

# DNA language models in biomedical research

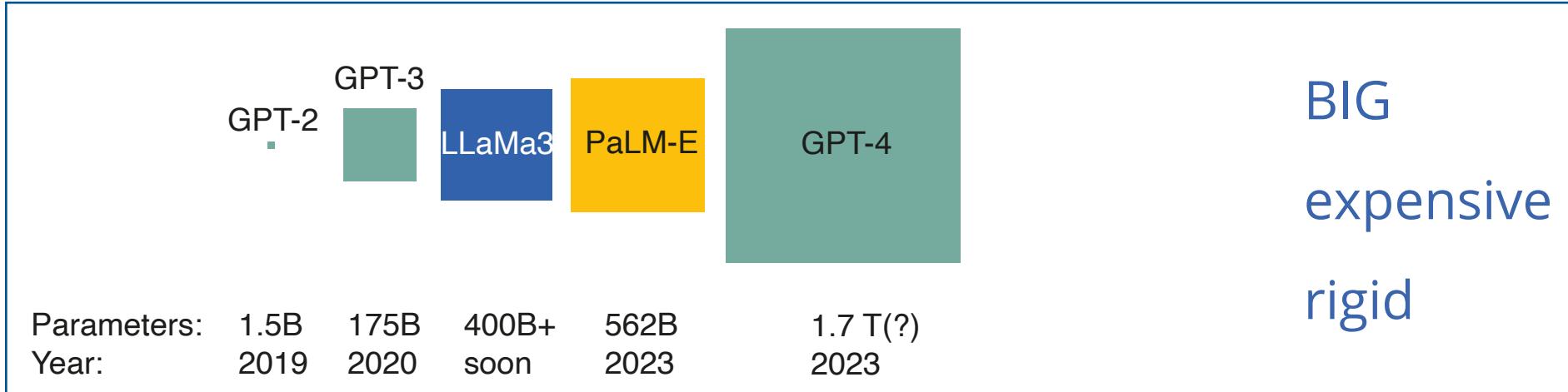
Anna Poetsch

Research Group „Biomedical Genomics“, Biotechnology Center TU Dresden, NCT Dresden, and CSBD

- What are Large Language Models?
- A short glimpse into a transformer foundation model
- DNA language models
- The inner workings of GROVER
- applications in biomedical research

# Pre-training and fine-tuning

Foundation models train language on large corpora of data



Fine-tuning

- assistant
- image generation
- translation
- speech recognition
- ...

"Dalle2, please give me an illustration of damaged DNA in comic style":



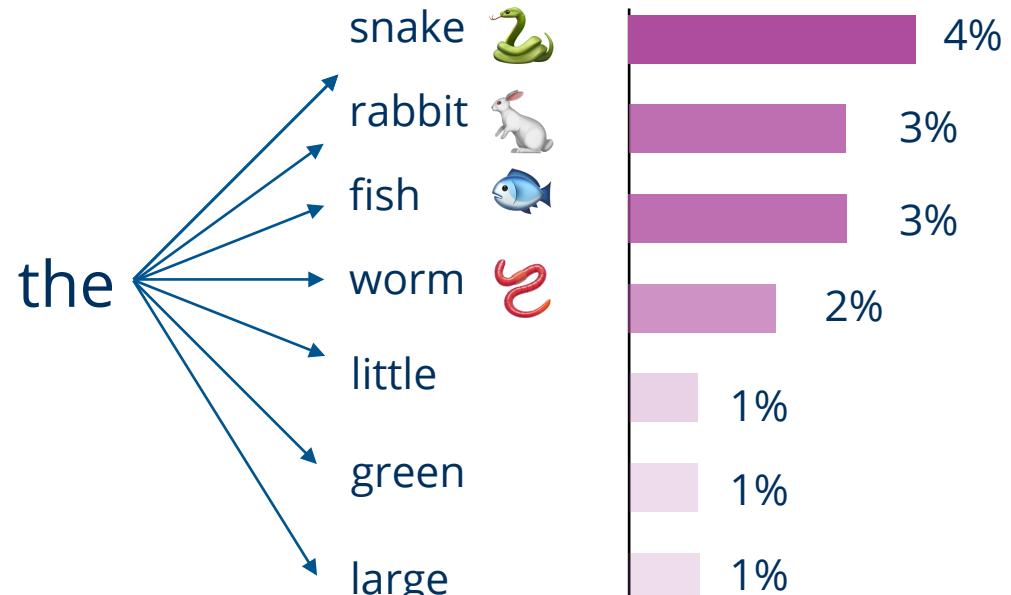
efficient  
flexible  
cheap

# Large Language Models are trained on next-token prediction

Predicting the next word...

