

LLM модели в биологии

DNABERT, DNABERT-2, GENA-LM, ESM-1v

DNABERT

Bioinformatics, 37(15), 2021, 2112–2120

doi: 10.1093/bioinformatics/btab083

Advance Access Publication Date: 4 February 2021

Original Paper

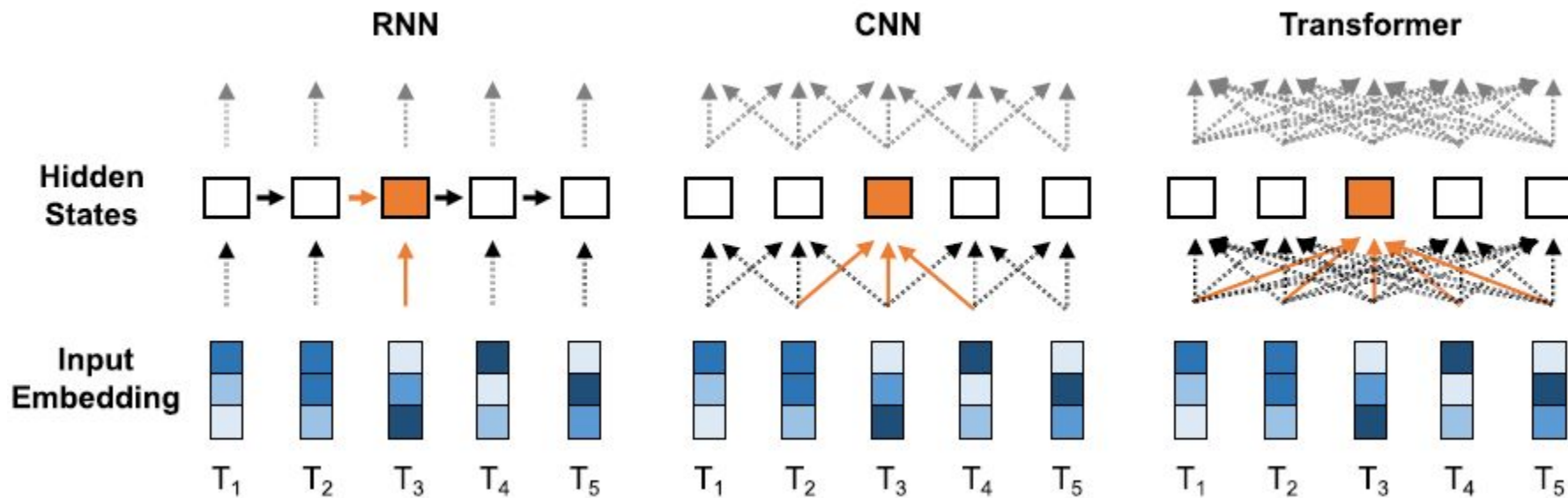
OXFORD

Genome analysis

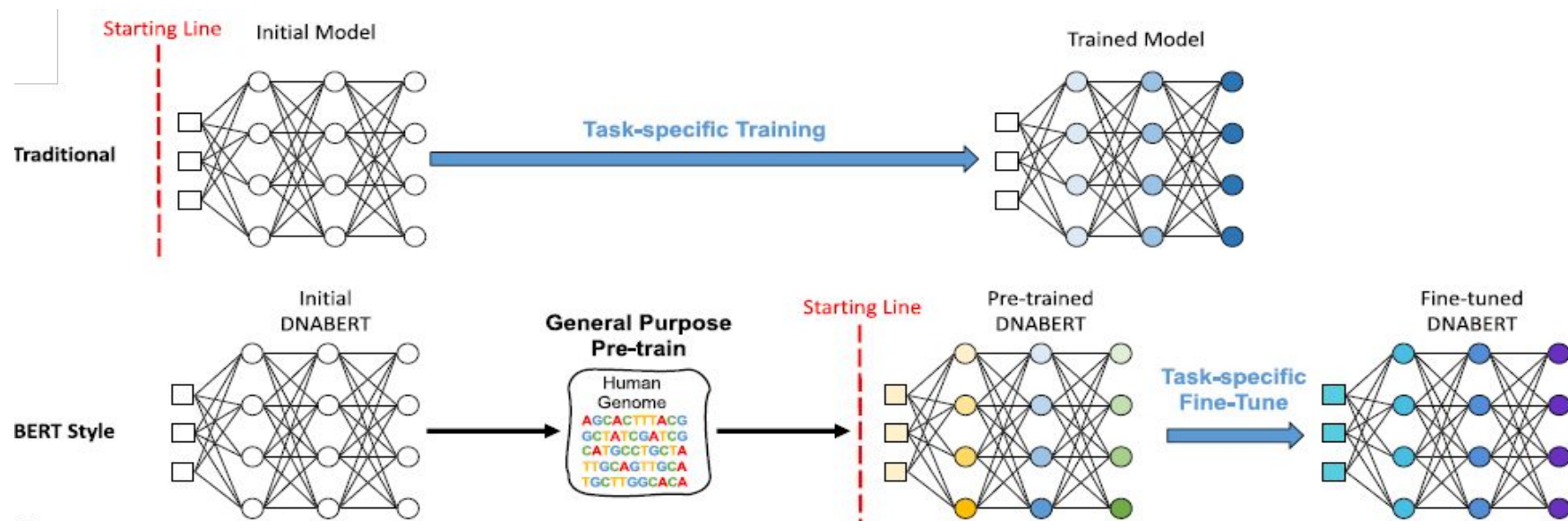
DNABERT: pre-trained Bidirectional Encoder Representations from Transformers model for DNA-language in genome

