

DNAChunker:

Learnable Tokenization for DNA Language Models

Taewon Kim, Jihwan Shin, Hyomin Kim,
Youngmok Jung, Jonghoon Lee, Won-Chul Lee, Sungsoo Ahn*, Insu Han*

GLMs are essentially representation encoders

Genetic Domain



\mathcal{Z}



OpenAI CLIP



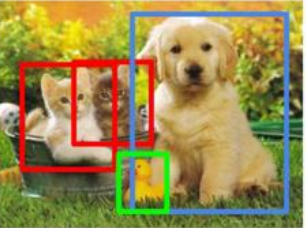


DINO v2, v3



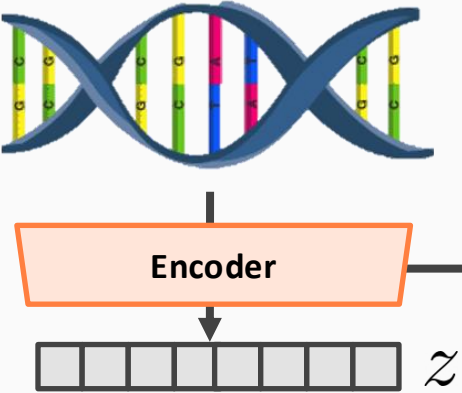
BEiT, Florence

Diverse SSL approaches have enabled
Compression of images to vectors, which can solve:

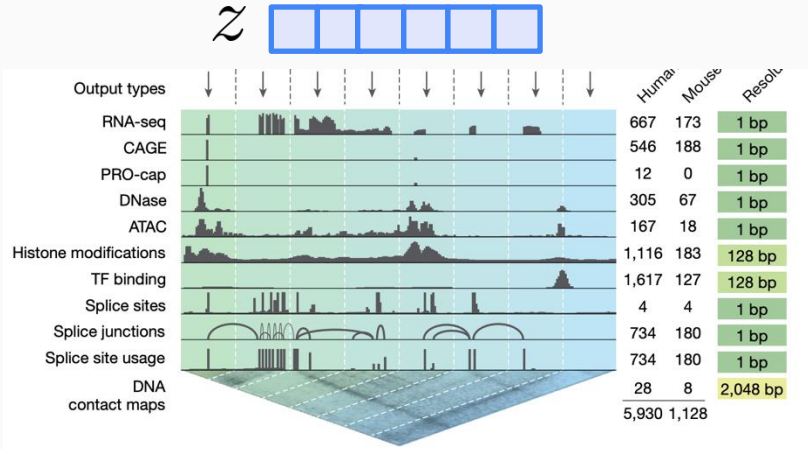
Classification	Classification + Localization	Object Detection
		
CAT	CAT	CAT, DOG, DUCK

GLMs are essentially representation encoders

Genetic Domain



Input: Genomic Sequence -> Output: Representation
Either trained with SSL (autoregressive, MLM) or SL (label guidance)



Can solve DNA-specific tasks such as DNA contact map prediction, etc.

In this work, we focus on tokenization

Natural Language : I / love / genomic / language / models

Many meaningful separation exists – ex. **spaces already divide input in a semantically-aware manner**

Genomic Language: ATGCATTACATGCATTACATGCATTAC...



No natural separation exists

Thus, people have naively used NLP-style tokenizations

Fixed K-Mer (6)

Original A A G A G G | A T T A A A | T C C T G A | A G G T T A | G G G A T G

Divide input sequence into 6 bp-length chunks, each with unique numbering.

