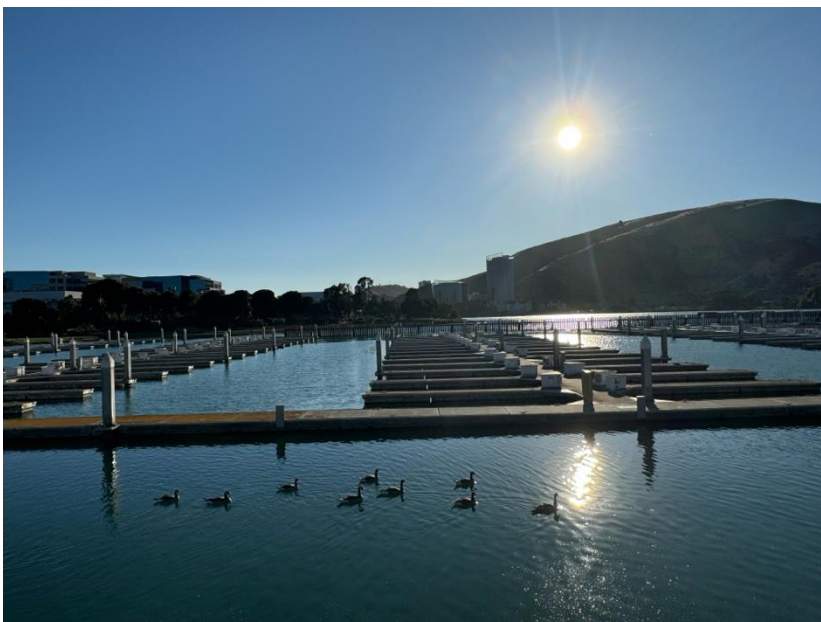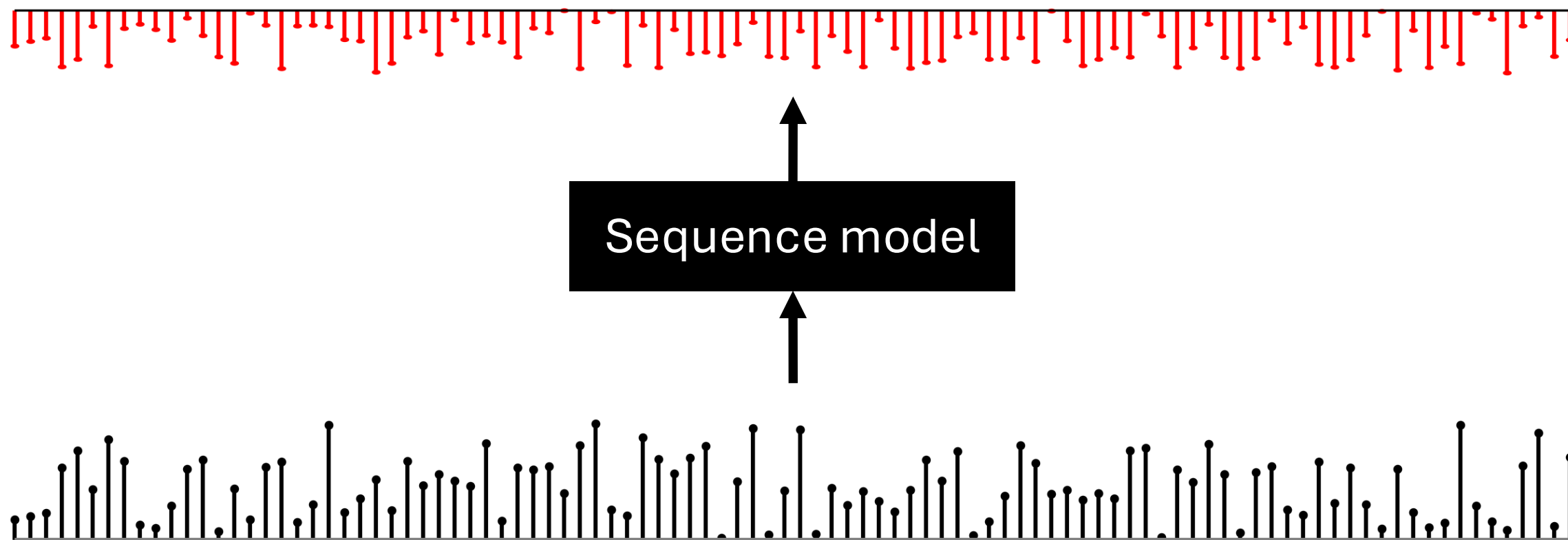Photo with you!