Yanrong Ji^{1,†}, Zhihan Zhou^{2,†}, Han Liu^{2,*} and Ramana V. Davuluri ^{3,*}

DNABERT



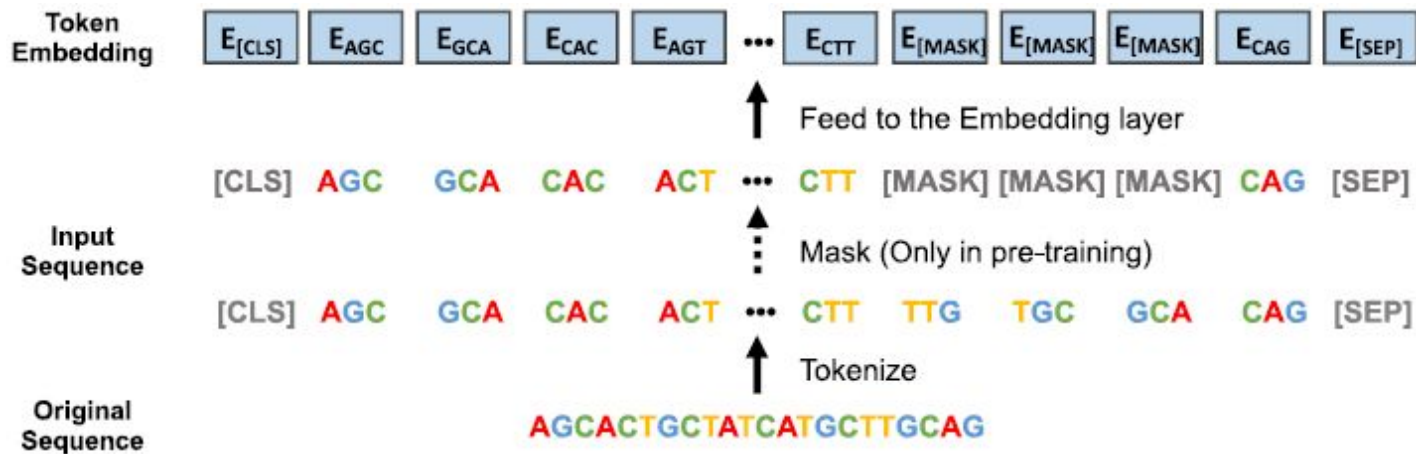
DNABERT



DNABERT

Модель обучалась на геноме человека с длиной последовательностей от 5 до 510.

Использовалась k-mer токенизация:



DNABERT

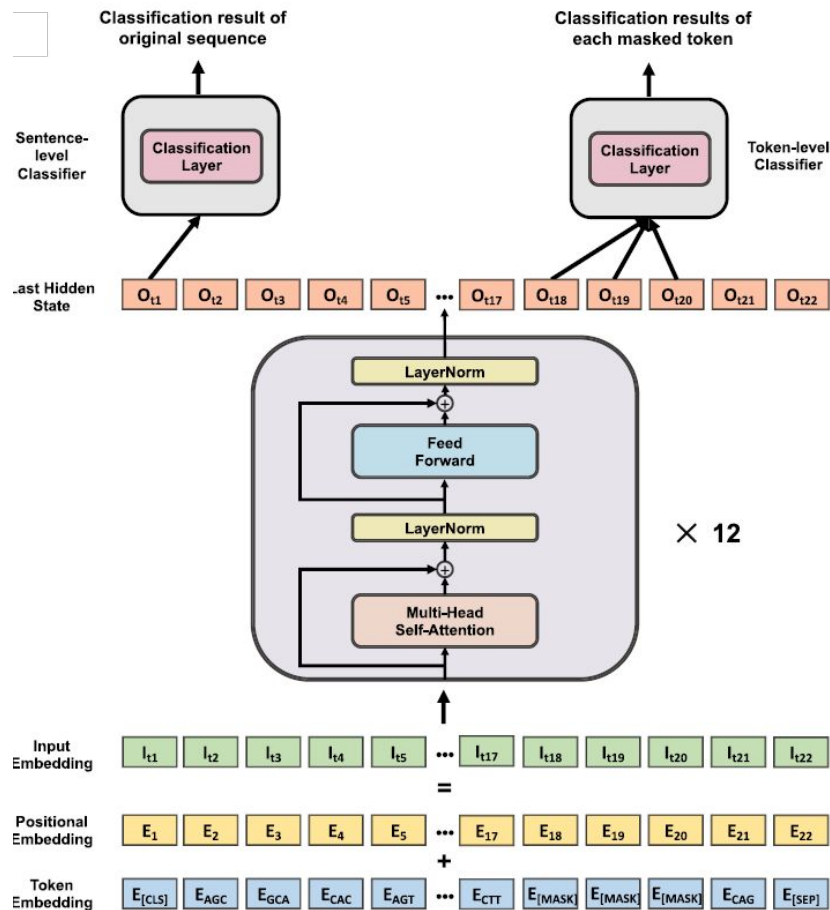
DNABERT captures contextual information by performing the multi-head self-attention mechanism on M :

$$MultiHead(M) = Concat(head_1, \dots, head_b) W^O$$

where

$$head_i = softmax\left(\frac{M W_i^O M W_i^K T}{\sqrt{d_k}}\right) \cdot M W_i^V$$

W^O and $W_i^O, W_i^K, W_i^V \{W_i^O, W_i^K, W_i^V\}_{i=0}^b$ are learned parameters



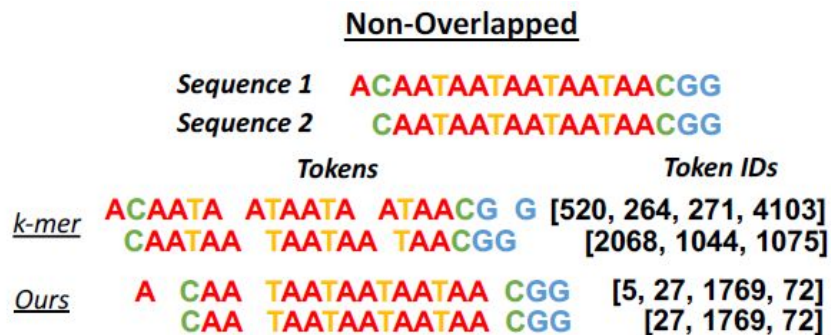
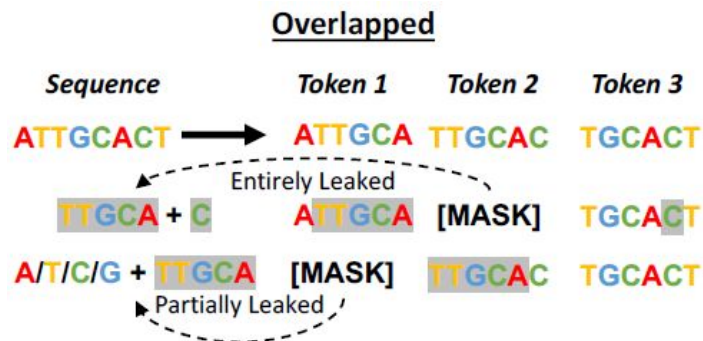
DNABERT-2

Модель обучалась на геноме 4-х видов с длиной последовательностей от 70 до 10000.

