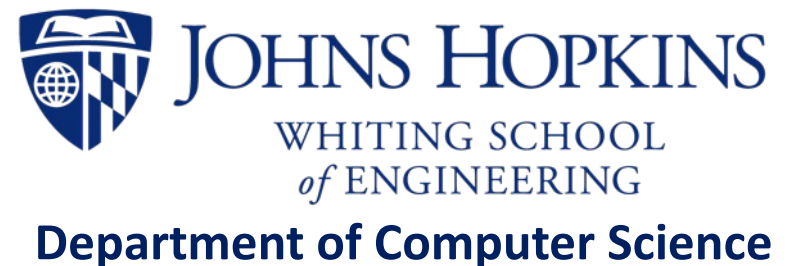# Predicting splice sites in DNA sequences with sequence models

2024.05.15
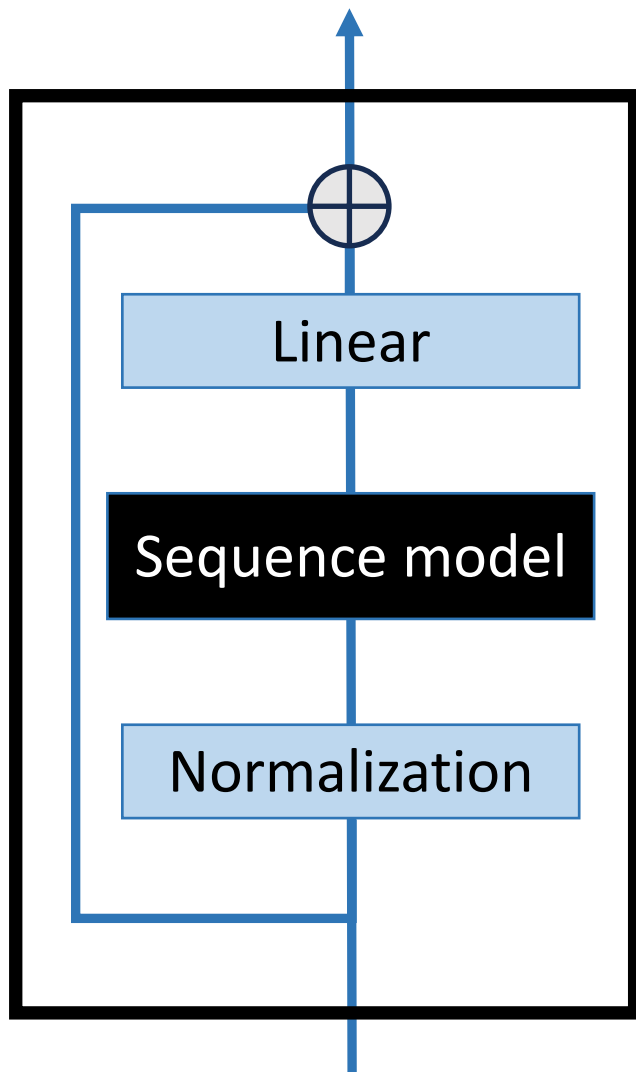
Kuan-Hao Chao

🌐 **khchao.com**      𝕏 **@KuanHaoChao**      **Kuanhao-Chao**

JOHNS HOPKINS UNIVERSITY
CENTER FOR COMPUTATIONAL BIOLOGY
**CCB**

JOHNS HOPKINS
WHITING SCHOOL
*of* ENGINEERING
**Department of Computer Science**

Sequence models map a sequence to a sequence
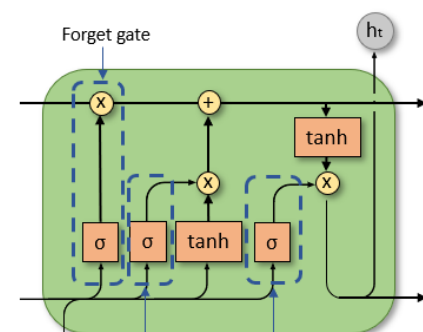
(batch, length, dim)