the eagle

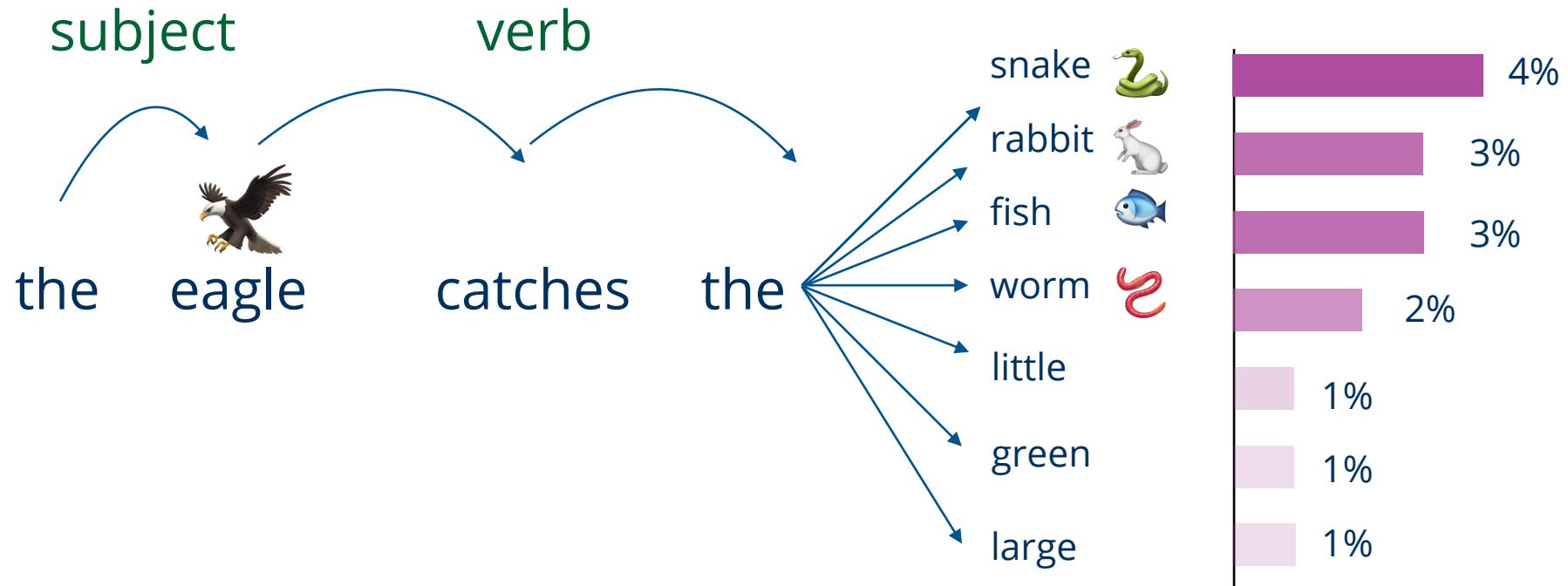
catches



# Large Language Models

Predicting the next word...  
...requires context

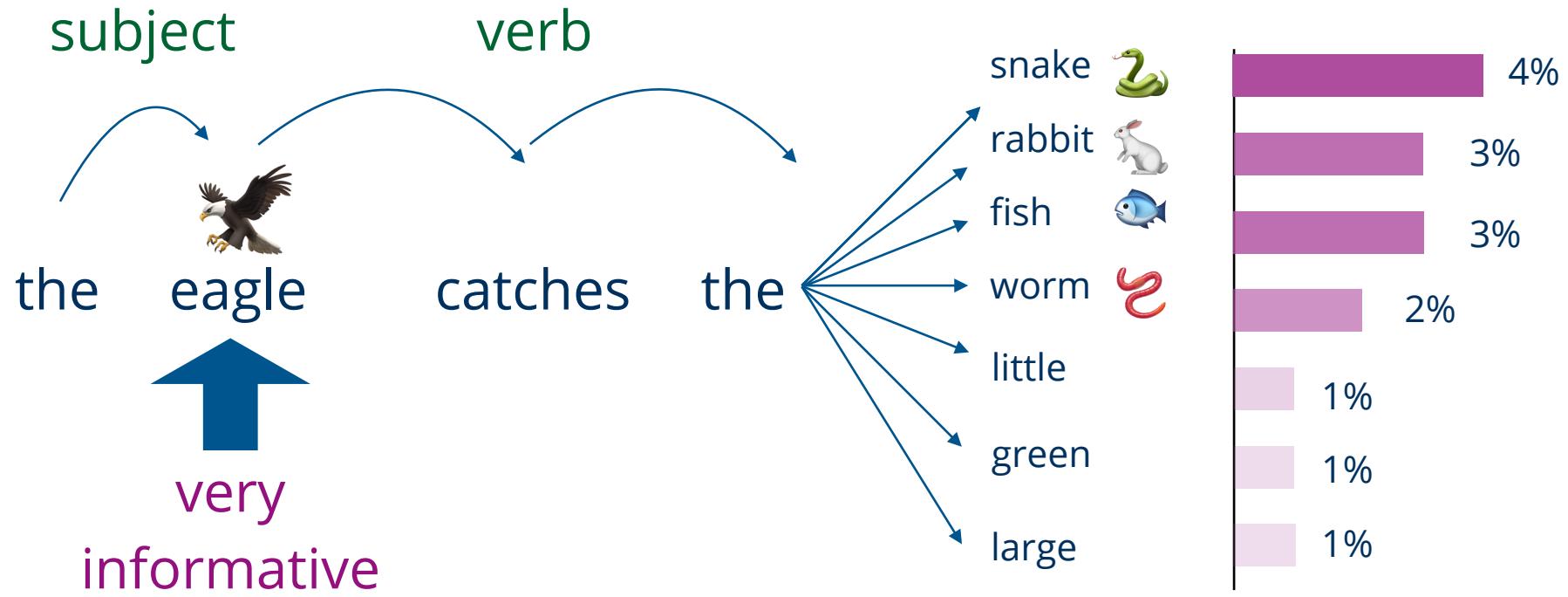
object  
noun?  
adjective?



