

Caduceus: Bi-Directional Equivariant Long-Range DNA Sequence Modeling

Advanced Biomedicine Seminar

28/11/2024 - John Shahla

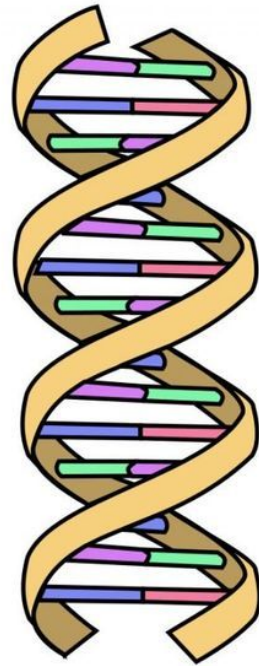
Roadmap

- Definitions
- Motivation
- Models
- Results



Definitions

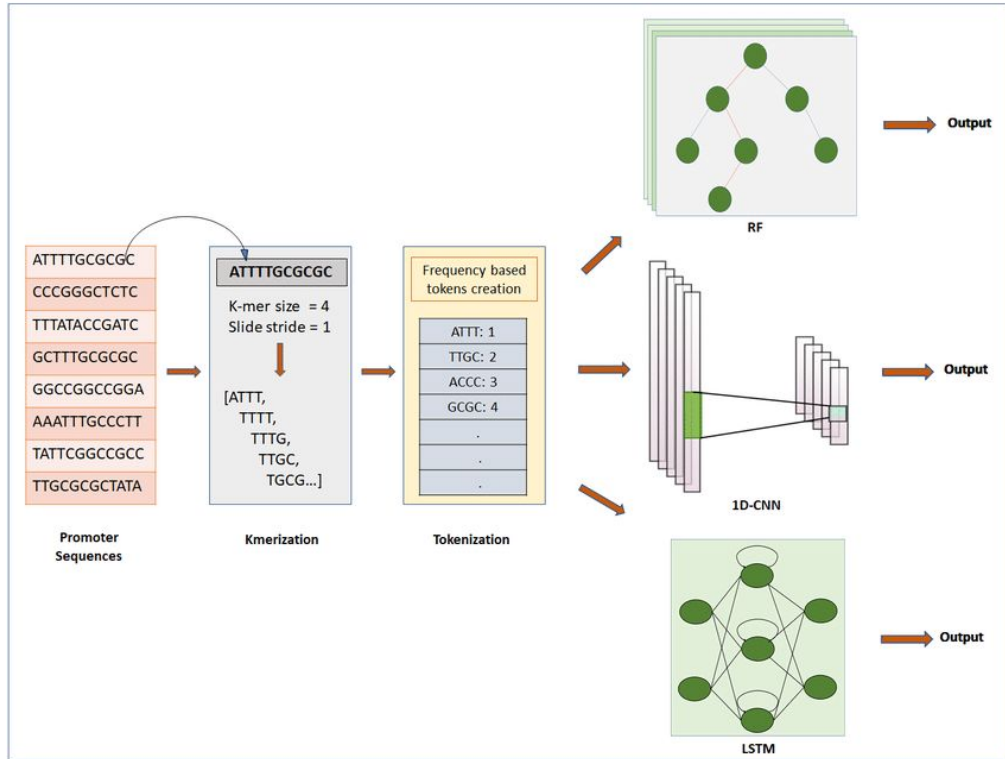
DNA Structure



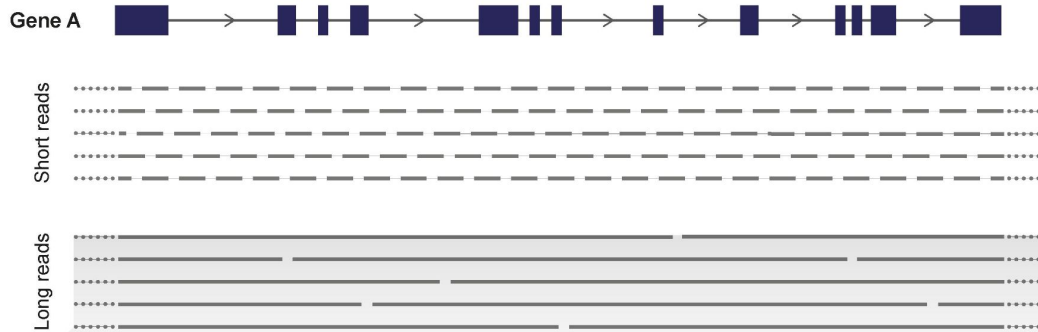
DNA

-  = Adenine
-  = Thymine
-  = Cytosine
-  = Guanine
-  = Phosphate backbone

DNA Sequence Modelling

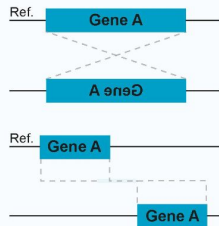


Long-Range DNA Sequencing



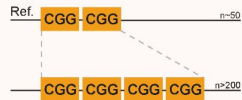
1. Structural variation

e.g. *PRKAR1A*, *G6PC*, *BBS9*, *ARGHEF9*, *TAF1*



