

An introduction to Large Language Models

Anna Poetsch

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

Organisation

25.6.2024: Dimensionality reduction

2.7.2024: Machine Learning

**9.7.2024: Machine Learning exercise
Deep Learning, Large Language Models**

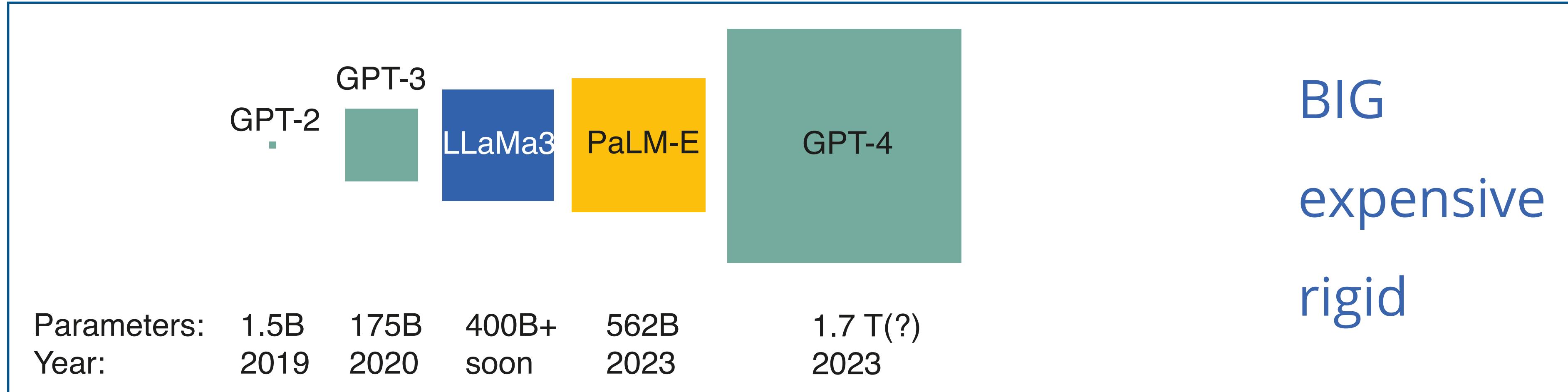
16.7.2024: Deep Learning exercise
Summary, recap, Q&A

Outline

- What are Large Language Models?
- What do they learn?
- Token embeddings
- A short glimpse into a transformer foundation model
- Hallucination and bias
- Large Language Models with language-like data

Pre-training and fine-tuning

Foundation models train language on large corpora of data



Fine-tuning

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

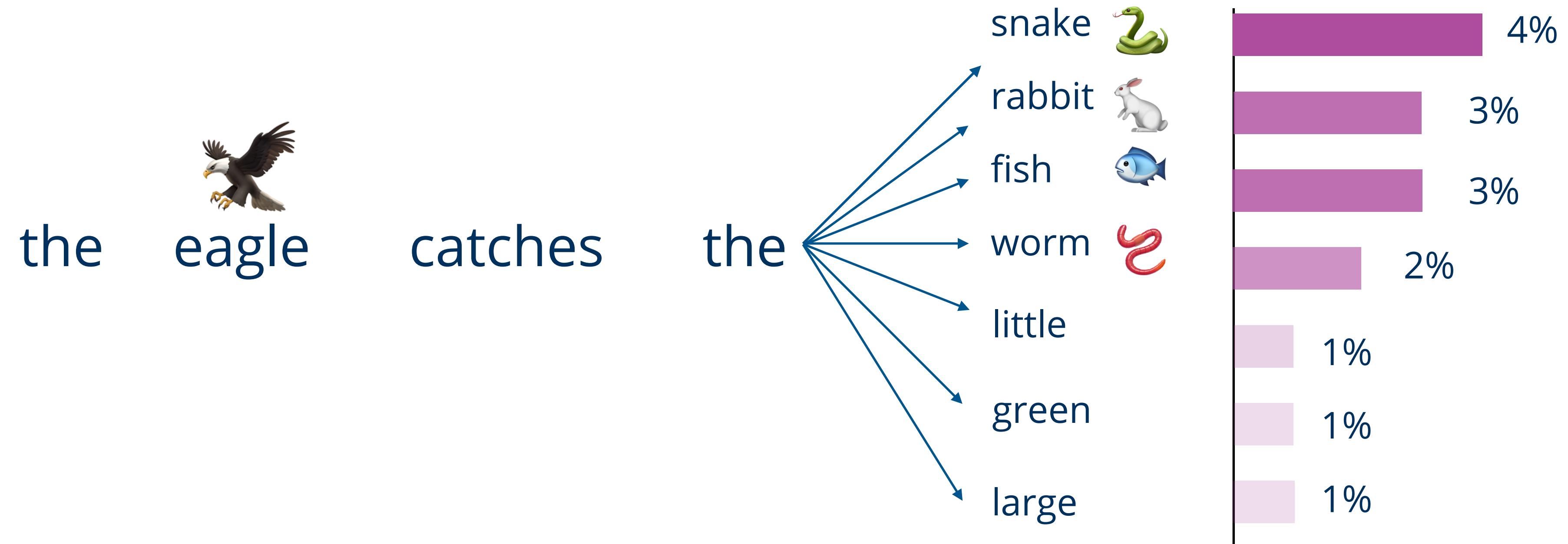
“DallE2, please give me an illustration of damaged DNA in comic style”:



efficient
flexible
cheap

Large Language Models

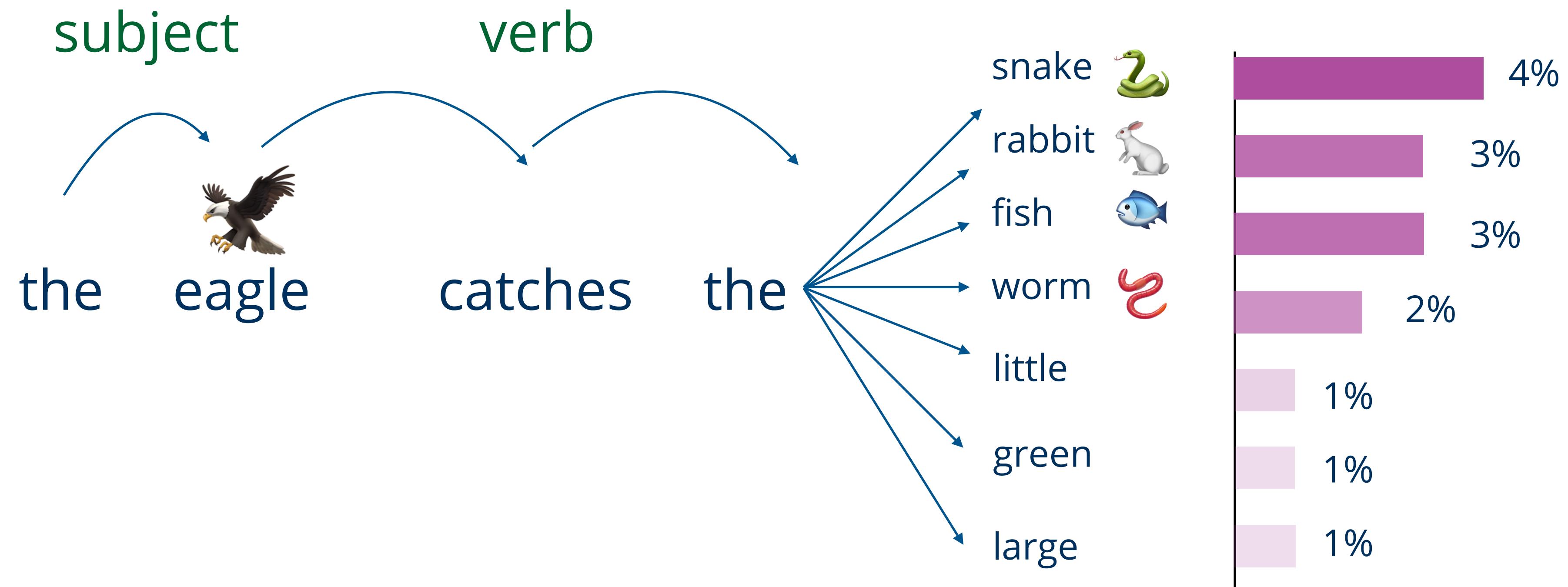
Predicting the next word...



