



Northwestern  
University



Stony Brook University

# DNABERT-2: EFFICIENT AND EFFECTIVE FOUNDATION MODEL FOR MULTI-SPECIES GENOME

Zhihan Zhou, Yanrong Ji, Weijian Li, Pratik Dutta,  
Ramana V. Davuluri, Han Liu

# DNABERT

Understand genome is a **fundamental** task in biology research.

Various important sub-task:

- Promoter Detection
- Splice Site Prediction
- Transcription Factor Prediction
- Epigenetic Marks Prediction
- ...

# Pre-train of DNABERT

# DNABERT

Human Genome

ATGTC AATGTC AATGTC AATGTC AACTGTC AATTACTGTC AATTACTGTC AATTGCACTGTC AGACTGTC AACTGTC AATT  
CAAATGTC AATTGCACTGTC AGACTGTC AATTACTGTC AATTACTGTC AATTGCACTGTC AGACTGTC AACTGTC AATTTT  
GTC AATTACTGTC AATTACTGTC AATTGCACTGTC AGACTAATGTC AATTGCACTGTC AGACTGTC AATTACTGTC AGACTG

DNA sequence

ATTGCACTGTCAG

k-mer sequence

ATTGCA TTGCAC TGC ACT GCACTG CACTGT ACTGTC CTGTCA TGTCAG

masked k-mer sequence

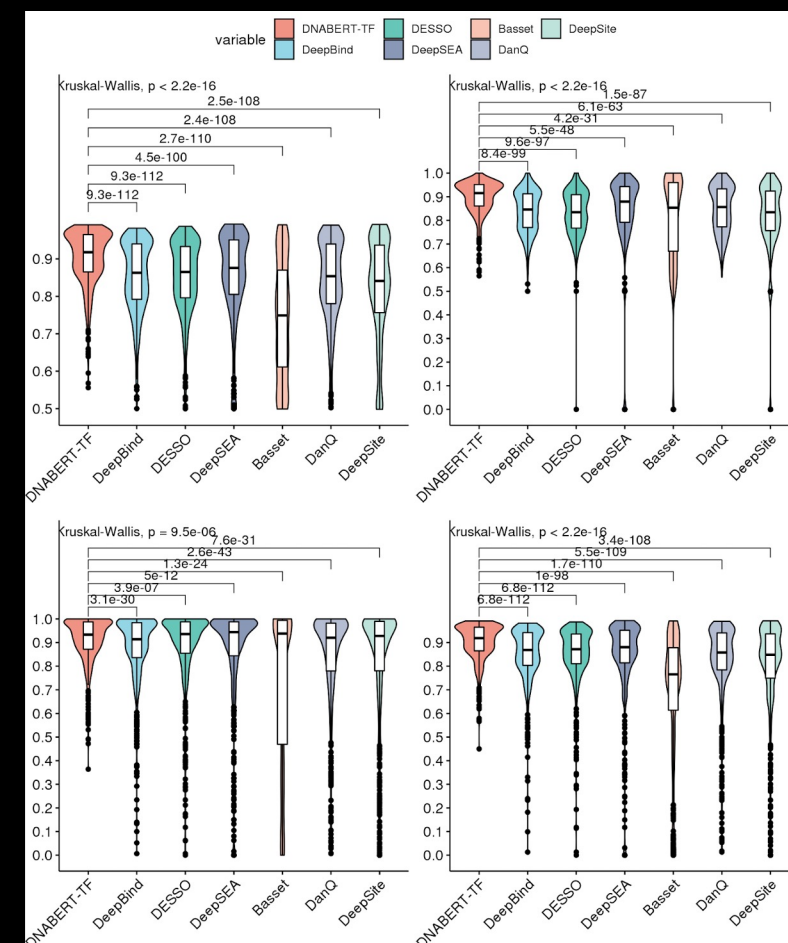
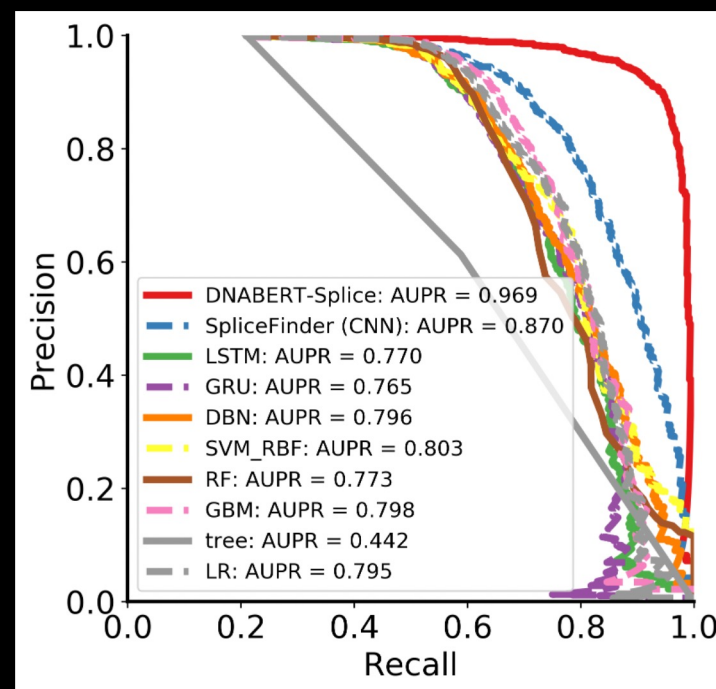
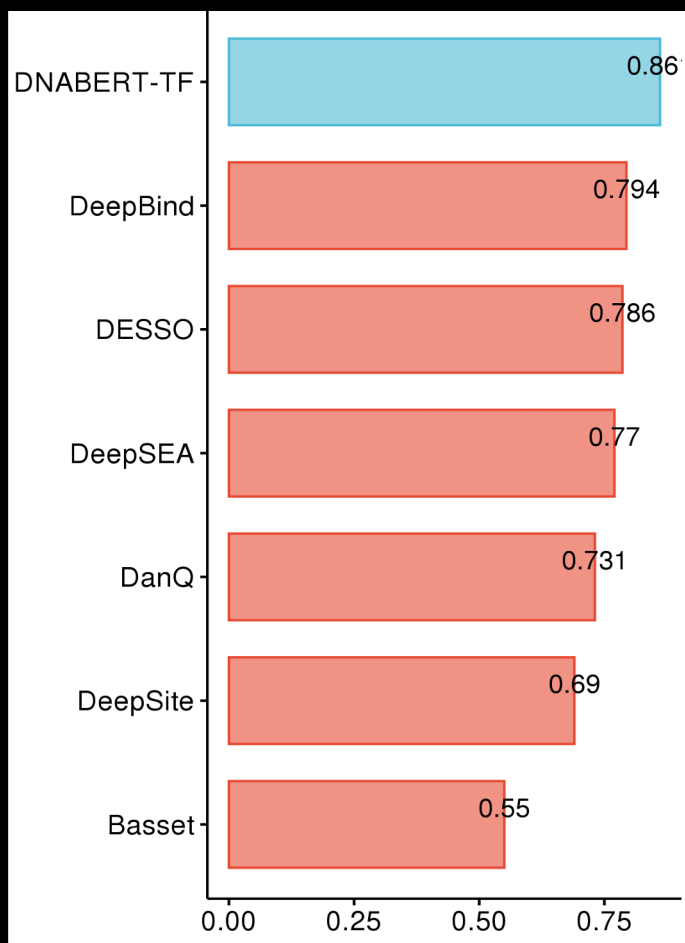
ATTGCA TTGCAC [MASK] [MASK] [MASK] [MASK] CTGTCA TGTCAG