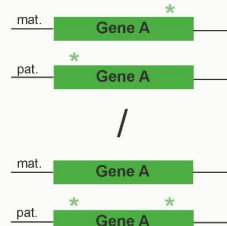
2. Repeat expansion

e.g. *FMR1*, *DMPK*, *ATXN10*, *HTT*



3. Phasing

e.g. Compound heterozygosity, Parental origin of de novo mutations, Mosaicism

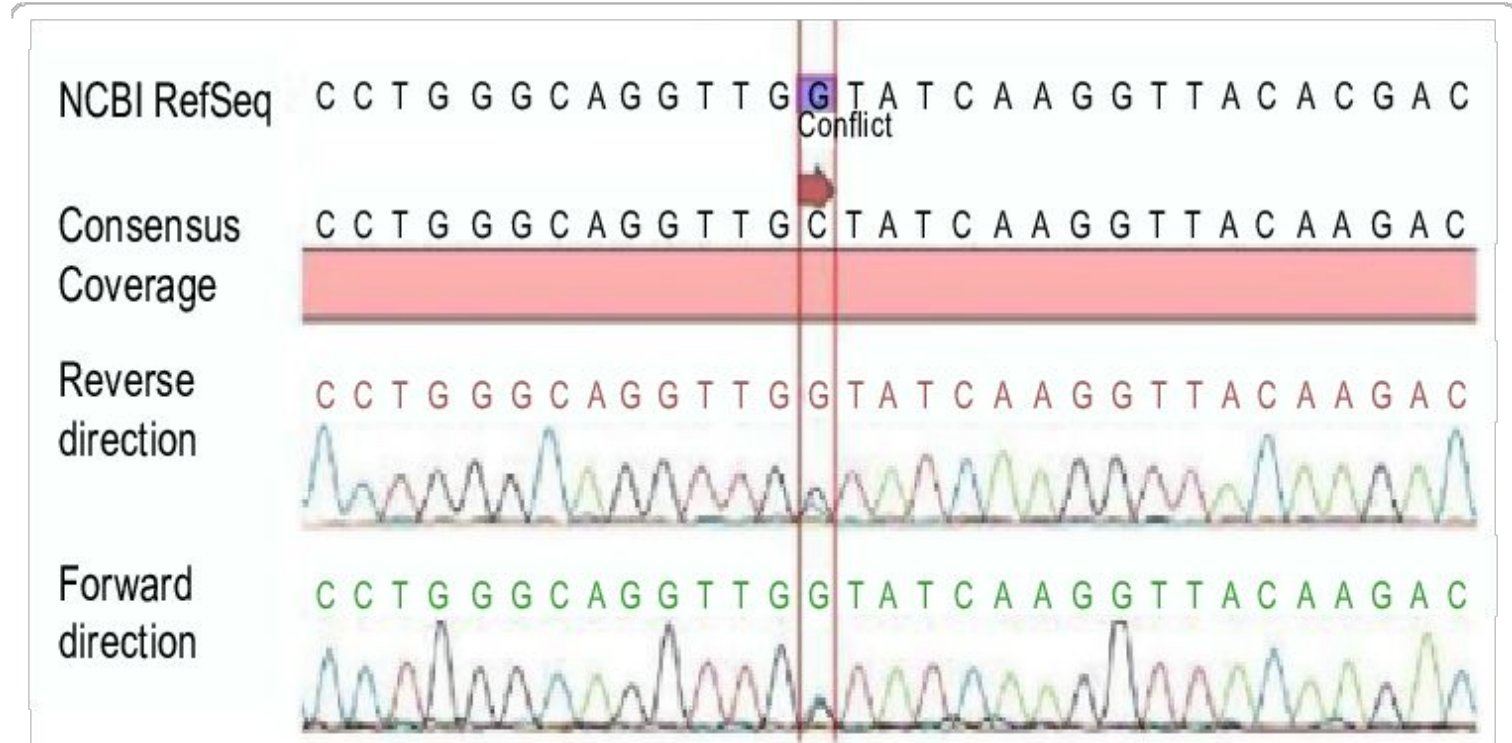


4. Pseudogenes

e.g. *PMS2*, *CYP2D6*, *CHEK2*, *SMN1*, *PKD1*



Bi-Directional Sequencing



Reverse Complementarity (RC)

Writing sequences

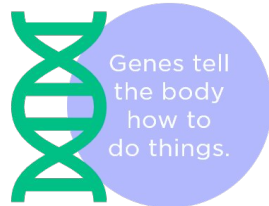
- ❑ Written 5'-3'
 - ❑ ATGGGTAGCGGTCATGATAC
- ❑ Complement
 - ❑ TACCCATCGCCAGTACTATG
- ❑ Reverse (inverse)
 - ❑ CATAGTACTGGCGATGGGTA
- ❑ Reverse complement
 - ❑ GTATCATGACCGCTACCCAT