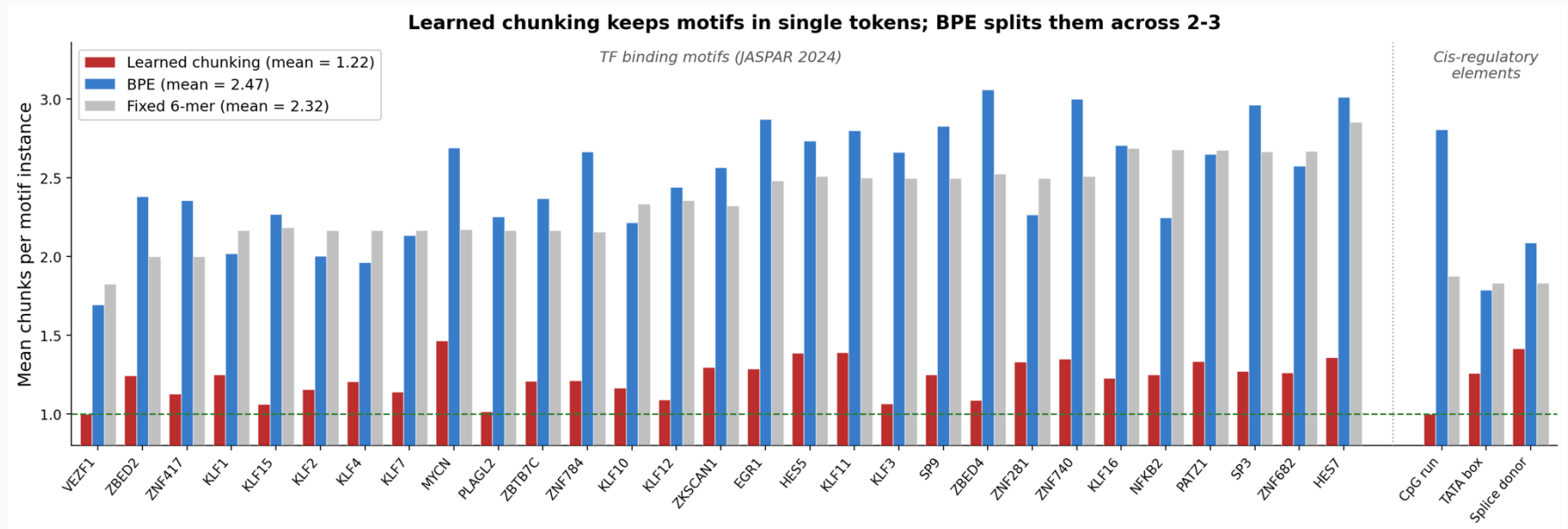
Byte Pair Encoding (BPE)

Original A A G A G G A T | T A A A | T C C T G A A G | G T T A G | G G A T G

Predefine a set of chunks that are statistically – merging frequent pairs of characters.

Prior approaches are ‘semantically blind’ of DNA

Question: Does prior tokenization preserve the known “Grammar (motif)” of DNA?

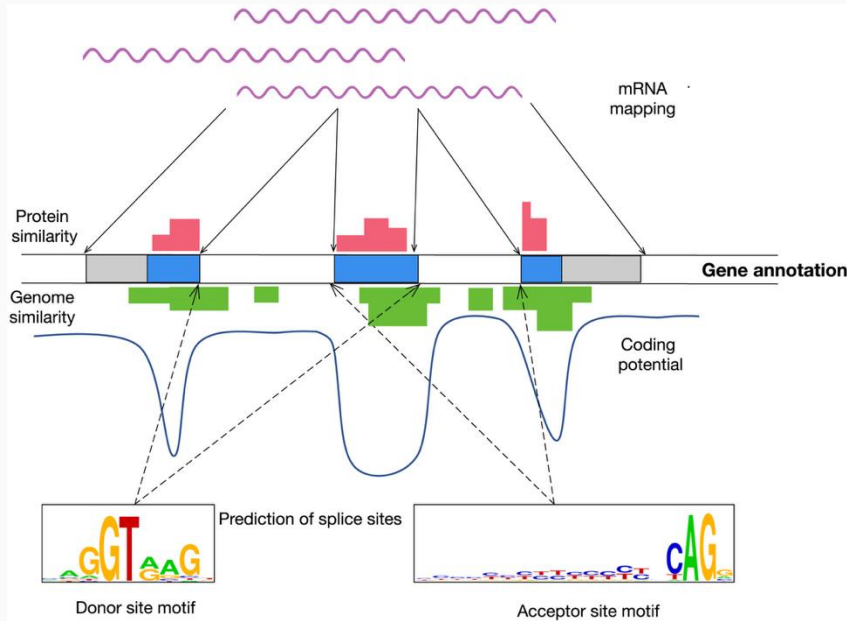


Answer: NO! in fact, it ACTIVELY fragments it!

Single Motif -> BPE divides it into 2.47, K-Mer into 2.32 tokens on average

Why Not Use Pre-annotated Gene Elements as Guidance?

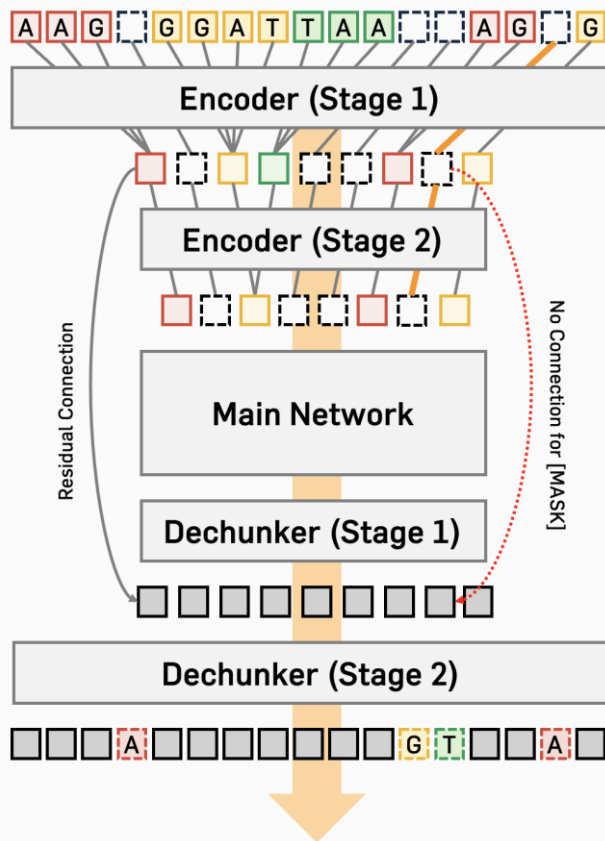
TL;DR: Lack of annotation across genomes