Linear

Sequence model

Normalization

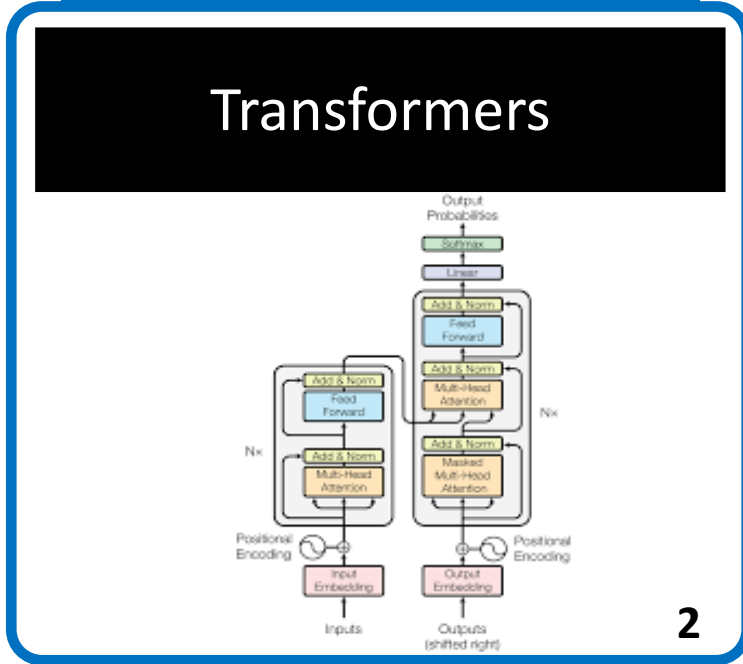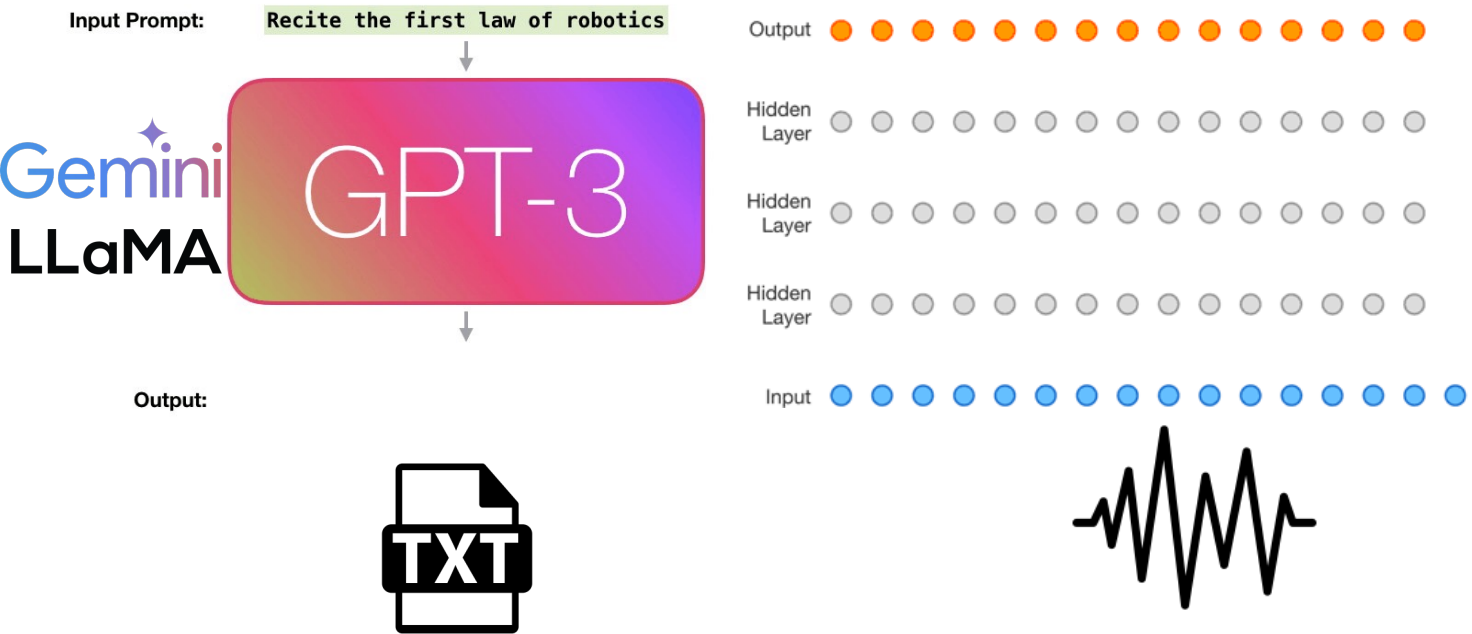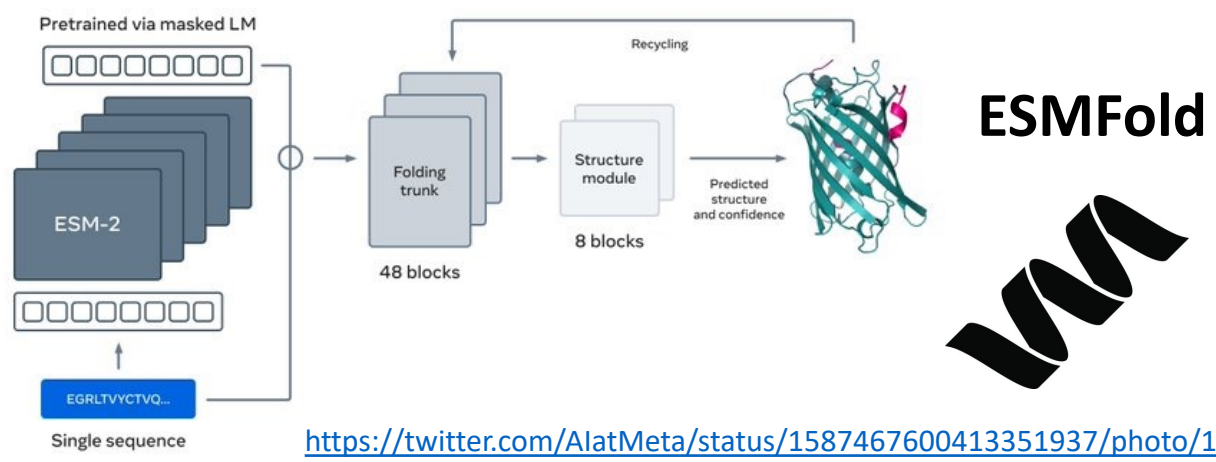(batch, length, dim)