# Large Language Models

Predicting the next word...  
...requires context

object  
noun?  
adjective?



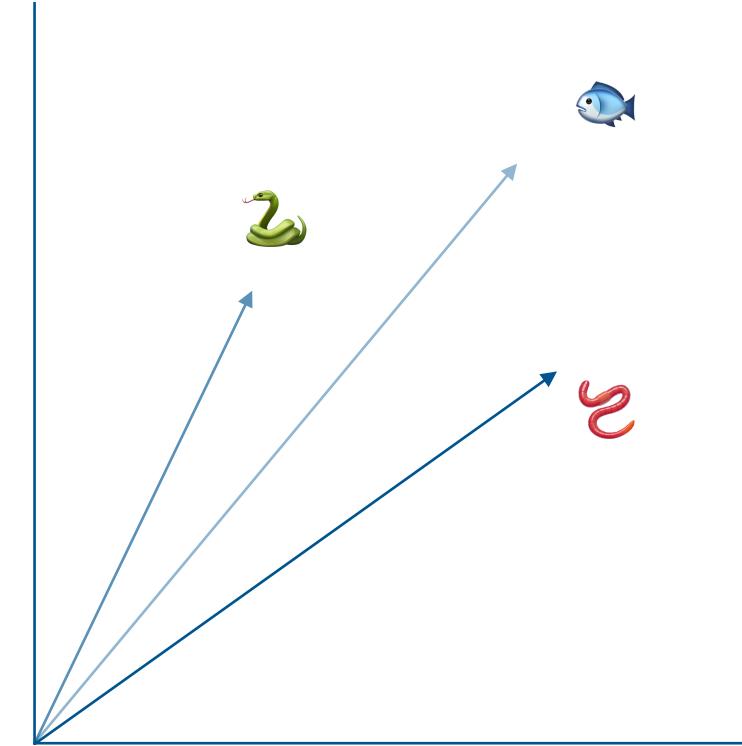
# The embedding of words/ tokens

Tokens are assigned a vector, which places them into a multi-dimensional space

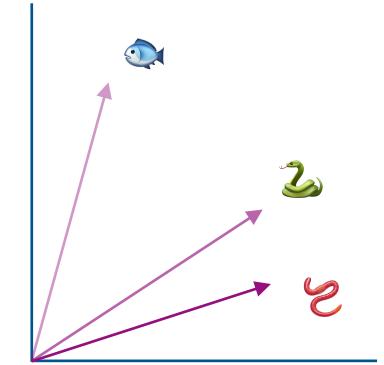
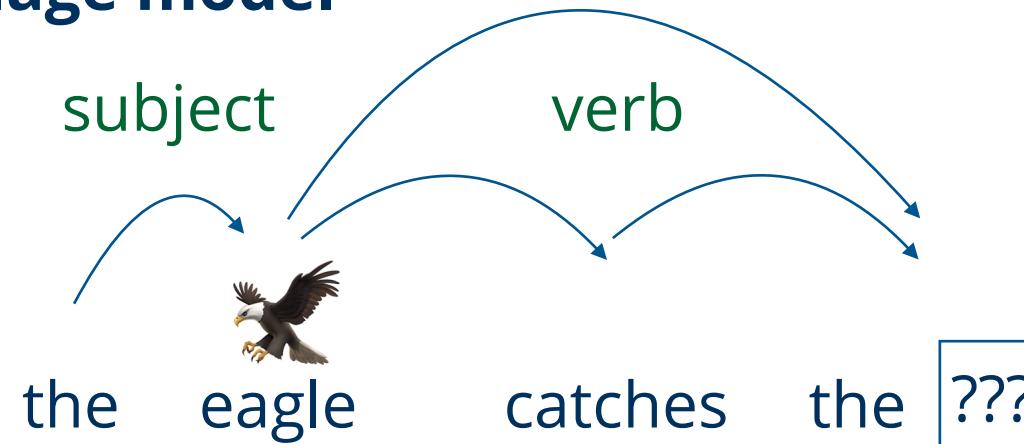
the eagle catches the

???

5.4	3.2	0.2	5.4
2.3	1.4	0.3	2.3
1.2	0.4	5.6	1.2
.	.	.	.
.	.	.	.
.	.	.	.
6.4	1.2	3.2	6.4
0.4	8.6	4.1	0.4



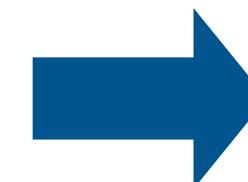
# Training a large language model



Input embedding

5.4	3.2	0.2	5.4
2.3	1.4	0.3	2.3
1.2	0.4	5.6	1.2
.	.	.	.
.	.	.	.
6.4	1.2	3.2	6.4
0.4	8.6	4.1	0.4

