

# JanusDNA: A Powerful Bi-directional Hybrid DNA Foundation Model

Qihao Duan, Bingding Huang, Zhenqiao Song, Irina Lehmann, Lei Gu, Roland Eils, Benjamin Wild

qihao.duan, benjamin.wild@bih-charite.de

BIH Berlin Institute of Health @Charité

復旦大學 FUDAN UNIVERSITY

MAX PLANCK INSTITUTE FOR HEART AND LUNG RESEARCH

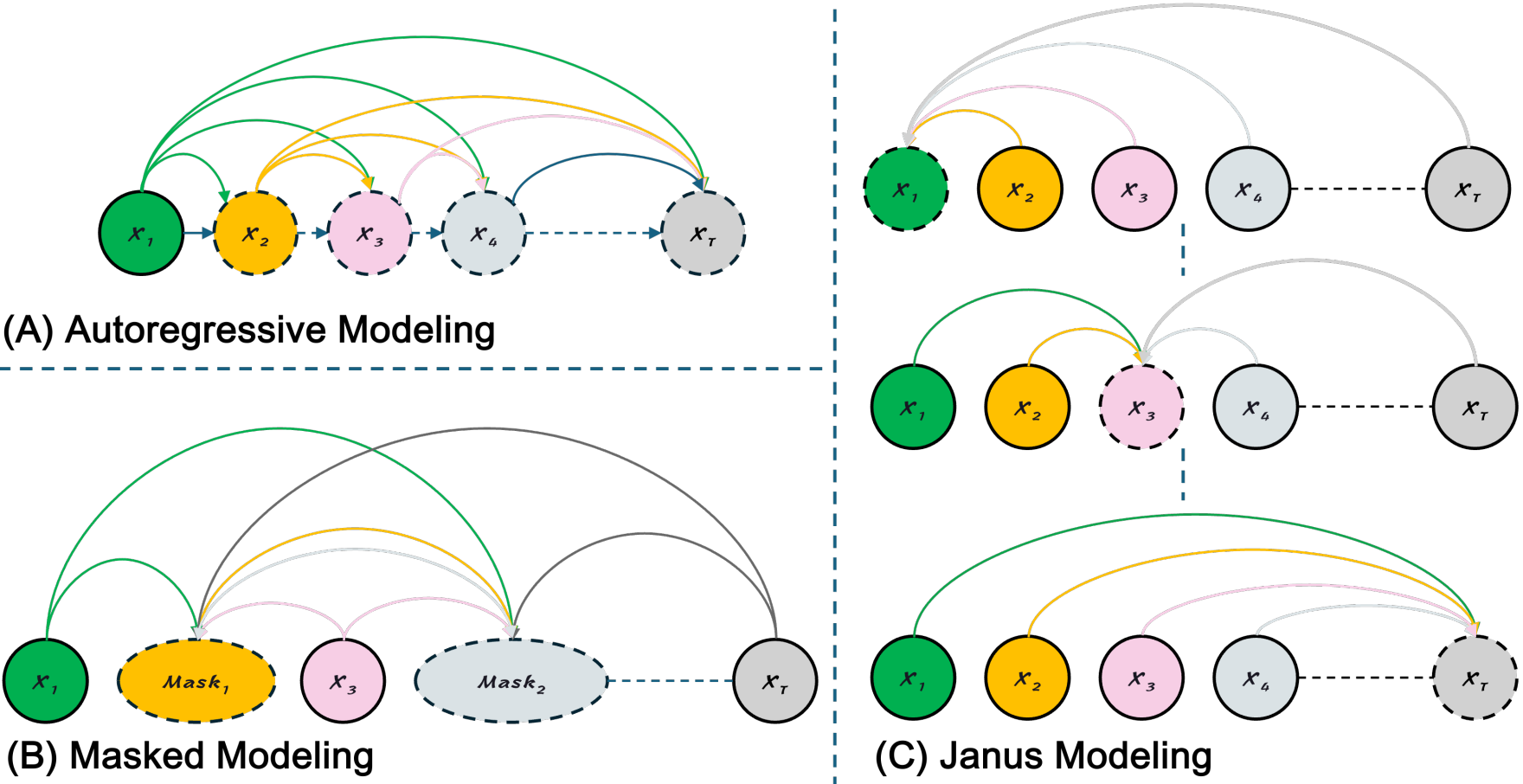
Carnegie Mellon University

SZTU 深圳技术大学