1. Conservation-based reasoning
 - ✓ Highly conservative == Important
 - ✓ **Lack of lineage-specific functionals**
2. Association-based reasoning
 - ✓ Statistically associated with disease
 - ✓ **Correlation != Causal**
3. Experimental Assays ...
 - ✓ **Hard to obtain...**

Furthermore, we do not know
WHICH are best for tokenization

DNACChunker: Unsupervised, Data-driven tokenization strategy

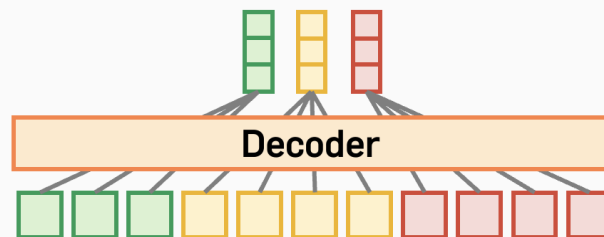
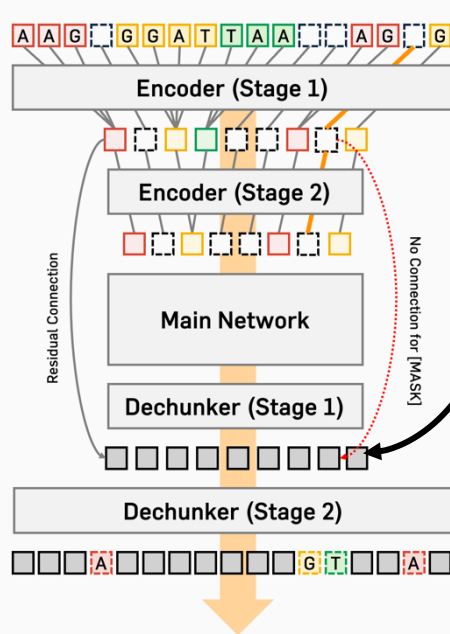


TL;DR: Learn the tokenization strategy from data!

- ✓ First approach to use **learnable tokenization in gLMs**
- ✓ Enforce bidirectional nature of DNA sequences
- ✓ Almost **NO inductive bias**
- ✓ **SOTA performance** on Benchmark Tasks

Decoding chunks back to original resolution

TL;DR: Tricks need be in place for gradient propagation

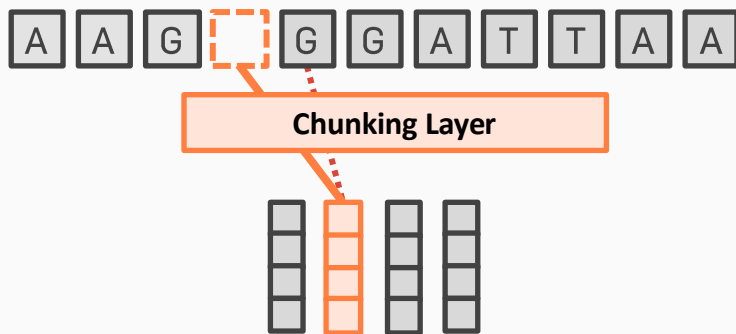


$$z_t^{(s+1)} = \frac{1}{2} (\text{SCAN}_{\rightarrow}(\tilde{z}^{(s+1)}, p)_t + \text{SCAN}_{\leftarrow}(\tilde{z}^{(s+1)}, p)_t),$$

1. Embeddings are upsampled to original resolution
2. **Bidirectional EMA scan** is applied
Trick to ensure gradient propagation to the chunking module (via mixing between boundaries)
3. + Residual connections (like U-Net) to original embedding

[MASK] are different from other letters

- ✓ [MASK] is *NOT seem in test data*. -> Train vs. Test gap exits
- ✓ Learning a chunker WITH *[MASK] makes the tokenization mask-dependent*