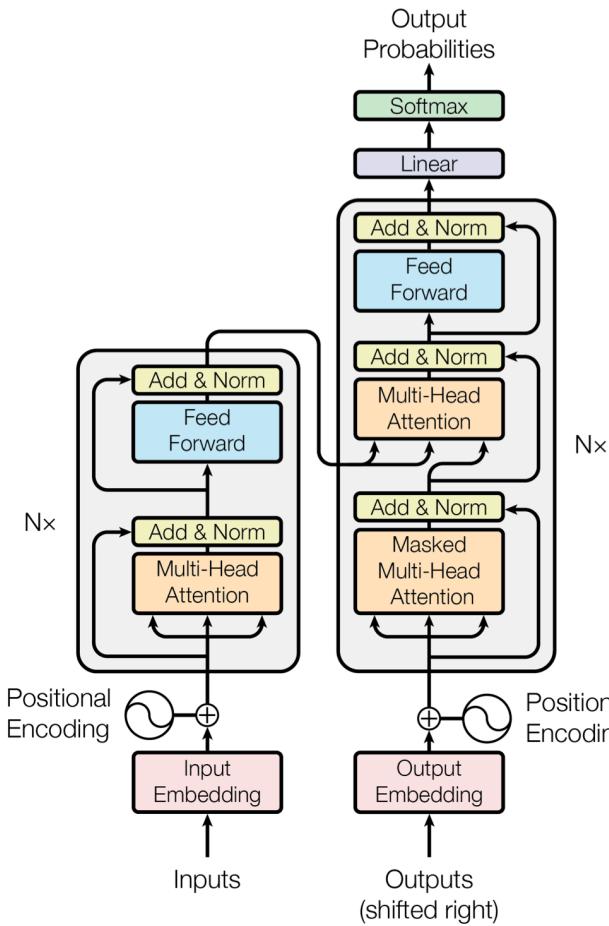
context



Trained embedding

3.2	6.3	0.6	3.7
4.6	0.4	1.7	5.3
5.2	8.2	4.8	7.2
.	.	.	.
.	.	.	.
5.9	3.2	3.9	2.1
0.3	0.4	6.2	0.9

# Context is grasped through transformers: “Attention is all you need”

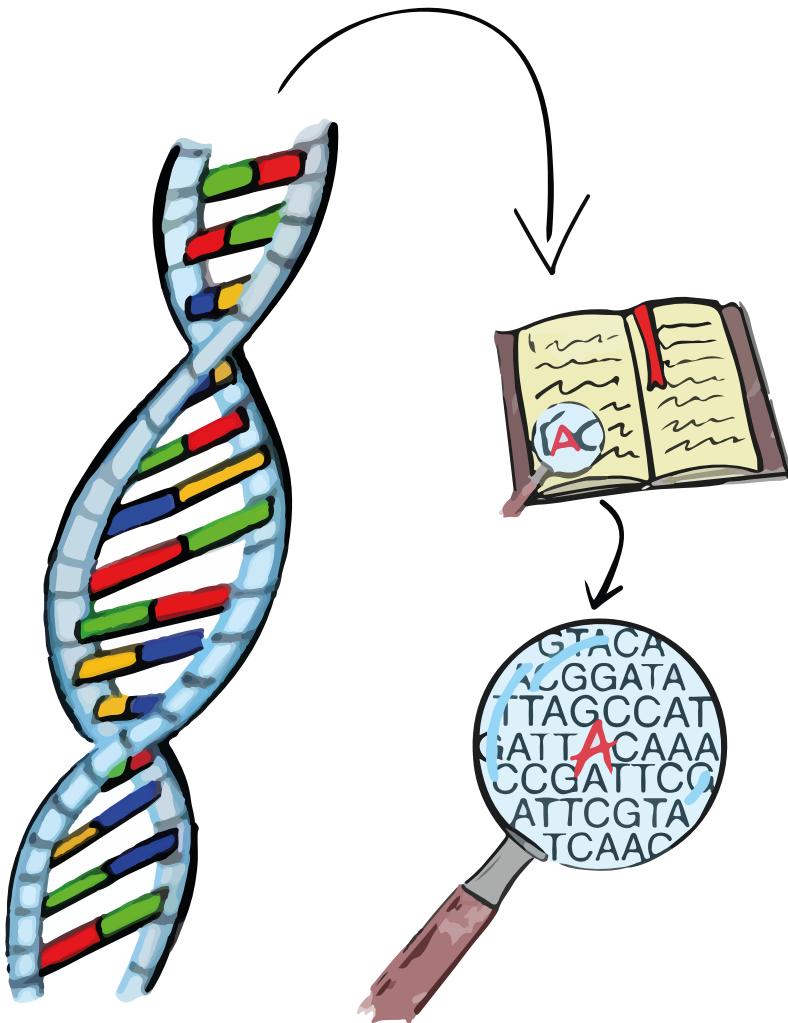


		GPT-3				Total parameters:				
		12,288	50,257			175,181,291,520				
attention	Embedding	$d_{\text{embed}}$	$\ast$	$n_{\text{vocab}}$		=	617,558,016			
	Key	128	$\ast$	12,288	96	96	$14,495,514,624$			
	Query	$d_{\text{query}}$	$\ast$	$d_{\text{embed}}$	$\ast$	$n_{\text{heads}}$	$\ast$	$n_{\text{layers}}$	=	14,495,514,624
	Value	128	$\ast$	12,288	96	96	$=$	14,495,514,624		
	Output	$d_{\text{query}}$	$\ast$	$d_{\text{embed}}$	$\ast$	$n_{\text{heads}}$	$\ast$	$n_{\text{layers}}$	=	14,495,514,624
	Up-projection	12,288	$\ast$	128	96	96	$=$	14,495,514,624		
	Down-projection	$d_{\text{embed}}$	$\ast$	$d_{\text{value}}$	$\ast$	$n_{\text{heads}}$	$\ast$	$n_{\text{layers}}$	=	57,982,058,496
	Unembedding	49,152	$\ast$	12,288	96	$n_{\text{neurons}} \ast d_{\text{embed}}$	$\ast$	$n_{\text{layers}}$	$=$	57,982,058,496
	12,288	$\ast$	49,152	96	$d_{\text{embed}} \ast n_{\text{neurons}}$	$\ast$	$n_{\text{layers}}$	$=$	617,558,016	
	50,257	$\ast$	12,288	$n_{\text{vocab}} \ast d_{\text{embed}}$						