Neural ODEs



RNN



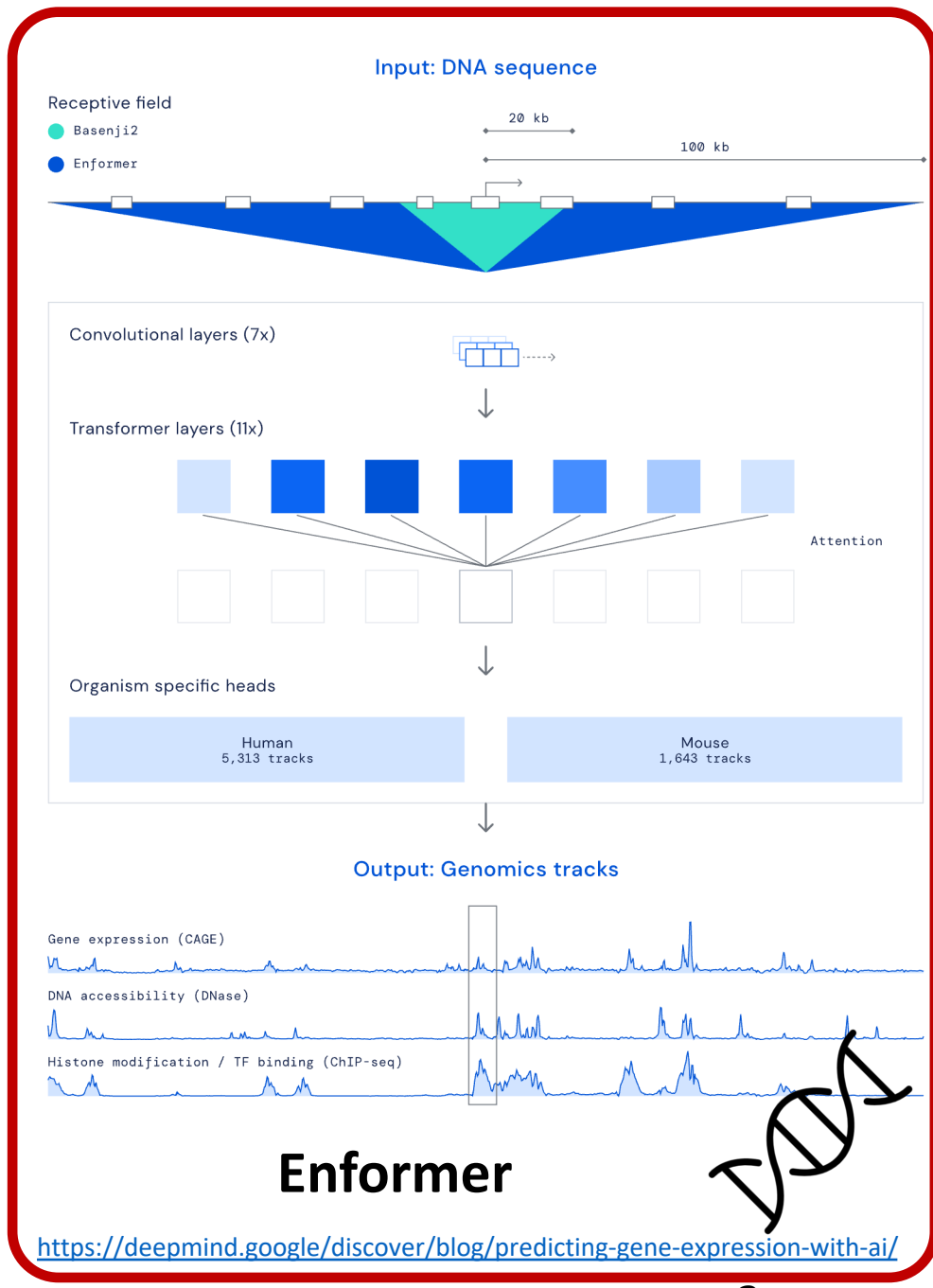**Future work**

CNNs



Transformers



Figure 1: The Transformer - model architecture.

2

Introduction    SpliceAI-toolkit    Splam    Future work

Input Prompt: Recite the first law of robotics

Gemini
LLaMA
GPT-3

Output:

Output
Hidden Layer
Hidden Layer
Hidden Layer
Input

TXT

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

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

Pretrained via masked LM
ESM-2
Single sequence
EGRLTVYCTVQ...
Folding trunk
48 blocks
Structure module
8 blocks
Predicted structure and confidence
Recycling

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-ai/

3

# Spectrum of Sequential Data



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

Text　　　Graph　　　DNA　　　Video　　Sound signal　　Time-series data

# Why Convolutional Neural Network to DNA ❓

# Why Convolutional Neural Network to DNA ❓

2015   2016   2017   2018   2019   2020   2021   2022   2023   2024

# Why Convolutional Neural Network to DNA ❓

# Why Convolutional Neural Network to DNA ❓

# "All models are wrong, but some are useful"

## - George Box, 1978

Box, G. E., Jenkins, G. M., Reinsel, G. C., & Ljung, G. M. (2015).
*Time series analysis: forecasting and control*. John Wiley & Sons.

# SpliceAI-toolkit: splice site predictor

Steven Salzberg

Mihaela Pertea

Anqi Liu

**Pre-mRNA**

**Chao, K. H.**, Mao, A., Liu, Anqi, Salzberg, S. L., & Pertea, M. (<u>2024</u>). SpliceAI-toolkit. Manuscript in preparation. 📕 **https://ccb.jhu.edu/spliceai-toolkit/**

**Chao, K. H.**, Mao, A., Salzberg, S. L., & Pertea, M. (<u>2023</u>). Splam: a deep-learning-based splice site predictor that improves spliced alignments. **bioRxiv.** 📕 **https://ccb.jhu.edu/splam/**

# Can we predict splice sites using only DNA?
## Yes!



X    AGACTCAGCCCCGGAGACTTAGTTAGAGGAAGAAAAAGGTAGGACAGAAGAAAAAGGCAGGACATACAAGGTGCTGGCCCAGGGCGG

Y    0000000000000000000[2]0000000[1]0000000[2]00000000000[1]00000000000000000000000000000000

▲ Donor: 2          ◆ Acceptor: 1          Neither: 0

# SpliceAI: splite site predictor

$$\begin{bmatrix} A \\ C \\ G \\ T \\ N \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} Acceptor \\ Donor \\ Neither \\ Padding \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

**Input sequence (len: L)**

Flanking sequence (len: 5000)     Flanking sequence (len: 5000)

A T G          C

Dimension:

**X:** $(L+10000) * 4$

How do we train the model given that gene sequences vary?