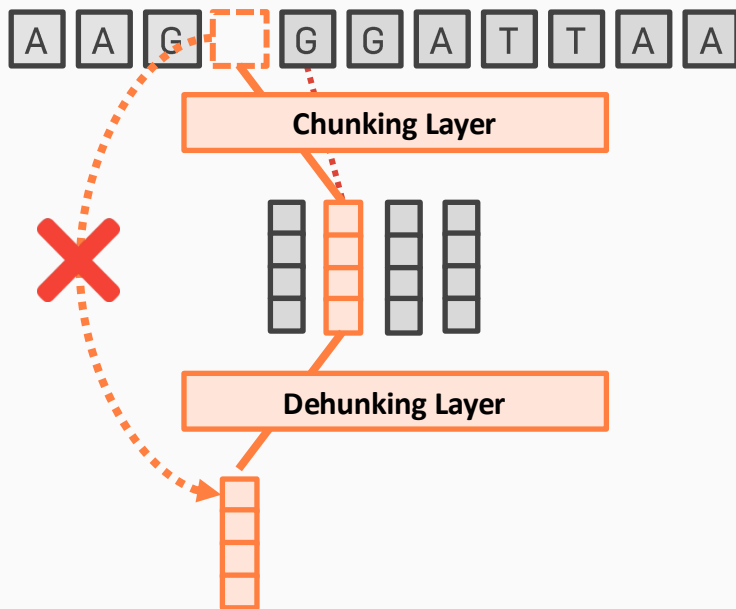


Trick#1: [MASK] Protection Mechanism

- [MASK] should not be merged with other tokens (A, T, G, C, N).

[MASK] are different from other letters

- ✓ [MASK] is *NOT seem in test data*. -> Train vs. Test gap exits
- ✓ Learning a chunker WITH *[MASK] makes the tokenization mask-dependent*



Trick#1: [MASK] Protection Mechanism

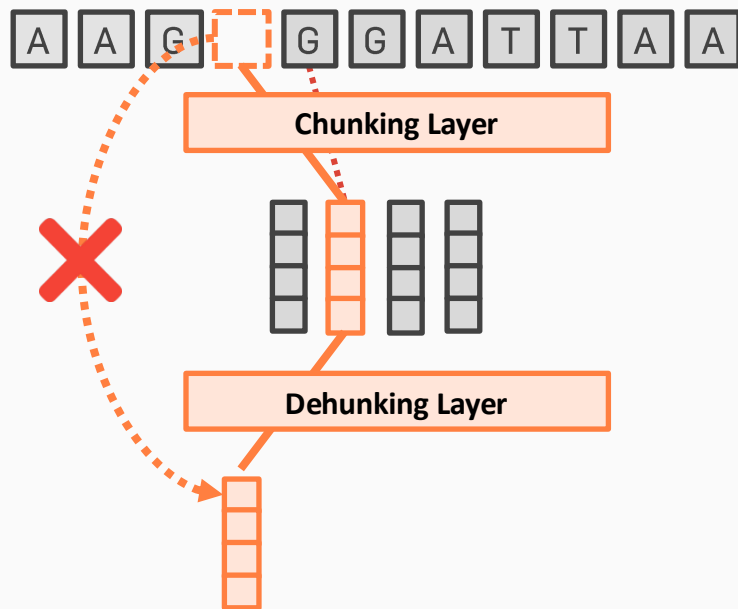
- [MASK] should not be merged with other tokens (A, T, G, C, N).

Trick#2: Residual Gating of [MASK]

- If residual connection is ALSO given to the mask, encoder leaks information to decoder

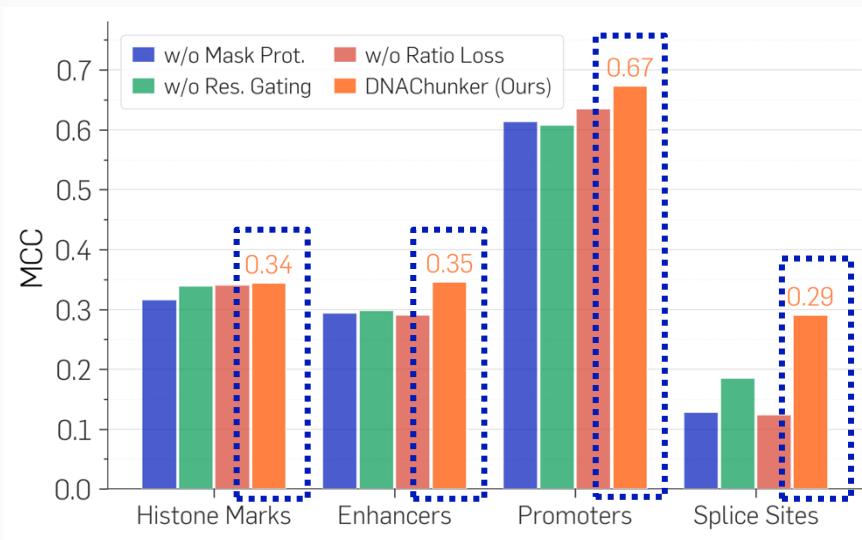
[MASK] are different from other letters

- ✓ [MASK] is *NOT seem in test data*. -> Train vs. Test gap exits
- ✓ Learning a chunker WITH *[MASK] makes the tokenization mask-dependent*