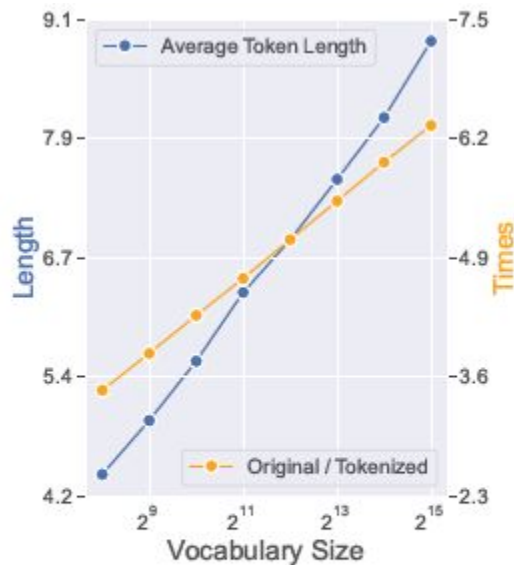
Для токенизации последовательностей ДНК использовался алгоритм Byte Pair Encoding (BPE):

<i>Iteration</i>	<i>Corpus</i>	<i>Vocabulary</i>
0	AACGCACTATATA	{A,T,C,G}
1	A AC G C AC TA TA TA	{A,T,C,G,TA}
2	A AC G C AC TA TA TA	{A,T,C,G,TA,AC}
3	A AC G C AC TA TA TA

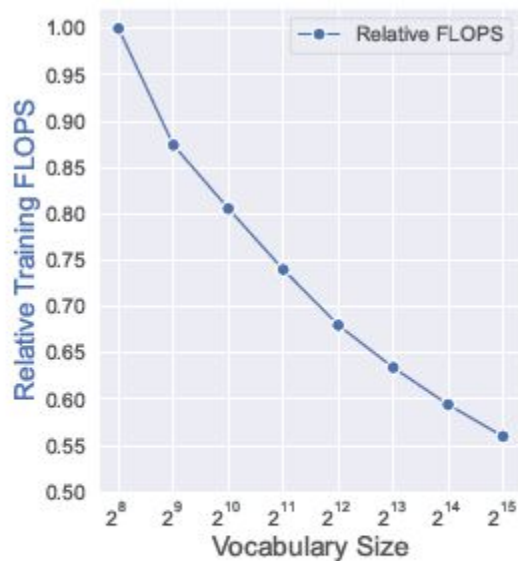
DNABERT-2



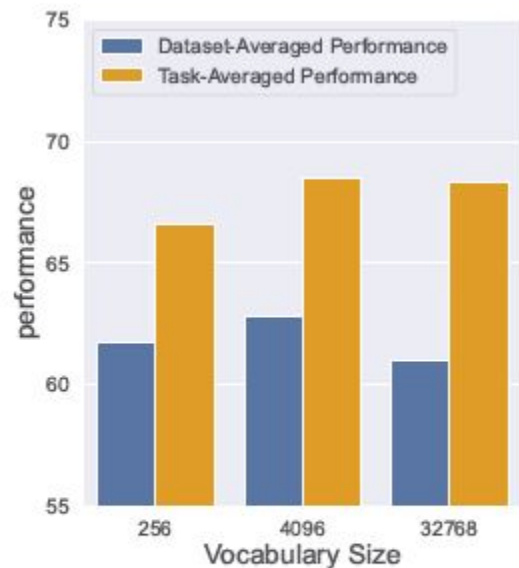
DNABERT-2



(a) Average token length and the length ratio of original sequence v.s. tokenized sequence.



(b) Training FLOPs on 500-length sequences compared to model with 2^8 vocabulary.



(c) Model performance averaged over each tasks (macro) and individual dataset (micro).

DNABERT-2, низкоранговая адаптация (LoRA)

Пусть $W_0, W_1 \in \mathbb{R}^{m \times n}$ определяют одну и ту же весовую матрицу до и после тонкой настройки для конкретной задачи

Представим W_1 как $W_1 = W_0 + \Delta W$ где $\Delta W \in \mathbb{R}^{m \times n}$

А ΔW как перемножение 2-х низкоранговых матриц:

$\Delta W = BA$, где $B \in \mathbb{R}^{m \times r}$, $A \in \mathbb{R}^{r \times n}$ и $r \ll m, r \ll n$

Низкоранговая декомпозиция уменьшает количество обучаемых параметров с $m \times n$ до $r \times (m + n)$

DNABERT-2 vs DNABERT/Nucleotide Transformer

Статистика и производительность моделей:

Model	Params. ↓	FLOPs ↓	Trn. Tokens	Num. Top-2 ↑	Ave. Scores ↑
DNABERT (3-mer)	86M	3.27	122B	2 0	61.62
DNABERT (4-mer)	86M	3.26	122B	0 1	61.14
DNABERT (5-mer)	87M	3.26	122B	0 1	60.05
DNABERT (6-mer)	89M	3.25	122B	0 1	60.51
NT-500M-human	480M	3.19	50B	0 0	55.43
NT-500M-1000g	480M	3.19	50B	0 1	58.23
NT-2500M-1000g	2537M	19.44	300B	0 1	61.41
NT-2500M-multi	2537M	19.44	300B	7 9	<u>66.93</u>
DNABERT-2	117M	1.00	262B	8 4	66.80
DNABERT-2♦	117M	1.00	263B	11 10	67.77

В пяти столбцах представлено количество параметров модели, относительные FLOP по сравнению с DNABERT-2, количество токенов, использованных при предварительном обучении, количество попаданий в топ-2 среди всех моделей (1-й // 2-й), а также средние оценки на 28 наборах данных теста GUE.

♦ предварительное обучение выполнено на обучающих наборах эталонного теста GUE.

NT-500M-human, NT-500M-1000g, NT-2500M-1000g и NT-2500M-multi, где human, 1000g и multi соответственно относятся к эталонному геному человека GRCh38/hg38, набор из 3202 геномов человека с высоким охватом из проекта «1000 геномов» (Byrka-Bishop et al., 2021) и геномы 850 различных видов.

