

# 1 **A long-context language model for the generation of bacteriophage genomes**

2 Bin Shao<sup>1,2\*</sup>

3

4 <sup>1</sup>Department of Molecular and Cellular Biology, Harvard University, Cambridge, MA, USA

5 <sup>2</sup>Present address: Klarman Cell Observatory, Broad Institute of MIT and Harvard, Cambridge, MA, USA

6 \*Correspondence should be addressed to Bin Shao ([shaobin@broadinstitute.org](mailto:shaobin@broadinstitute.org))

7

8 **Abstract:** Generative pre-trained transformers (GPTs) have revolutionized the field of natural language  
9 processing. Inspired by this success, we develop a long-context generative model for genomes. Our  
10 multiscale transformer model was pre-trained on unannotated bacteriophage genomes with byte-level  
11 tokenization. It generates *de novo* sequences up to 96K with functional genomic structure, including  
12 regulatory elements and novel proteins with phage-related functions. Our work paves the way for the *de*  
13 *novo* design of the whole genome.

14

15

16

17

18

19

20

21

22

23

24

25 Large pre-trained language models have drastically transformed the natural language processing (NLP)  
26 field<sup>1,2</sup>. Drawing on the similarity of natural language and genome sequences, genomic language models  
27 have been developed. These models were trained on large scale genomic datasets, and they effectively  
28 predict regulatory elements, uncover co-regulation patterns in protein and identify genome-wide  
29 variant effects<sup>3-7</sup>. However, it remains an open question whether language models can be tailored to  
30 generate genome-scale sequence with functional structure, which holds potential for the rational design  
31 of the whole genome.

32 Most of the current models used masked language modeling like BERT (Bidirectional Encoder  
33 Representations from Transformers)<sup>1</sup> which is not ideal for tasks that involve generating new content. In  
34 addition, current models face technical constraints such as short context size and aggregation of  
35 sequences in k-mer tokenization. These limitations hinder their ability to efficiently learn from genome  
36 scale data while maintaining the resolution needed for precise design of functional elements.

37 In this work, we developed a long-context language model that can generate *de novo* sequence with  
38 functional genomic structure. Our model draws inspiration from the GPT model<sup>2</sup>, which is renowned for  
39 its proficiency in generating long and coherent texts. We utilized a multiscale transformer structure<sup>8</sup>  
40 that enables us to train the model on unannotated bacteriophage genomes up to 96K bp at the single  
41 nucleotide-level. The trained model generates sequences that share similar genomic structure with the  
42 natural bacteriophage genomes. We found functional promoter and ribosome binding sites (RBS) in the  
43 5' untranslated regions (5'UTR) of the predicted genes. In addition, the proteins from the generated  
44 sequences are predicted to be structurally plausible and span a wide variety of functional families. We  
45 believe our model is a timely advance that paves the way for DNA design at the genome scale. The  
46 model is available from GitHub: <https://github.com/lingxusb/megaDNA>

47 To construct the training dataset, we collected bacteriophage genomes with high confidence from three  
48 sources including NCBI genebank, Metagenomic Gut Virus (MGV) catalogue<sup>9</sup> and Gut Phage Database  
49 (GPD)<sup>10</sup> (Supplementary Fig. 1). After data cleaning, we constructed a dataset with 99.7K bacteriophage  
50 genomes up to 96K bp, which was used to pre-train our model (Methods). The training data was byte-  
51 level tokenized. We employed the multi-scale transformer structure with long-context length from Yu et  
52 al.<sup>8</sup> and our model is named megaDNA, correspondingly. In model inference, we generated a total of  
53 1,024 sequences longer than 1K bp. Then geNomad<sup>11</sup> was used for functional annotation of the  
54 generated sequences. Among all these sequences, 607 of them have a virus score larger than zero and  
55 we focused our analysis on them. Their mean sequence length is 43K bp and the mean number of

56 predicted genes is 67, both are similar to the training dataset (Mean length: 48K bp, mean number of  
57 predicted genes: 68). The gene length distribution is close to that of the training dataset (Fig. 1b,  
58 average gene length: 558 bp vs 646 bp), while the predicted gene numbers show wider spread (Fig. 1c).  
59 The median virus score of these generated sequences is 0.84 and the maximum virus score is 0.97,  
60 comparable to the virus scores for natural bacteriophage genomes which range from 0.70 to 0.98 (Fig.  
61 1c). 223/607 (37%) of the generated sequences are predicted to be Caudoviricetes by geNomad (Fig.  
62 1d). As a comparison, 98% of the training dataset was classified as Caudoviricetes.