[MASK] handling is crucial for performance

Linear probing upon NT-revised Benchmark in MCC



Loss Function

Standard MLM loss + Repeat region down-weighting + Ratio Loss

$$\mathcal{L}_{\text{MLM}} = \sum_{t \in M} w_t \mathcal{L}_{\text{CE}}(t) \quad w_t = \begin{cases} 0.1 & \text{if position } t \text{ is in a repetitive region} \\ 1.0 & \text{otherwise} \end{cases}$$

- Following formulation of Evo2[1], down-weight repeat regions
- Repeat regions lack biological significance

Loss Function

Standard MLM loss + Repeat region down-weighting + Ratio Loss [1]

$$\mathcal{L}_{\text{ratio}}^{(s)} = \frac{\bar{b}^{(s)}\bar{p}^{(s)}}{\alpha^{(s)}} + \frac{(1 - \bar{b}^{(s)})(1 - \bar{p}^{(s)})}{1 - \alpha^{(s)}}, \quad \bar{b}^{(s)} = \frac{1}{T} \sum_{t=1}^T b_t^{(s)}, \quad \bar{p}^{(s)} = \frac{1}{T} \sum_{t=1}^T p_t^{(s)},$$

- To avoid collapse to trivial solutions (i.e. assigning each character to a separate chunk)
- Loss attains minimum when the fraction of selected boundaries ($\bar{b}^{(s)}$) and boundary probability ($\bar{p}^{(s)}$) becomes the target compression ratio ($\alpha^{(s)}$)

State-of-the-art upon Benchmarks

Table 1. Nucleotide Transformer Benchmark. The reported values represent the Matthews Correlation Coefficient (MCC; mean \pm standard error) averaged over 10-fold cross-validation. Best results are **bold**; second best are underlined.