GENA-LM

GENA-LM: A Family of Open-Source Foundational DNA Language Models for Long Sequences

Veniamin Fishman^{1,2*†}, Yuri Kuratov^{1,3†}, Maxim Petrov¹,
Aleksei Shmelev^{1,4}, Denis Shepelin¹, Nikolay Chekanov¹,
Olga Kardymon^{1,4*}, Mikhail Burtsev^{5*}

¹AIRI, Moscow, Russia.

²Institute of Cytology and Genetics, Novosibirsk, Russia.

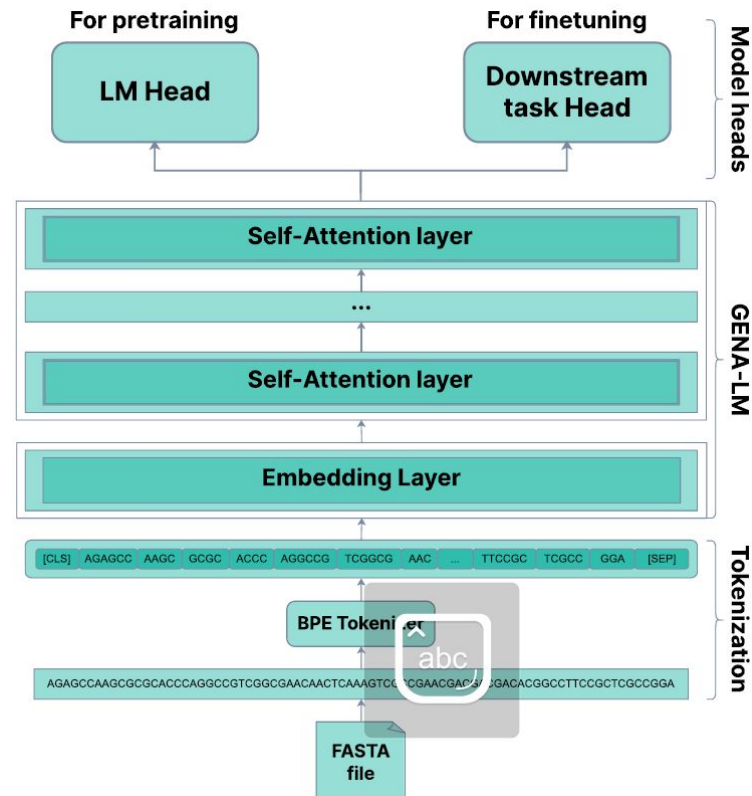
³Moscow Institute of Physics and Technology, Dolgoprudny, Russia.

⁴HSE University, Moscow, Russia.

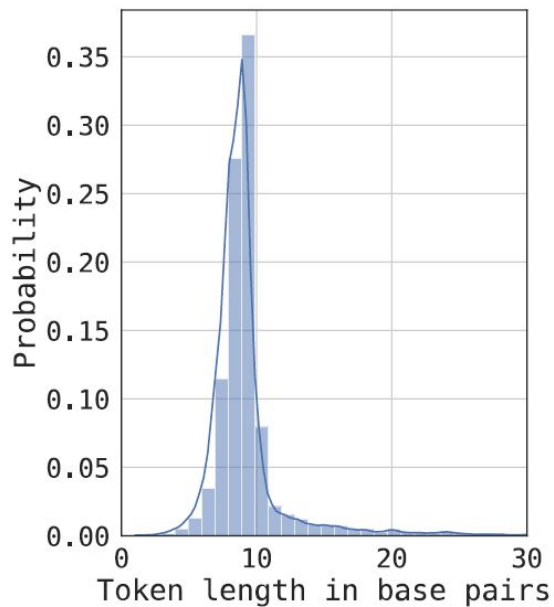
⁵London Institute for Mathematical Sciences, London, UK.

GENA-LM

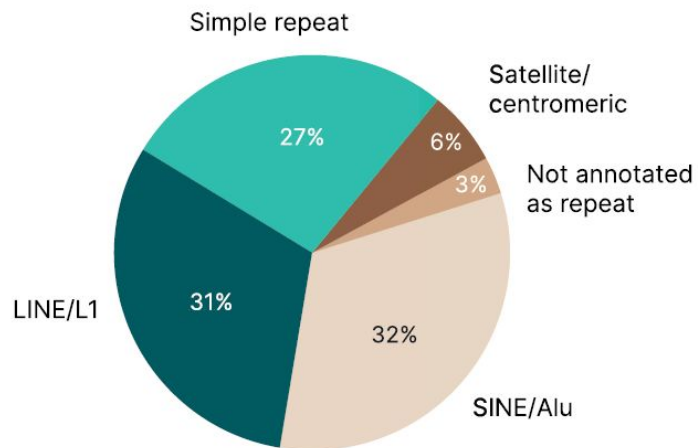
“В этой работе мы демонстрируем успешное применение передовых нейронных сетей на основе трансформеров для прогнозного анализа различных функциональных геномных элементов в последовательностях ДНК, включая активность промотора, сплайсинг, сайты полиаденилирования, аннотации энхансеров и профили хроматина. Мы вносим свой вклад в исследовательское сообщество, представляя GENA-LM, семейство моделей с открытым исходным кодом, доступных на GitHub, и предварительно обученные модели (с префиксом gena-lm-) на <https://huggingface.co/AIRI-Institute>. Мы эмпирически показываем, что точная настройка наших моделей превосходит результаты, полученные с помощью современных архитектур. Более того, мы доказываем, что дополнение GENA-LM архитектурой трансформера рекуррентной памяти (RMT) способствует увеличению длины входной последовательности, увеличивая качество решения сложных биологических задач.”