Upstream and downstream

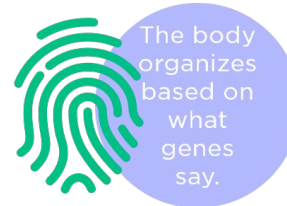


Phenotypes

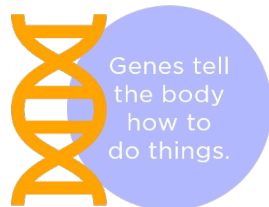
Genotype: DNA



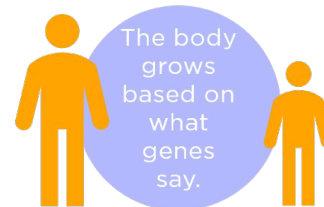
Phenotype: Thumbprint



Genotype: DNA



Phenotype: Height



Caduceus use

Variant effect prediction: a task to detect whether a genetic mutation influences a phenotype

Motivation

- Bi-Directionality
- Reverse Complementarity
- Long-Range Dependencies
- Limitations of Existing Models



Models

Mamba

- Processes DNA in long-range sequences.
- Recognizes reverse complement strands.

Selective State Space Model

with Hardware-aware State Expansion

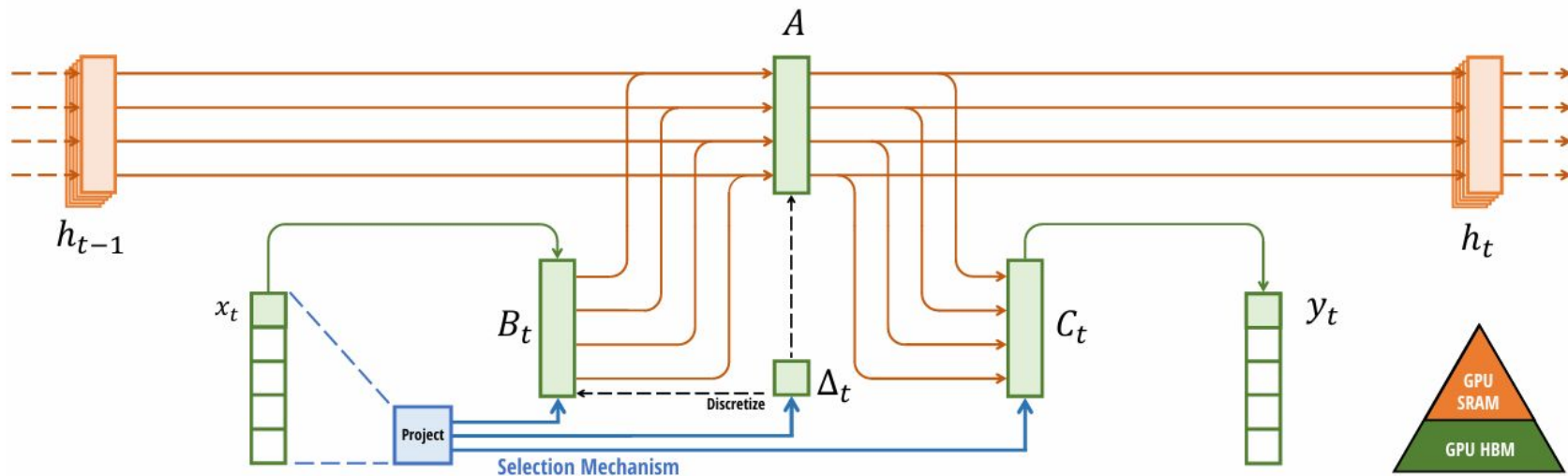
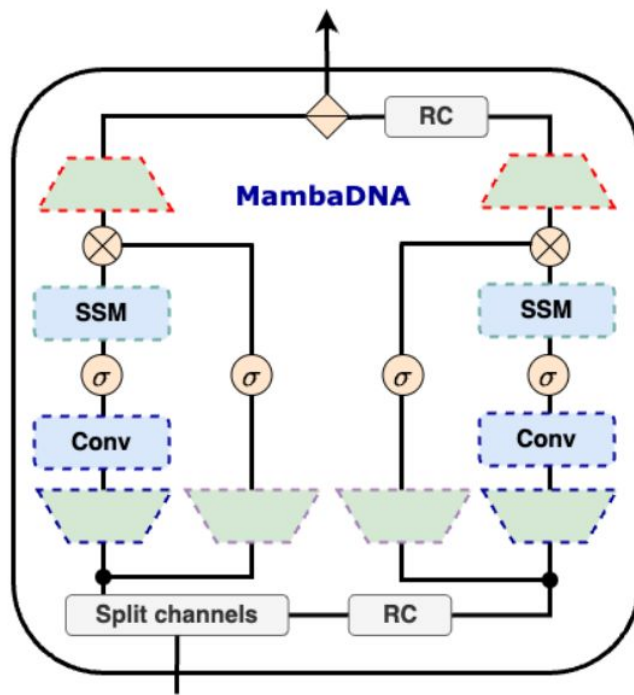