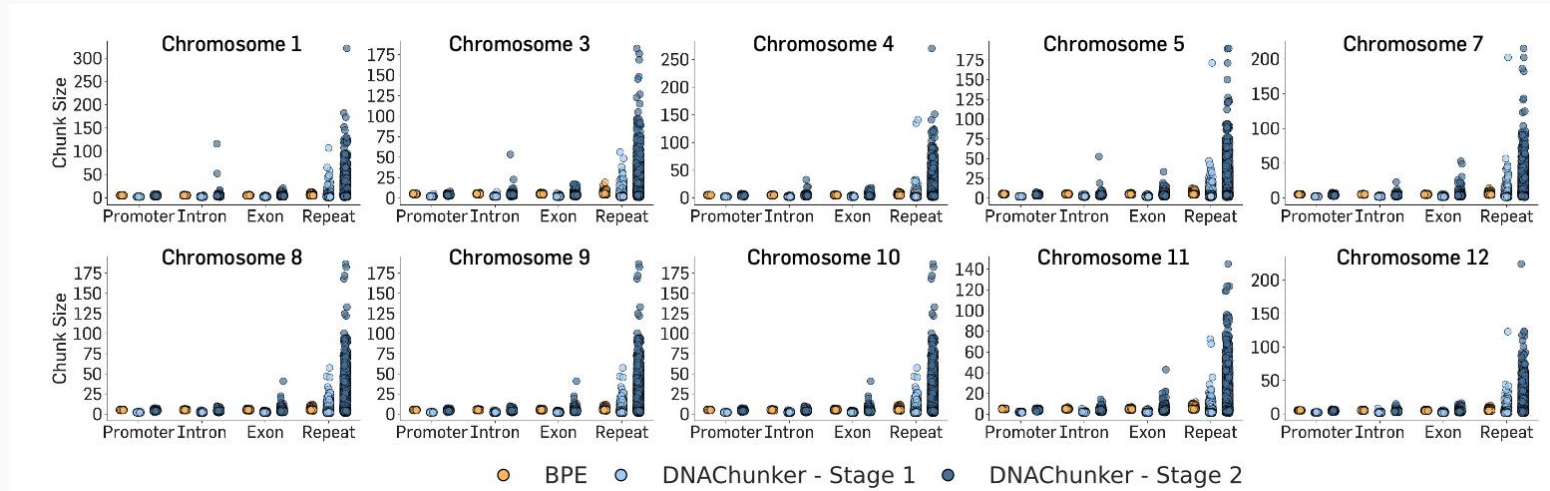
	Enformer (252M)	DNABERT-2 (117M)	HyenaDNA (55M)	NT-multi (2.5B)	NT-v2 (500M)	Caduceus-Ph (8M)	Caduceus-PS (8M)	GROVER (87M)	GENERator (1.2B)	DNACHUNKER (172M)
Histone Markers										
H3	0.724 \pm 0.018	<u>0.785</u> \pm 0.012	0.781 \pm 0.015	0.793 \pm 0.013	0.788 \pm 0.010	0.794 \pm 0.012	0.772 \pm 0.022	0.768 \pm 0.008	0.806 \pm 0.005	0.817 \pm 0.011
H3K14ac	0.284 \pm 0.024	0.515 \pm 0.009	0.608 \pm 0.020	0.538 \pm 0.009	0.538 \pm 0.015	0.564 \pm 0.033	0.596 \pm 0.038	0.548 \pm 0.020	0.605 \pm 0.008	0.711 \pm 0.021
H3K36me3	0.345 \pm 0.019	0.591 \pm 0.005	0.614 \pm 0.014	0.618 \pm 0.011	0.618 \pm 0.015	0.590 \pm 0.018	0.611 \pm 0.048	0.563 \pm 0.017	<u>0.657</u> \pm 0.007	0.677 \pm 0.003
H3K4me1	0.291 \pm 0.016	0.512 \pm 0.008	0.512 \pm 0.008	0.541 \pm 0.005	0.544 \pm 0.009	0.468 \pm 0.015	0.487 \pm 0.029	0.461 \pm 0.018	0.553 \pm 0.009	0.631 \pm 0.009
H3K4me2	0.207 \pm 0.021	0.333 \pm 0.013	<u>0.455</u> \pm 0.028	0.324 \pm 0.014	0.302 \pm 0.020	0.332 \pm 0.034	0.431 \pm 0.016	0.403 \pm 0.042	0.424 \pm 0.013	0.599 \pm 0.011
H3K4me3	0.156 \pm 0.022	0.353 \pm 0.021	<u>0.550</u> \pm 0.015	0.408 \pm 0.011	0.437 \pm 0.028	0.490 \pm 0.042	0.528 \pm 0.033	0.458 \pm 0.022	0.512 \pm 0.009	0.660 \pm 0.045
H3K79me3	0.498 \pm 0.013	0.615 \pm 0.010	0.669 \pm 0.014	0.623 \pm 0.010	0.621 \pm 0.012	0.641 \pm 0.028	0.682 \pm 0.018	0.626 \pm 0.026	0.670 \pm 0.011	0.731 \pm 0.012
H3K9ac	0.415 \pm 0.020	0.545 \pm 0.009	0.586 \pm 0.021	0.547 \pm 0.011	0.567 \pm 0.020	0.575 \pm 0.024	0.564 \pm 0.018	0.581 \pm 0.015	0.612 \pm 0.006	0.678 \pm 0.007
H4	0.735 \pm 0.023	0.797 \pm 0.008	0.763 \pm 0.012	0.808 \pm 0.007	0.795 \pm 0.008	0.788 \pm 0.010	0.799 \pm 0.010	0.769 \pm 0.017	0.815 \pm 0.008	0.813 \pm 0.012
H4ac	0.275 \pm 0.022	0.465 \pm 0.013	0.564 \pm 0.011	0.492 \pm 0.014	0.502 \pm 0.025	0.548 \pm 0.027	0.585 \pm 0.018	0.530 \pm 0.017	<u>0.592</u> \pm 0.015	0.687 \pm 0.027
Average MCC (\uparrow)	0.393	0.551	0.610	0.569	0.571	0.579	0.606	0.571	<u>0.625</u>	0.701
Regulatory Annotation										
Enhancer	0.454 \pm 0.029	0.525 \pm 0.026	0.520 \pm 0.031	0.545 \pm 0.028	<u>0.561</u> \pm 0.029	0.522 \pm 0.024	0.511 \pm 0.026	0.516 \pm 0.018	0.580 \pm 0.015	0.558 \pm 0.011
Enhancer Type	0.312 \pm 0.043	0.423 \pm 0.018	0.403 \pm 0.056	0.444 \pm 0.022	0.444 \pm 0.036	0.403 \pm 0.028	0.410 \pm 0.026	0.433 \pm 0.029	0.477 \pm 0.017	0.519 \pm 0.005
Promoter All	0.910 \pm 0.004	0.945 \pm 0.003	0.919 \pm 0.003	0.951 \pm 0.004	0.952 \pm 0.002	0.937 \pm 0.002	0.941 \pm 0.003	0.926 \pm 0.004	0.952 \pm 0.002	0.967 \pm 0.013
Promoter NonTATA	0.910 \pm 0.006	0.944 \pm 0.003	0.919 \pm 0.004	0.969 \pm 0.003	0.952 \pm 0.003	0.935 \pm 0.007	0.940 \pm 0.002	0.925 \pm 0.006	0.952 \pm 0.001	0.971 \pm 0.007
Promoter TATA	0.920 \pm 0.012	0.911 \pm 0.011	0.881 \pm 0.020	0.919 \pm 0.008	0.933 \pm 0.009	0.895 \pm 0.010	0.903 \pm 0.010	0.891 \pm 0.009	0.948 \pm 0.008	0.961 \pm 0.015
Average MCC (\uparrow)	0.701	0.750	0.728	0.766	0.768	0.738	0.741	0.738	<u>0.786</u>	0.796
Splice Site Annotation										
Splice Acceptor	0.772 \pm 0.007	0.909 \pm 0.004	0.935 \pm 0.005	0.973 \pm 0.002	0.973 \pm 0.004	0.918 \pm 0.017	0.907 \pm 0.015	0.912 \pm 0.010	0.981 \pm 0.002	0.969 \pm 0.013
Splice Site All	0.831 \pm 0.012	0.950 \pm 0.003	0.917 \pm 0.006	0.974 \pm 0.004	0.975 \pm 0.002	0.935 \pm 0.011	0.953 \pm 0.005	0.919 \pm 0.009	0.976 \pm 0.011	0.968 \pm 0.030
Splice Donor	0.813 \pm 0.015	0.927 \pm 0.003	0.894 \pm 0.013	0.974 \pm 0.002	0.977 \pm 0.007	0.912 \pm 0.009	0.930 \pm 0.010	0.888 \pm 0.012	0.978 \pm 0.001	0.960 \pm 0.007
Average MCC (\uparrow)	0.805	0.929	0.915	0.974	<u>0.975</u>	0.922	0.930	0.906	0.979	0.965
Total Average MCC (\uparrow)	0.547	0.669	0.694	0.690	0.693	0.680	0.697	0.673	0.728	0.772
Total Average Rank (\downarrow)	9.67	6.72	6.00	4.83	4.56	6.33	5.61	7.22	2.06	1.67