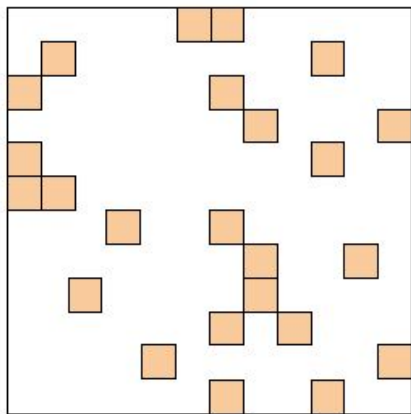
GENA-LM



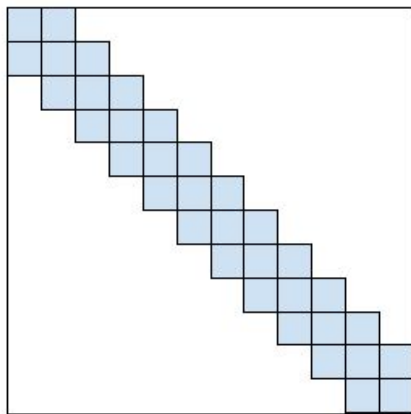
Bases annotated as repeats among top 100 longest tokens



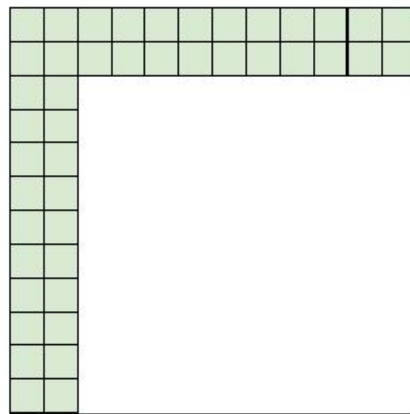
GENA-LM



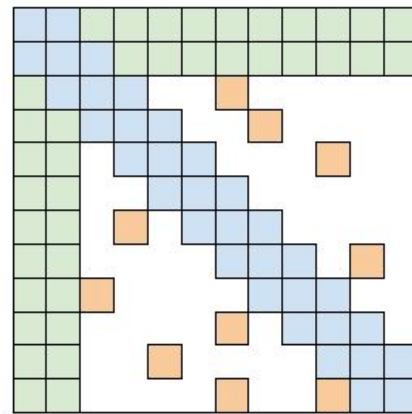
(a) Random attention



(b) Window attention



(c) Global Attention



(d) BIGBIRD

Figure 1: Building blocks of the attention mechanism used in BIGBIRD. White color indicates absence of attention. (a) random attention with $r = 2$, (b) sliding window attention with $w = 3$ (c) global attention with $g = 2$. (d) the combined BIGBIRD model.

GENA-LM

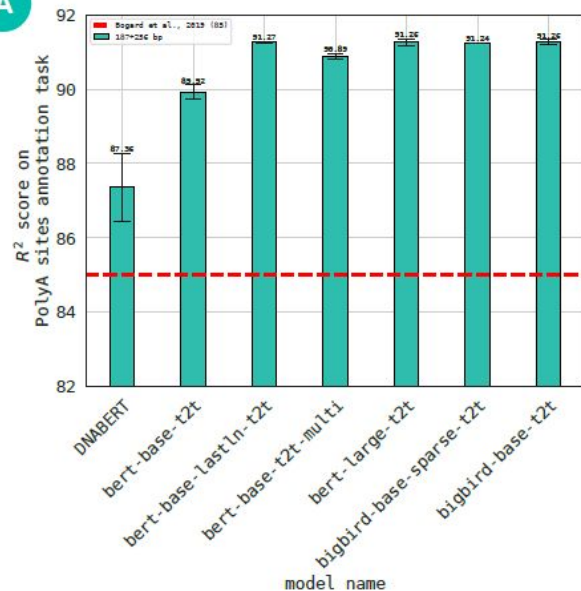
Model	Architecture	Maximum seq len, tokens (\approx bp)	Tokenizer data	Training data
DNABERT	BERT-12L	512 (512)	3,4,5,6-mer	GRCh38.p13
GENA-LM models:				
<i>bert-base</i>	BERT-12L	512 (4,500)	T2T split v1	T2T split v1
<i>bert-base-t2t</i>	BERT-12L	512 (4,500)	T2T+1KG+M	T2T+1KG
<i>bert-base-lastln-t2t</i>	BERT-12L	512 (4,500)	T2T+1KG+M	T2T+1KG
<i>bert-base-t2t-multi</i>	BERT-12L	512 (4,500)	T2T+1KG+M	T2T+1KG+M
<i>bert-large-t2t</i>	BERT-24L	512 (4,500)	T2T+1KG+M	T2T+1KG
<i>bigbird-base-sparse</i>	BERT-12L, RoPE DS Sparse Att	4,096 (36,000)	T2T split v1	T2T split v1
<i>bigbird-base-sparse-t2t</i>	BERT-12L, RoPE DS Sparse Att	4,096 (36,000)	T2T+1KG+M	T2T+1KG
<i>bigbird-base-t2t</i>	BERT-12L HF Sparse Attention	4,096 (36,000)	T2T+1KG+M	T2T+1KG

Характеристики моделей GENA-LM. Выделены различия в данных предварительного обучения, количестве слоев, типе attention и длине последовательности. «T2T-сплит v1» - словарь токенов построен на разбиении сборки генома человека T2T, «1KG» является сокращением 1000G, «М» - включение данных нескольких видов.

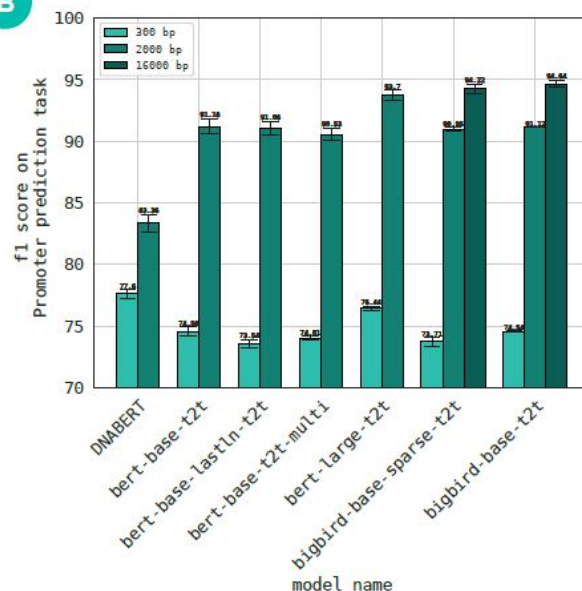
Обозначения «DS Sparse» и «HF Sparse» относятся к реализациям DeepSpeed разреженного внимания и HuggingFace BigBird соответственно. Аббревиатура «RoPE» означает использование rotary position embeddings в качестве альтернативы абсолютного PE в BERT. Модели были структурированы либо из 12 (обозначаемых BERT-12L), либо из 24 (обозначаемых BERT-24L) слоев, содержащих 110M и 336M параметров соответственно.