Sequence models map a sequence to a sequence

**Introduction** | Self-supervised LM | Supervised model | Fine-tuning LM

(batch, length, dim)



Linear

Sequence model

Normalization

(batch, length, dim)

## Neural ODEs



## RNN



LSTM

## CNNs



## Transformers



2

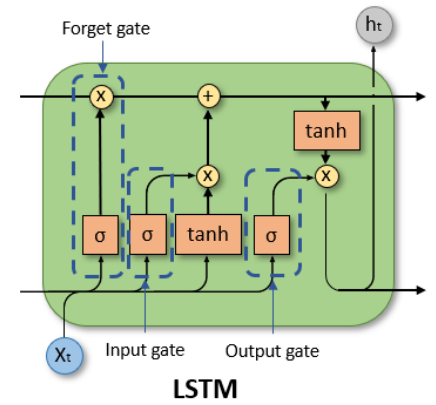**Introduction**     Self-supervised LM     Supervised model     Fine-tuning LM

Input Prompt: Recite the first law of robotics

Output:

TXT

https://jalammar.github.io/how-gpt3-works-visualizations-animations/
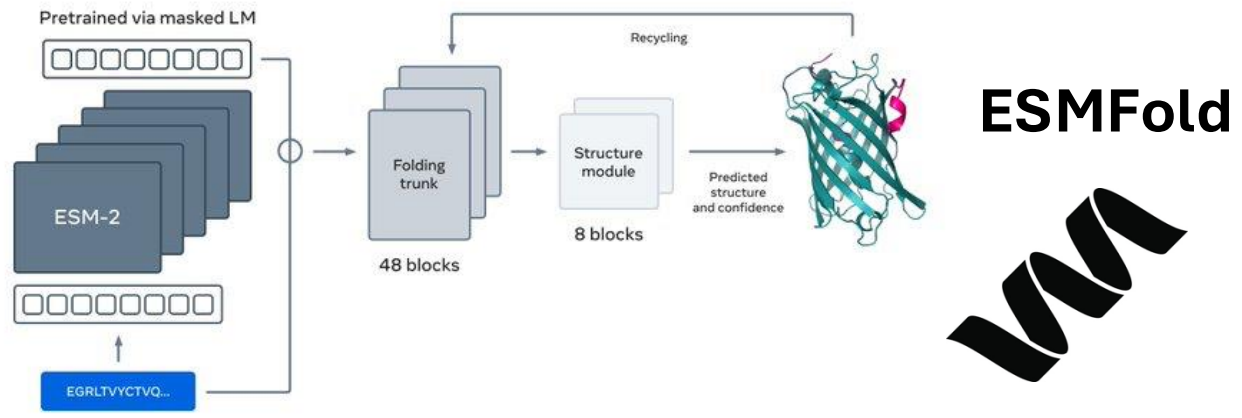
Output

Hidden Layer

Hidden Layer

Hidden Layer

Input

https://deepmind.google/discover/blog/wavenet-a-generative-model-for-raw-audio/

Pretrained via masked LM

ESM-2

Single sequence

EGRLTVYCTVQ...

Recycling

Folding trunk
48 blocks

Structure module
8 blocks

Predicted structure and confidence

ESMFold

https://twitter.com/AIatMeta/status/1587467600413351937/photo/1

Input: DNA sequence

Receptive field
● Basenji2
● Enformer

20 kb
100 kb

Convolutional layers (7x)

Transformer layers (11x)

Attention

Organism specific heads

Human
5,313 tracks

Mouse
1,643 tracks

Output: Genomics tracks

Gene expression (CAGE)

DNA accessibility (DNase)

Histone modification / TF binding (ChIP-seq)

Enformer

https://deepmind.google/discover/blog/predicting-gene-expression-with-a

Introduction    Self-supervised LM    Supervised model    Fine-tuning LM

# Spectrum of Sequential Data



Discrete ←————————————————————→ Continuous

| Text | Graph | DNA | Video | Sound signal | Time-series data |

# Why Deep learning sequence models to DNA ?

# Foundation model

- GPT-3, GPT-4 by OpenAI
- Gemini by Google
- Claude by Anthropic
- Llama 3 by Meta

- Stanford researchers called transformers "foundation models" in an August 2021 paper because they see them driving a paradigm shift in AI.



Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., ... & Liang, P. (2021). On the opportunities and risks of foundation models. arXiv preprint arXiv:2108.07258.

**Introduction** | Self-supervised LM | Supervised model | Fine-tuning LM

# Foundation model

- **Versatility**: wide range of downstream tasks

- **Transfer learning**: learn general representation of data. Task-specific is limited

- **Efficiency**: computational efficiency of fine-tuning models

- **Generalization**: "zero-shot" or "few-shot"

- **Emergent abilities**:
  - basic arithmetic
  - simple programming tasks
  - summarization, translation, or question-answering.

# Goals

- Building an interpretable fungi LLM to help Calico

  construct gene regulatory networks (GRN) in the future.

- Predicting ChIP-exo, histone marks, and RNA-Seq

- Does fine-tuning a pretrained LM outperform training a new model

  from scratch under the exact model architecture?

# Why yeast?

- Simple Eukaryotic Model

- Rapid Growth and Easy Culturing

- Genetic Manipulability

- Well-Characterized Genome

- Conserved Regulatory Mechanisms

# Part I

## Fungi Language Model

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

- Q: What is the quality of the annotation? Coding vs non-coding regions

- Q: How repetitive are the genomes?

# Why building a Fungi Language Model?

- 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.