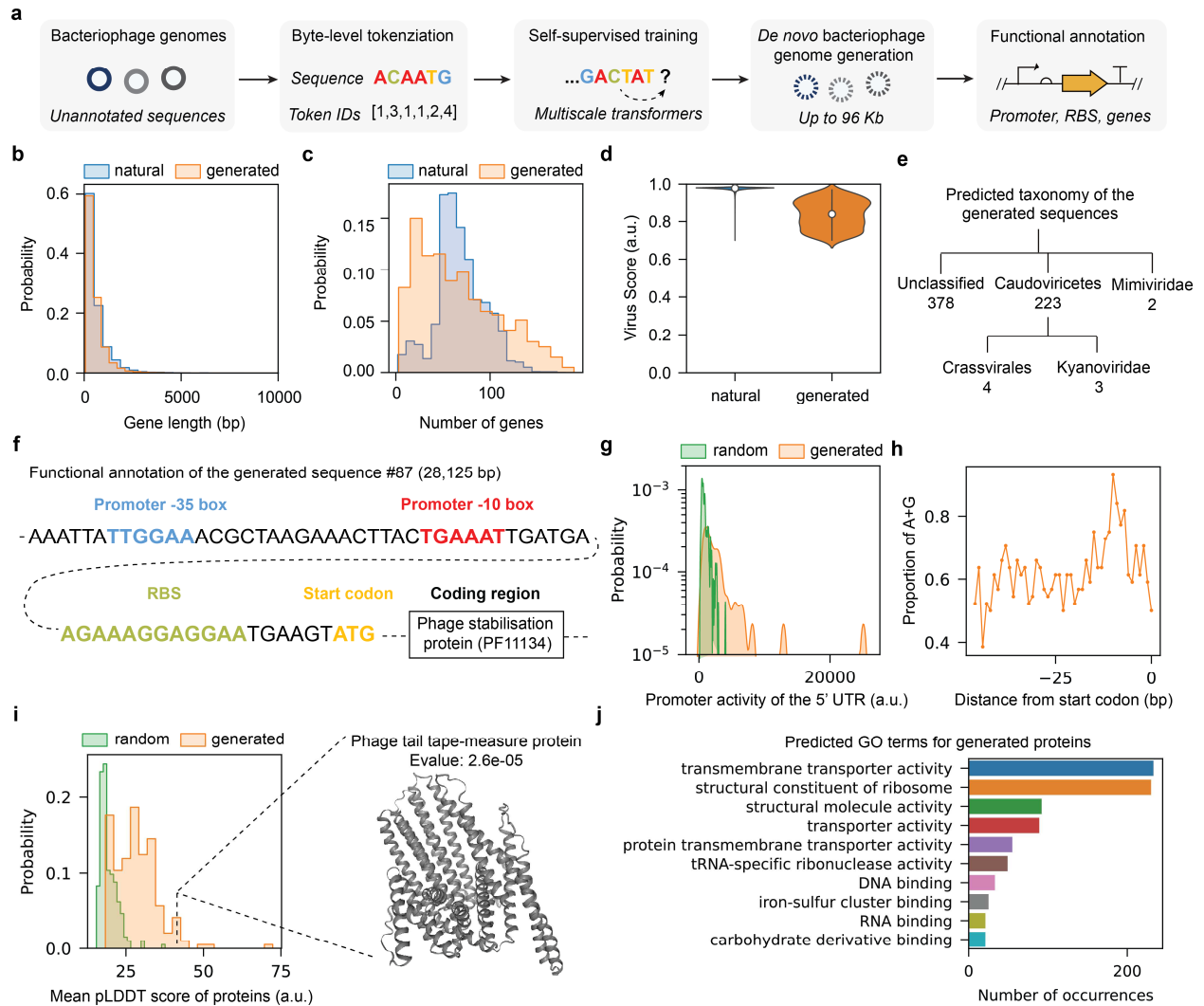
63 We then examined the 5'UTR region of the annotated genes in the generated sequences to see whether  
64 they contain functional regulatory elements like promoters and RBS to initiate transcription and  
65 translation. We chose generated sequence #87 for further analysis due to its high predicted virus score  
66 (0.96) and relatively small size (28K bp). Using a machine learning tool (Promoter Calculator)<sup>12</sup>, we  
67 identified the -35 box and -10 box of the promoter within the 5'UTR of the predicted gene. Notably, they  
68 closely aligned with the established consensus motifs: TTGACA and TATAAT (Fig. 1f). Prior to the start  
69 codon of the predicted phage stabilization protein, we observed a region enriched in adenine (A) and  
70 guanine (G) nucleotides, which is a motif characteristic of functional ribosome binding sites (Fig. 1f).  
71 Analyzing all 5'UTR sequences of the predicted genes from this sequence, we observed significantly  
72 higher mean promoter activity compared to random sequences of the same length. (Fig. 1g).  
73 Intriguingly, the proportion of A and G nucleotides peaked around 10 bp upstream of the start codon,  
74 close to the optimal position for RBS to drive translation initiation (Fig. 1h). This trend of A/G enrichment  
75 is also consistent across all the generated sequences (Supplementary Fig. 2). In short, our generated  
76 sequences harbor functional regulatory sequences that would enable expression of the predicted genes.