DNABERT

Predict Masked k-mers

TGC ACT GCACTG CACTGT ACTGTC

# DNABERT



# DNABERT-2

We identify 3 key limitations of DNABERT:

1. The k-mer tokenization is data and computational inefficient
2. The model is trained on human genome only
3. The model architecture is less optimal
  - Input length limitation
  - Weak representation capability
  - ...

# DNABERT-2

## Limitation 1: Inefficiency of k-mer tokenization



[MASK]: starts with TGCAC and ends with TTGCA

Data  
Inefficiency



starts with TTGCA. Search space 4096 → 4

## Solution

Use BPE to replace k-mer tokenization.



High-frequently co-occur segments

# DNABERT-2

## Limitation 1: Inefficiency of k-mer tokenization

### Computational Inefficiency

With k-mer tokenization, a DNA sequence with length  $N$  will end up with  $(N-5)$  tokens.

It's computationally heavy considering the  $O(n^2)$  computational efficiency of Transformer

### Solution

Use **BPE** to replace k-mer tokenization.

- Reduce sequence length by **5** times on average
- **No information leakage** in pre-training

# DNABERT-2

## Limitation 2: Only trained on human genome

### Solution

Train the model on the combination of genome from **135** different species.

### Solution

## Limitation 3: Less optimal architecture

- Replace positional embedding with ALiBi. Now support **unlimited input length**.
- Incorporate Flash Attention to improve **efficiency**.
- Use new activation function and other tricks to improve **effectiveness**.



## Genome Evaluation Benchmark

Species	Task	Num. Datasets	Num. Classes	Sequence Length
<b>Human</b>	Core Promoter Detection	3	2	70
	Transcription Factor Prediction	5	2	100
	Promoter Detection	3	2	300
	Splice Site Detection	1	3	400
<b>Mouse</b>	Transcription Factor Prediction	5	2	100
<b>Yeast</b>	Epigenetic Marks Prediction	10	2	500
<b>Virus</b>	Covid Variant Classification	1	9	1000

Table 1: Summarization of the Genome Understanding Evaluation (GUE) benchmark.

Species	Task	Num. Datasets	Num. Classes	Sequence Length
<b>Human</b>	Enhancer Promoter Interaction	6	2	5000
<b>Fungi</b>	Species Classification	1	25	5000
<b>Virus</b>	Species Classification	1	20	10000

Table 2: Summarization of the Genome Understanding Evaluation Plus (GUE<sup>+</sup>) benchmark.