Large Language Models

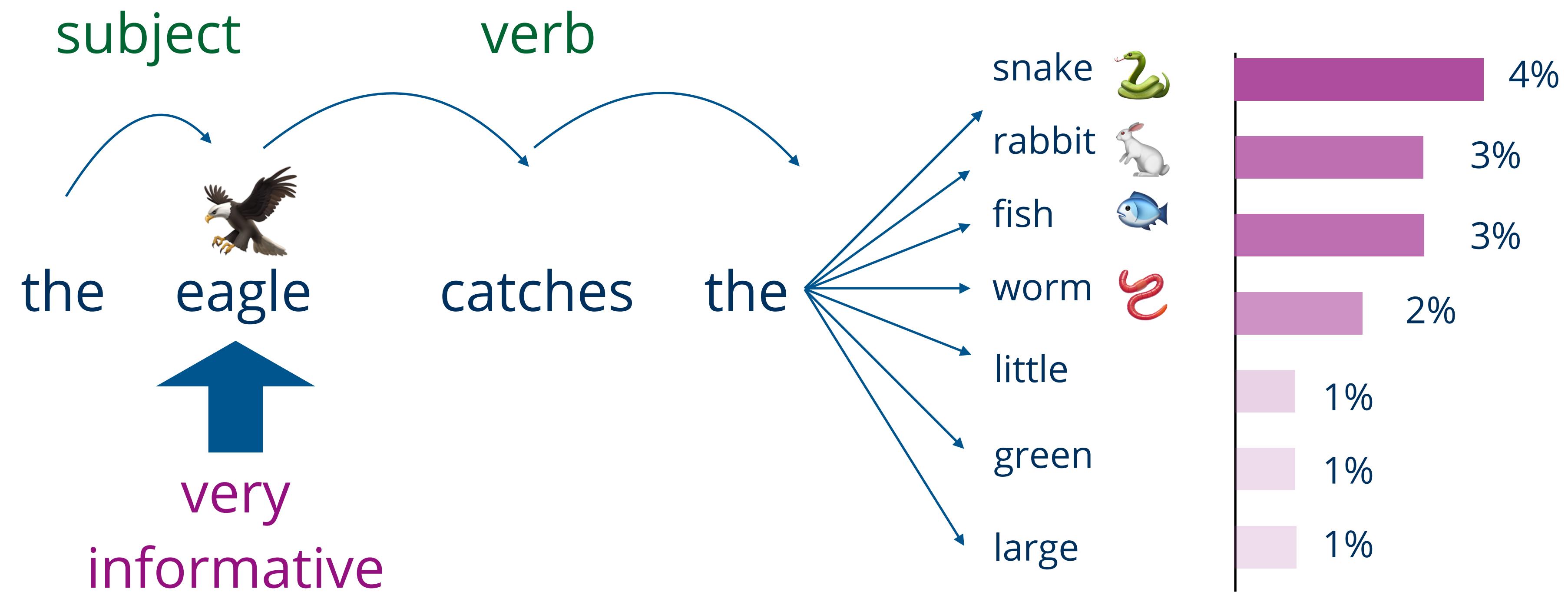
Predicting the next word...
...requires context

object
noun?
adjective?



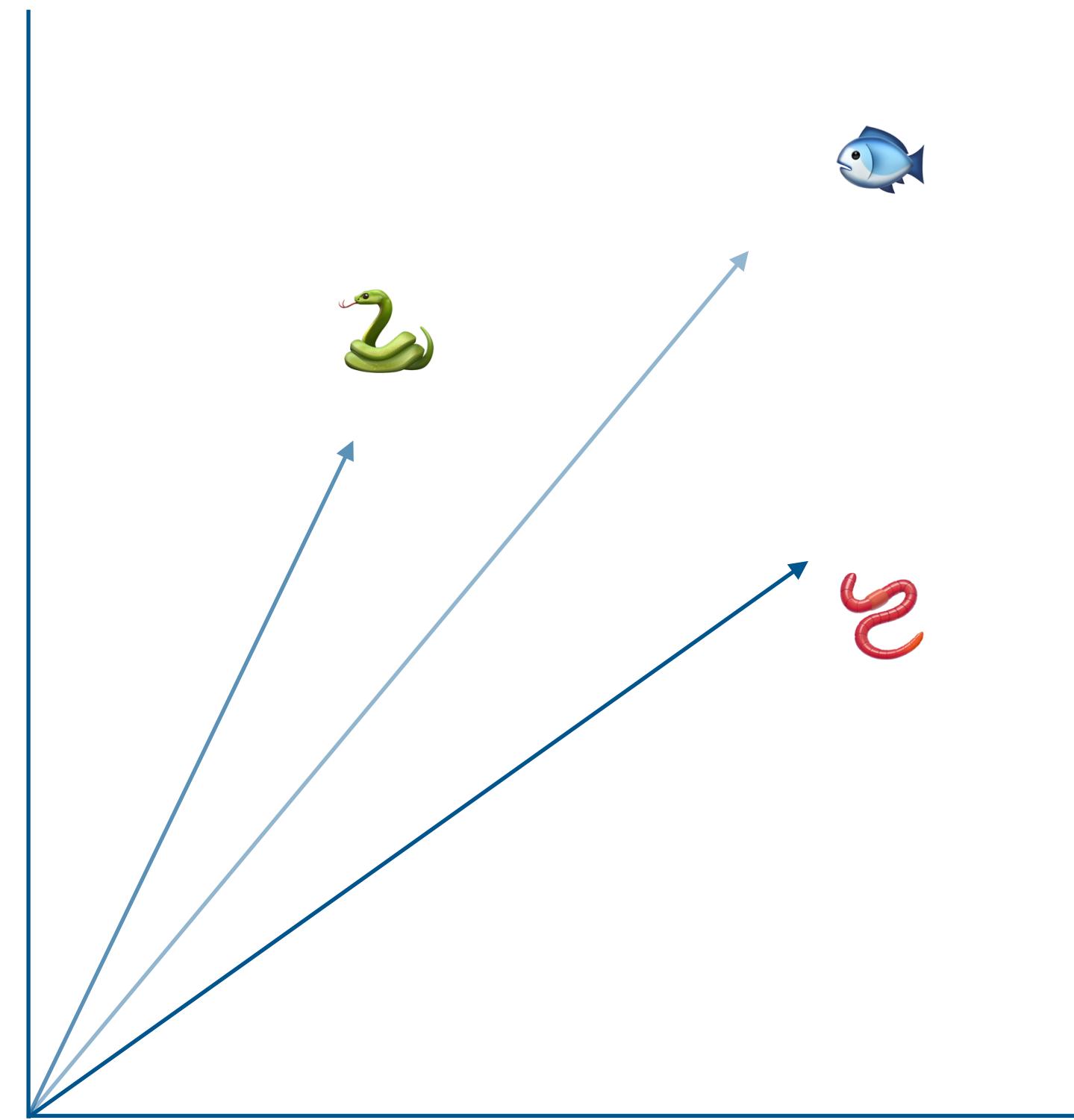
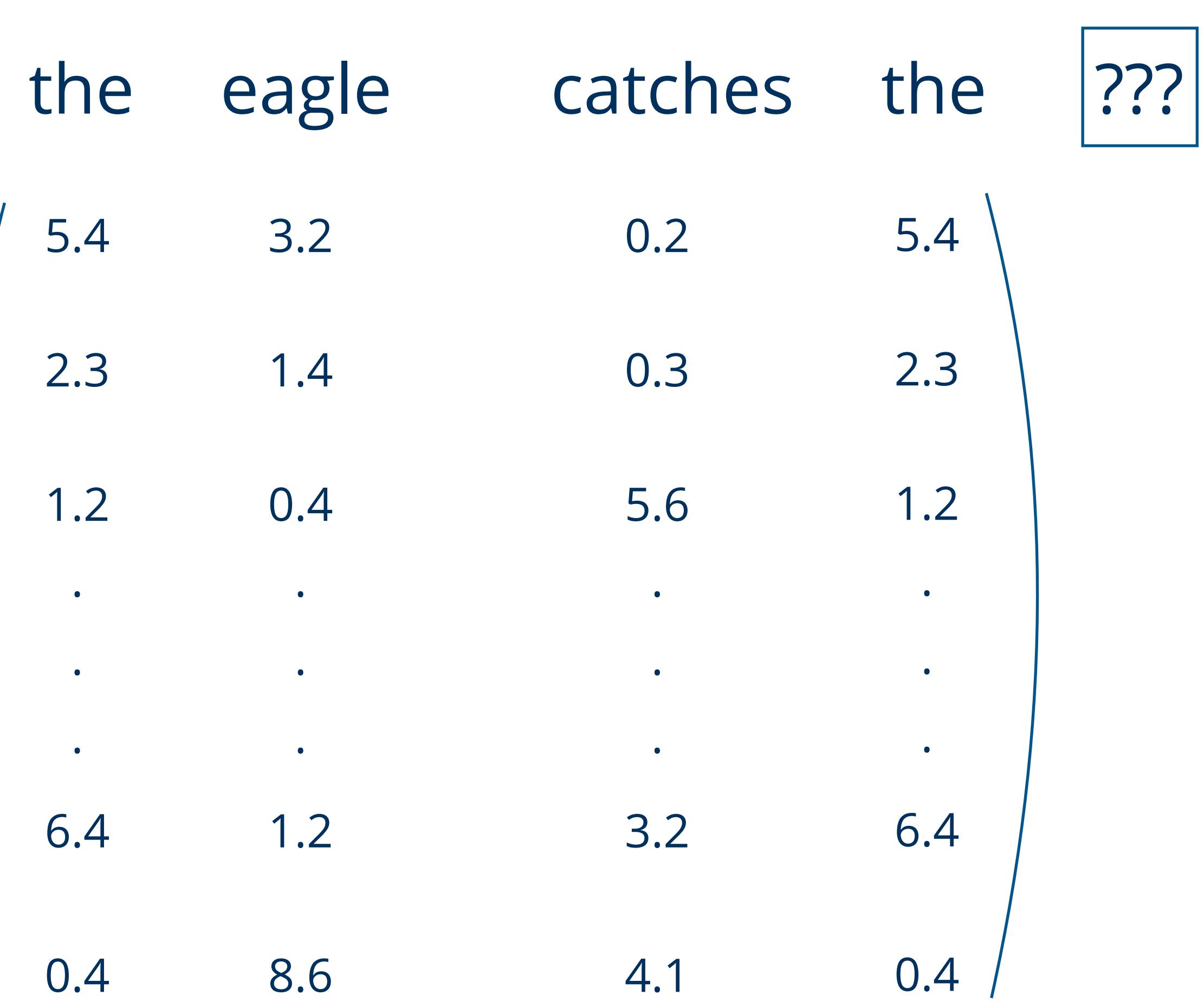
Large Language Models