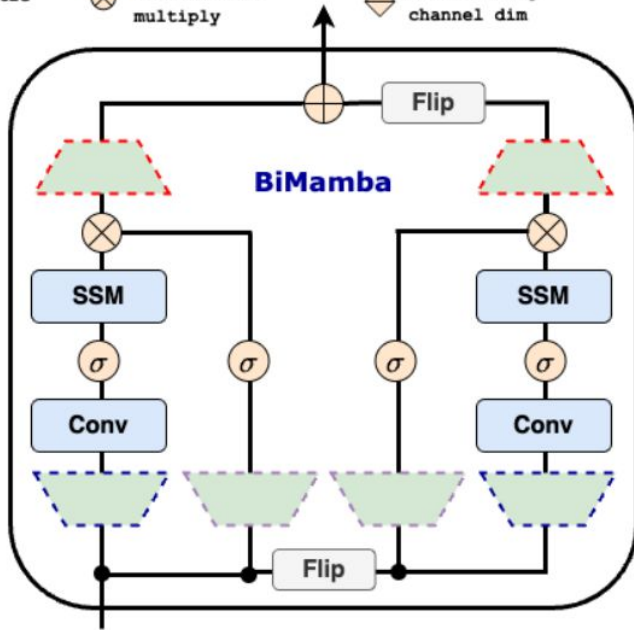
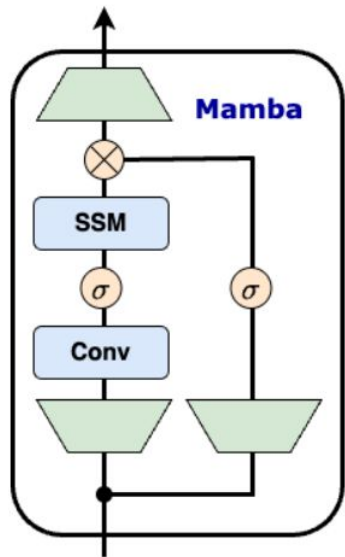
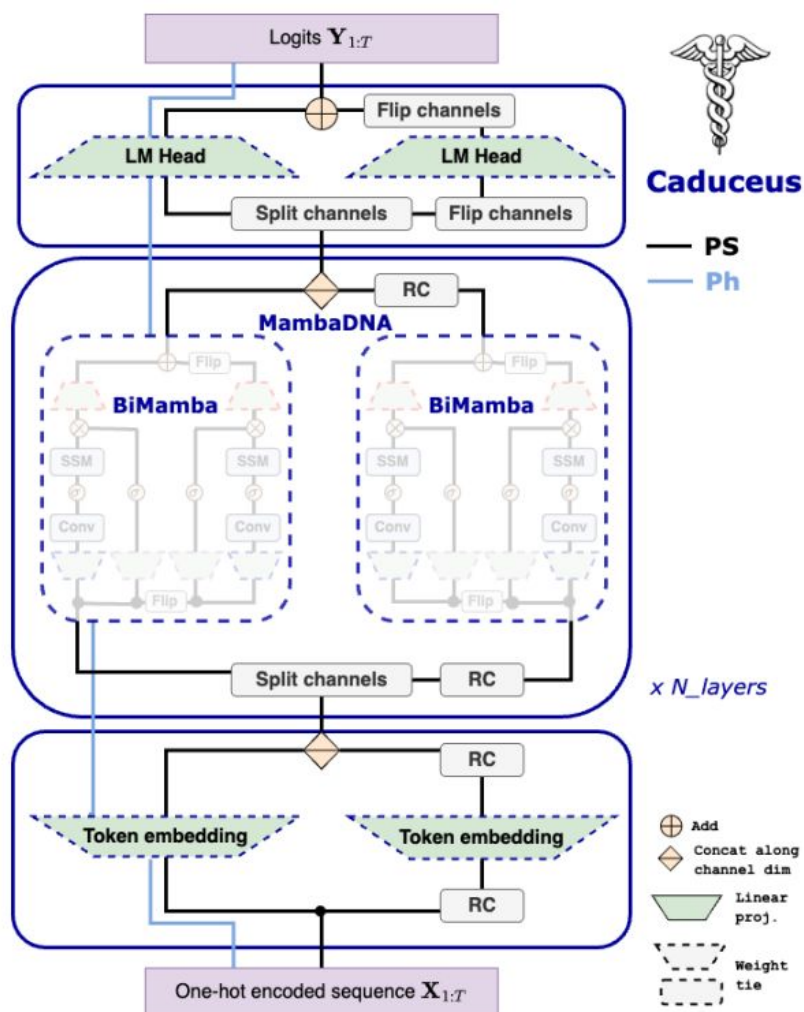


Figure 1: (**Overview.**) Structured SSMs independently map each channel (e.g. $D = 5$) of an input x to output y through a higher dimensional latent state h (e.g. $N = 4$). Prior SSMs avoid materializing this large effective state (DN , times batch size B and sequence length L) through clever alternate computation paths requiring time-invariance: the (Δ, A, B, C) parameters are constant across time. Our selection mechanism adds back input-dependent dynamics, which also requires a careful hardware-aware algorithm to only materialize the expanded states in more efficient levels of the GPU memory hierarchy.