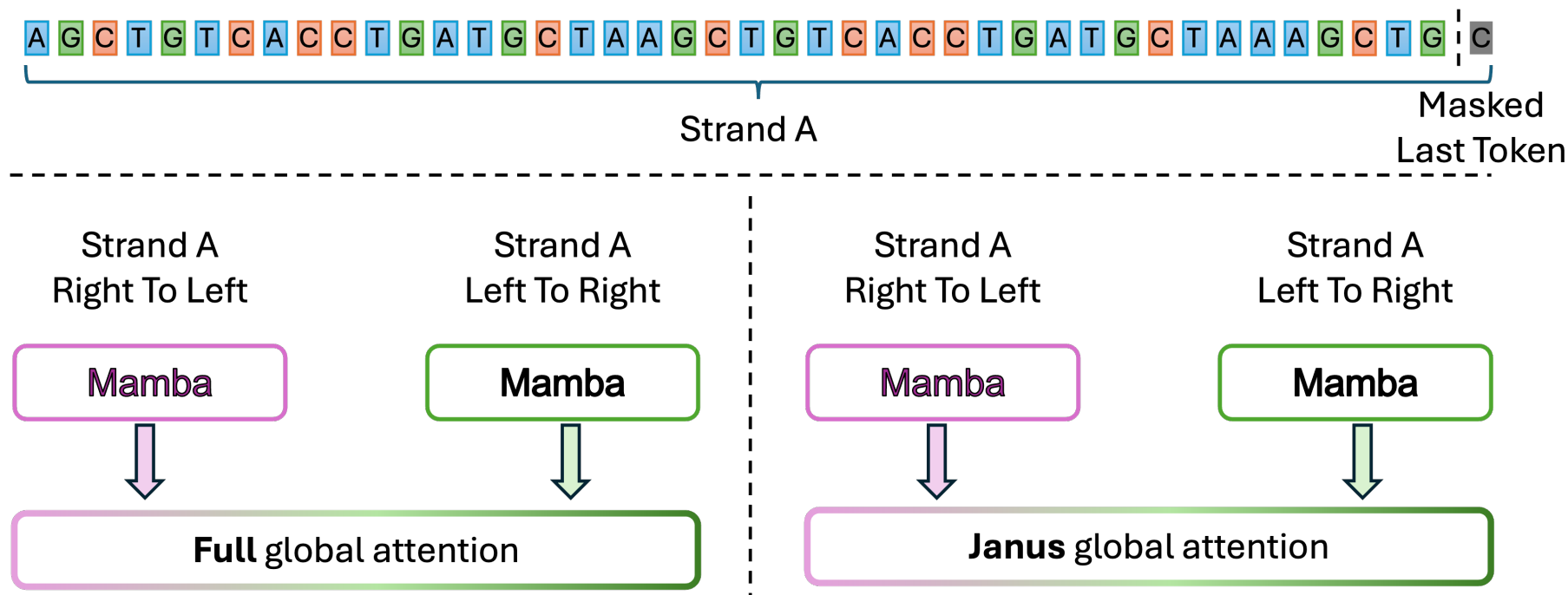
NEURAL INFORMATION PROCESSING SYSTEMS

## Janus Modeling

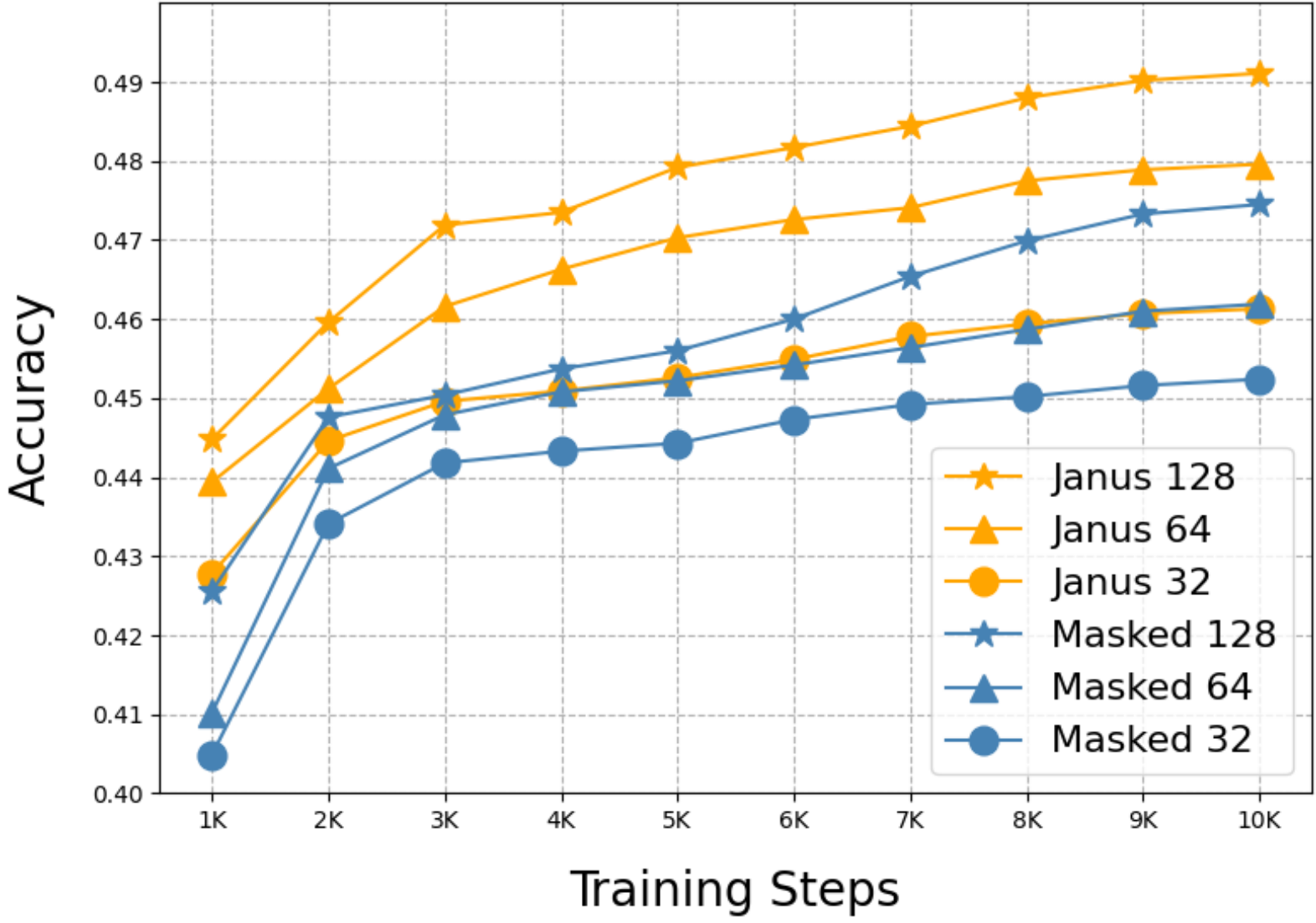
Efficiently learn from every token with bidirectional understanding.