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

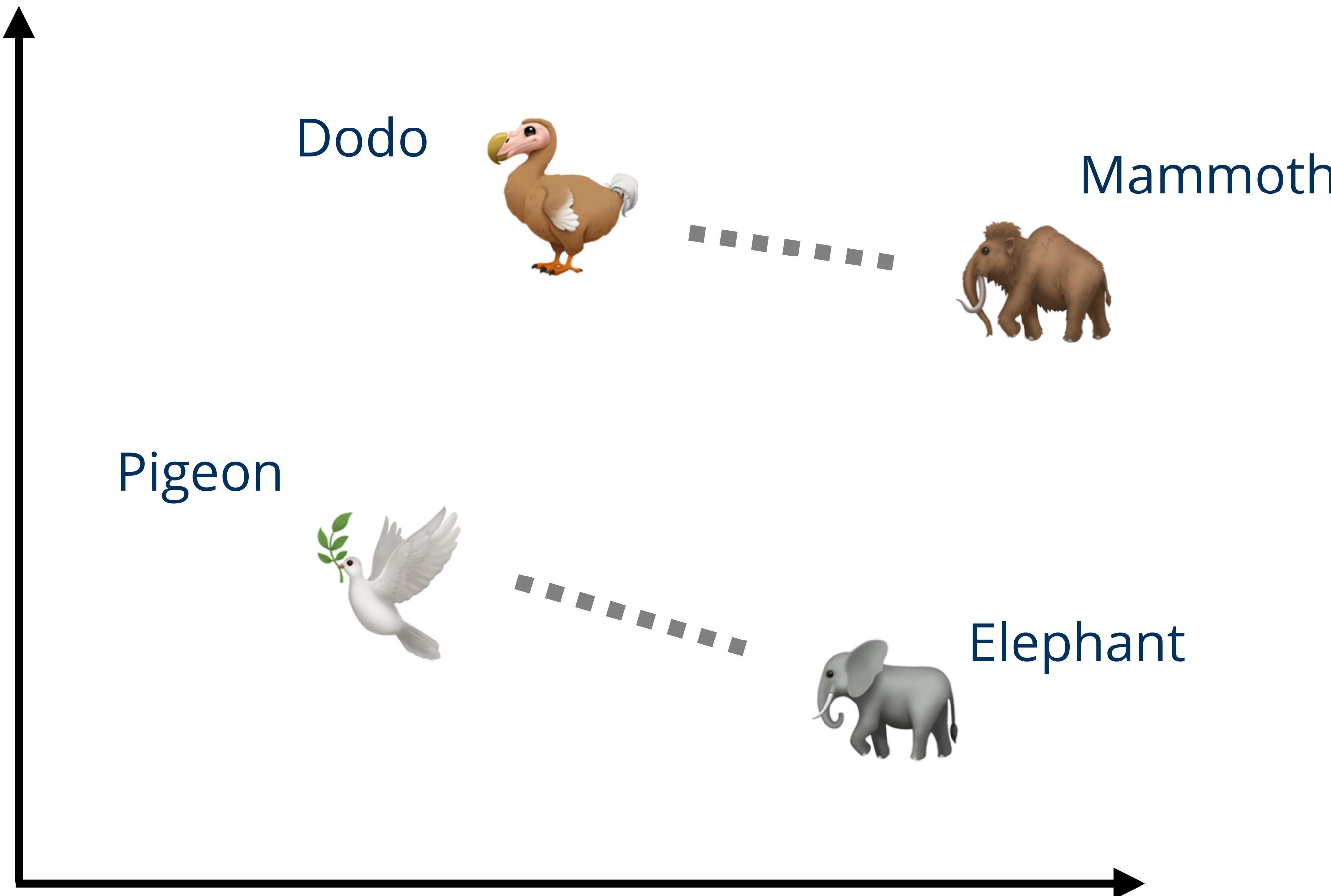


The embedding of words/ tokens

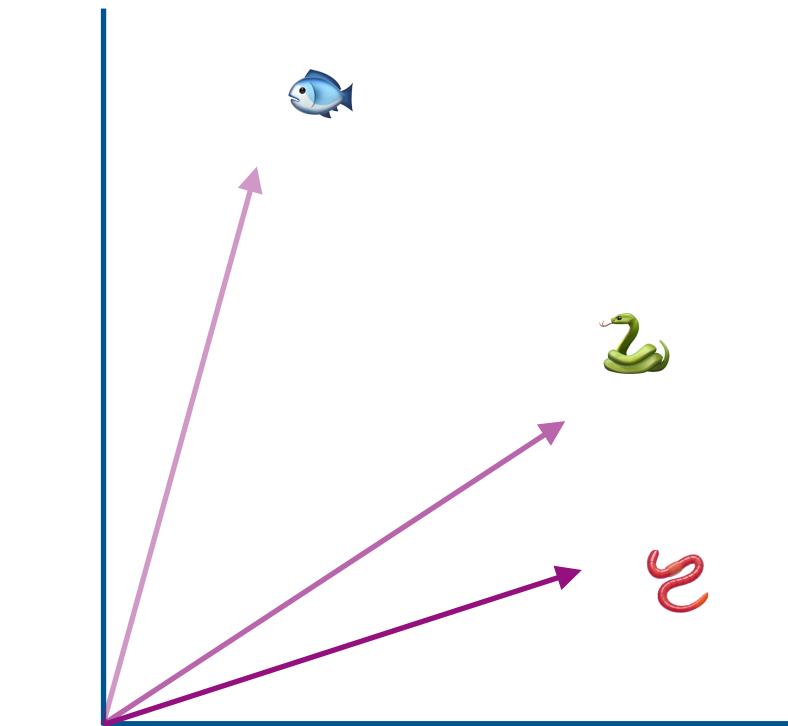
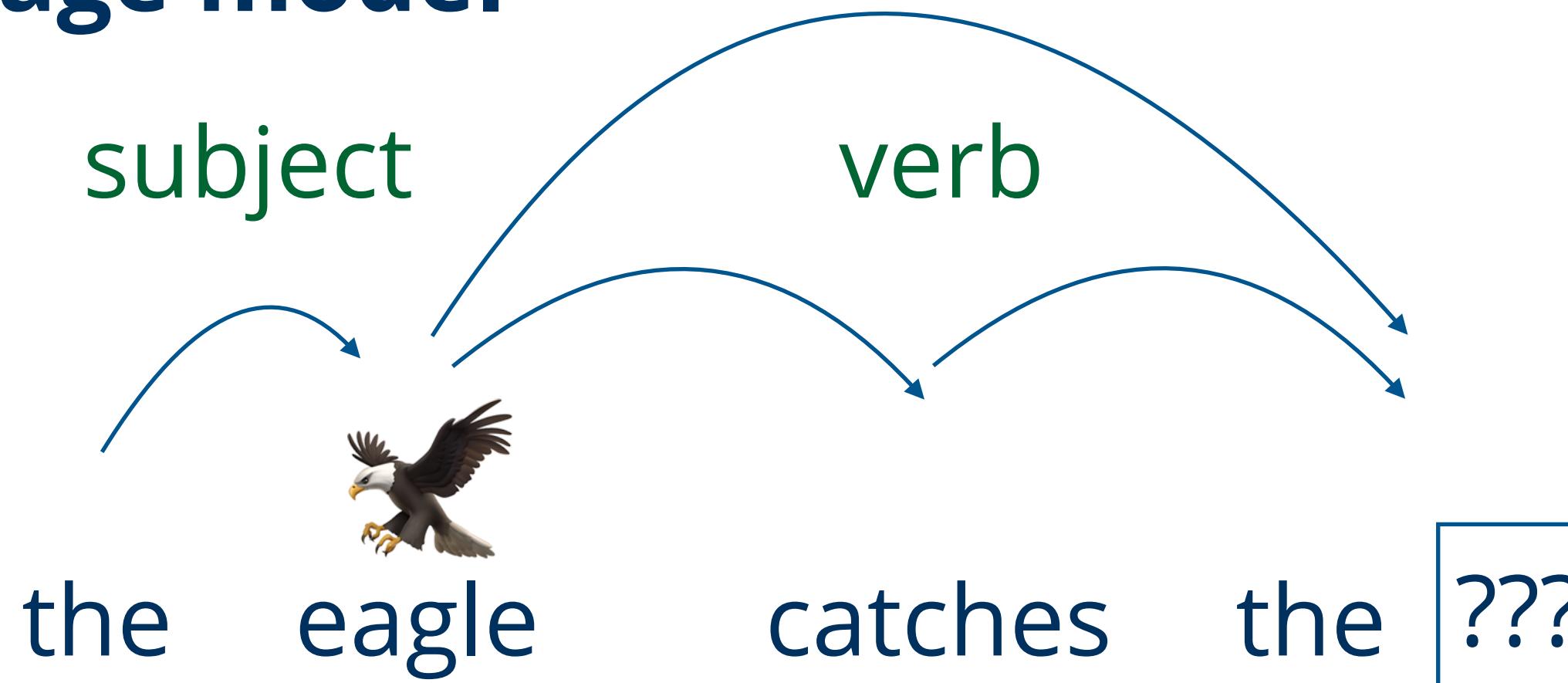
Tokens are assigned a vector, which places them into a multi-dimensional space



Recap: PCA in the embedding of language