# Data preprocessing

# Data preprocessing



Repeat regions

Coding regions

16384

4096

~ 7 genes per window

# Data preprocessing

Repeat regions

Coding regions

**16384**

**4096**

~ 7 genes per window

# Data preprocessing



Repeat regions
Coding regions

7% repeat threshold

Training

Validation
(chrXI, chrXIII, chrXV)

Testing
(chrXII, chrXIV, chrXVI)

# 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

**Same species, Different strains**

**Order level**

**Kingdom level**

**Dataset 1**

**Dataset 2**

**Dataset 3**

**Dataset 4**

R64 reference yeast

80 strains of yeasts

165 Saccharomycetales

1361 Fungus genomes

**Q1: Diversity of strains?**

**Q2: Diversity of species?**

**Q3: Even more diverse?**

# Genome distance evaluation

**R64 Reference Yeast** | **80 strains of yeasts** | **165 Saccharomycetales**

Introduction | **Self-supervised LM** | Supervised model | Fine-tuning LM

**Q2:** What is the quality of the annotation?
Coding vs non-coding regions?

# Genome annotation completeness evaluation

**R64 Reference Yeast**



**Conclusion:**
**~95% completeness**

**80 strains of yeasts**



**165 Sachramonycetales**



*Mosè Manni, Matthew R Berkeley, Mathieu Seppey, Felipe A Simão, Evgeny M Zdobnov, BUSCO Update: Novel and Streamlined Workflows along with Broader and Deeper Phylogenetic Coverage for Scoring of Eukaryotic, Prokaryotic, and Viral Genomes. Molecular Biology and Evolution, Volume 38, Issue 10, October 2021, Pages 4647–4654*

The following protocol covers the various BUSCO running modes and workflows, BUSCO setup, guidelines to interpret the results, and additional analyses, e.g., for building phylogenomic trees and visualizing syntenies using BUSCO results:

*Manni, M., Berkeley, M. R., Seppey, M., & Zdobnov, E. M. (2021). BUSCO: Assessing genomic data quality and beyond. Current Protocols, 1, e323. doi: 10.1002/cpz1.323*

Introduction | **Self-supervised LM** | Supervised model | Fine-tuning LM

# Genome evaluation – coding / noncoding regions

**R64 Reference Yeast** *72.46% coding regions*



**80 strains of yeasts**



**165 Sachramonycetales**

# Genome evaluation – # genes per window

## R64 Reference Yeast    *Median: 9.0;  Mean:  8.98*

## 80 strains of yeasts



## 165  Sachramonycetales

Gene Locus 1    Gene Locus 2    Gene Locus 3

# **Q3:** How repetitive are the genomes?

7% repeat threshold

# Genome evaluation – repeat regions

## R64 Reference Yeast    *7.39% repeat regions*



## 80 strains of yeasts



## 165 Sachramonvcetales

# Repeats Detection

- **RepeatModeler**: Identifies de novo transposable element (TE) families.

  - BuildDatabase

  - RepeatModeler

- **RepeatMasker**: Screens DNA sequences for interspersed repeats and low

  complexity DNA sequences using Dfam (or RepBase, **30K**💰) database.

- **Dust**: Masks low-compexity regions

# Repeats masking evaluation



Scatter Plot of Precision vs Recall for fungi_gtf Samples

Introduction    **Self-supervised LM**    Supervised model    Fine-tuning LM

# Data cleaning – repeats removal

a 7% threshold removes
~10% of the sequences.

**Train**                **Test**                **Validation**

Training

Validation (chrXI, chrXIII, chrXV)

Testing (chrXII, chrXIV, chrXVI)

**Q4:** How many homologous sequences are there between training and testing?

Training

Validation
(chrXI, chrXIII, chrXV)

Testing
(chrXII, chrXIV, chrXVI)

Introduction    **Self-supervised LM**    Supervised model    Fine-tuning LM

Training

Validation
(chrXI, chrXIII, chrXV)

Testing
(chrXII, chrXIV, chrXVI)

Detect homologous sequence using **DNA sequence aligner**

# Homology sequence removal

- Nucmer:

  - minimum length of maximal exact matches (MEMs) (20) MEMs shorter than this length will be ignored.

  - A cluster is a group of MEMs that are close to each other and are used to build the alignment (65) Smaller clusters will be ignored

- Minimap2: minimap2 -x asm20

  - - asm5/asm10/asm20: - asm-to-ref mapping, for ~0.1/1/**5%** sequence divergence

# Homology sequence removal evaluation

Introduction    **Self-supervised LM**    Supervised model    Fine-tuning LM

# Homology sequence removal evaluation (minimap2)

**Train - Test**

**Train - Validation**

r64



Strains



**Is it good enough?**

Saccharomycetales

Introduction | **Self-supervised LM** | Supervised model | Fine-tuning LM

# Homology sequence removal evaluation

Introduction | **Self-supervised LM** | Supervised model | Fine-tuning LM

# Homology sequence removal evaluation (minimap2)

# Final sequence for training / testing / validation

**r64**       **80 strains**      **165 Saccharomycetales**

**Before cleaning**

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Train | : | 1440 | | Train | : | 108960 | | Train | : | 404608 |
| Test | : | 608 | | Test | : | 608 | | Test | : | 608 |
| Validation | : | 576 | | Validation | : | 576 | | Validation | : | 576 |