77 Among the annotated genes in our generated sequences, 343 of them were predicted to match  
78 geNomad markers. We employed ESMfold<sup>13</sup> to predict structures for these genes and calculated the  
79 average predicted local distance difference test (pLDDT) score. This score reflects the confidence of  
80 ESMfold on the predicted structures. The median pLDDT score for these proteins is higher than the  
81 pLDDT of random peptide sequence of the same length (28 vs 18). We also randomly sampled 10K  
82 annotated genes from the generated sequences and found a high pLDDT score for them (median value  
83 of 36, Supplementary Fig. 3), suggesting that these generated proteins are more likely to adopt a stable  
84 conformation. We further used deepFRI<sup>14</sup> for functional annotation of all generated proteins and we  
85 only retained proteins with high scores (> 0.5). Our analysis reveals several large protein families with  
86 functional roles, including the transporter activity and the structural molecule activity (Fig. 1j).

87 Interestingly, we identified several proteins with DNA-binding activity, and the predicted structure  
88 resembles the canonical helix-turn-helix (HTH) domain in this protein family (Supplementary Fig. 4).

89 To the best of our knowledge, our work presents the first long-context generative model for genomic  
90 sequences. Our language model effectively learns the high-order genomic grammar via a single step of  
91 self-supervised training on unannotated whole genomes. The generated sequences match the length of  
92 natural bacteriophage genomes and display functional genomic architecture. With further scaling up, we  
93 envision that the generative genomic models will pave the way for the *de novo* design of the genome  
94 sequence, which offers opportunities for breakthroughs in medicine, agriculture, and environmental  
95 science. This field also faces ongoing challenges in ethical considerations, biosafety, and regulatory  
96 frameworks, which are critical for the responsible advancement of generative modeling in synthetic  
97 biology.

98 **Figures**



99

100 **Figure 1. Language model generates sequences with functional genomic structures.**

101 **a)** the workflow schematic. **b)** comparison of gene length distributions between predicted genes in  
 102 generated sequences (n = 40,399) and a randomly sampled subset of genes from the training dataset (n  
 103 = 10,000). **c)** distributions of the number of predicted genes for the generated sequences (n = 607) and  
 104 the training dataset (n = 99,673). **d)** comparison of the predicted virus scores for generated sequences  
 105 and the training dataset. **e)** predicted taxonomy for the generated sequences. Only taxonomies with > 1  
 106 sequence are shown. **f)** functional annotation of a selected sequence fragment (generated sequence  
 107 #87). **g)** predicted promoter activity for all the 5'UTRs in the generated sequence #87 (n = 44), along  
 108 with the promoter activity of the random sequences with the same length. **h)** proportions of adenine (A)  
 109 and guanine (G) nucleotides preceding the start codon of all the predicted genes in the generated

110 sequence #87. **i)** mean predicted pLDDT scores for proteins with geNomad markers from generated  
111 sequences (sample size: n = 343; median value: 28) against random peptide sequences of the same  
112 lengths (sample size: n = 343; median value: 18). A sample generated protein is shown on the right. **j)**  
113 top 10 predicted functions of proteins derived from the generated sequences.