Training a large language model

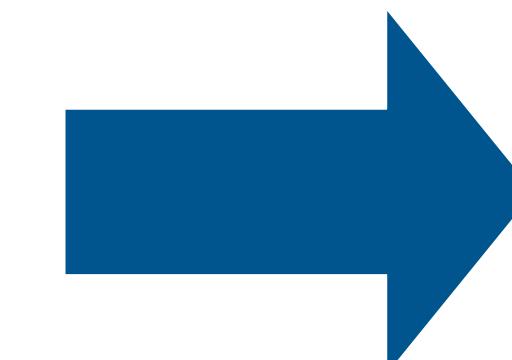


Input embedding

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

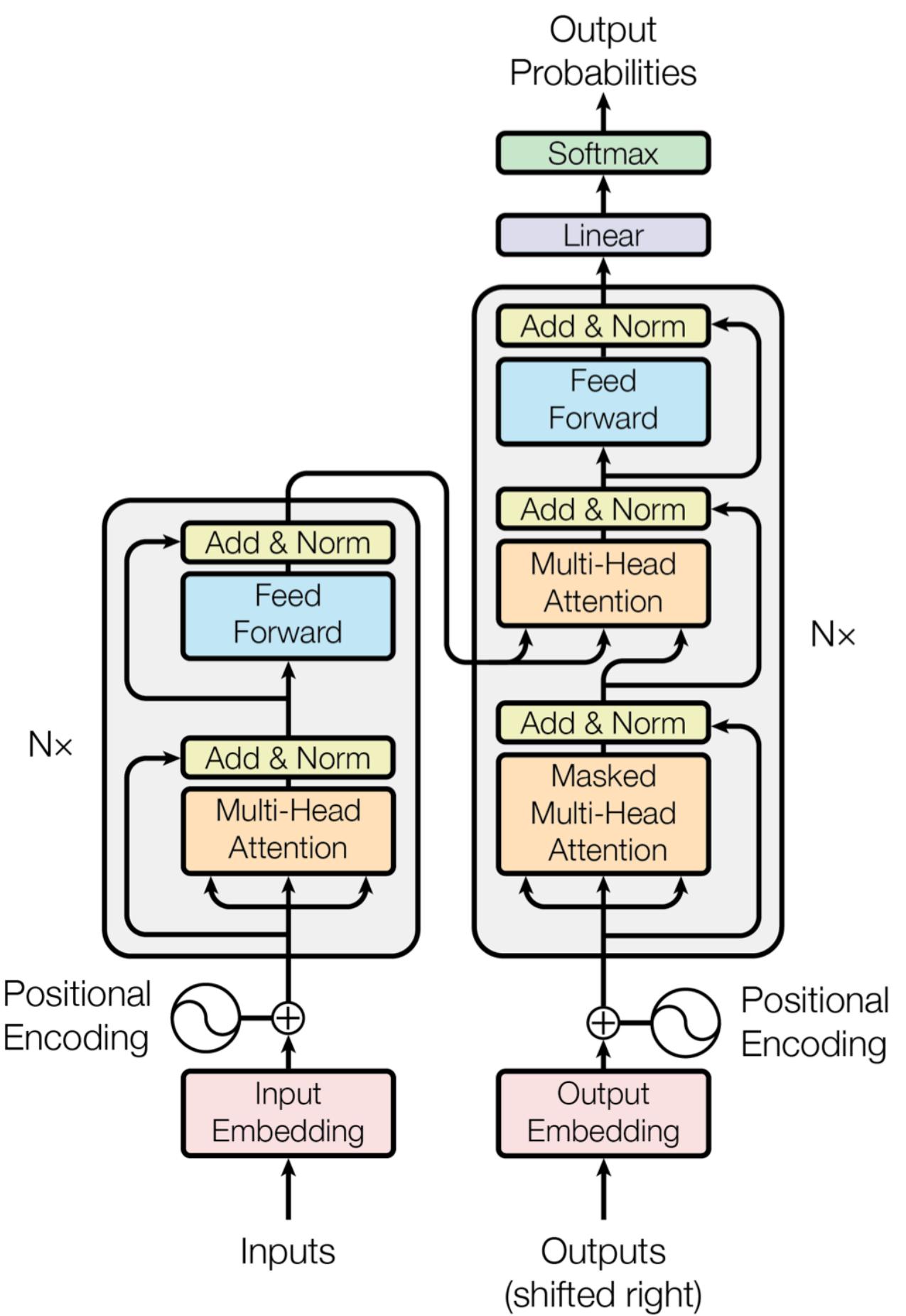
learned through
attention blocks



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

“Attention is all you need”