GENA-LM

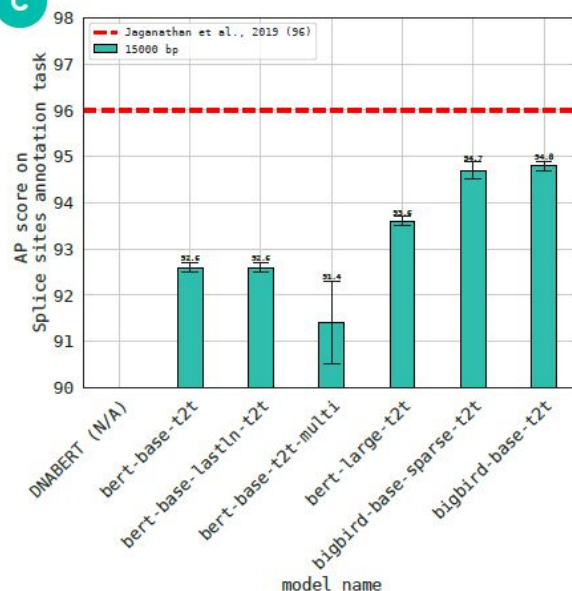
A



B

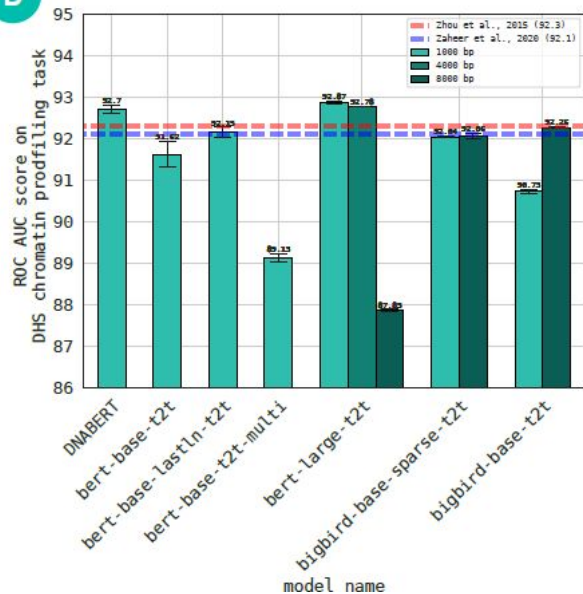


C

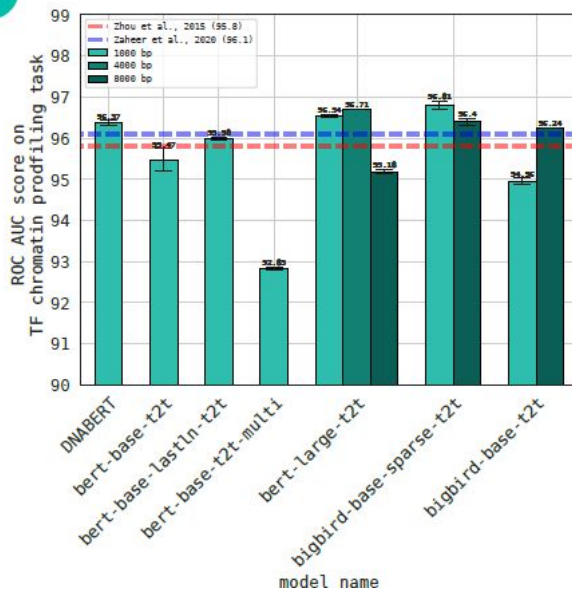


GENA-LM

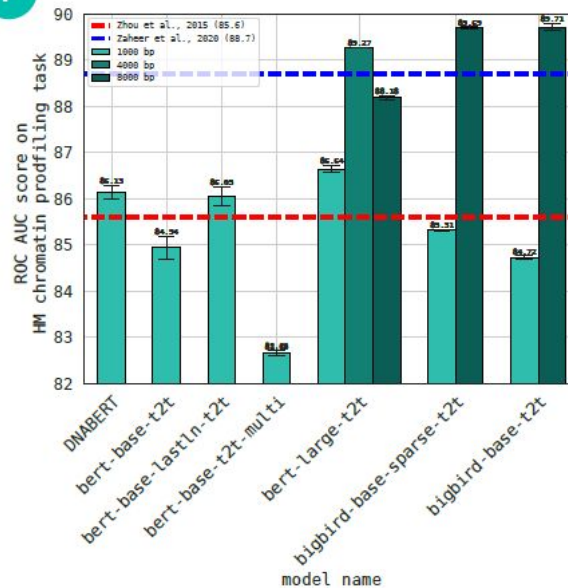
D



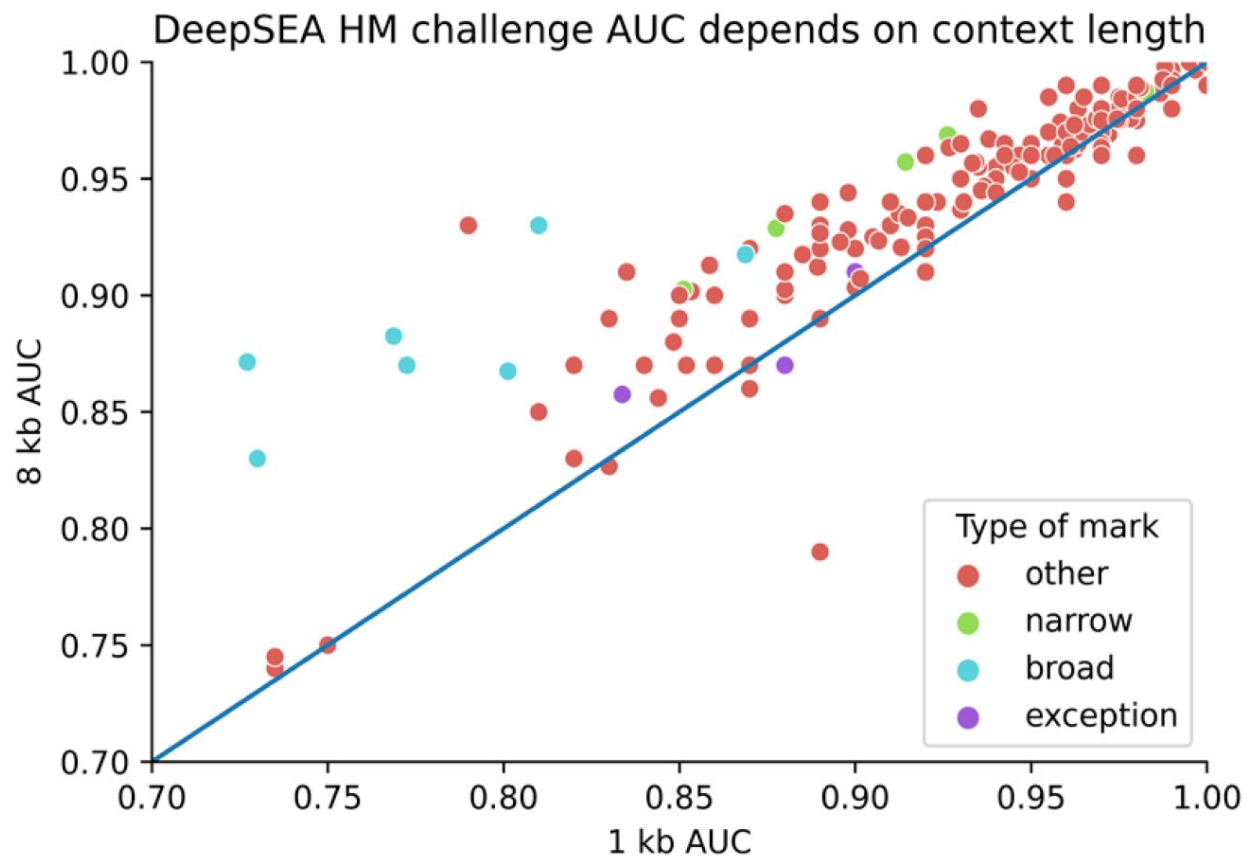
E



F

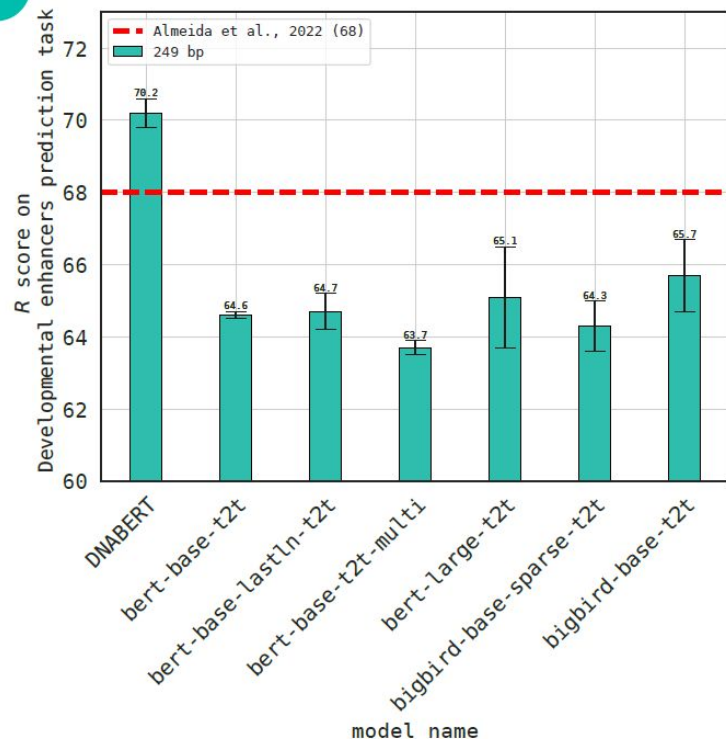


GENA-LM

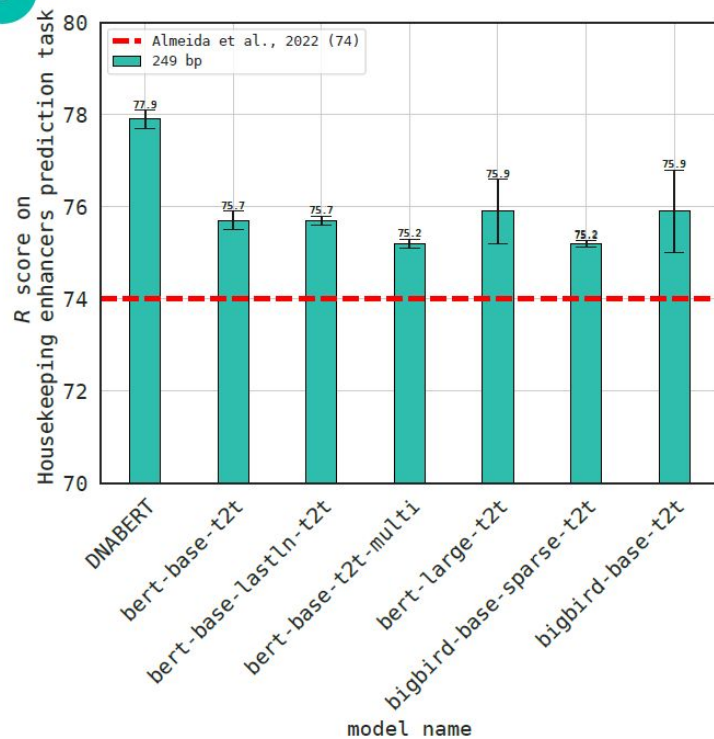


GENA-LM

H

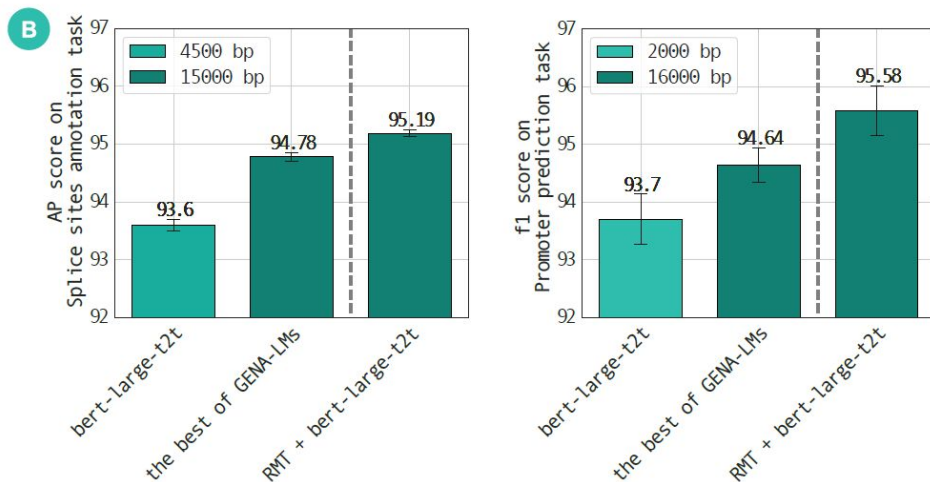
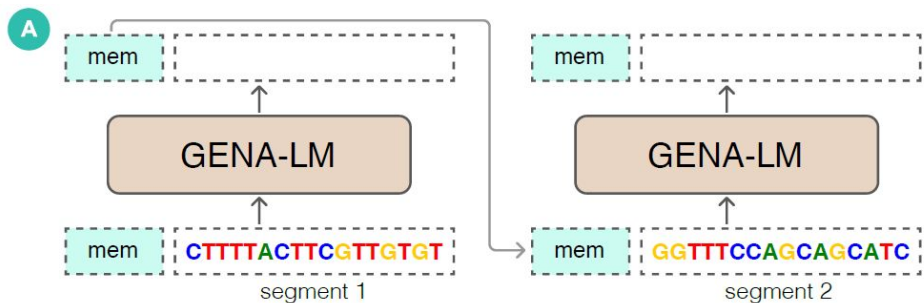