Table 2. Genomics Benchmark. The reported values represent the Top-1 Accuracy (Acc; mean \pm standard error) averaged over 10-fold cross-validation. Best results are **bold**; second best are underlined.

	DNABERT-2 (117M)	HyenaDNA (55M)	NT-v2 (500M)	Caduceus-Ph (8M)	Caduceus-PS (8M)	GROVER (87M)	GENERator (1.2B)	GENERator-All (1.2B)	DNACHUNKER (172M)
Coding vs. Intergenic	0.951 \pm 0.002	0.902 \pm 0.004	0.955 \pm 0.001	0.933 \pm 0.001	0.944 \pm 0.002	0.919 \pm 0.002	0.963 \pm 0.000	0.959 \pm 0.001	0.955 \pm 0.012
Drosophila Enhancers Stark	0.774 \pm 0.011	0.770 \pm 0.016	0.797 \pm 0.009	0.827 \pm 0.010	0.816 \pm 0.015	0.761 \pm 0.011	0.821 \pm 0.005	0.768 \pm 0.015	0.779 \pm 0.021
Human Enhancers Cohn	0.758 \pm 0.005	0.725 \pm 0.009	0.756 \pm 0.006	0.747 \pm 0.003	0.749 \pm 0.003	0.738 \pm 0.003	0.763 \pm 0.002	0.754 \pm 0.006	0.761 \pm 0.011
Human Enhancers Ensembl	0.918 \pm 0.003	0.901 \pm 0.003	0.921 \pm 0.004	0.924 \pm 0.002	0.923 \pm 0.002	0.911 \pm 0.004	0.917 \pm 0.002	0.912 \pm 0.002	0.922 \pm 0.007
Human Ensembl Regulatory	0.874 \pm 0.007	0.932 \pm 0.001	0.941 \pm 0.001	0.938 \pm 0.004	0.941 \pm 0.002	0.897 \pm 0.001	0.928 \pm 0.001	0.926 \pm 0.001	0.935 \pm 0.005
Human non-TATA Promoters	0.957 \pm 0.008	0.894 \pm 0.023	0.932 \pm 0.006	0.961 \pm 0.003	0.961 \pm 0.002	0.950 \pm 0.005	0.958 \pm 0.001	0.955 \pm 0.005	0.962 \pm 0.001
Human OCR Ensembl	0.806 \pm 0.003	0.774 \pm 0.004	0.813 \pm 0.001	0.825 \pm 0.004	0.826 \pm 0.003	0.789 \pm 0.002	0.823 \pm 0.002	0.812 \pm 0.003	0.810 \pm 0.007
Human vs. Worm	0.977 \pm 0.001	0.958 \pm 0.004	0.976 \pm 0.001	0.975 \pm 0.001	0.976 \pm 0.001	0.966 \pm 0.001	0.980 \pm 0.000	0.978 \pm 0.001	0.969 \pm 0.001
Mouse Enhancers Ensembl	0.865 \pm 0.014	0.756 \pm 0.030	0.855 \pm 0.018	0.788 \pm 0.028	0.826 \pm 0.021	0.742 \pm 0.025	<u>0.871</u> \pm 0.015	0.784 \pm 0.027	0.874 \pm 0.020
Average Acc (\uparrow)	0.876	0.846	0.883	0.880	0.885	0.853	0.892	0.872	0.885
Average Rank (\downarrow)	5.11	8.22	4.17	3.89	3.33	8.11	2.89	5.44	3.29