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

<https://medium.com/@tech-gumptions/transformer-architecture-simplified-3fb501d461c8#:~:text=The%20transformer%20architecture%20consists%20of,translation%20or%20a%20text%20continuation.>

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

“Attention is all you need”

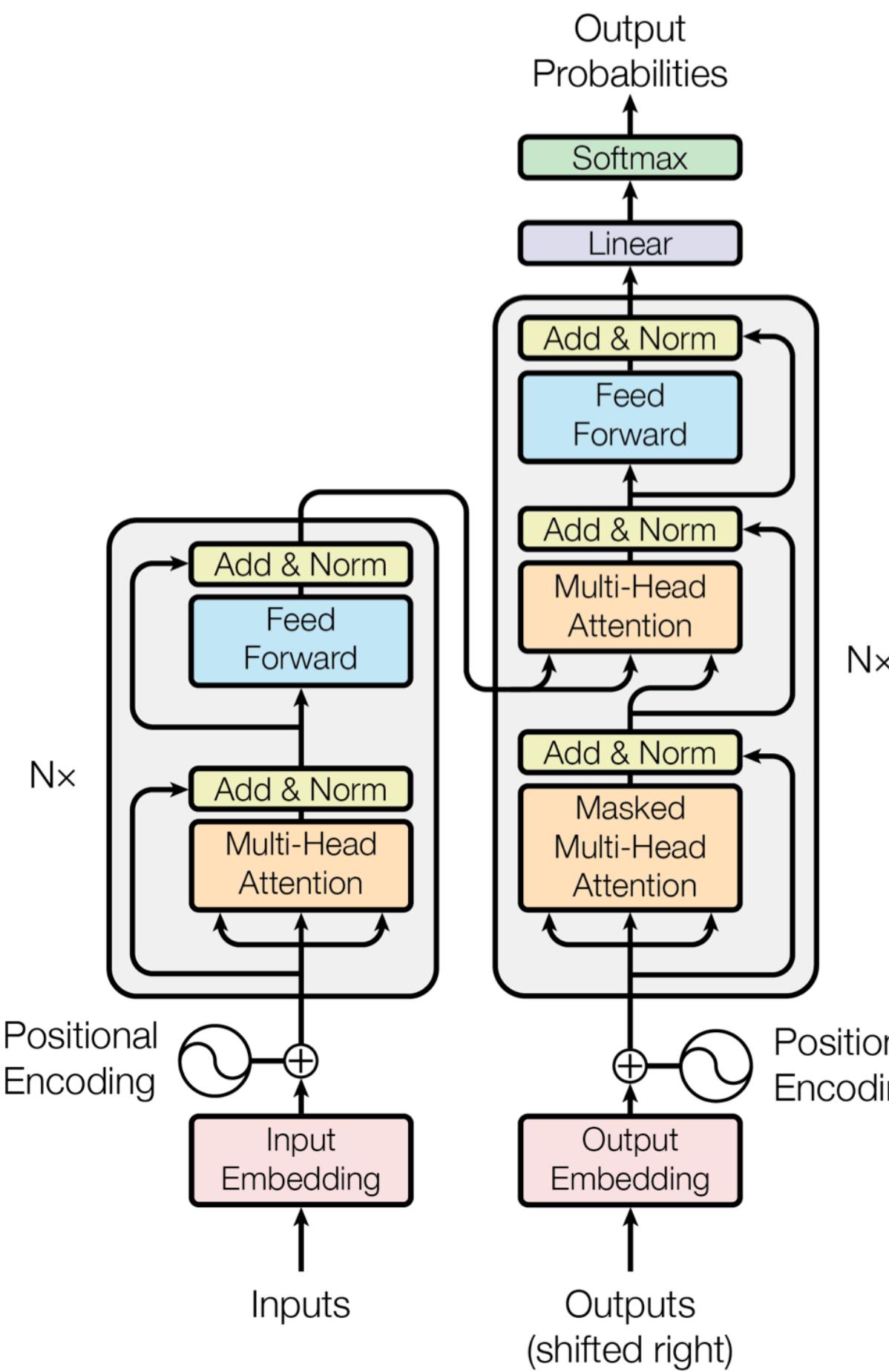
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Query (Q): Represents the word we are currently focusing on. Each word has its corresponding Query vector.

Key (K): Represents all words we want to pull information from to help determine the relevance of each word to the Query.

Value (V): Contains the information of words that we'll extract when there's a match between Query and Key.

“Attention is all you need”



GPT-3

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

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

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

Hallucination, “making things up”

Anna Poetsch is a Computational Biologist of the TU Dresden.

She has received a grant from the European Research Council (ERC) and is one of the winners of the “Cluster of Excellence: Inflammation at Interfaces”. The Cluster of Excellence led by the TU Dresden and the University of Leipzig has been funded with more than 40 million euros by the German Federal Ministry of Education and Research (BMBF) for a period of ten years.

Out of context conversations



ChatGPT

I apologize if it seems like I am asking the same questions over and over again. I am a large language model trained by OpenAI and do not have personal experiences or memories. As a result, I am unable to remember past conversations or specific questions that I have asked. I generate responses based on the prompts given to me and do not have the ability to browse the web or access other information.

Bias

A model is only as good as the data used to train it.

Trained it with racism and sexism, it will return just that.

Compensating for bias is hard.

Google Gemini: Adi Robertson / The Verge

Other text-like data

Proteins

DNA

Sound

Music

•

•

•

Summary, take-home messages, and further information

- Large Language models are models trained self-supervised on next- or masked-token prediction
- They learn a sense of grammar and syntax, as well as language context
- They can be fine-tuned for a myriad of tasks (e.g. assistants and image generators)
- They are at risk of hallucination and amplifying bias
- Any data that resembles “language” can be used for training a model like this