114 **References**

- 115 1. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. Bert: Pre-training of deep bidirectional  
116 transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- 117 2. Brown, T. *et al.* Language models are few-shot learners. *Adv Neural Inf Process Syst* **33**, 1877–  
118 1901 (2020).
- 119 3. Ji, Y., Zhou, Z., Liu, H. & Davuluri, R. V. DNABERT: pre-trained Bidirectional Encoder  
120 Representations from Transformers model for DNA-language in genome. *Bioinformatics* **37**,  
121 2112–2120 (2021).
- 122 4. Dalla-Torre, H. *et al.* The nucleotide transformer: Building and evaluating robust foundation  
123 models for human genomics. *bioRxiv* 2021–2023 (2023).
- 124 5. Benegas, G., Batra, S. S. & Song, Y. S. DNA language models are powerful predictors of genome-  
125 wide variant effects. *Proceedings of the National Academy of Sciences* **120**, e2311219120 (2023).
- 126 6. Hwang, Y., Cornman, A. L., Kellogg, E. H., Ovchinnikov, S. & Girguis, P. R. Genomic language  
127 model predicts protein co-regulation and function. *bioRxiv* 2023–2024 (2023).
- 128 7. Nguyen, E. *et al.* Hyenadna: Long-range genomic sequence modeling at single nucleotide  
129 resolution. *arXiv preprint arXiv:2306.15794* (2023).
- 130 8. Yu, L. *et al.* Megabyte: Predicting million-byte sequences with multiscale transformers. *arXiv*  
131 *preprint arXiv:2305.07185* (2023).
- 132 9. Nayfach, S. *et al.* Metagenomic compendium of 189,680 DNA viruses from the human gut  
133 microbiome. *Nat Microbiol* **6**, 960–970 (2021).
- 134 10. Camarillo-Guerrero, L. F., Almeida, A., Rangel-Pineros, G., Finn, R. D. & Lawley, T. D. Massive  
135 expansion of human gut bacteriophage diversity. *Cell* **184**, 1098–1109 (2021).
- 136 11. Camargo, A. P. *et al.* Identification of mobile genetic elements with geNomad. *Nat Biotechnol* 1–  
137 10 (2023).
- 138 12. LaFleur, T. L., Hossain, A. & Salis, H. M. Automated model-predictive design of synthetic  
139 promoters to control transcriptional profiles in bacteria. *Nat Commun* **13**, 5159 (2022).
- 140 13. Lin, Z. *et al.* Language models of protein sequences at the scale of evolution enable accurate  
141 structure prediction. *BioRxiv* **2022**, 500902 (2022).
- 142 14. Gligorijević, V. *et al.* Structure-based protein function prediction using graph convolutional  
143 networks. *Nat Commun* **12**, 3168 (2021).

144

145

146

## 147 **Methods**

### 148 **Training dataset**

149 Our training dataset was curated from three sources. Firstly, we downloaded all the complete virus  
150 genomes from NCBI genebank, retaining only those with "phage" in the organism's name. Secondly, the  
151 phage genomes from MGV were downloaded, and we only included genomes with a completeness  
152 score larger than 95% and classified under the order Caudovirales. Our third source was GPD, and we  
153 kept all the genomes with a completeness score above 0.95. Following the initial collection, we  
154 undertook an additional round of filtering. We used geNomad to predict the taxonomy of these  
155 genomes and then deleted all the genomes whose predicted host is not a unicellular organism. All  
156 genomes smaller than 96K bp were used to construct the final training dataset.

157

### 158 **Model training and inference**

159 Our megaDNA model utilized a three-layer transformer structure<sup>8</sup>. Each layer had a depth of 8 and  
160 progressively fewer dimensions (512, 256, 196), capturing local-to-global information. The sequence  
161 length for three layers is 128, 64, 16. The model contains 145M parameters in total. We assigned  
162 numerical tokens (1, 2, 3, and 4) to the nucleotides A, T, C, and G, respectively. For model training, we  
163 used a batch size of 1 and set the learning rate at 0.0002. The learning rate was progressively increased  
164 during the initial 50,000 steps as part of a warmup schedule. We utilized the Adam optimizer and  
165 applied gradient clipping with a norm of 0.5 to prevent gradient explosion.

166 We generated sequences from the trained model using a predefined set of parameters. Specifically, we  
167 adjusted the temperature to 0.95 to ensure a balance between variety and coherence in the sequences  
168 and kept the filter threshold at 0.0 to avoid limiting the range of token probabilities. For model training  
169 and inference, we utilized Nvidia's A100 GPU (40GB) and 3090 Ti GPU (24GB) and used the PyTorch  
170 version 2.1.1 software package.