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

<https://arxiv.org/abs/1706.03762>

# Can we treat DNA as if it were text?



**DNABERT-2**



**GROVER**



**HyenaDNA**

**Nucleotide  
Transformer**

**Proformer**

**Enformer**

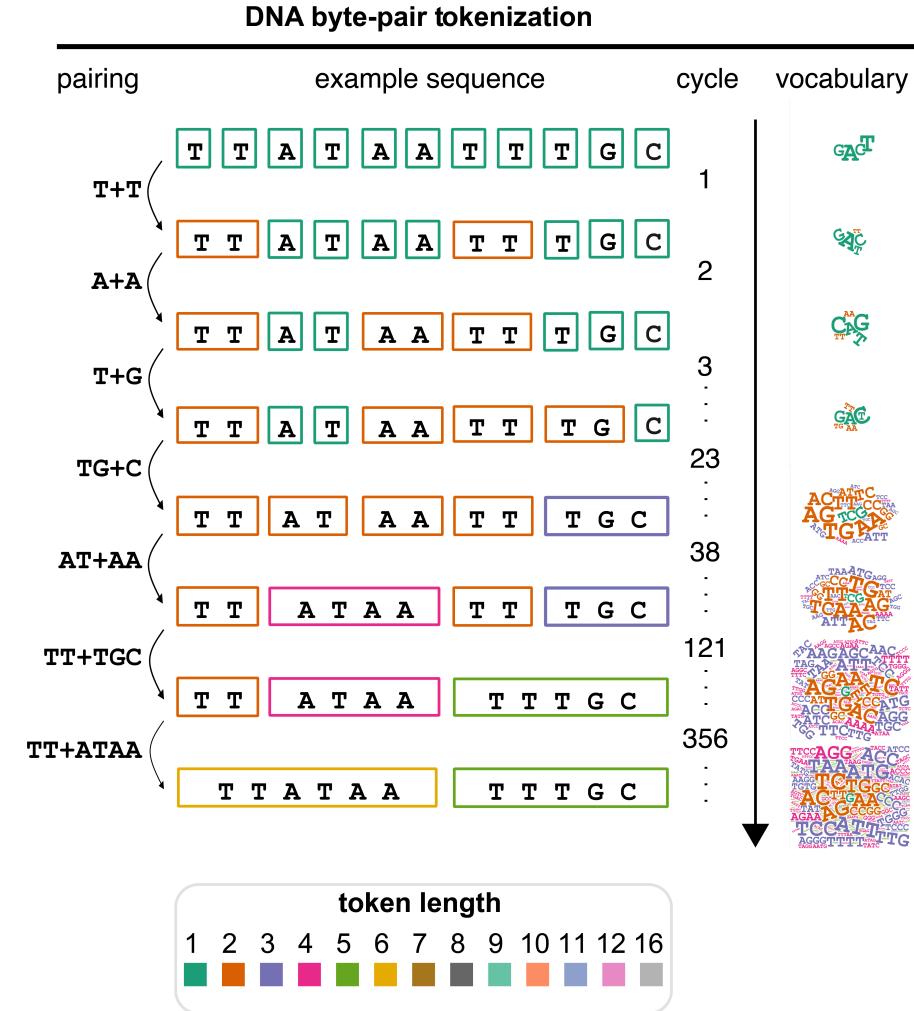
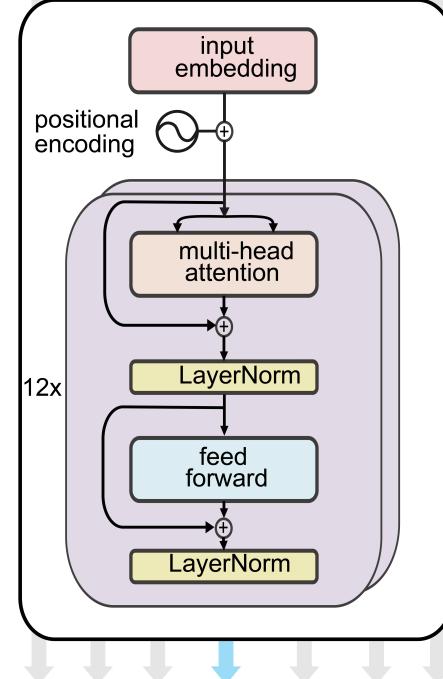
:



**GROVER** (Genome Rules Obtained Via Extracted Representations)

# trains on masked token prediction

TAGA GAC GAGG TTTG ATCAT GTT ACCC  
TAGA GAC GAGG [MASK]ATCAT GTT ACCC

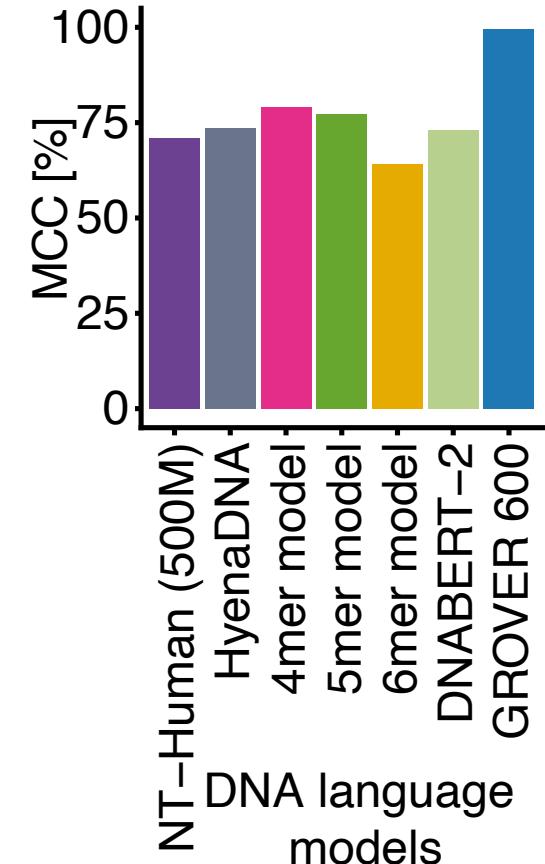
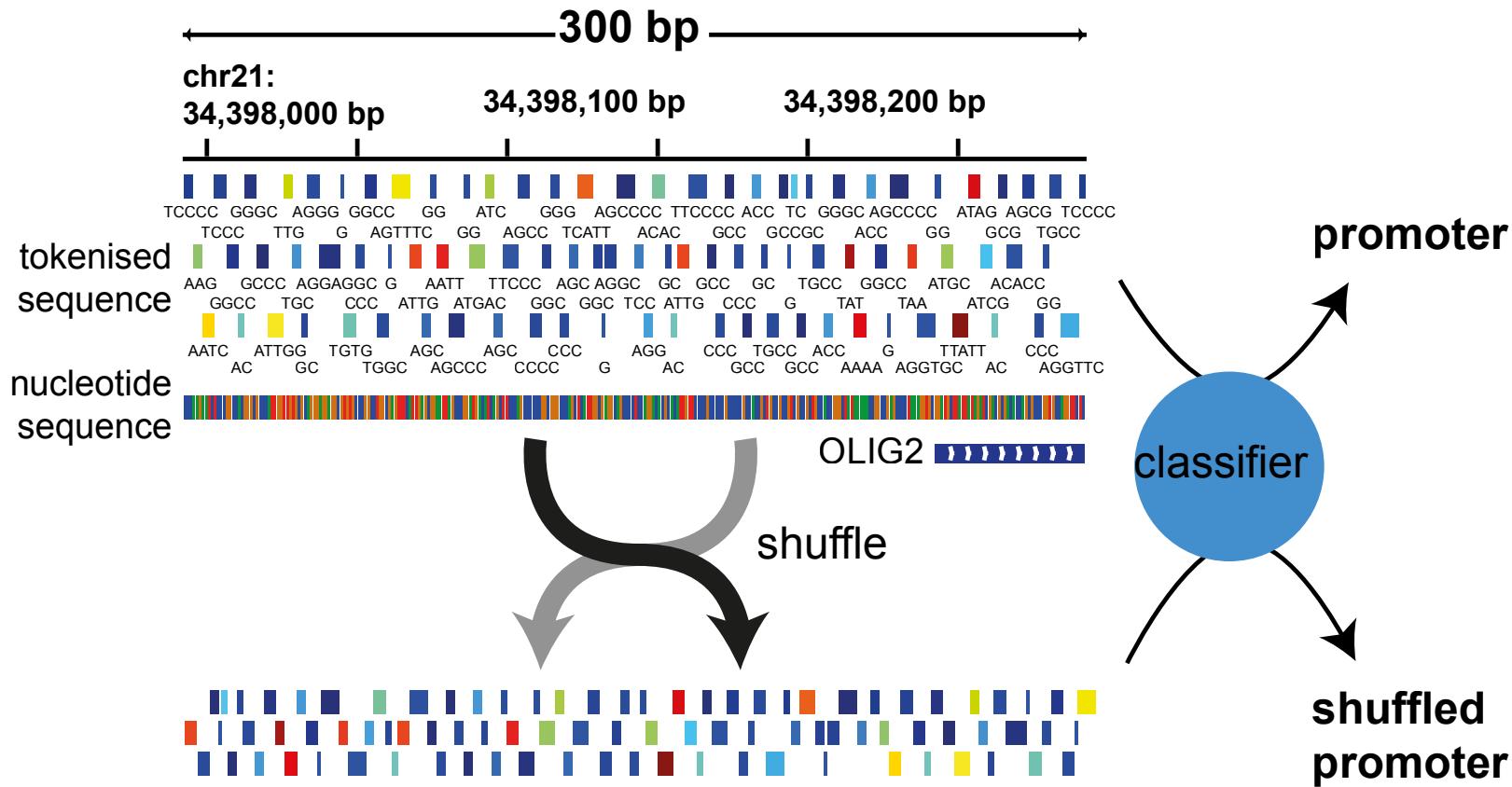


Sanabria *et al* Nature Machine Intelligence (accepted) & Sanabria *et al* BMC Bioinformatics (accepted)