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>

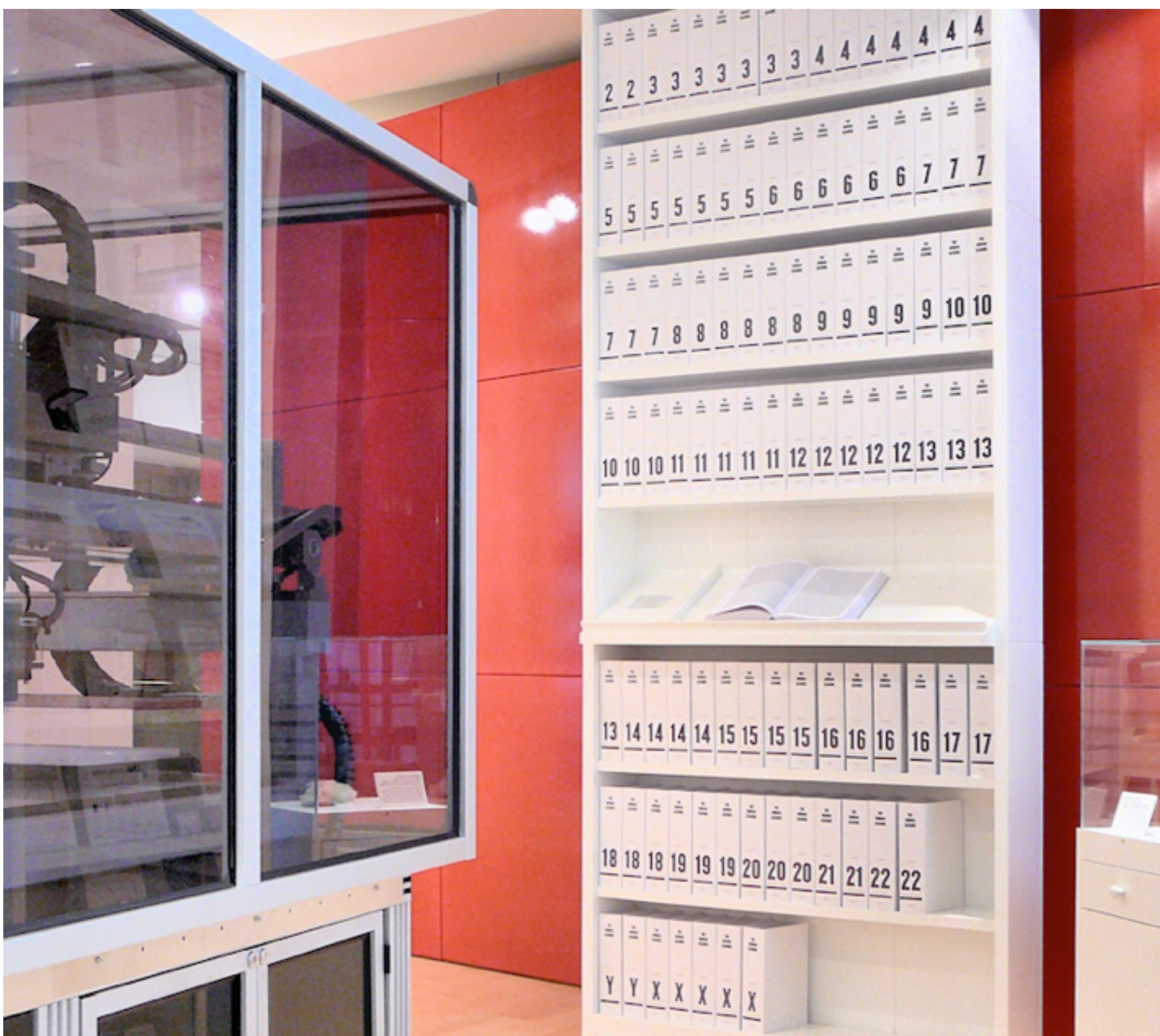
Attention in transformers, visually explained: <https://www.youtube.com/watch?v=eMlx5fFNoYc>

Using DNA language models to understand different layers of code in genomes

Anna Poetsch

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

The Human Genome

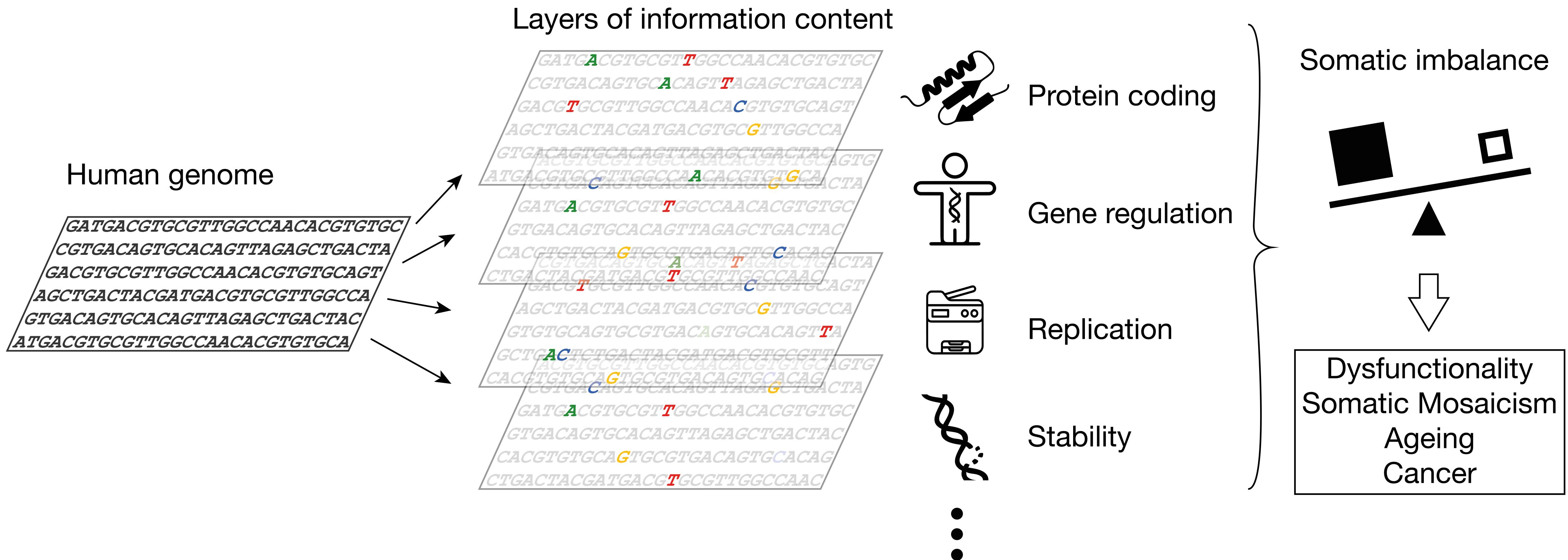


**3.2 Billion letters
119 books
20.000 genes (1-2 %)**

Wellcome Trust Collection in London

[Image source](#): Russ London at English Wikipedia, licensed [CC-BY-SA 3.0](#)

Information content of DNA



Can we treat DNA as if it were text?



DNABERT-2



GROVER



HyenaDNA

**Nucleotide
Transformer**

Proformer

Enformer

⋮



DNA-BERT

Yi *et al*: <https://academic.oup.com/bioinformatics/article/37/15/2112/6128680>
Sanabria *et al*: <https://www.biorxiv.org/content/10.1101/2023.07.11.548593v1>