## Ablation Experiments



## Masked vs. Janus training



## Contributions

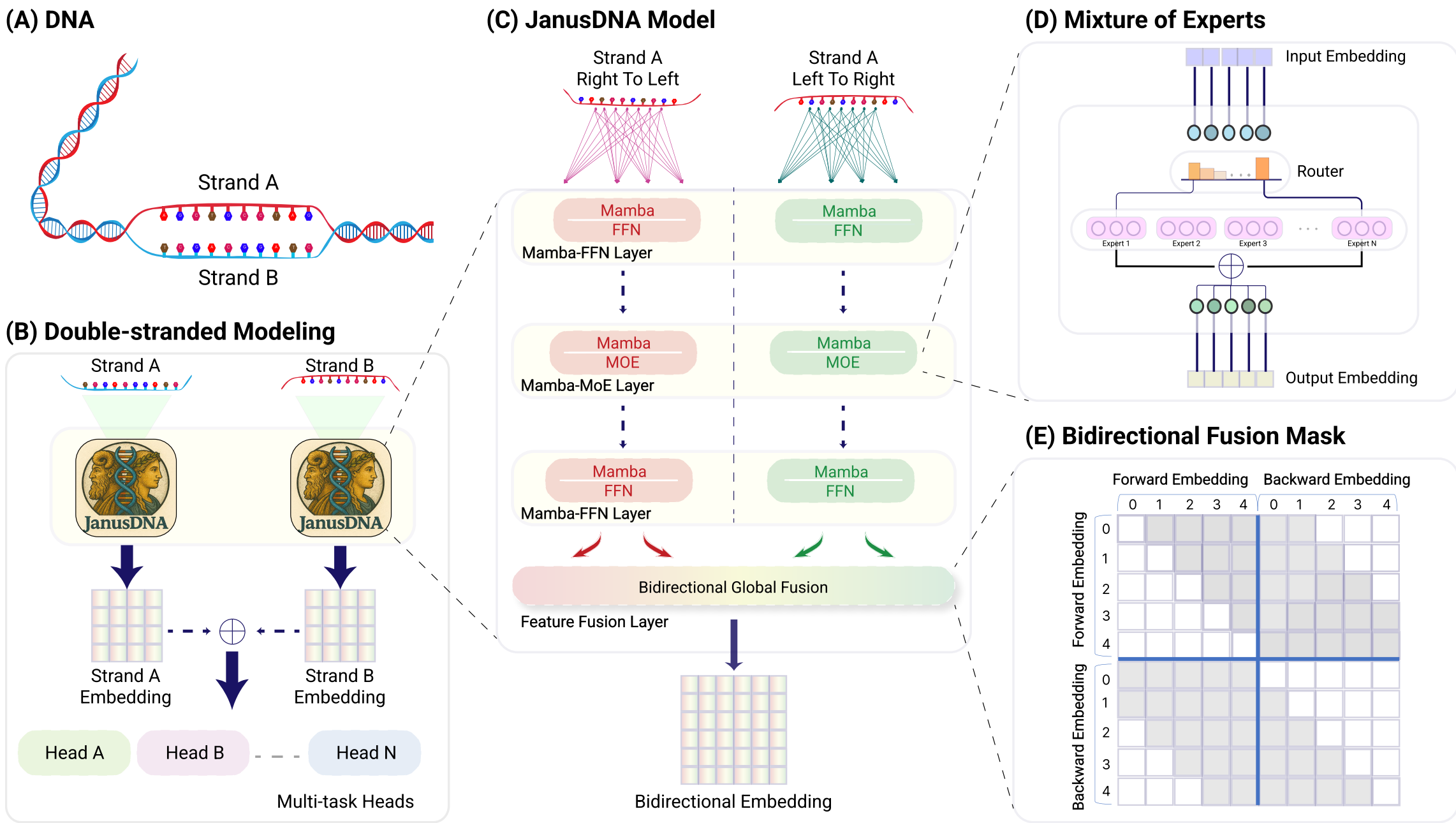
- JanusDNA introduces the first bidirectional DNA foundation model, leveraging a novel **Janus Modeling** paradigm that integrates autoregressive optimization efficiency with masked modeling's bidirectional comprehension.
- Its **Hybrid Mamba–Attention–MoE** architecture enables **global long-range dependency modeling** and **single-nucleotide resolution** processing of sequences **up to 1 million base pairs on a single 80 GB GPU**.
- JanusDNA outperforms models up to **250×** larger and achieves state-of-the-art results across diverse short- and long-range genomic tasks.

## Motivation

- Global attention is stronger in modeling long-range dependencies in DNA, but it is memory-intensive → **leading to limited length coverage and low resolution.**
- DNA functions bidirectionally → **while unidirectional understanding introduces knowledge bias.**
- Traditional masked modeling learns from only about 15% of input tokens per training step → **resulting in low learning efficiency.**

## JanusDNA Architecture and Workflow

Integrating *bidirectional understanding* into training while integrating *bi-strand understanding* into inference with *global attention*.



## Short-range Benchmarks Genomic Benchmark, Nucleotide Transformer

- Achieves performance comparable to or exceeding all previous SOTA models of similar scale.
- Surpasses models with **250×** more activated parameters on 12 out of 18 NT benchmark tasks.

## Long-range Benchmarks DNALONGBENCH

- Outperforms models **30×** larger, including expert and SOTA models.

