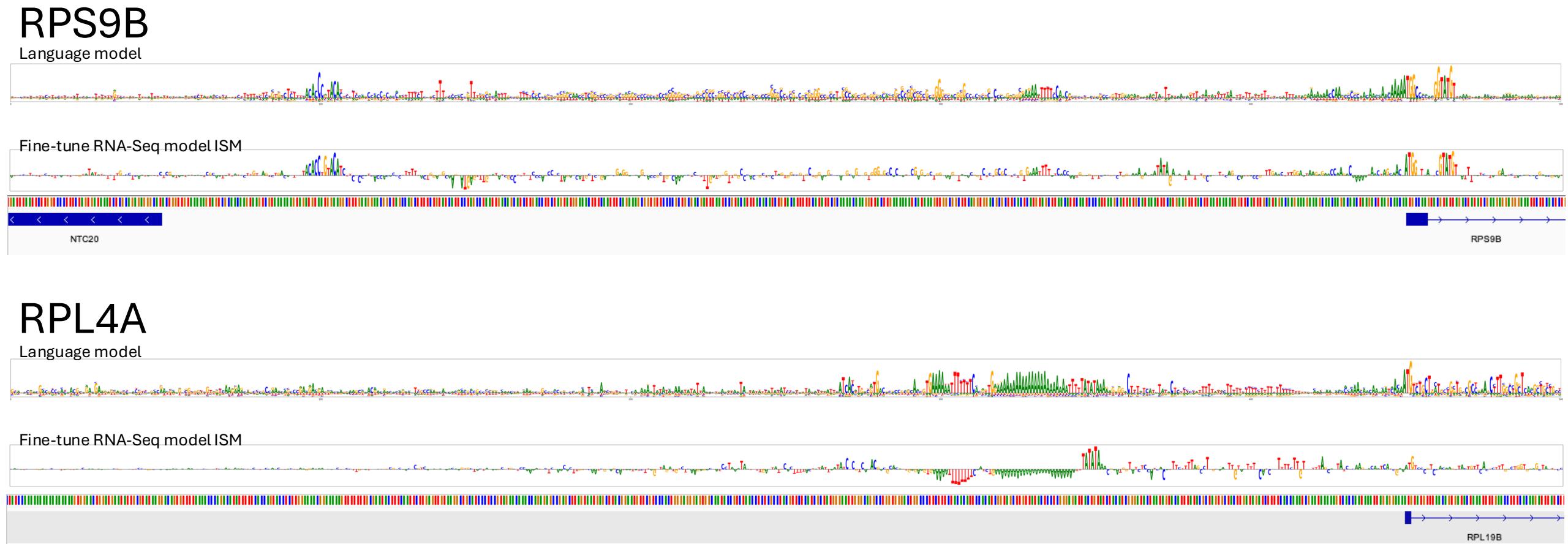
Scratch-trained model vs.

Fine-tuned Fungal LM:

Motif usage



Fine-tuned vs Scratch-Trained



Part III

Applications

1. Assessing the variant effects on eQTLs and negatively selected eQTLs
2. Predicting the influence of distal regulatory elements (i.e. enhancers) on gene expression.
3. MPRA mutation effect prediction

Caudal, É., Loegler, V., Dutreux, F., Vakirlis, N., Teyssonnière, É., Caradec, C., ... & Schacherer, J. (2024). Pan-transcriptome reveals a large accessory genome contribution to gene expression variation in yeast. *Nature Genetics*, 1-10.