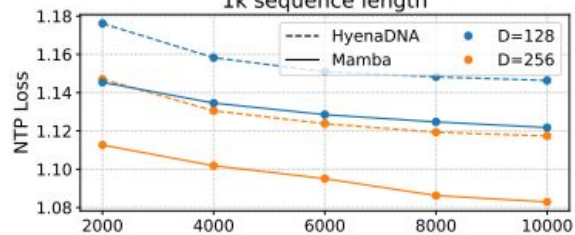
Caduceus



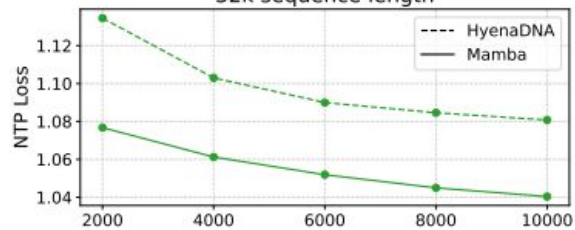


Results

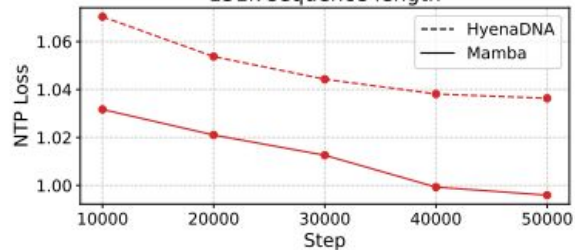
Mamba vs. HyenaDNA
1k sequence length



32k sequence length

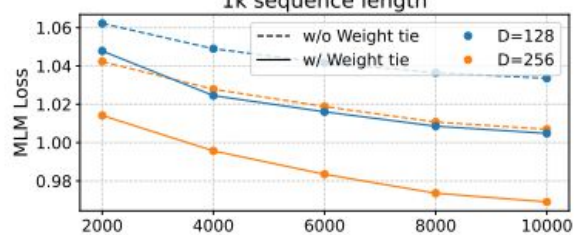


131k sequence length

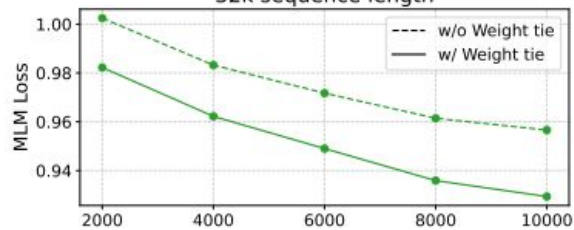


(a)

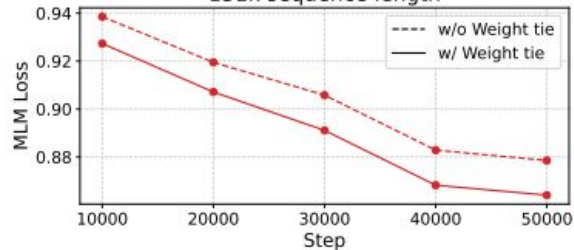
Effect of Weight Tying on Pre-training
1k sequence length



32k sequence length

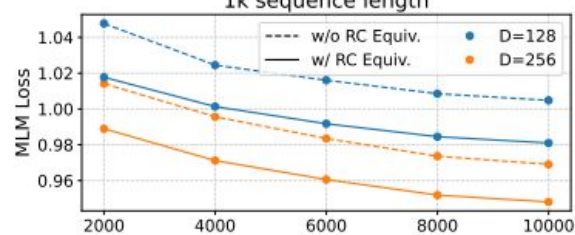


131k sequence length

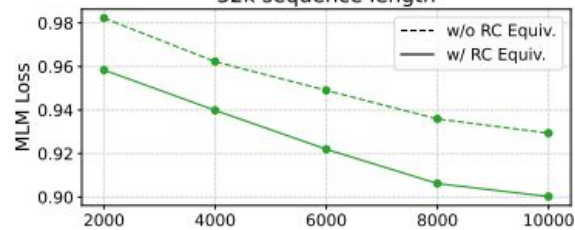


(b)

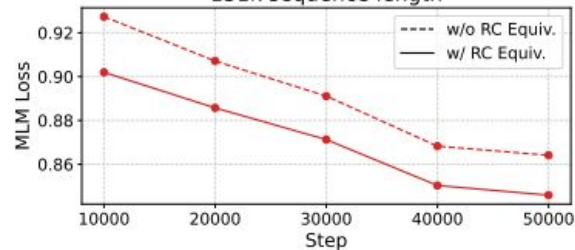
Effect of RC Equiv. on Pre-training
1k sequence length