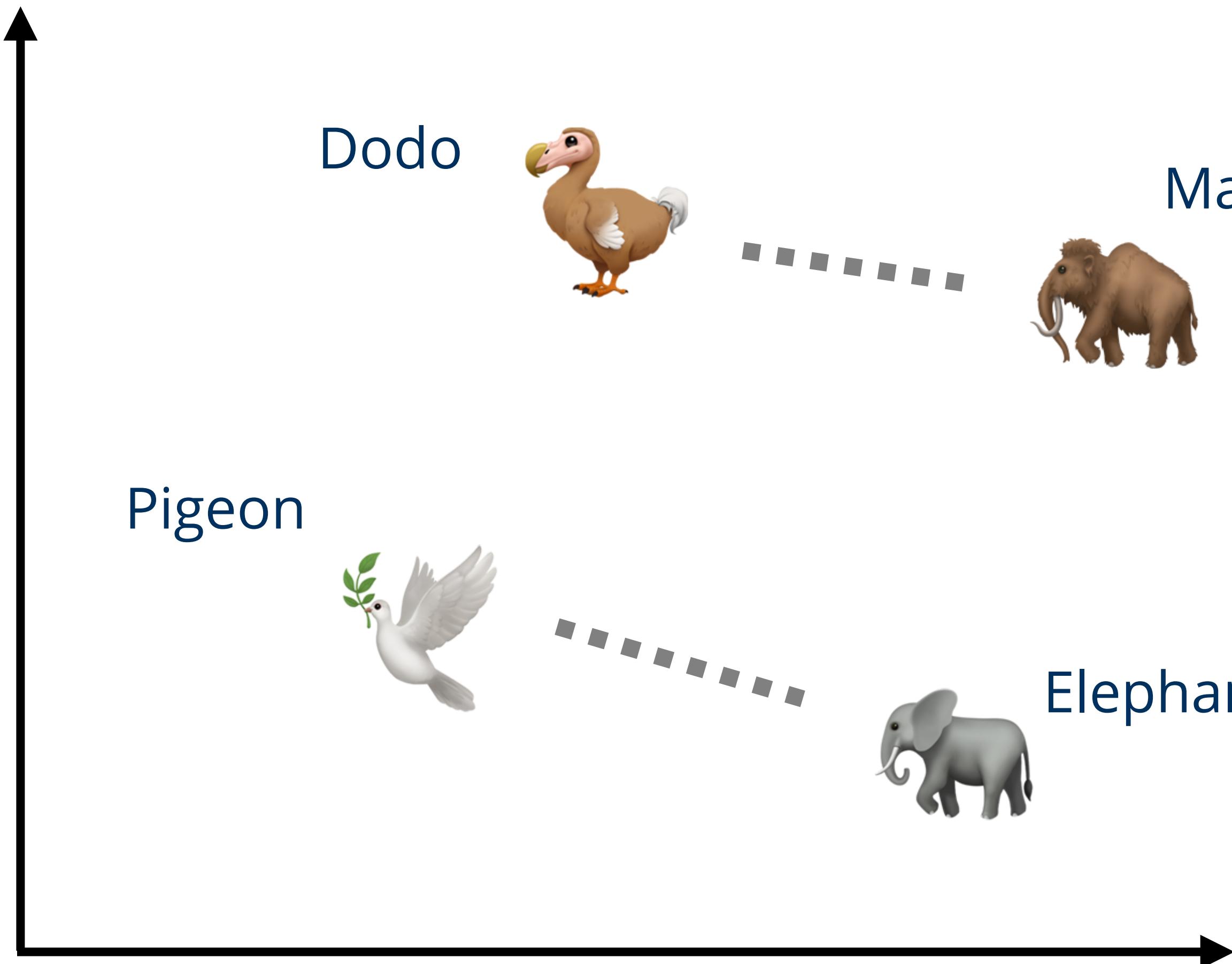


DNA-BERT

Sanabria *et al*: <https://www.biorxiv.org/content/10.1101/2023.07.11.548593v1>

Word2Vec embedding

Algorithm for semantic word associations through static embeddings



- shallow neural networks
- input is large corpus of text
- Each unique word is assigned a position in the vector space
- Different methods can be used, e.g. continuous bag-of-words (CBOW) “fill in the blank” task.



Using next-k-mer prediction to compare models

Sanabria *et al*: <https://www.biorxiv.org/content/10.1101/2023.07.11.548593v1>

**Overlaps lead to learning how to spell
How can we avoid this?**

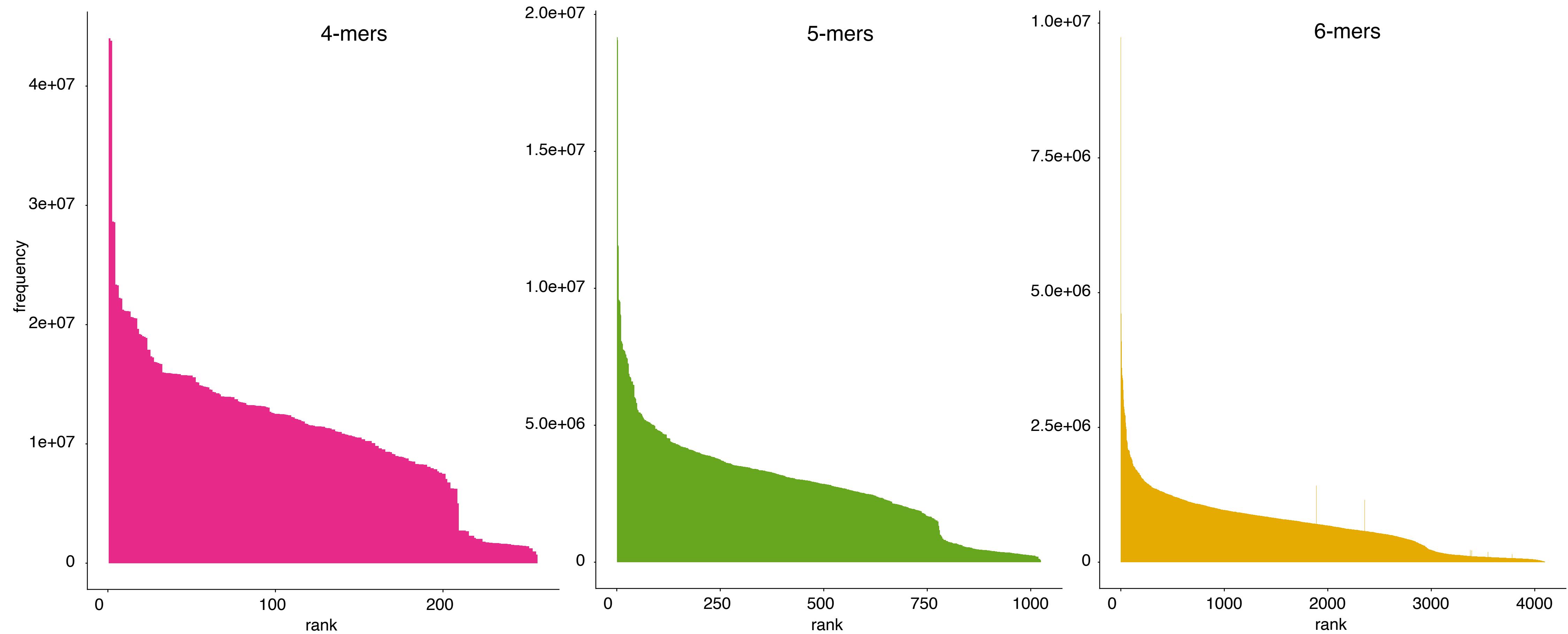


GROVER: Genome Rules Obtained Via Extracted Representations

Sanabria *et al* BioRxiv

What is a word?