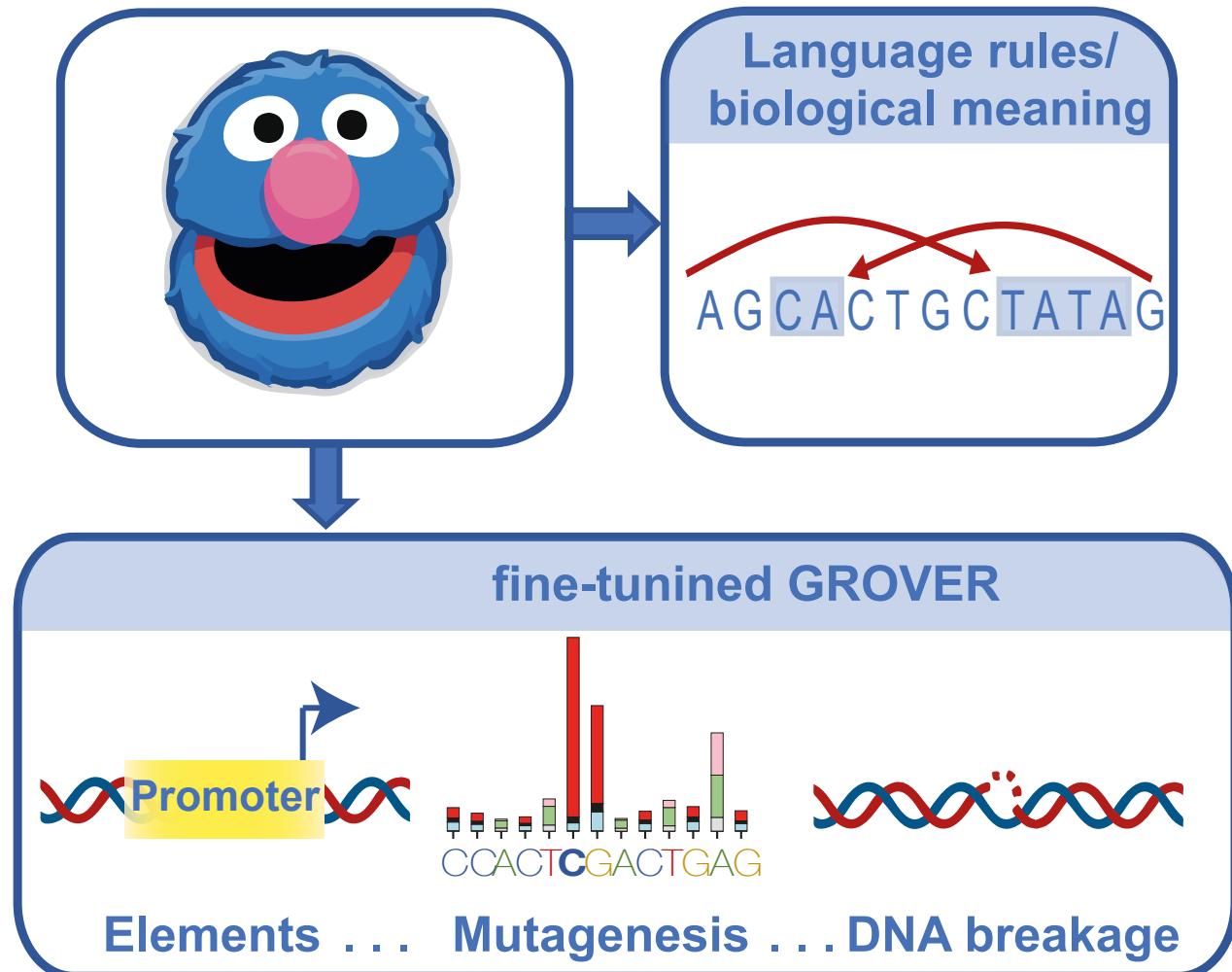
# Fine-tuning tasks to understand genome biology

## Prom300



Sanabria et al Nature Machine Intelligence (accepted)

# What can we use GROVER for?



The model:

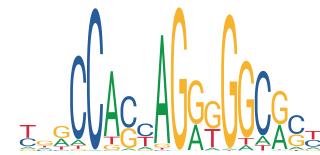


<https://huggingface.co/PoetschLab/GROVER>

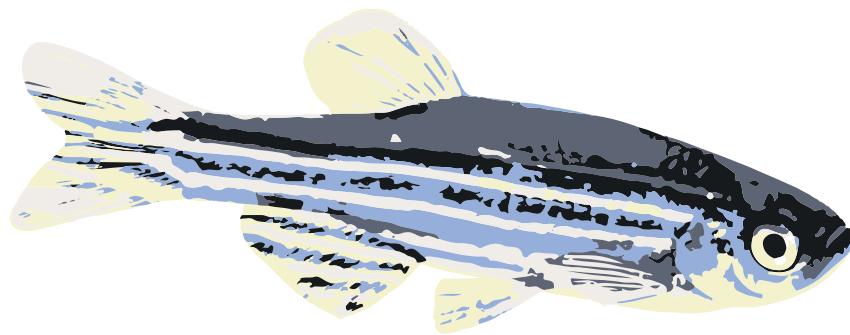
Tutorial for GROVER-based CTCF binding prediction:



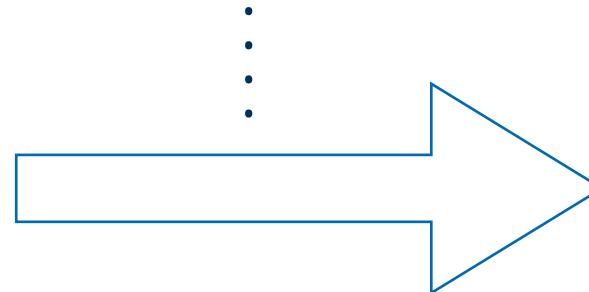
<https://doi.org/10.5281/zenodo.8373158>



