

## Benchmarking Diagnostic Performance of Genomic Language Model

Woojong Sim<sup>1</sup>, Jiho Lee<sup>1</sup>, Jayoung Park<sup>2</sup>, Yujin Kim<sup>3,4</sup>, Soowhee Kim<sup>3,4</sup>, Yeojin Ryu<sup>3,4</sup>, Hwanseok

Sim<sup>5</sup>, Joon-Yong An<sup>3,4,5,\*</sup> and Minji Jeon<sup>1,6,\*\*</sup>

<sup>1</sup>*Department of Biomedical Sciences, Korea University College of Medicine*

<sup>2</sup>*Department of Biomedical Informatics, Korea University College of Medicine*

<sup>3</sup>*Department of Integrated Biomedical and Life Science, Korea University*

<sup>4</sup>*L-HOPE Program for Community-Based Total Learning Health Systems*

<sup>5</sup>*School of Biosystems and Biomedical Sciences, Korea University College of Health Sciences*

<sup>6</sup>*Department of Biomedical Informatics, Biomedical Research Center, Korea University Anam Hospital*

\*Corresponding author: [joonan30@korea.ac.kr](mailto:joonan30@korea.ac.kr)

\*\*Corresponding author: [mjeon@korea.ac.kr](mailto:mjeon@korea.ac.kr)

Genomic language models (gLMS) have emerged as powerful tools for interpreting noncoding DNA, beyond the limitations of traditional variant effect predictors that rely on fixed annotations or sequence alignments. While these deep learning models have shown promise for recognition of functional regions through learning on regulatory grammar from raw sequence, recent reviews and benchmarks suggest a significant performance gap exists between the zero-shot capabilities of pre-trained gLMS and highly specialized biological tasks such as variant effect prediction and gene expression regulation by non-coding elements. This disparity arises due to the fact that language modeling is a separate field from understanding biological roles of genomic elements, highlighting the need for rigorous benchmarking on clinically relevant tasks. Here, we present the first systematic benchmarking of genomic language models, including DNABERT, DNABERT-2, Nucleotide Transformer, PhyloGPN, HyenaDNA and Evo 2 for whole-genome disease prediction. Using whole-genome sequencing data from 14,606 individuals across three major ASD cohorts, we trained a Set Transformer on variant-centered embeddings from each gLM and predict ASD diagnosis. DNABERT-2 consistently outperformed other models(AUROC = 0.5568), though overall predictive power remained modest. Fine-tuning marginally improved model performance, while integrating variant-level functional annotations such as conservation scores and regulatory impact predictions significantly enhanced prediction. Notably, high-impact variants prioritized by the models were enriched in ultra-conserved noncoding regions, implicating regulatory disruption as a

mechanism in ASD pathogenesis. These findings provide the systematic evaluation of gLMS for complex disease prediction and emphasize the need for integrative approaches to realize their full potential in complex trait diagnosis.