**SpliceAI**

| .1 | 0 | .6 | | .2 |
|----|---|----|--|----|
| .1 | .1 | .2 | | .3 |
| .8 | .9 | .2 | | .5 |

=> **Donor**
=> **Acceptor**
=> **Neither**

**Y:**      $L * 3$

**Predicting Splicing from Primary Sequence with Deep Learning**

Kishore Jaganathan,[1,6] Sofia Kyriazopoulou Panagiotopoulou,[1,6] Jeremy F. McRae,[1,6] Siavash Fazel Darbandi,[2] David Knowles,[3] Yang I. Li,[3] Jack A. Kosmicki,[1,4] Juan Arbelaez,[2] Wenwu Cui,[1] Grace B. Schwartz,[2] Eric D. Chow,[5] Efstathios Kanterakis,[1] Hong Gao,[1] Amirali Kia,[1] Serafim Batzoglou,[1] Stephan J. Sanders,[2] and Kyle Kai-How Farh[1,7,*]

# SpliceAI: data preprocessing

# SpliceAI: data preprocessing

# SpliceAI: data preprocessing

# SpliceAI: data preprocessing

X  NN … NN**ATGTCGTGTCGAGTTGTCGTGTTCAGGTCAGTCAGGTCAGTAAGTAGAGCTCA**NN … NN

Y  0000000**2**00000**1**0000000**2**00000**1**00000**2**000**1**00000**2**0000000**1**000000

Shape: $\left\lceil \dfrac{L}{W} \right\rceil * \left( \dfrac{F}{2} + W + \dfrac{F}{2} \right)$

Shape: $\left\lceil \dfrac{L}{W} \right\rceil * (W)$

$$X = \begin{bmatrix} \text{NN … NN} \textbf{ATGTCGTGTCGAGTTG} \\ \textbf{TCGTGTCGAGTTGTCGTGTTCAG} \\ \textbf{TTGTCGTGTTCAGGTCAGTCAGG} \\ \vdots \\ \textbf{AAGTAGAGCTCANNN} \text{NN … NN} \end{bmatrix}$$

$$Y = \begin{bmatrix} 0000000\textbf{2}00 \\ 000\textbf{1}000000 \\ \textbf{2}000001\textbf{0}00 \\ \vdots \\ 0000000000 \end{bmatrix}$$

# SpliceAI-toolkit : better than SpliceAI!

# SpliceAI-toolkit : retrain on different species



**A** Splice site prediction metrics for mouse

**B** Splice site prediction metrics for zebrafish

**C** Splice site prediction metrics for arabadop

**D** Splice site prediction metrics for bee

# SpliceAI-toolkit : new concept – Calibration

**Input sequence (len: L)**

Flanking sequence (len: 5000)  Flanking sequence (len: 5000)

| A | T | G | | C |

X

**SpliceAI-MANE**

| .1 | 0 | .6 | | .2 |
| .1 | .1 | .2 | | .3 |
| .8 | .9 | .2 | | .5 |

=> **Donor**

=> **Acceptor**

=> **Neither**

Y

- What do SpliceAI-MANE scores signify?

- Do the model's predicted probabilities accurately represent the true likelihood of an event occurring?"

# SpliceAI-toolkit : new concept – Calibration



Model → 0.99

0          0.2         0.4         0.6         0.8        1.0

Model predicted probability

15

# SpliceAI-toolkit : new concept – Calibration



**Model** → 0.00



0     0.2     0.4     0.6     0.8     1.0

Model predicted probability

# SpliceAI-toolkit : new concept – Calibration



**Model** → 0.60



| 0 | 0.2 | 0.4 | 0.6 | 0.8 | 1.0 |
|---|---|---|---|---|---|

Model predicted probability

15

# SpliceAI-toolkit : new concept – Calibration



Model → 0.62



| 0 | 0.2 | 0.4 | 0.6 | 0.8 | 1.0 |

Model predicted probability

15

# SpliceAI-toolkit : new concept – Calibration

# SpliceAI-toolkit : new concept – Calibration



Empirical probability (fraction of chihuahua) vs Model predicted probability

# SpliceAI-toolkit : new concept – Calibration

# SpliceAI-toolkit : new concept – Calibration

# SpliceAI-toolkit : new concept – Calibration

Empirical probability
(fraction of chihuahua)

Model predicted probability

15

# SpliceAI-toolkit : new concept – Calibration

Introduction    **SpliceAI-toolkit**    Splam    Future work

SpliceAI-toolkit : new concept – Calibration

# SpliceAI-toolkit : new concept – Calibration

15

# SpliceAI-toolkit : new concept – [calibr]ation

SpliceAI-80nt

Input

Conv(32, 1, 1)

Conv(32, 1, 1)

RB(32, 11, 1)

RB(32, 11, 1)

RB(32, 11, 1)

RB(32, 11, 1)

Conv(32, 1, 1)

+

Conv(3, 1, 1)

- Platt scaling (Temperature scaling)

$$P_i = \frac{e^{\frac{y_i}{T}}}{\sum_{k=1}^{n} e^{\frac{y_k}{T}}}$$

**Logits($y_i$)**

**Calibration variable $T$**

Softmax

**Probability ($P_i$)**

Output

Platt, John. "Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods." Advances in large margin classifiers 10.3 (1999): 61-74.

$i \in \{Neither, Acceptor, Donor\}$

**16**

# SpliceAI-toolkit : new concept – Calibration

$$conf(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \hat{p}_i \qquad acc(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} 1(\hat{y}_i = y_i)$$

Before calibration - **NLL** : 0.13310300, **ECE**: 0.00001282

Optimal temperature: 1.28049

After calibration - **NLL**: 0.11896934, **ECE** : 0.00000079

$$\mathcal{L}_{NLL} = -\sum_{i=1}^{n} \log(\hat{\pi}(y_i|x_i))$$

$$\mathcal{L}_{ECE} = \sum_{m=1}^{M} \frac{|B_m|}{n} |acc(B_m) - conf(B_m)|$$



Guo, Chuan, et al. "On calibration of modern neural networks." International conference on machine learning. PMLR, 2017.

Introduction | **SpliceAI-toolkit** | Splam | Future work

# SpliceAI-toolkit : Summary

1. Data preprocessing: sliding window chunking

2. Easy-to-run framework to train your own SpliceAI

3. Pretrained SpliceAI-MANE

4. Pretrained SpliceAI on different species

5. Predict genetic variants' effect on splice sites    **Problem solved?**

6. Model calibration: temperature scaling

**Chao, K. H.**, Mao, A., Liu, Anqi, Salzberg, S. L., & Pertea, M. (**2024**). SpliceAI-toolkit. Manuscript in preparation. 📕 **https://ccb.jhu.edu/spliceai-toolkit/ (in preparation)**

🚨 Is canonical labelling approach correct?