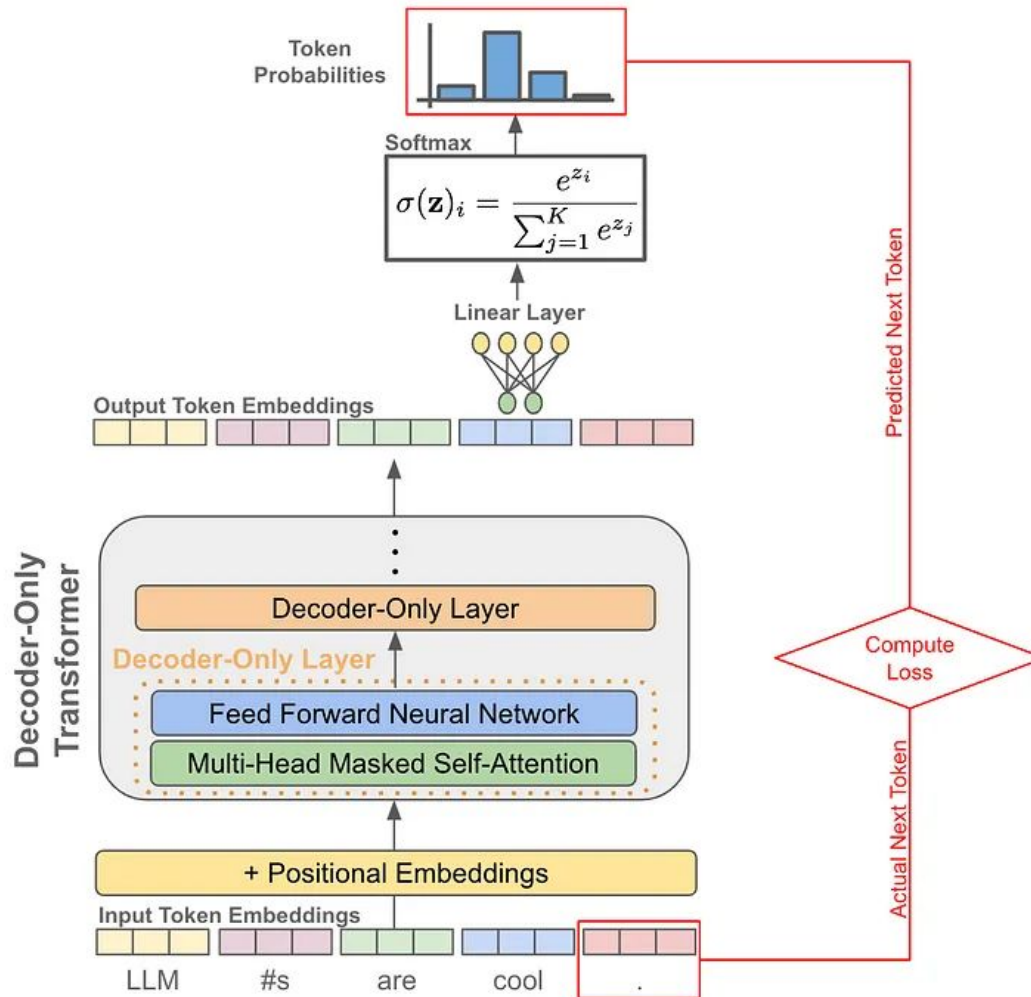
32k sequence length

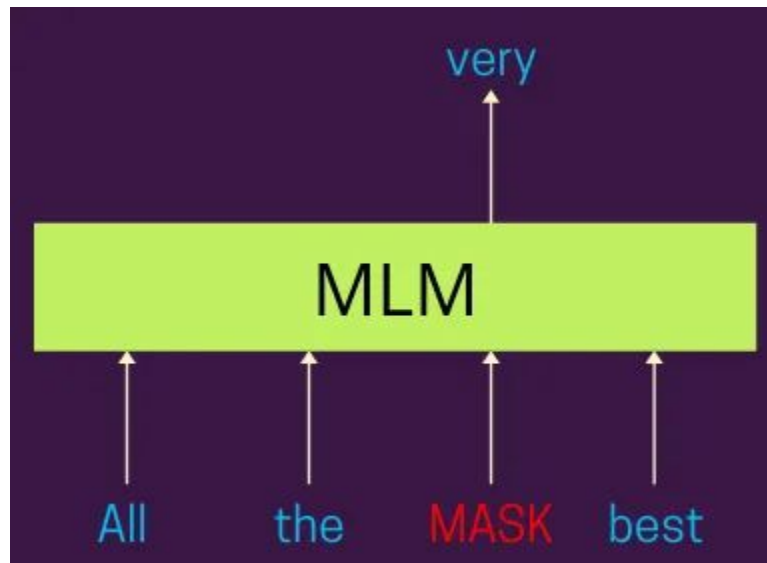


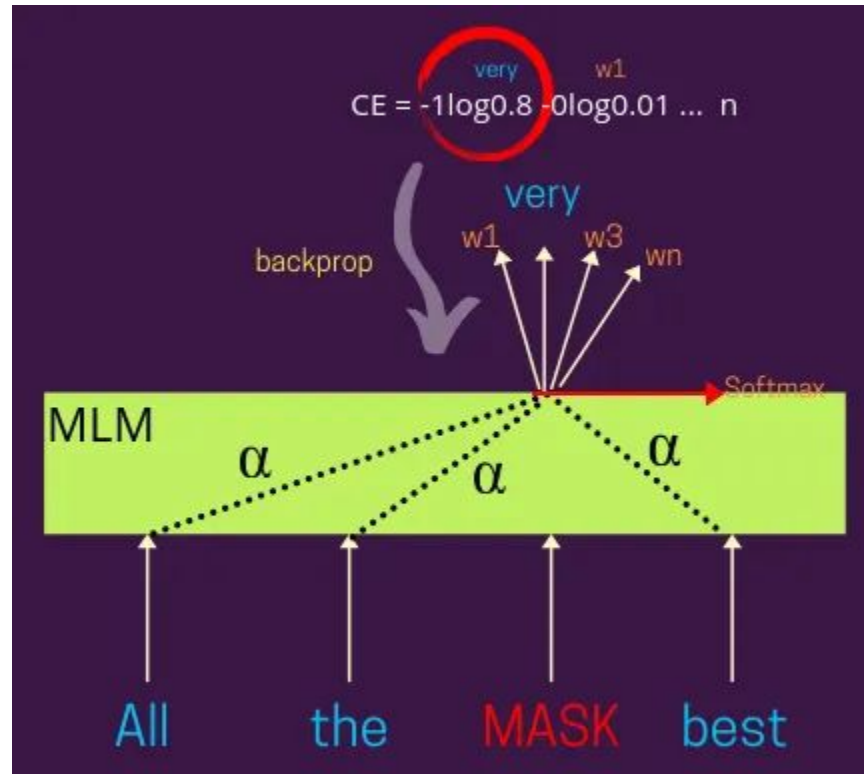
131k sequence length



(c)









Genomics Benchmarks

	CNN (264k)	HYENADNA (436k)	MAMBA (468k)	CADUCEUS w/o EQUIV. (470k)	CADUCEUS-PH (470k)	CADUCEUS-PS (470k)
MOUSE ENHANCERS	0.715 \pm 0.087	<i>0.780</i> \pm 0.025	0.743 \pm 0.054	0.770 \pm 0.058	0.754 \pm 0.074	0.793 \pm 0.058
CODING VS. INTERGENOMIC	0.892 \pm 0.008	0.904 \pm 0.005	0.904 \pm 0.004	0.908 \pm 0.003	0.915 \pm 0.003	<i>0.910</i> \pm 0.003
HUMAN VS. WORM	0.942 \pm 0.002	0.964 \pm 0.002	0.967 \pm 0.002	<i>0.970</i> \pm 0.003	0.973 \pm 0.001	0.968 \pm 0.002
HUMAN ENHANCERS COHN	0.702 \pm 0.021	0.729 \pm 0.014	0.732 \pm 0.029	0.741 \pm 0.008	0.747 \pm 0.004	<i>0.745</i> \pm 0.007
HUMAN ENHANCER ENSEMBL	0.744 \pm 0.122	0.849 \pm 0.006	0.862 \pm 0.008	0.883 \pm 0.002	<i>0.893</i> \pm 0.008	0.900 \pm 0.006
HUMAN REGULATORY	0.872 \pm 0.005	0.869 \pm 0.012	0.814 \pm 0.211	0.871 \pm 0.007	<i>0.872</i> \pm 0.011	0.873 \pm 0.007
HUMAN OCR ENSEMBL	0.698 \pm 0.013	0.783 \pm 0.007	0.815 \pm 0.002	0.818 \pm 0.003	0.828 \pm 0.006	<i>0.818</i> \pm 0.006
HUMAN NONTATA PROMOTERS	0.861 \pm 0.009	0.944 \pm 0.002	0.933 \pm 0.007	0.933 \pm 0.006	0.946 \pm 0.007	<i>0.945</i> \pm 0.010