**-597 (-41.4%)**      **-40447 (-37.1%)**      **-65442 (-16.2%)**

**After cleaning**

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Train | : | 843 | | Train | : | 68513 | | Train | : | 339166 |
| Test | : | 507 | | Test | : | 507 | | Test | : | 507 |
| Validation | : | 488 | | Validation | : | 488 | | Validation | : | 488 |

Introduction      **Self-supervised LM**      Supervised model      Fine-tuning LM

# Fungi 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)

**Masked language modeling loss**

$$L_{MLM}^{(x)} = -\frac{1}{|M_x|} \sum_{i \in M_x} log P(x_i / x_{\setminus M_x})$$

where:

$x_{\setminus M_x}$ represents masked version of x

$M_x$ represents set of masked token positions in x

# Self-supervised Fungi LM

# Language Model Results

# Model comparison

# Dataset comparison



Validation Losses

LM R64 U-Net small (valid); loss = 8.9725 (8.9652)
LM R64 U-Net big (valid); loss = 8.9815 (8.9737)
LM strains U-Net small (valid); loss = 8.8472 (8.8390)
LM strains U-Net big (valid); loss = 8.8948 (8.8825)
LM saccharomycetales U-Net small (valid); loss = 8.7766 (8.7555)
LM saccharomycetales U-Net big (valid); loss = 8.7340 (8.7086)

# Different resolutions of input to transformer blocks

**128bp res**

**32bp res**

Transform er Blocks (11x)



**Both resolutions reach the similar loss**

Validation Losses

- - - LM saccharomycetales U-Net small ( 32 res) (valid); loss = 8.7724 (8.7609)
- - - LM saccharomycetales U-Net big ( 32 res) (valid); loss = 8.7598 (8.7382)
- - - LM saccharomycetales U-Net small (128 res) (valid); loss = 8.7658 (8.7586)
- - - LM saccharomycetales U-Net big (128 res) (valid); loss = 8.7391 (8.7243)

# Fungi LM Language Model

# Motif inference

# Constructing PWM from Fungi LM

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

Testing
(chrXII, chrXIV, chrXVI)



|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **A** | .8 | .8 | .0 | .0 | .8 | .0 | .0 | .0 | .0 | .0 | .8 | .0 | .0 | .8 | .0 | .8 |
| **C** | .1 | .1 | .0 | .7 | .1 | .7 | .0 | .0 | .0 | .7 | .1 | .0 | .7 | .1 | .7 | .1 |
| **G** | .1 | .1 | .9 | .1 | .1 | .1 | .9 | .9 | .9 | .1 | .1 | .9 | .1 | .1 | .1 | .1 |
| **T** | .0 | .0 | .1 | .2 | .0 | .2 | .1 | .1 | .1 | .2 | .0 | .1 | .2 | .0 | .2 | .0 |

PolyA tracks

YLR057W; chrXII:255305-257855 (+)

TATA box

YLR438W; chrXII:1012500-1013775 (+)    PolyA tracks

TATA box

YLR056W; chrXII:253860-254958 (+)

TATA box

YLR134W; chrXII:410722-412414 (+)

Initiator (Inr)

YLR015W; chrXII:175226-176744(+)

34

YLR057W; chrXII:255305-257855 (+)

PolyA tracks

YLR438W; chrXII:1012500-1013775 (+)

YLR056W; chrXII:253860-254958 (+)

YLR134W; chrXII:410722-412414 (+)

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

YLR015W; chrXII:175226-176744(+)

# Fungi LM: Summary

1. Fungi language model: The **Saccharomycetales order** is a good evolutionary distance, offering good species diversity.

2. Orthologous gene annotations are **95%** complete.

3. Coding regions make up **50% - 75%** of the genome (**72.46%** in r64). Down-weighting is important!

4. A window size of 16,384 captures approximately **5-10 genes** (**9** in r64).

5. Repetitive regions account for **~2% - 15%** of the genome (**7.39%** in r64). Down-weighting is important!

6. Homologous sequence removal between train-test/validation is crucial (**40% / 60% / 16%**)

7. Transformer-based U-Net architecture overfits in r64 but generalizes best in Saccharomycetales.

8. Self-supervised learning is able to capture cis-regulatory motifs (preliminary results)

# Part II

Supervised ChiP-exo, histone marks, RNA-Seq prediction

Linder, J., Srivastava, D., Yuan, H., Agarwal, V., & Kelley, D. R. (2023). Predicting RNA-seq coverage from DNA sequence as a unifying model of gene regulation. Biorxiv, 2023-08.

# Label data introduction & preprocessing

# ChiP-exo + Histone Marks

- ChIP-exo provides high res view of DNA binding

- Dataset includes 800 ChIP-exo experiments:

- Epigenetic regulators, DNA replication, centromeres, subtelomeres, transposons, RNA polymerase I/II/III

- 161 matched TF ChIP-exo from IDEA 1.0

- Histone Mods MNase-ChIP-seq



Rossi, Matthew et al. Nature. 2021.

# RNA-Seq

- Genome-scale perturbation dynamics propagate signals across regulatory networks

- Measuring dynamics allows events to be ordered

- Aggregating dynamics across many time-courses enables disambiguation of cause > effect relationships

# RNA-Seq

- IDEA (the Induction Dynamics gene Expression Atlas)

# Supervised model architecture

# Basenji Model Training

- Divide genome into 8 folds.