# SPLAM : Data Processing

Donor: 400bp

Acceptor: 400bp

Intron          Intron          Intron          Intron

Exon        Exon        Exon        Exon        Exon

Splice junction

Donor          Acceptor

**Chao, K. H.**, Mao, A., Salzberg, S. L., & Pertea, M. (2023). Splam: a deep-learning-based splice site predictor that improves spliced alignments. **bioRxiv.** **https://ccb.jhu.edu/splam/**

# : Splam Model Architecture



**Chao, K. H.**, Mao, A., Salzberg, S. L., & Pertea, M. (2023). Splam: a deep-learning-based splice site predictor that improves spliced alignments. **bioRxiv.** **https://ccb.jhu.edu/splam/**

# SPL✂M : deep-learning splice site predictor



**Score stability**

**Transcriptome assembly improvement**

Chao, K. H., Mao, A., Salzberg, S. L., & Pertea, M. (2023). Splam: a deep-learning-based splice site predictor that improves spliced alignments. **bioRxiv.** 📕 **https://ccb.jhu.edu/splam/**

# SPLAM : deep-learning splice site predictor

Interpretability: ablation study

Interpretability: input sequence



Chao, K. H., Mao, A., Salzberg, S. L., & Pertea, M. (2023). Splam: a deep-learning-based splice site predictor that improves spliced alignments. **bioRxiv.** **https://ccb.jhu.edu/splam/**

# SPL✂M : Summary

- Better than Spli...

- Generalize to no...

$$Loss_{CLE} = - \sum_{class \in \{donor, acceptor, neither\}} I_{class} \times \log(P_{class}) \quad (2)$$

$$Loss_{FL} = - \sum_{class \in \{donor, acceptor, neither\}} I_{class} \times (1 - P_{class})^{\gamma} \times \log(P_{class}) \quad (3)$$

## Technical takeaways

- Focal loss improves cross entropy loss

- Learing rate warm up + sinusoidal decay

- Residual connection is powerful

- Grouped convolution helps (cardinality)
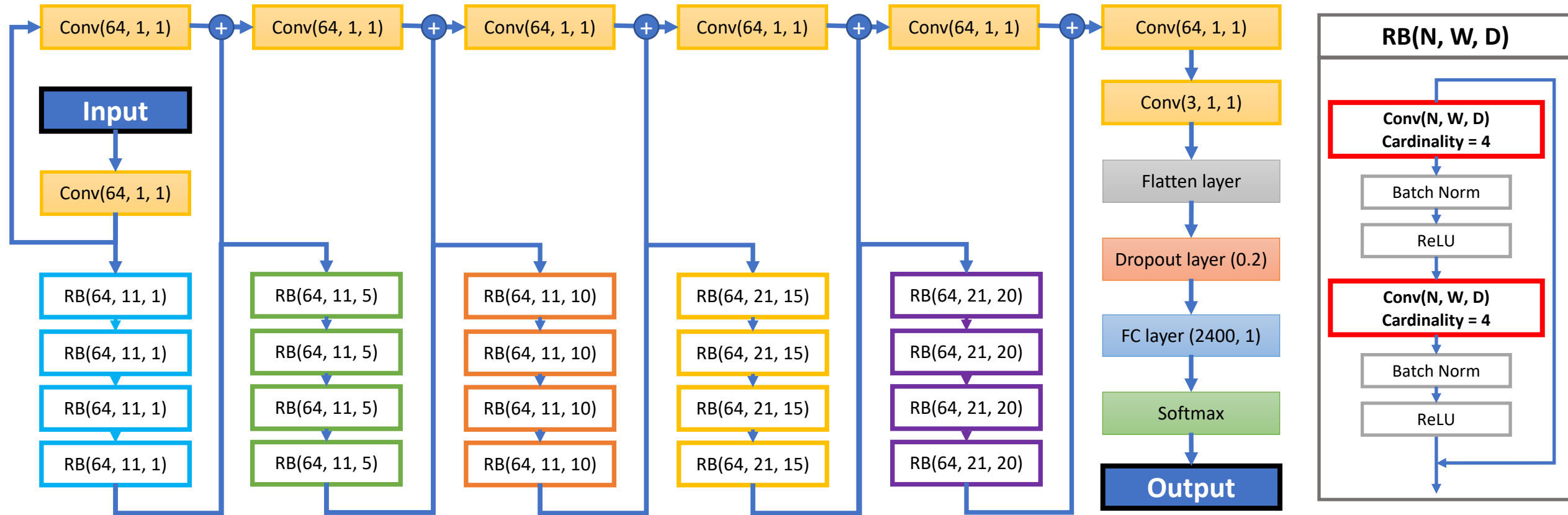