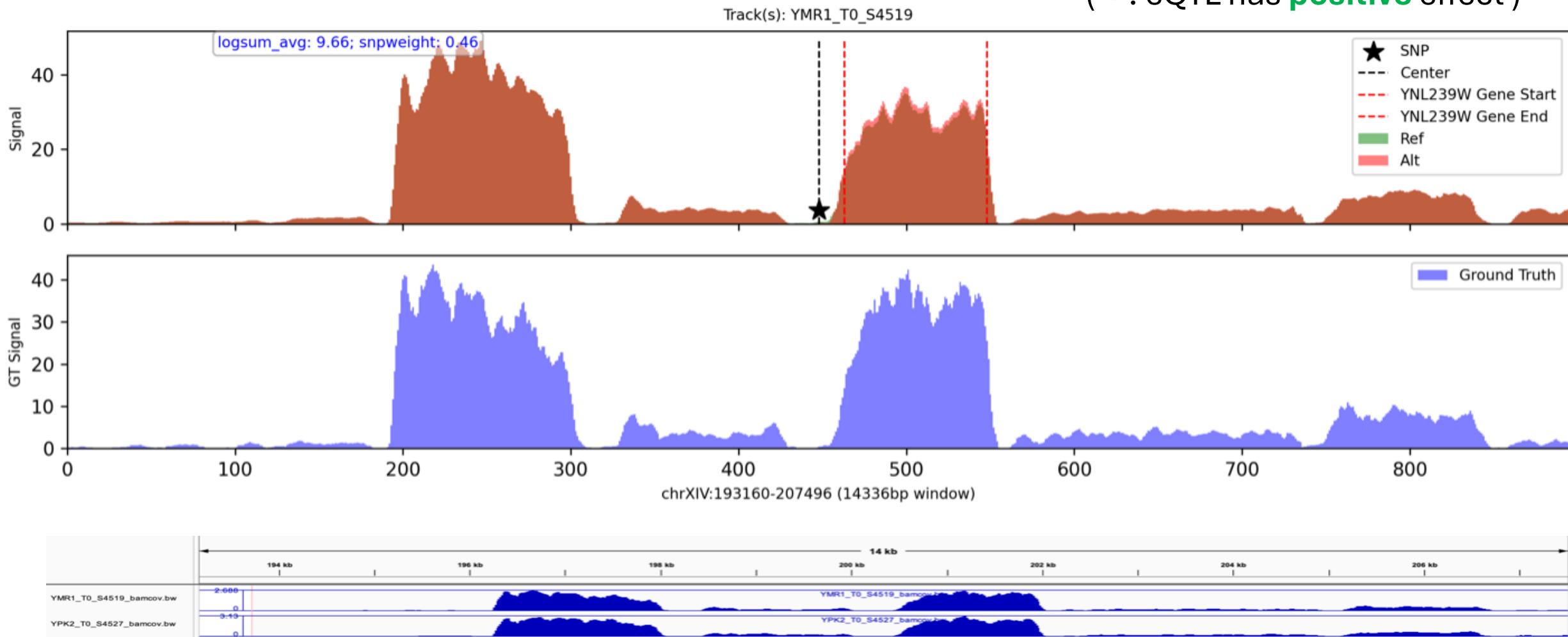
Peter, J., De Chiara, M., Friedrich, A., Yue, J. X., Pflieger, D., Bergström, A., ... & Schacherer, J. (2018). Genome evolution across 1,011 *Saccharomyces cerevisiae* isolates. *Nature*, 556(7701), 339-344.

Predicting eQTLs

SNPWeight from GWAS (+)