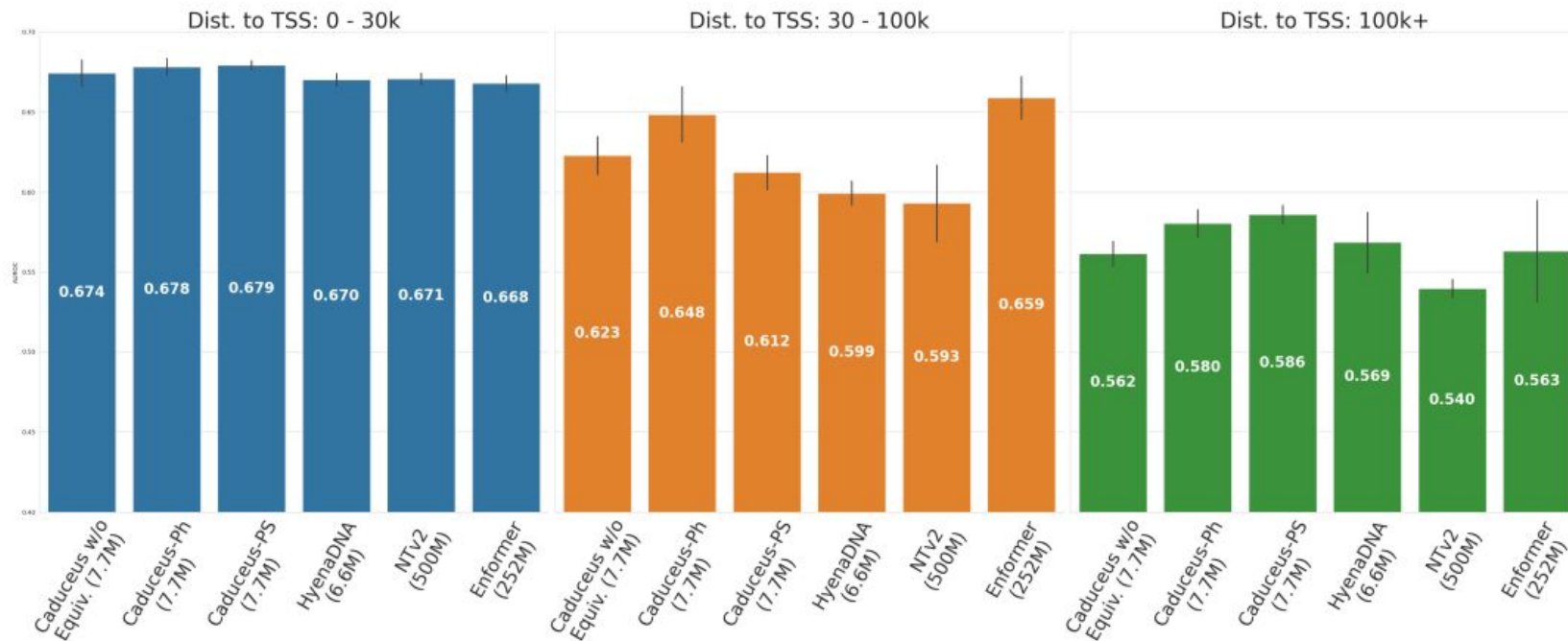
Nucleotide Transformer Tasks

	> 100M PARAM. MODELS			< 2M PARAM. MODELS		
	ENFORMER (252M)	DNABERT-2 (117M)	NT-v2 (500M)	HYENADNA (1.6M)	CADUCEUS-PH (1.9M)	CADUCEUS-PS (1.9M)
<i>Histone Markers</i>						
H3	0.719±0.048	0.785±0.033	0.784±0.047	0.779±0.037	0.815 ±0.048	0.799±0.029
H3K14AC	0.288±0.077	0.516±0.028	0.551±0.021	0.612±0.065	0.631 ±0.026	0.541±0.212
H3K36ME3	0.344±0.055	0.591±0.020	0.625 ±0.013	0.613±0.041	0.601±0.129	0.609±0.109
H3K4ME1	0.291±0.061	0.511±0.028	0.550 ±0.021	0.512±0.024	0.523±0.039	0.488±0.102
H3K4ME2	0.211±0.069	0.336±0.040	0.319±0.045	0.455±0.095	0.487 ±0.170	0.388±0.101
H3K4ME3	0.158±0.072	0.352±0.077	0.410±0.033	0.549 ±0.056	0.544±0.045	0.440±0.202
H3K79ME3	0.496±0.042	0.613±0.030	0.626±0.026	0.672±0.048	0.697 ±0.077	0.676±0.026
H3K9AC	0.420±0.063	0.542±0.029	0.562±0.040	0.581±0.061	0.622 ±0.030	0.604±0.048
H4	0.732±0.076	0.796±0.027	0.799±0.025	0.763±0.044	0.811 ±0.022	0.789±0.020
H4AC	0.273±0.063	0.463±0.041	0.495±0.032	0.564±0.038	0.621 ±0.054	0.525±0.240
<i>Regulatory Annotation</i>						
ENHANCER	0.451±0.108	0.516±0.098	0.548 ±0.144	0.517±0.117	0.546±0.073	0.491±0.066
ENHANCER TYPES	0.309±0.134	0.423±0.051	0.424±0.132	0.386±0.185	0.439 ±0.054	0.416±0.095
PROMOTER: ALL	0.954±0.006	0.971±0.006	0.976 ±0.006	0.960±0.005	0.970±0.004	0.967±0.004
NONTATA	0.955±0.010	0.972±0.005	0.976 ±0.005	0.959±0.008	0.969±0.011	0.968±0.006
TATA	0.960±0.023	0.955±0.021	0.966 ±0.013	0.944±0.040	0.953±0.016	0.957±0.015
<i>Splice Site Annotation</i>						
ALL	0.848±0.019	0.939±0.009	0.983 ±0.008	0.956±0.011	0.940±0.027	0.927±0.021
ACCEPTOR	0.914±0.028	0.975±0.006	0.981 ±0.011	0.958±0.010	0.937±0.033	0.936±0.077
DONOR	0.906±0.027	0.963±0.006	0.985 ±0.022	0.949±0.024	0.948±0.025	0.874±0.289



Gene Expression

Predicting Effects of Variants on Gene Expression





Any Questions?



Thank you