**Chao, K. H.**, Mao, A., Salzberg, S. L., & Pertea, M. (2023). Splam: a deep-learning-based splice site predictor that improves spliced alignments. **bioRxiv.** 📕**https://ccb.jhu.edu/splam/**
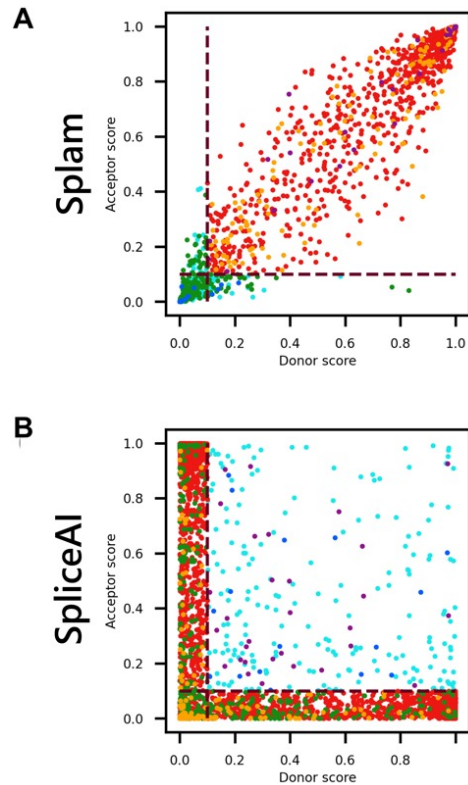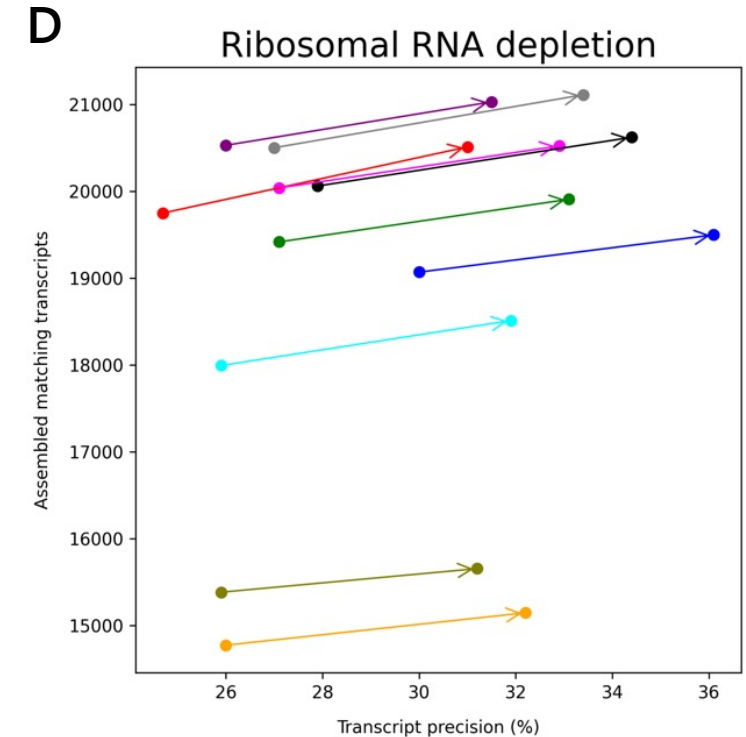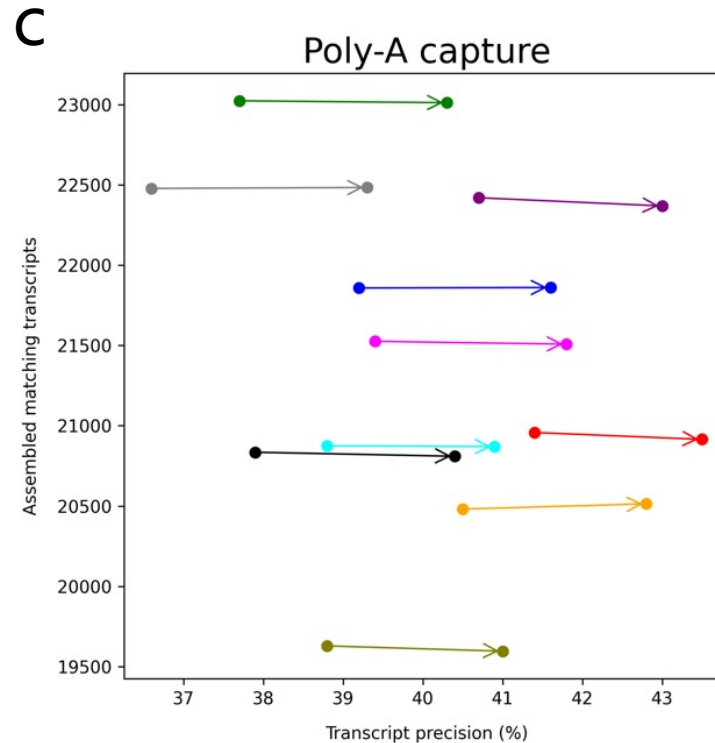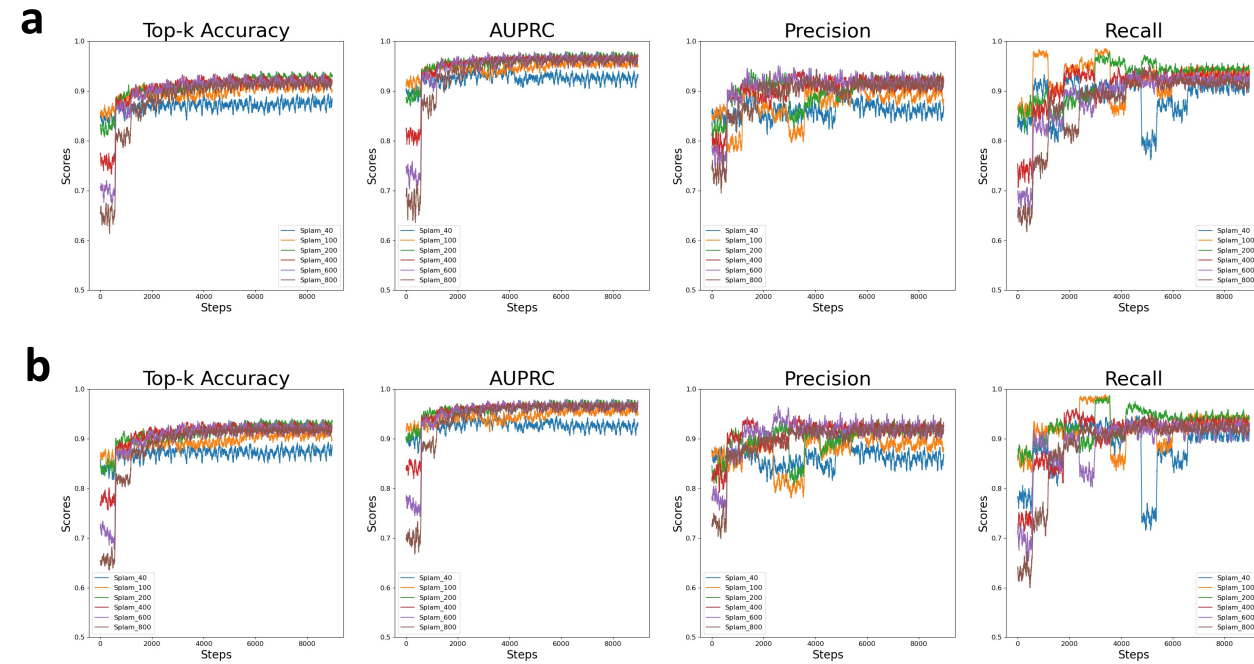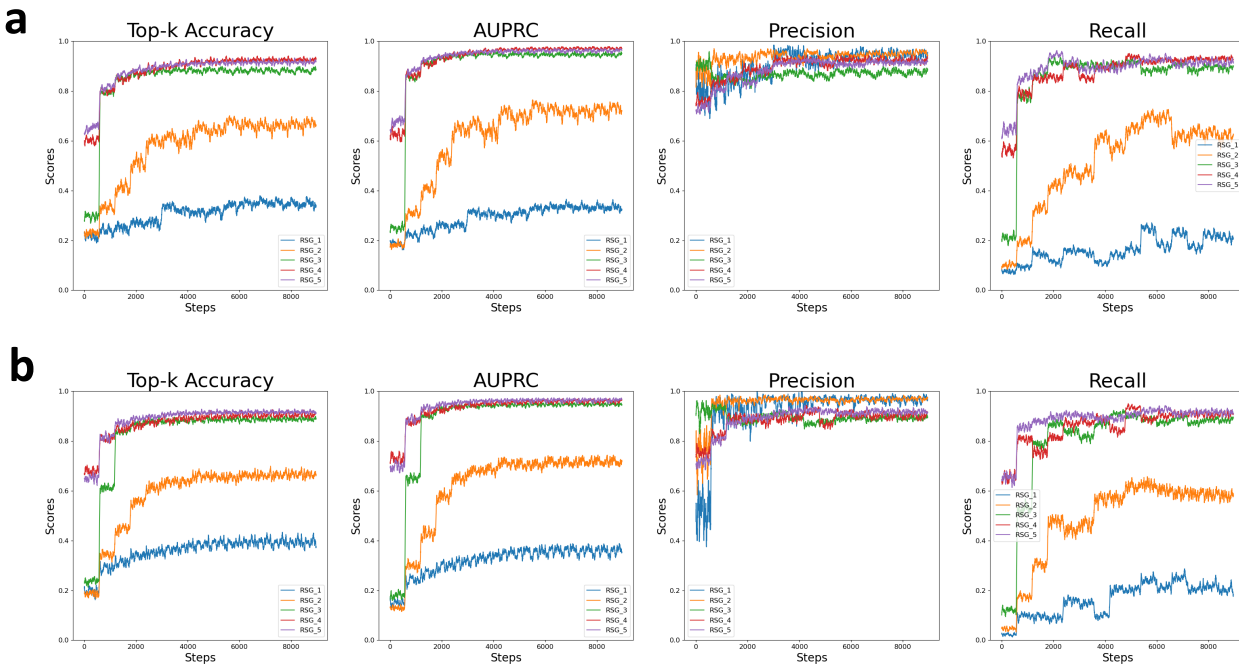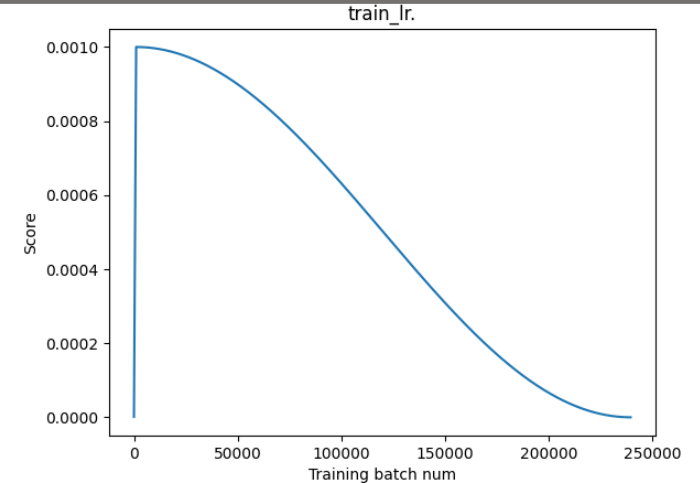
**23**