## Short-range Benchmarks 200 – 4,776 base pairs

MODELS ACTIVATED PARAM	CNN (264K)	HYENADNA (436K)	MAMBA (468K)	CADUCEUS PH (470K)	CADUCEUS PS (470K)	CONVNOVA (386K)	JANUSDNA MLP w/ MID-ATTN (426K)	JANUSDNA MLP w/o MID-ATTN (431K)
MOUSE ENHANCERS	0.715±0.087	0.780±0.025	0.743±0.054	0.754±0.074	<b>0.793</b> ±0.058	0.784±0.009	0.770±0.048	0.769±0.029
CODING VS. INTERGENOMIC	0.892±0.008	0.904±0.005	0.904±0.004	0.915±0.003	0.910±0.003	<b>0.943</b> ±0.001	0.912±0.003	0.911±0.001
HUMAN VS. WORM	0.942±0.002	0.964±0.002	0.967±0.002	<b>0.973</b> ±0.001	0.968±0.002	0.967±0.002	0.971±0.001	0.971±0.001
HUMAN ENHANCERS COHN	0.702±0.021	0.729±0.014	0.732±0.029	<b>0.747</b> ±0.004	0.745±0.007	0.743±0.005	0.741±0.005	0.742±0.006
HUMAN ENHANCER ENSEMBL	0.744±0.122	0.849±0.006	0.862±0.008	0.893±0.008	<b>0.900</b> ±0.006	<b>0.900</b> ±0.004	0.897±0.004	0.899±0.004
HUMAN REGULATORY	0.872±0.005	0.869±0.012	0.814±0.211	0.872±0.011	0.873±0.007	0.873±0.002	<b>0.877</b> ±0.005	0.868±0.008
HUMAN OCR ENSEMBL	0.698±0.013	0.783±0.007	0.815±0.002	<b>0.828</b> ±0.006	0.818±0.006	0.793±0.004	0.822±0.003	0.824±0.001
HUMAN NONTATA PROMOTERS	0.861±0.009	0.944±0.002	0.933±0.007	0.946±0.007	0.945±0.010	0.951±0.003	<b>0.957</b> ±0.004	0.954±0.010

Genomic Benchmark

> 100M ACTIVATED PARAM. MODELS				< 2M ACTIVATED PARAM. MODELS			
ENFORMER (252M)	DNABERT-2 (117M)	NT-v2 (500M)		HYENADNA (1.6M)	CADUCEUS-PH (1.9M)	CONVNOVA (1.7M)	JANUSDNA MLP w/ MIDATTN (2.001M)

Histone Markers								
H3	0.719±0.048	0.785±0.033	0.784±0.047	0.779±0.037	0.815±0.048	0.812±0.017	<b>0.835</b> ±0.009	0.831±0.023
H3K14AC	0.288±0.077	0.516±0.028	0.551±0.021	0.612±0.065	0.631±0.026	0.644±0.009	<b>0.729</b> ±0.022	0.718±0.026
H3K36ME3	0.344±0.055	0.591±0.020	0.625±0.013	0.613±0.041	0.601±0.129	0.661±0.019	<b>0.702</b> ±0.015	0.699±0.025
H3K4ME1	0.291±0.061	0.511±0.028	0.550±0.021	0.512±0.024	0.523±0.039	0.554±0.023	0.615±0.035	<b>0.616</b> ±0.018
H3K4ME2	0.211±0.069	0.336±0.040	0.319±0.045	0.455±0.095	0.487±0.170	0.485±0.032	<b>0.589</b> ±0.023	0.586±0.019
H3K4ME3	0.158±0.072	0.352±0.077	0.410±0.033	0.549±0.056	0.544±0.045	0.566±0.027	<b>0.688</b> ±0.026	0.675±0.014
H3K79ME3	0.496±0.042	0.613±0.030	0.626±0.026	0.672±0.048	0.697±0.077	0.700±0.007	<b>0.747</b> ±0.013	0.743±0.009
H3K9AC	0.420±0.063	0.542±0.029	0.562±0.040	0.581±0.061	0.622±0.030	0.658±0.011	<b>0.673</b> ±0.014	0.661±0.027
H4	0.732±0.076	0.796±0.027	0.799±0.025	0.763±0.044	0.811±0.022	0.808±0.008	0.812±0.011	<b>0.813</b> ±0.013
H4AC	0.273±0.063	0.463±0.041	0.495±0.032	0.564±0.038	0.621±0.054	0.636±0.011	0.698±0.013	<b>0.705</b> ±0.023