- Train 8 models with distinct

  validation and test folds.



Kristy May for Pets Advisor

Introduction    Self-supervised LM    **Supervised model**    Fine-tuning LM

# Part III

## Fine-tuning Fungi Language Model

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

# Supervised Fungi model VS Fine-tuning Language Model

**16384bp**

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

**16384 * 4**

|   | A | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 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

**Transformer Blocks (8x)**

128bp res 128bp res

16384bp

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

1bp res

...

16bp res

32bp res

64bp res

128bp res

Transformer Blocks (8x)

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

**Coverage Tacks**

16bp res

32bp res

64bp res

128bp res

Introduction  Self-supervised LM  Supervised model  **Fine-tuning LM**

**Supervised**

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

**Input**

16,384 * 4

A
C
G
T

**Fine-tuning LM**

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

16,384 * ( 4 + 1 + 165)

A
C
G
T

Masked encoding

Species encoding
(r64 : **109**)

...

Model

Model

# Fine-tuning vs Training from Scratch (16 bp resolution)



Validation Losses

Legend:
- supervised 16bp unet big f0c0 (valid); loss = 27.1389 (27.0946)
- supervised 16bp unet big f1c0 (valid); loss = 27.2156 (27.1461)
- supervised 16bp unet big f2c0 (valid); loss = 21.3880 (21.3216)
- supervised 16bp unet big f3c0 (valid); loss = 22.3028 (22.2175)
- supervised 16bp unet big f4c0 (valid); loss = 24.9767 (24.8778)
- supervised 16bp unet big f5c0 (valid); loss = 25.4824 (25.3507)
- supervised 16bp unet big f6c0 (valid); loss = 22.0502 (21.9841)
- supervised 16bp unet big f7c0 (valid); loss = 23.5928 (23.4970)
- Fine-tuned LM 16bp unet big f0c0 (valid); loss = 25.8804 (25.8370)
- Fine-tuned LM 16bp unet big f1c0 (valid); loss = 25.9723 (25.9049)
- Fine-tuned LM 16bp unet big f2c0 (valid); loss = 20.3451 (20.3153)
- Fine-tuned LM 16bp unet big f3c0 (valid); loss = 20.9475 (20.9285)
- Fine-tuned LM 16bp unet big f4c0 (valid); loss = 23.8470 (23.8168)
- Fine-tuned LM 16bp unet big f5c0 (valid); loss = 24.1303 (24.1008)
- Fine-tuned LM 16bp unet big f6c0 (valid); loss = 21.2125 (21.1887)
- Fine-tuned LM 16bp unet big f7c0 (valid); loss = 22.4292 (22.4001)

# Fine-tuning vs Training from Scratch (4 bp resolution)



Validation Losses

Introduction   Self-supervised LM   Supervised model   **Fine-tuning LM**

Validation Losses

supervised 16bp U-Net small F0 (valid); loss = 28.6248 (28.4843)
supervised 16bp U-Net small F1 (valid); loss = 28.2678 (28.2117)
supervised 16bp U-Net small F2 (valid); loss = 22.0851 (22.0750)
supervised 16bp U-Net small F3 (valid); loss = 23.1358 (23.0719)
supervised 16bp U-Net small F4 (valid); loss = 25.8020 (25.6840)
supervised 16bp U-Net small F5 (valid); loss = 26.6529 (26.5296)
supervised 16bp U-Net small F6 (valid); loss = 23.0085 (22.9141)
supervised 16bp U-Net small F7 (valid); loss = 24.3076 (24.2682)
Fine-tuned LM 16bp U-Net small F0 (valid); loss = 25.9486 (25.9145)
Fine-tuned LM 16bp U-Net small F1 (valid); loss = 26.2580 (26.1771)
Fine-tuned LM 16bp U-Net small F2 (valid); loss = 20.4224 (20.3980)
Fine-tuned LM 16bp U-Net small F3 (valid); loss = 21.0006 (20.9930)
Fine-tuned LM 16bp U-Net small F4 (valid); loss = 23.9728 (23.9450)
Fine-tuned LM 16bp U-Net small F5 (valid); loss = 24.3522 (24.2979)
Fine-tuned LM 16bp U-Net small F6 (valid); loss = 21.3005 (21.2772)
Fine-tuned LM 16bp U-Net small F7 (valid); loss = 22.5353 (22.5008)

supervised 16bp U-Net big F0 (valid); loss = 27.1389 (27.0946)
supervised 16bp U-Net big F1 (valid); loss = 27.2156 (27.1461)
supervised 16bp U-Net big F2 (valid); loss = 21.3880 (21.3216)
supervised 16bp U-Net big F3 (valid); loss = 22.3028 (22.2175)
supervised 16bp U-Net big F4 (valid); loss = 24.9767 (24.8778)
supervised 16bp U-Net big F5 (valid); loss = 25.4824 (25.3507)
supervised 16bp U-Net big F6 (valid); loss = 22.0502 (21.9841)
supervised 16bp U-Net big F7 (valid); loss = 23.5928 (23.4970)
Fine-tuned LM 16bp U-Net big F0 (valid); loss = 25.7919 (25.7356)
Fine-tuned LM 16bp U-Net big F1 (valid); loss = 25.8828 (25.8564)
Fine-tuned LM 16bp U-Net big F2 (valid); loss = 20.2898 (20.2702)
Fine-tuned LM 16bp U-Net big F3 (valid); loss = 20.9088 (20.9180)
Fine-tuned LM 16bp U-Net big F4 (valid); loss = 23.8441 (23.8013)
Fine-tuned LM 16bp U-Net big F5 (valid); loss = 24.1333 (24.0781)
Fine-tuned LM 16bp U-Net big F6 (valid); loss = 21.2041 (21.1740)
Fine-tuned LM 16bp U-Net big F7 (valid); loss = 22.4327 (22.4049)