I



GENA-LM

Токены памяти добавляются в каждый сегмент для передачи информации между последовательными сегментами, что позволяет им использовать информацию из всех предыдущих сегментов. Таким образом, весь предварительно обученный Трансформер эффективно функционирует как единая рекуррентная единица.

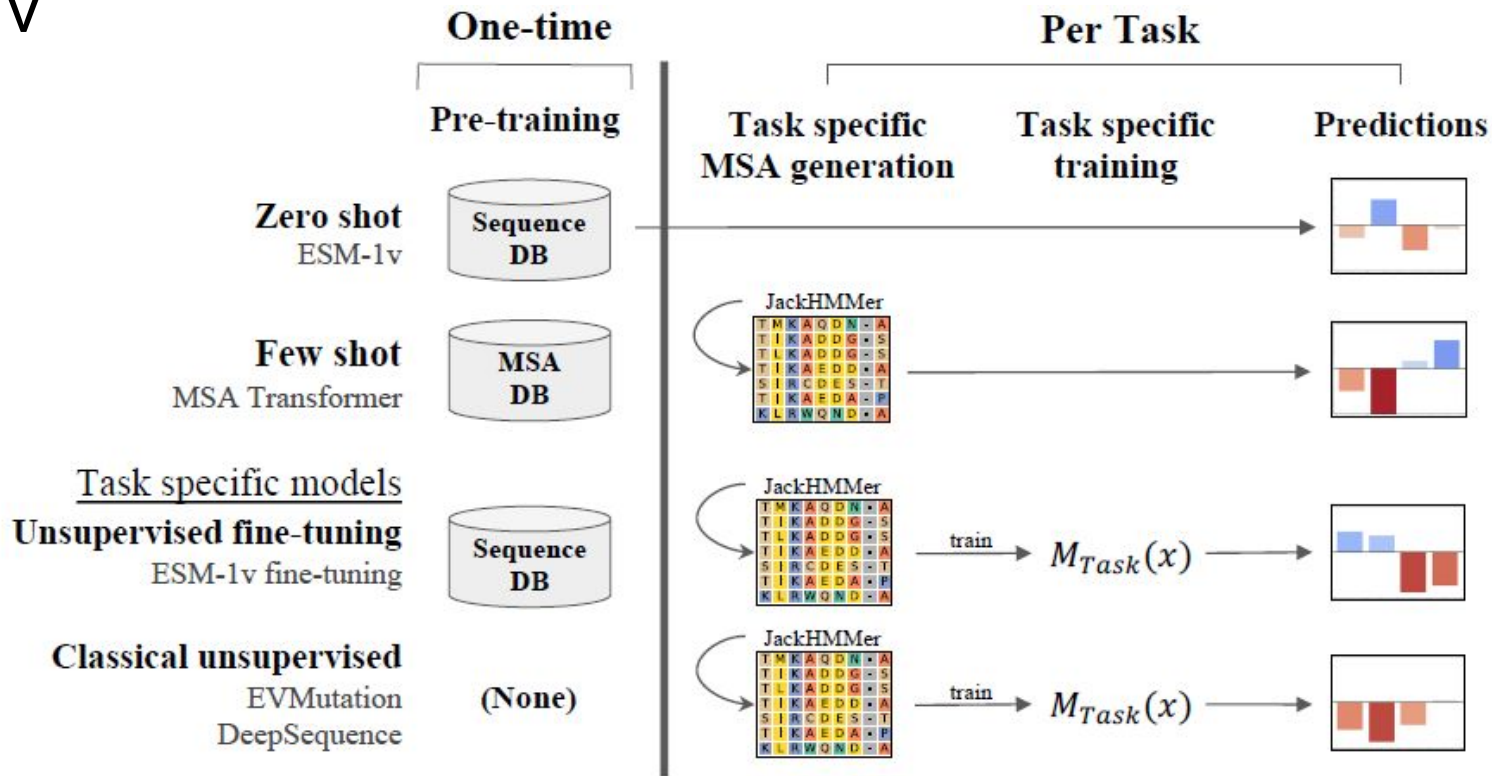


ESM-1v

Language models enable zero-shot prediction of the effects of mutations on protein function

Joshua Meier^{1 2} Roshan Rao³ Robert Verkuil¹ Jason Liu¹
Tom Sercu¹ Alexander Rives^{1 2}

ESM-1v



ESM-1v

Мы оцениваем мутации, используя логарифмическое отношение вероятностей в мутированной позиции, предполагая аддитивную модель, когда несколько мутаций T существуют в одной и той же последовательности:

$$\sum_{t \in T} \log p(x_t = x_t^{mt} | x_{\setminus T}) - \log p(x_t = x_t^{wt} | x_{\setminus T})$$