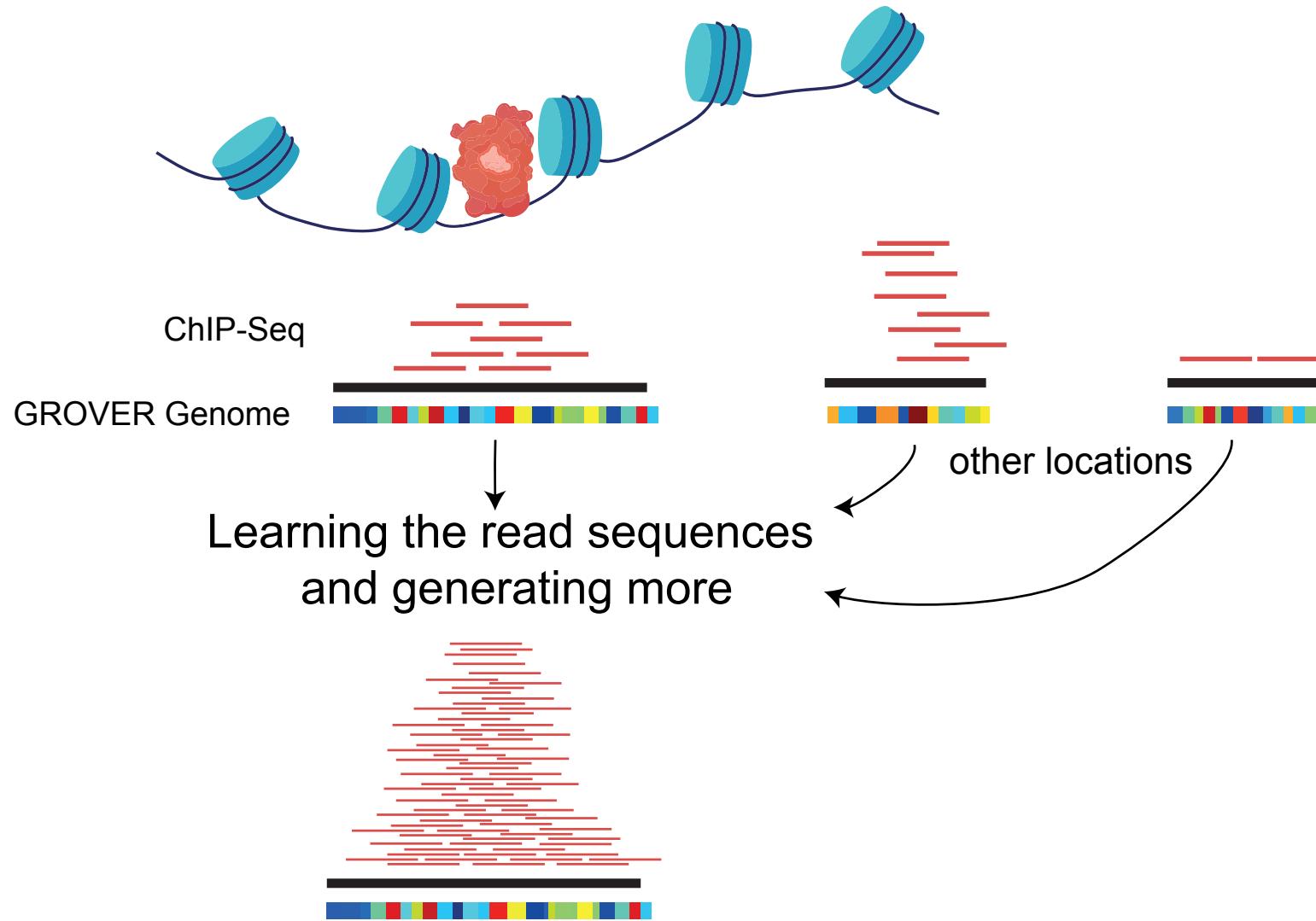
# Translation AI - translation of genomics data across species



annotation  
regulatory genomics data



# Generative AI - Boosting of genomics data



# TextToImage - GenotypeToPhenotype

“DallE2, please give me an illustration of damaged DNA in comic style”:



“...ACGCGTAAAATCGATTAGCGATTGCAA...”:



<https://www.ebi.ac.uk/training/online/courses/human-genetic-variation-introduction/what-is-genetic-variation/genotype-or-phenotype/>

# Summary, take-home messages, and further information

- Large Language models are models trained on next- or masked-token prediction
- They learn a sense of grammar and syntax, as well as language context
- DNA resembles “language” and can be trained analogously
- DNA language models can be used to understand genomes
- DNA language models can be fine-tuned for a myriad of tasks



Attention is all you need: <https://arxiv.org/abs/1706.03762>

How Transformers Work: A Detailed Exploration of Transformer Architecture:  
<https://www.datacamp.com/tutorial/how-transformers-work>

But what is a GPT? Visual intro to transformers: <https://www.youtube.com/watch?v=wjZofJX0v4M>

Attention in transformers, visually explained: <https://www.youtube.com/watch?v=eMlx5fFNoYc>