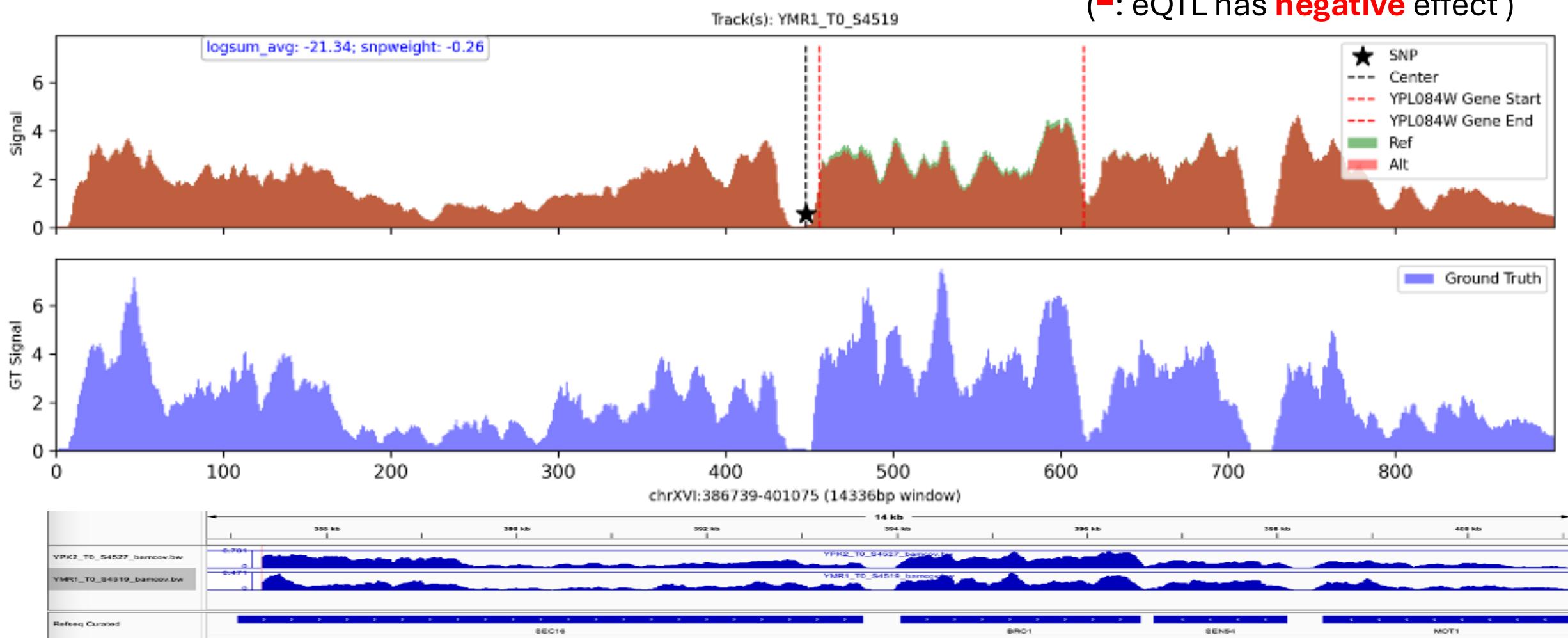
Yeast LM prediction (+)

(+: eQTL has **positive** effect)



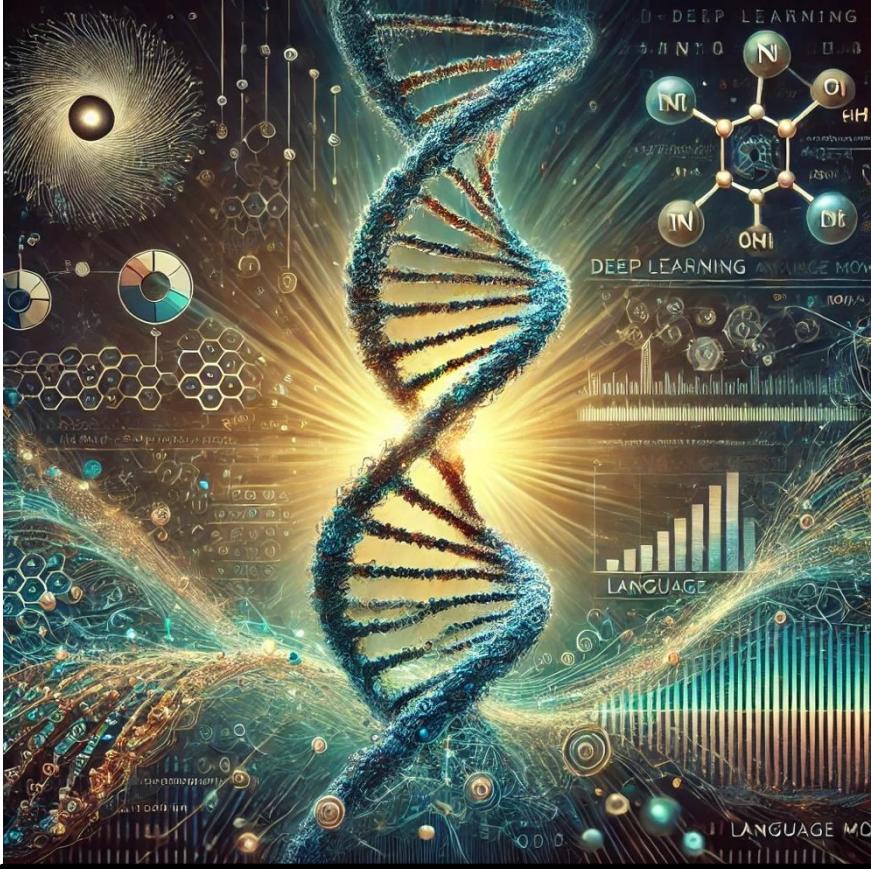
Predicting eQTLs

SNPWeight from GWAS (-)
Yeast LM prediction (-)
(-: eQTL has **negative** effect)