# Future sequence models in genomics?

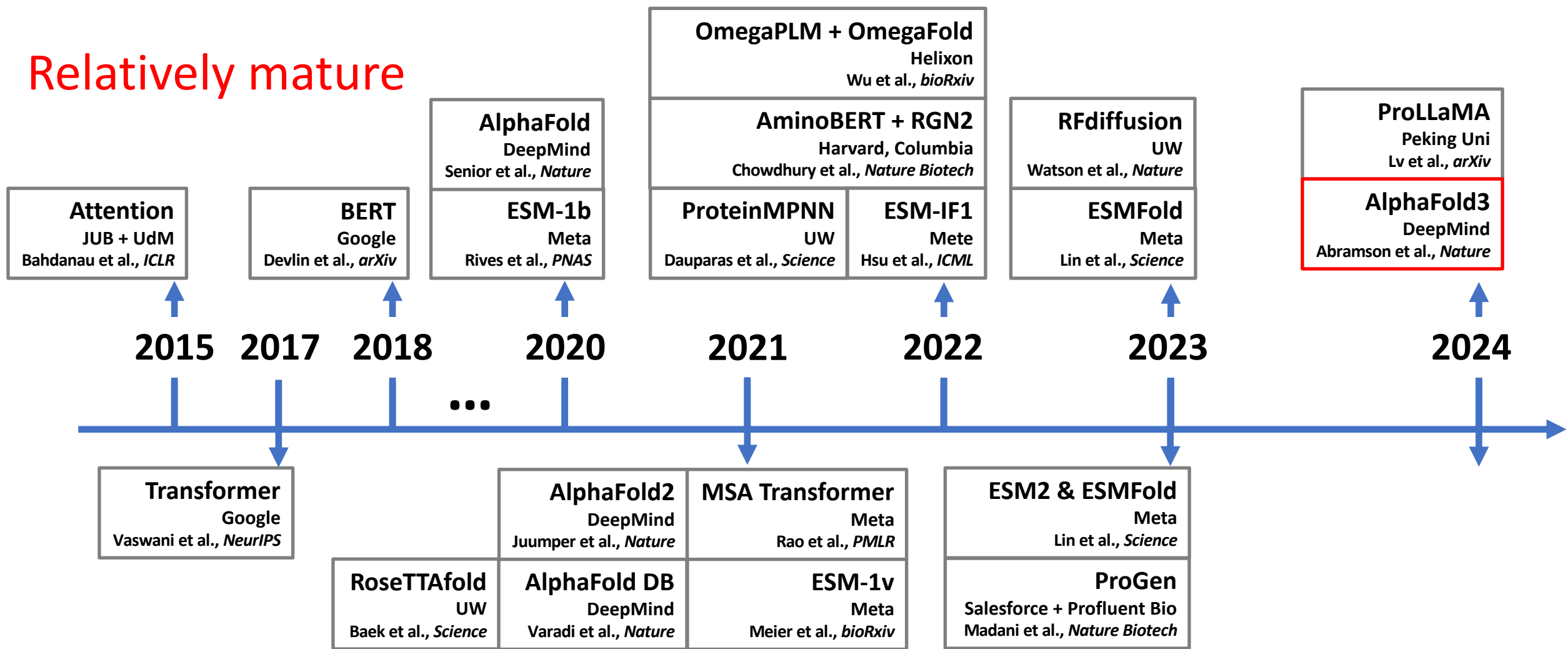# CNN or ?

# Future sequence models in genomics?

# CNN or /and Transformer?

# Future? – Protein transformer-based models

Relatively mature

**OmegaPLM + OmegaFold**
Helixon
Wu et al., *bioRxiv*

**AlphaFold**
DeepMind
Senior et al., *Nature*

**AminoBERT + RGN2**
Harvard, Columbia
Chowdhury et al., *Nature Biotech*

**RFdiffusion**
UW
Watson et al., *Nature*

**ProLLaMA**
Peking Uni
Lv et al., *arXiv*

**Attention**
JUB + UdM
Bahdanau et al., *ICLR*

**BERT**
Google
Devlin et al., *arXiv*

**ESM-1b**
Meta
Rives et al., *PNAS*

**ProteinMPNN**
UW
Dauparas et al., *Science*

**ESM-IF1**
Mete
Hsu et al., *ICML*

**ESMFold**
Meta
Lin et al., *Science*

**AlphaFold3**
DeepMind
Abramson et al., *Nature*

2015  2017  2018  2020  2021  2022  2023  2024

...

**Transformer**
Google
Vaswani et al., *NeurIPS*

**RoseTTAfold**
UW
Baek et al., *Science*

**AlphaFold2**
DeepMind
Juumper et al., *Nature*

**AlphaFold DB**
DeepMind
Varadi et al., *Nature*

**MSA Transformer**
Meta
Rao et al., *PMLR*

**ESM-1v**
Meta
Meier et al., *bioRxiv*

**ESM2 & ESMFold**
Meta
Lin et al., *Science*

**ProGen**
Salesforce + Profluent Bio
Madani et al., *Nature Biotech*

# Application❓ – Genome annotation

**ANNUAL REVIEW OF GENOMICS AND HUMAN GENETICS**

## Deep Learning Sequence Models for Transcriptional Regulation

Ksenia Sokolova[1], Kathleen M. Chen[1], Yun Hao[2], Jian Zhou[3], and Olga G. Troyanskaya[1,2,4]

**SegmentNT: annotating the genome at single-nucleotide resolution with DNA foundation models**

Bernardo P. de Almeida, Hugo Dalla-Torre, Guillaume Richard, Christopher Blum, Lorenz Hexemer, Maxence Gélard, Priyanka Pandey, Stefan Laurent, Alexandre Laterre, Maren Lang, Uğur Şahin, Karim Beguir, Thomas Pierrot
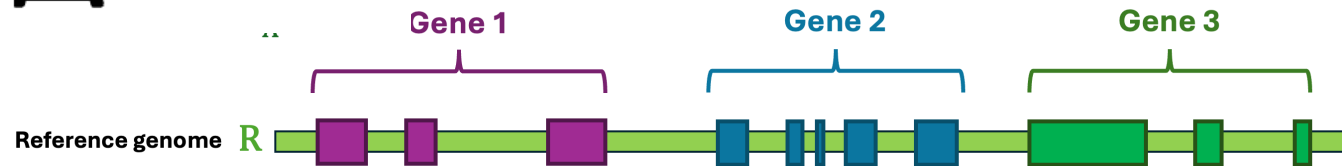
**github.com/Kuanhao-Chao/LiftOn**

**ccb.jhu.edu/lifton**    Preprint coming soon!

27

Introduction    SpliceAI-toolkit    Splam    Future work

# Acknowledge

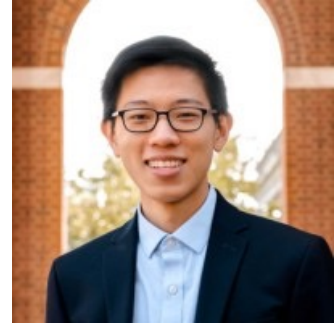Steven Salzberg    Mihaela Pertea    Anqi Liu    Alaina Shumate    Jakob Heinz    Celine Hoh    Alan Mao

- All members in Salzberg lab, Pertea lab
- All friends at Malone & CCB
- All friends at JHU Computational Biology Group

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

Demis Hassabis, CEO of DeepMind, 2022