Regulatory Annotation								
ENHANCER	0.451±0.108	0.516±0.098	0.548±0.144	0.517±0.117	0.546±0.073	<b>0.586</b> ±0.038	0.559±0.042	0.542±0.044
ENHANCER TYPES	0.309±0.134	0.423±0.051	0.424±0.132	0.386±0.185	0.439±0.054	0.500±0.018	<b>0.503</b> ±0.038	0.492±0.096
PROMOTER: ALL	0.954±0.006	0.971±0.006	<b>0.976</b> ±0.006	0.960±0.005	0.970±0.004	0.967±0.001	0.970±0.002	0.970±0.003
NONTATA	0.955±0.010	0.972±0.005	<b>0.976</b> ±0.005	0.959±0.008	0.969±0.011	0.968±0.003	0.971±0.004	0.971±0.003
TATA	0.960±0.023	0.955±0.021	0.966±0.013	0.944±0.040	0.953±0.016	<b>0.969</b> ±0.003	0.958±0.007	0.960±0.008

Splice Site Annotation								
ALL	0.848±0.019	0.939±0.009	<b>0.983</b> ±0.008	0.956±0.011	0.940±0.027	0.965±0.004	0.967±0.005	0.943±0.020
ACCEPTOR	0.914±0.028	0.975±0.006	<b>0.981</b> ±0.011	0.958±0.010	0.937±0.033	0.971±0.003	0.957±0.012	0.961±0.009
DONOR	0.906±0.027	0.963±0.006	<b>0.985</b> ±0.022	0.949±0.024	0.948±0.025	<b>0.965</b> ±0.003	0.948±0.008	0.935±0.016

Nucleotide Transformer

## Long-range Benchmarks 450,000 base pairs

MODELS ACTIVATED PARAM	EXPERT MODEL (252M)	CADUCEUS-PH (7.7M)	JANUSDNA MLP w/o MID-ATTN (7.662M)	JANUSDNA MLP w/o MID-ATTN (7.745M)
ARTERY TIBIAL	0.741	0.690	0.803	<b>0.852</b>
ADIPOSE SUBCUTANEOUS	0.736	0.759	0.741	<b>0.769</b>
CELLS CULTURED FIBROBLASTS	0.639	0.690	0.771	<b>0.802</b>
MUSCLE SKELETAL	0.621	0.789	0.803	<b>0.864</b>
NERVE TIBIAL	0.683	0.842	0.877	<b>0.914</b>
SKIN NOT SUN EXPOSED SUPRAPUBIC	0.710	0.812	0.875	<b>0.903</b>
SKIN SUN EXPOSED LOWER LEG	0.700	0.692	0.706	<b>0.846</b>
THYROID	0.612	0.703	0.752	<b>0.793</b>
WHOLE BLOOD	0.689	0.769	0.794	<b>0.821</b>

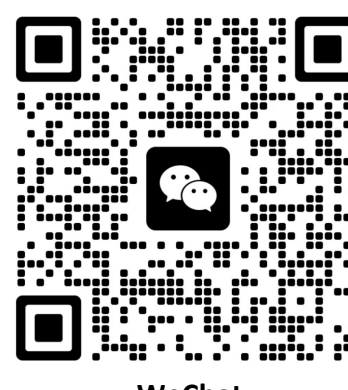
DNALONGBENCH

## References

- Yair Schiff, et al. Caduceus: Bi-directional equivariant long-range dna sequence modeling, ICML 2024.
- Hugo Dalla-Torre, et al. Nucleotide transformer: building and evaluating robust foundation models for human genomics. Nature Methods, 22(2):287–297, 2025.
- Nguyen, E., et al. HyenaDNA: long-range genomic sequence modeling at single nucleotide resolution. Neurips, 2024.
- Žiga Avsec, et al. Effective gene expression prediction from sequence by integrating long-range interactions. Nature methods, 18(10):1196–1203, 2021.
- Gu, A. and T. Dao. Mamba: Linear-time sequence modeling with selective state spaces. arXiv preprint arXiv:2312.00752.
- Zhou, Z., et al. Dnabert-2: Efficient foundation model and benchmark for multi-species genome. arXiv preprint arXiv:2306.15006.
- Cheng, Wenduo, et al. "Dnalongbench: a benchmark suite for long-range dna prediction tasks." bioRxiv (2025).



Code and Paper



WeChat