# DNABERT-2

## Results on GUE

Size

Efficiency

Effectiveness

Model	Num. Params. ↓	FLOPs ↓	Trn. Tokens	Num. Top-2 ↑	Ave. Scores ↑
<b>DNABERT (3-mer)</b>	86M	3.27	122B	2    0	61.62
<b>DNABERT (4-mer)</b>	86M	3.26	122B	0    1	61.14
<b>DNABERT (5-mer)</b>	87M	3.26	122B	0    1	60.05
<b>DNABERT (6-mer)</b>	89M	3.25	122B	0    1	60.51
<b>NT-500M-human</b>	480M	3.19	50B	0    0	55.43
<b>NT-500M-1000g</b>	480M	3.19	50B	0    1	58.23
<b>NT-2500M-1000g</b>	2537M	19.44	300B	0    1	61.41
<b>NT-2500M-multi</b>	2537M	19.44	300B	<u>7</u>    <u>9</u>	<u>66.93</u>
<b>DNABERT-2</b>	117M	1.00	262B	8    4	66.80
<b>DNABERT-2♦</b>	117M	1.00	263B	<b>11</b>    <b>10</b>	<b>67.77</b>

DNABERT-2 vs DNABERT: **3x faster**, **unlimited input length**, and **much better performance**.

DNABERT-2 vs NT-2.5B: **19x faster**, **92x less training cost**, and **similarly performance**.

# DNABERT-2

## Results on GUE+

<b>Task</b>	<b>SC (Fungi)</b>	<b>SC (Virus)</b>	<b>EPI (Human)</b>					
<b>Dataset</b>			<b>0</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>
<b>DNABERT (6-mer)</b>	89.29	44.51	-	-	-	-	-	-
<b>NT-2500M-multi</b>	92.85	45.00	61.91	72.15	73.13	79.49	86.48	68.64
<b>DNABERT-2</b>	<b>93.04</b>	<b>48.50</b>	<b>76.21</b>	<b>79.19</b>	<b>83.50</b>	<b>86.71</b>	<b>92.90</b>	<b>73.70</b>

DNABERT-2 vs NT-2.5B: better performance on long-sequence tasks

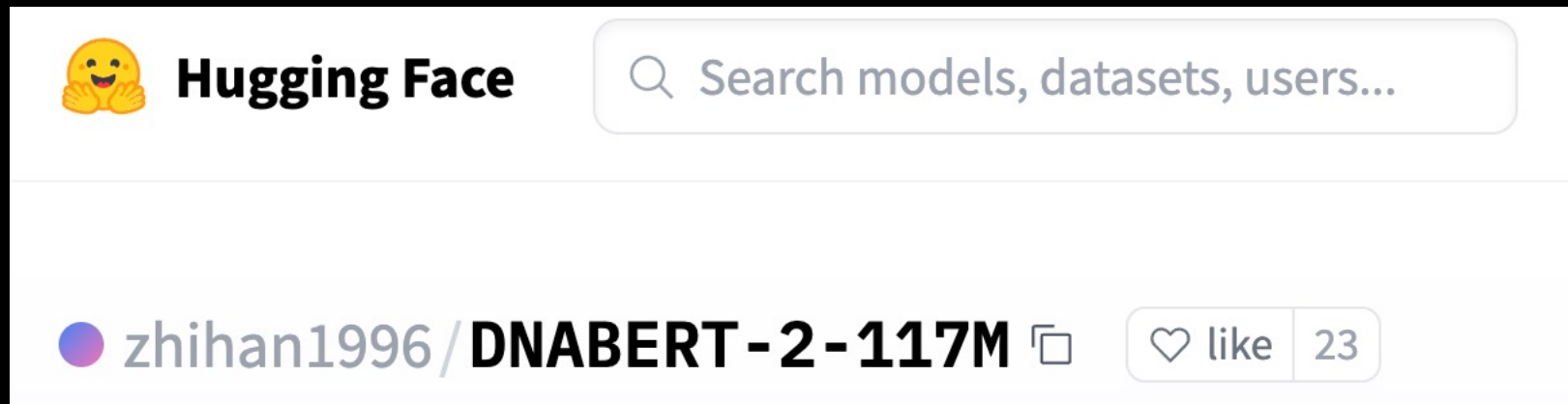
# DNABERT-2

## Code and Model

The code is publicly available at [https://github.com/MAGICS-LAB/DNABERT\\_2](https://github.com/MAGICS-LAB/DNABERT_2)

The trained models is available at HuggingFace:

<https://huggingface.co/zhihan1996/DNABERT-2-117M>



The screenshot shows the Hugging Face interface. At the top left is the Hugging Face logo (a yellow smiley face with hands) and the text "Hugging Face". To the right is a search bar with the placeholder text "Search models, datasets, users...". Below this is a horizontal line. At the bottom of the screenshot, the model name "zhihan1996 / DNABERT-2-117M" is displayed, followed by a folder icon. To the right of the model name is a "like" button with a heart icon and the number "23".