Validation Loss

# Training Batches

Introduction      Self-supervised LM      Supervised model      **Fine-tuning LM**

# RNA-Seq track visualization

Introduction   Self-supervised LM   Supervised model   **Fine-tuning LM**
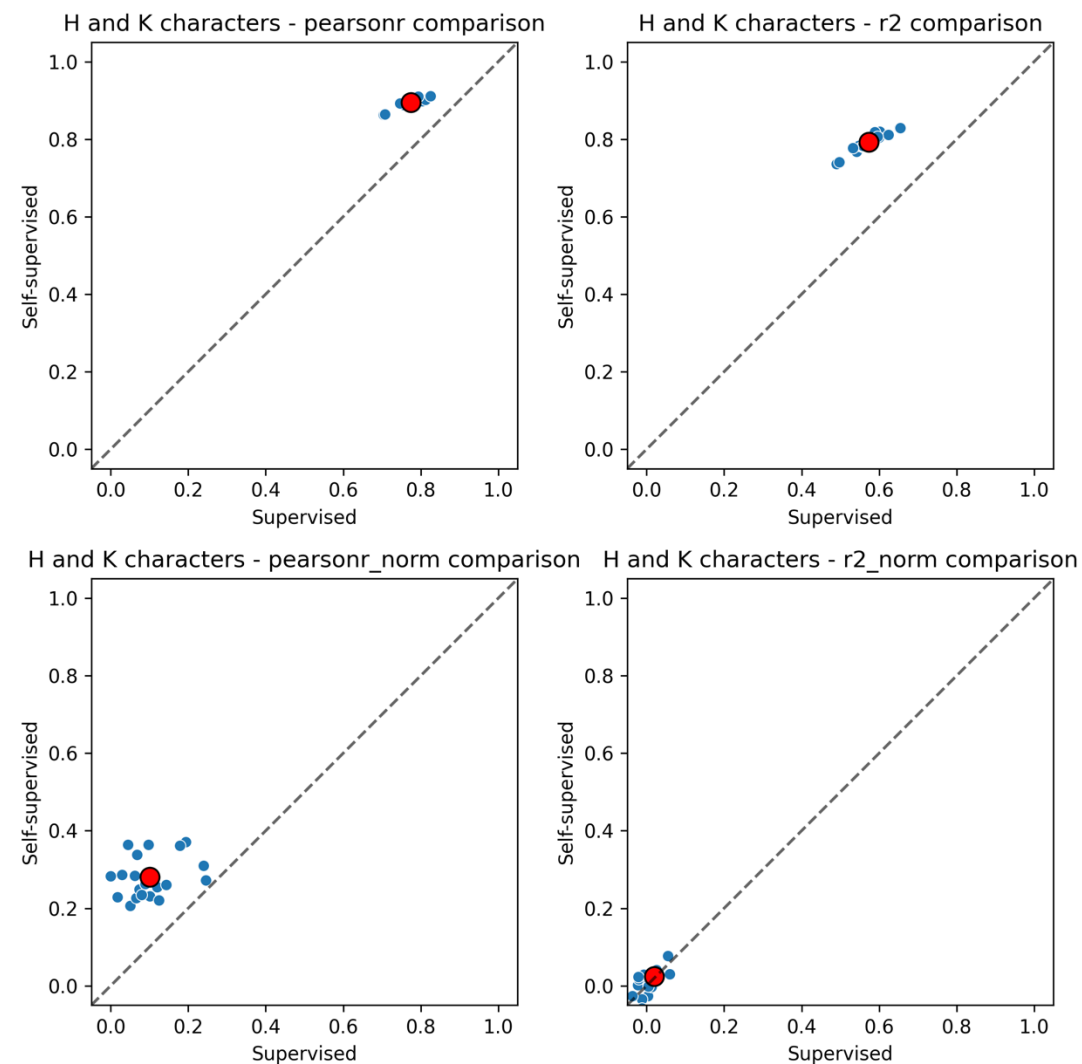
# Track level prediction evaluation