171

### 172 **Analysis of the generated sequence**

173 geNomad<sup>11</sup> was used for sequence annotation of all generated sequences. The 100 base pair regions  
174 preceding the start codon of each predicted gene was designated as the 5'UTR. We employed the



175 Promoter Calculator<sup>12</sup> to find the promoters in these regions. Only the promoter with the highest  
176 predicted activity was annotated. For protein structure prediction, we used the pretrained ESMfold  
177 model v1<sup>13</sup>. The chunk size of the model was set to be 64 for proteins longer than 700 AA and 128 for  
178 shorter proteins. We limited our structure calculation to proteins less than 1000 AA in length. Function  
179 prediction for these proteins was carried out using the default deepFRI model<sup>14</sup>, as available on GitHub  
180 (<https://github.com/flatironinstitute/DeepFRI>). We used a score cutoff of 0.5 which was reported to be  
181 significant in the original publication.

182

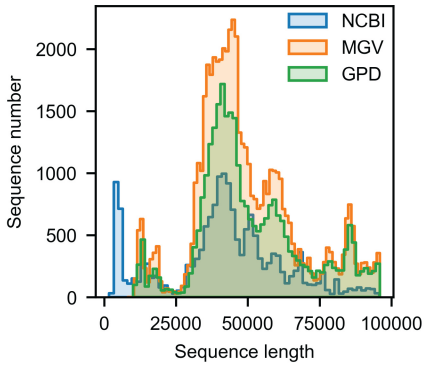
### 183 **Data availability**

184 The bacteriophage genomes were downloaded from public databases including NCBI genebank  
185 (<ftp://ftp.ncbi.nlm.nih.gov/genomes/genbank/>), MGV (<https://portal.nersc.gov/MGV>), and GPD  
186 (<https://www.sanger.ac.uk/data/gut-phage-database/>).

187

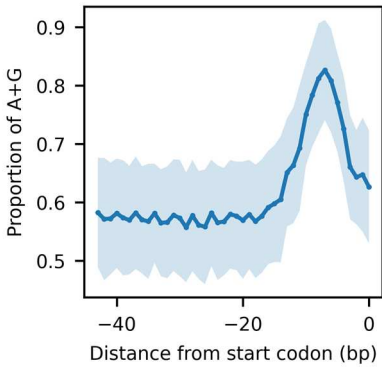
### 188 **Code availability**

189 Our trained model and codes for model inference are available from GitHub:  
190 <https://github.com/lingxusb/megaDNA>



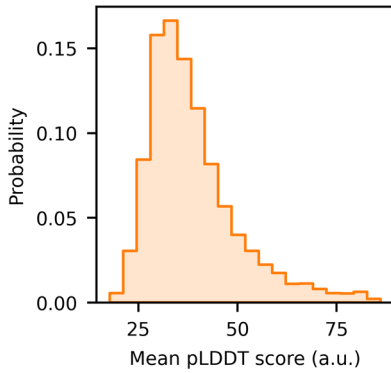
191

192 **Supplementary Figure 1: Distribution of genome sizes of three data sources.** Distributions of genome  
193 sizes within the training dataset: NCBI (n = 16,609), MGV (n = 53,032) and GPD (n = 30,032).



194

195 **Supplementary Figure 2: Proportions of adenine (A) and guanine (G) nucleotides preceding the start**  
196 **codon for all the generated sequences.** Blue line denotes the mean A+G nucleotides proportion profile  
197 for all the generated sequences (n = 607). The shaded region represents the standard derivation of all  
198 profiles.



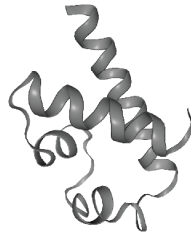
199

200 **Supplementary Figure 3: Mean pLDDT score for proteins derived from the generated sequences.** The  
201 distribution from a randomly sampled subset of the generated proteins is shown (sample size: n =  
202 10,000; median value: 36).

Iron-sulfur cluster binding  
GO:0051536  
score: 0.92683



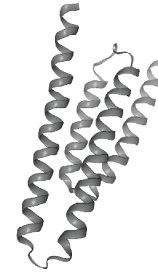
DNA binding  
GO:0003677  
score: 0.83517



Structural molecule activity  
GO:0005198  
score: 0.95713



Ion transmembrane transporter activity  
GO:0015075  
score: 0.70163



203

204 **Supplementary Figure 4: Representative proteins from the generated sequence with predicted**  
205 **functions and structures.** The protein structures were predicted using ESMfold<sup>13</sup> and the functions were  
206 annotated using deepFRI<sup>14</sup>. Predicted scores and GO terms from deepFRI are shown.