Conclusions

- **Fungal LM:**
 - Learned gene structure
 - Learned conserved regulatory motifs
- **Fine-tuning Fungal LM:**
 - Improved models training from scratch substantially
 - **0.7** Pearson's R in test set
- **Applications**
 - Assessing variant effects on eQTLs
 - Predicting distal regulatory elements influencing gene expression.
 - Predicting mutation effects with MPRA.



ChatGPT prompt:
Generate a figure about deep learning,
genomics, DNA, and language model



khchao.com



@KuanHaoChao



Kuanhao-Chao

Acknowledgement



Johannes Linder



Majed Mohamed Magzoub



David Kelley



Sean Hackett



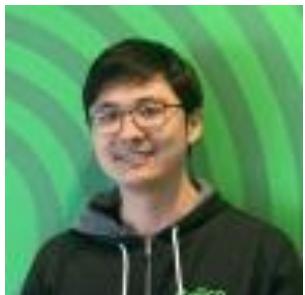
Steven Salzberg



Mihaela Pertea

Kelley Lab & Calico Computing Team

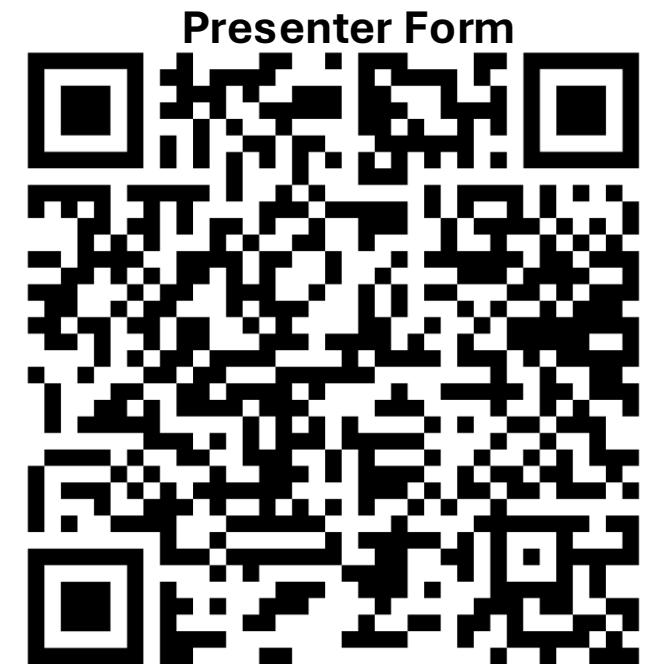
Great mentors, collaborators and good friends!





JHU Deep Learning in Genomics Study Group

- **Date and Time:** Every other Tuesday , 12:00 pm - 1:00 pm.
 - Next meeting 01/28. Celine's presenting "*A foundation model of transcription across human cell types*"
- **Location:** Room 228 at Malone or on Zoom
- **Slack Channel:** #deep-learning-reading-group
- Come join us!!



10/22	Kuan-Hao Chao	Predicting RNA-seq coverage from DNA sequence	https://doi.org/10.1038/s41588-024-02053-6
11/05	Mahler Revsine	HyenaDNA: Long-Range Genomic Sequence Model	https://doi.org/10.48550/arXiv.2306.15794
11/19	Cristina Martin Linares	Machine-guided design of cell-type-targeting cis-regulatory elements	https://doi.org/10.1038/s41586-024-08070-z
12/03	Eduarda Vaz	Effective gene expression prediction from sequencing data	https://doi.org/10.1038/s41592-021-01252-x
1/7	Gus Fridell	Applying interpretable machine learning in computational biology	https://doi.org/10.1038/s41592-024-02359-7
01/21	Stephen Hwang	Evolutionary-scale prediction of atomic-level protein conformations	https://doi.org/10.1126/science.adc2574
02/04	Celine Hoh	A foundation model of transcription across human cell types	https://doi.org/10.1038/s41586-024-08391-z