- On average, **BEST performance**
- Evaluated upon:
 - NT-Benchmark (Best avg. Perf. & Rank)
 - NT-revised Benchmark (Best avg. Rank)
 - Genomics Benchmark (2nd avg. Rank)
 - BEND Benchmark (Best avg. Rank)
 - DNALongBench (Best avg. Perf. & Rank)

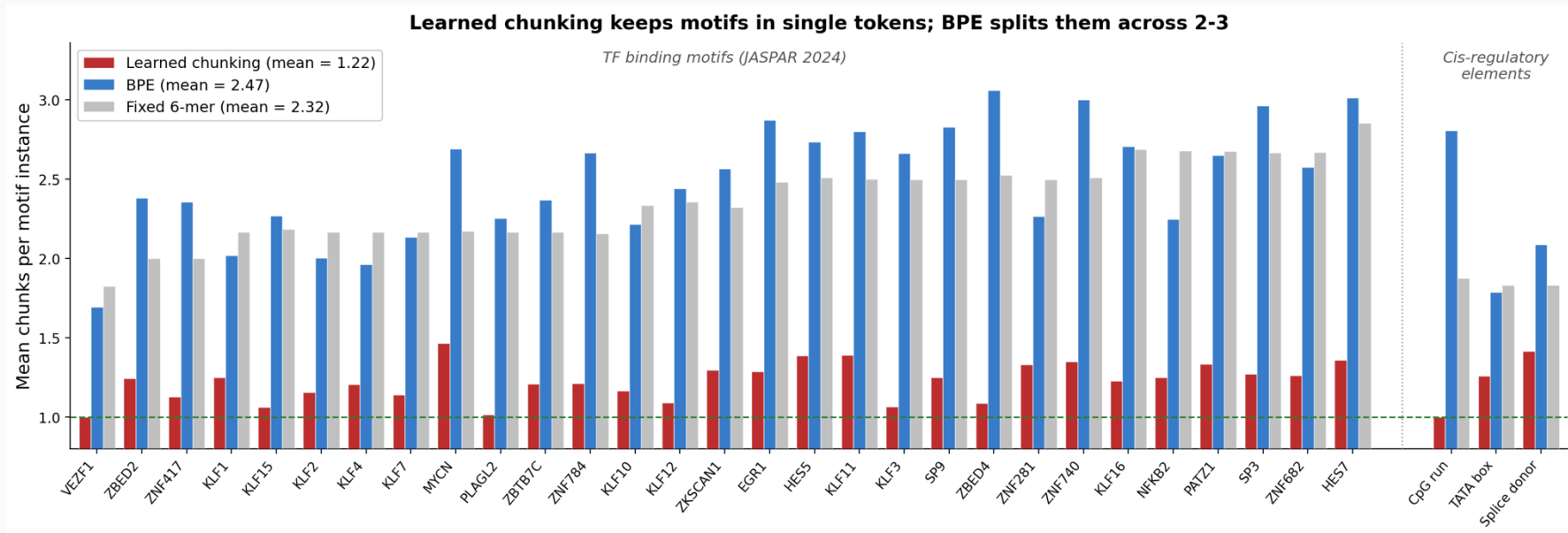
DNACunker learns biological grammar



- Scatter Plot upon *chunk size distribution* in chromosomes, *divided by annotation*.
- Large chunk sizes assigned to repeats – where chunk sizes are larger than 150+ bp

Revisiting 'semantic blindness' to biological grammar

Question: **DNACHunker** preserve the known "Grammar (motif)" of DNA?



Answer: Yes! in fact, it assigns in average, **1.22 tokens to motifs, ACTIVELY preserving it.**

Single Motif -> BPE divides it into 2.47, K-Mer into 2.32 tokens on average