# RNA-Seq

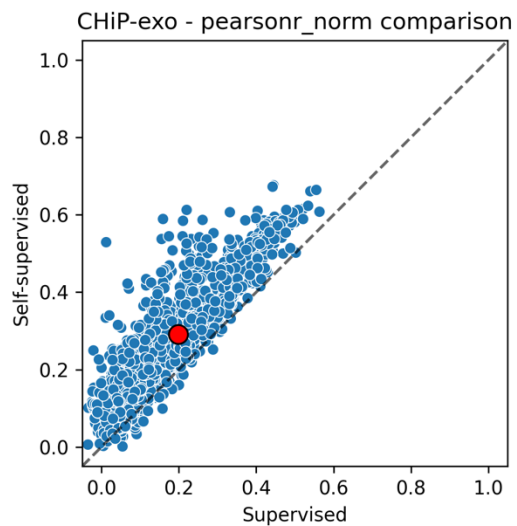# Histone Marks



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

Introduction | Self-supervised LM | Supervised model | **Fine-tuning LM**

# CHiP-exo

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



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

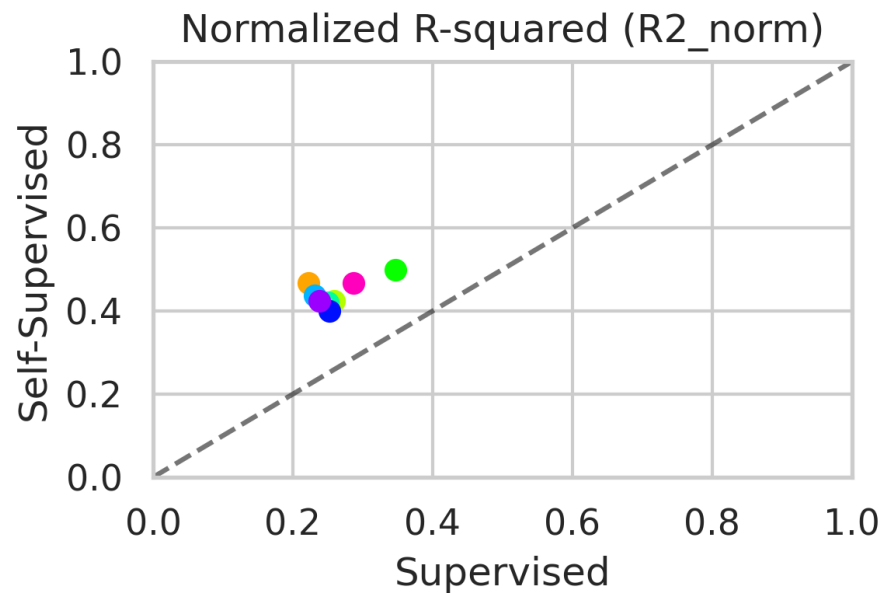Introduction | Self-supervised LM | Supervised model | **Fine-tuning LM**
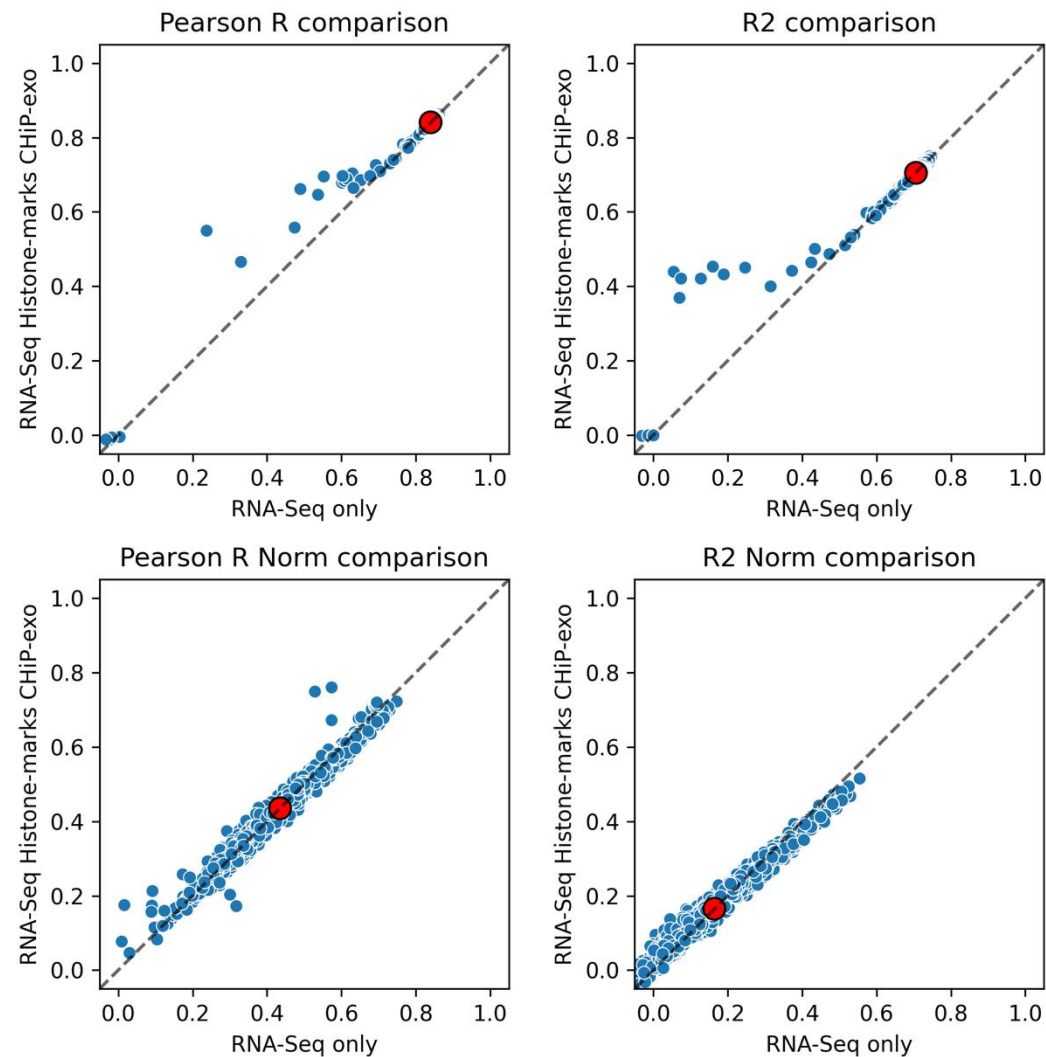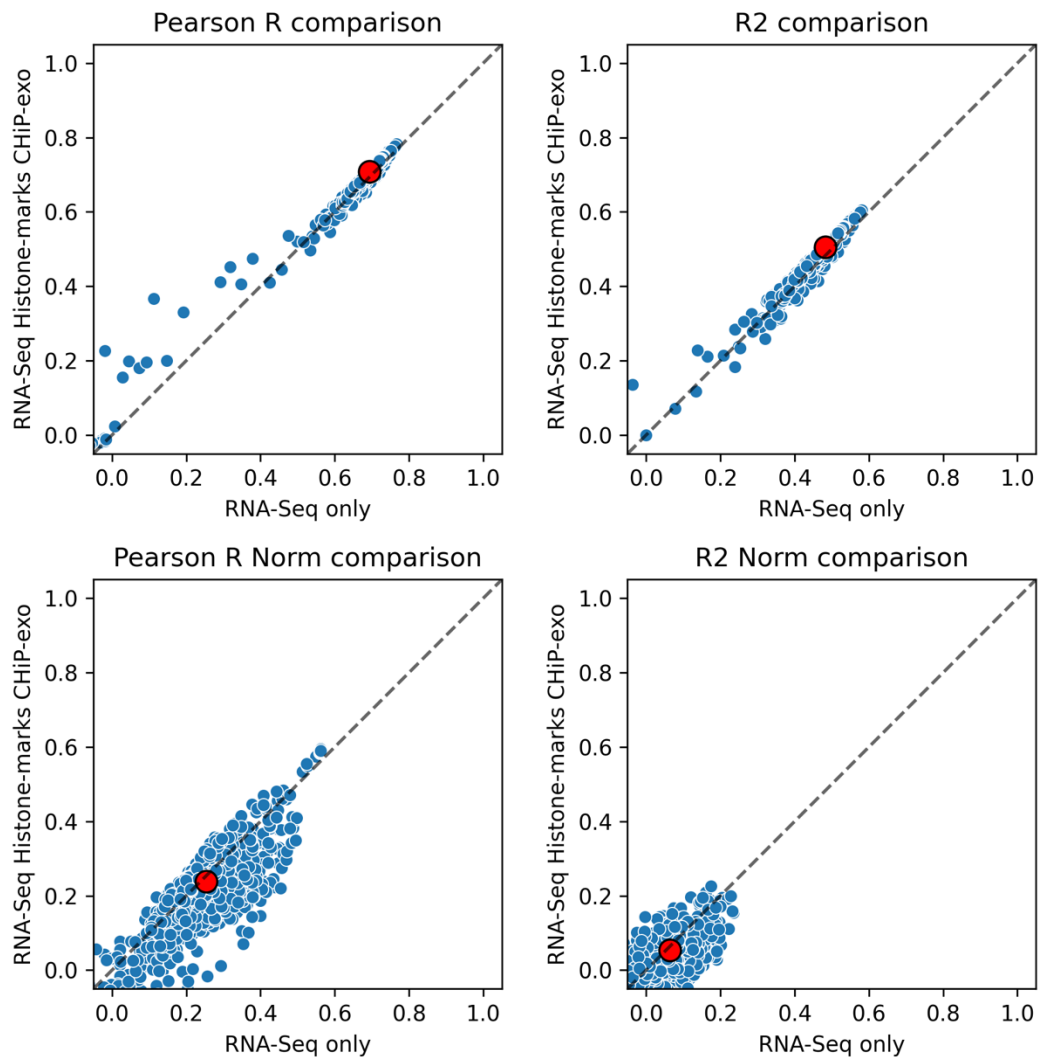
Average results across tracks. Each dot is a fold

RNA-Seq tracks alone  VS

RNA-Seq + Histone Marks +

CHiP-exo tracks

# Supervised trained models

# Self-supervised trained models



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

Introduction  Self-supervised LM  Supervised model  **Fine-tuning LM**

# Project Conclusion

1. Built the first fungi language model. The Saccharomycetales order is a good evolutionary distance, offering good species diversity. Processing 1361 fungus genomes.

2. Under the exact model architecture, pretrained LM weights & fine-tuning can outperform training a model from scratch.

   - Loss / gene level Pearson R / gene level $R^2$

# Acknowledgement



Johannes Linder   Majed Mohamed Magzoub   David Kelley        Sean Hackett

Kelley Lab & Calico Computing Team          **Great mentors, collaborators and good friends!**