TACCCATGTACCTGGAGGCCATGTTGCTTCCTCTAGCAAAAGATGCAGCAGATGCAGTTGGGGCTGCCACT



Byte Pair Tokenisation

T | **T** -> **TT**

A | A-> AA

T | G-> TG

A | G -> AG

A | G-> AG

C | C-> CC

G|TT|C|G|AG|A|C|T|AG|CC|TG|AA|A|CC|T|C|A|C|G|CC|TG|AA|A|CC|T|C|A|T|C|T|C|T|A|C|T|AA|AA|AA|AA|AA|AA|A|T|AA|
T|AG|C|TG|G|G|C|G|T|C|G|G|C|AG|G|TG|CC|TG|T|AA|T|CC|C|AG|C|T|A|C|T|C|AG|TG|G|G|C|TG|AG|G|C|AG|C|
C|AG|AA|G|G|C|G|G|AG|G|TT|G|C|AG|TG|AG|CC|G|AG|A|T|CC|C|G|CC|A|C|TG|C|A|C|T|CC|AG|CC|TG|G|G|TG|A|C|
C|T|CC|AA|AA|AA|AA|AA|AA|AA|AA|AG|C|AG|G|C|T|AG|G|C|T|AA|G|C|T|A|TG|A|TG|TT|CC|TT|AG|A|TT|AG|G|
C|AA|C|TT|A|C|AA|T|A|TT|TT|C|AA|C|T|G|A|C|G|AG|TT|T|A|T|C|AG|G|AA|G|T|AA|C|A|CC|A|T|C|G|T|AA|
T|A|T|C|AG|G|C|AA|AG|T|C|A|T|AG|AA|CC|A|T|TT|TT|C|A|TG|C|T|C|TT|T|AA|C|AA|TT|TT|C|TT|TT|TG|AA|AG|C|
AA|AA|T|A|T|A|T|A|TT|A|TG|G|T|A|T|AA|G|TT|G|G|TG|TT|C|TG|AA|G|TT|AG|C|T|A|C|AA|CC|AG|G|AG|CC|A|
T|C|A|C|TG|CC|C|CC|TG|A|TG|G|C|AA|A|TG|CC|C|C|AA|TT|G|C|AG|G|T|AA|AA|C|AG|T|C|AA|G|AA|AA|C|G|G|C|A|
A|C|TG|G|AA|A|C|TT|T|CC|A|C|TT|G|A|T|AA|G|AG|G|T|CC|C|AA|G|A|C|TT|AG|T|A|CC|TG|AA|G|G|G|TG|AA|A|T|A|
TT|T|C|TT|C|TT|TG|G|C|TG|G|G|AG|A|G|G|AG|C|T|TG|G|TG|TT|G|G|G|C|AG|TG|C|T|AG|G|AA|AG|A|G|G|C|AA|
G|TG|AA|T|C|TG|AG|G|C|AA|C|TG|C|A|CC|C|TT|G|G|T|C|CC|T|C|CC|A|CC|G|C|TT|C|TT|G|T|CC|TG|C|C|TT|G|C|
C|T|CC|C|TG|G|G|G|C|AG|C|T|C|G|TG|G|TG|AG|G|C|T|CC|C|C|TT|T|C|TT|G|C|G|AG|A|TT|C|T|C|TT|CC|T|C|TG|
T|CC|C|AG|G|A|C|AG|G|C|A|C|AA|A|C|A|C|G|C|A|CC|T|C|AA|AG|C|TG|TT|CC|G|T|CC|C|AG|T|AG|A|TT|A|CC|A|C|
A|AA|G|AG|AA|G|C|AA|G|AG|G|C|AG|T|AA|G|G|AA|A|T|C|AG|G|T|CC|T|A|CC|TG|T|CC|C|A|TT|T|AA|AA|AA|CC|AG|G|C|
C|AA|CC|A|CC|C|TT|G|T|CC|TT|T|C|TG|G|AG|CC|T|AA|G|C|T|CC|AG|C|T|CC|AG|G|T|AG|G|TG|G|AG|G|AG|AA|G|CC|A|C|
C|C|AA|AG|CC|AG|A|G|AA|AA|A|G|AA|AA|C|TG|AG|TG|G|G|AG|C|AG|T|AA|G|G|AG|A|TT|CC|C|C|G|CC|G|G|G|A|TG|T|G|A|
G|G|G|T|AG|T|AG|T|A|TG|G|AA|G|AA|A|T|C|G|G|T|AA|G|AG|G|TG|G|G|CC|C|AG|G|A|AA|G|C|AG|A|G|G|C|TG|G|G|C|A|
G|C|AG|G|CC|AG|TG|T|G|G|TG|G|C|AA|G|TG|CC|T|CC|TG|A|CC|TG|G|AG|T|C|TT|CC|AG|TG|T|G|A|TG|A|TG|G|TG|G|G|CC|
CC|G|G|TT|C|A|TG|CC|G|CC|C|A|TG|C|AG|G|AA|C|TG|TT|A|C|A|TG|T|AG|TT|G|T|AG|TG|G|A|TG|G|TG|G|T|A|C|AG|T|C|AG|A|
AG|G|AG|A|T|AA|C|A|C|AG|G|CC|C|AA|G|A|TG|AG|G|CC|TT|G|G|G|G|AG|A|CC|TG|T|G|G|C|AA|G|C|AG|G|G|G|CC|TT|TT|
TT|TT|TT|TT|G|AG|A|T|TG|G|AA|A|T|C|G|C|T|C|TG|T|C|G|CC|C|AG|G|C|TG|G|AG|TG|C|AG|TG|G|C|G|T|G|A|T|C|T|C|G|
A|G|C|T|CC|A|CC|G|CC|C|AG|G|TT|C|A|C|G|CC|A|TT|C|T|CC|TT|CC|T|C|AG|CC|T|C|G|AG|T|AG|C|T|G|G|A|C|T|A|C|AG|G|TG|CC|C|AG|C|
CC|A|C|G|CC|C|G|C|T|AA|TT|TT|TT|TG|T|A|TT|TT|C|AG|T|AG|A|G|A|C|G|G|G|TT|T|C|A|CC|G|TT|AG|CC|AG|G|A|TG|G|T|C|T|C|G|
T|C|T|CC|C|AA|CC|T|C|G|TG|A|T|CC|G|CC|TT|G|G|CC|T|CC|C|AA|AG|TG|C|TG|G|G|A|TT|A|C|AG|G|C|A|T|G|AG|CC|A|C|TG|C|G|CC|C|
CC|AA|G|C|AG|G|G|G|AG|G|CC|C|TT|AG|CC|T|C|TG|T|AA|G|G|C|TT|C|AG|TT|TT|C|AA|C|TG|T|G|C|AA|T|AG|TT|AA|A|CC|C|A|TT|T|A|C|
C|A|C|A|T|C|A|TG|G|G|TT|A|T|AG|G|G|AG|G|T|C|AA|A|T|AA|G|C|AG|C|AG|G|AG|AA|AG|CC|C|CC|C|T|A|C|TG|C|T|C|A|CC|TG|G|G|AG|G|
CC|A|C|TG|A|C|AA|CC|A|CC|C|TT|AA|CC|C|C|T|CC|C|AG|A|G|A|CC|C|C|AG|TT|G|G|C|AA|A|CC|T|C|AG|G|C|G|C|T|C|A|TG|G|TG|G|G|G|AG|C|
G|C|A|CC|A|CC|A|C|T|A|TG|T|C|G|AA|A|G|TG|TT|T|C|TG|T|C|A|T|CC|AA|A|T|A|C|TG|G|CC|A|TG|G|G|C|G|C|G|G|TG|G|CC|G|G|G|TG|G|C|
A|T|AA|G|A|TG|C|TG|AG|G|G|G|CC|AG|A|CC|T|AA|G|AG|C|AA|T|C|AG|TG|AG|G|AA|T|C|AG|A|G|G|CC|TG|G|G|A|CC|C|TG|G|G|C|AA|
AG|CC|C|TG|T|C|G|T|C|T|CC|AG|C|CC|C|AG|C|TG|C|T|C|A|CC|A|T|C|G|C|T|A|T|C|TG|AG|G|C|AG|C|G|C|T|C|A|TG|G|TG|G|G|G|AG|C|
CC|T|C|A|C|AA|CC|T|CC|G|T|C|A|TG|T|G|C|TG|T|G|A|C|TG|C|TT|G|T|AG|A|TG|G|CC|A|TG|G|G|C|G|C|G|G|TG|G|CC|G|G|G|TG|G|C|
G|G|TG|T|G|G|AA|T|C|AA|CC|C|A|C|AG|C|TG|C|A|C|AG|G|G|G|AG|C|TG|T|AA|G|AG|T|C|AG|G|G|C|A|TG|AA|C|AG|A|T|AA|AG|C|AA|C|TG|G|G|AA|
G|G|AG|T|A|C|TG|T|AG|G|AA|G|AG|G|AA|G|G|AG|A|C|AG|A|G|G|TT|G|AA|AG|T|C|AG|G|G|C|A|TG|AA|C|AG|A|T|AA|AG|C|AA|C|TG|G|G|AA|
A|C|G|G|C|AG|C|AA|AG|AA|A|C|AA|A|C|A|TG|C|G|T|AA|G|C|A|CC|T|CC|TG|C|AA|CC|C|A|C|T|AG|C|G|AG|C|T|AG|A|G|AG|A|G|TT|G|G|G|C|G|
T|A|C|A|CC|T|C|AG|G|AG|C|TT|TT|C|TT|TT|TT|TT|TT|TT|G|AG|A|T|AG|G|G|TG|T|C|TT|G|C|T|C|A|C|TG|T|C|A|C|AG|G|C|TG|G|AG|C|
C|AG|TG|G|TG|T|G|A|T|C|A|C|AG|C|T|C|A|C|TG|C|AG|CC|T|CC|TG|G|CC|T|C|AA|G|TG|A|T|C|TT|CC|C|A|C|AG|CC|T|C|A|CC|TG|C|T|C|A|

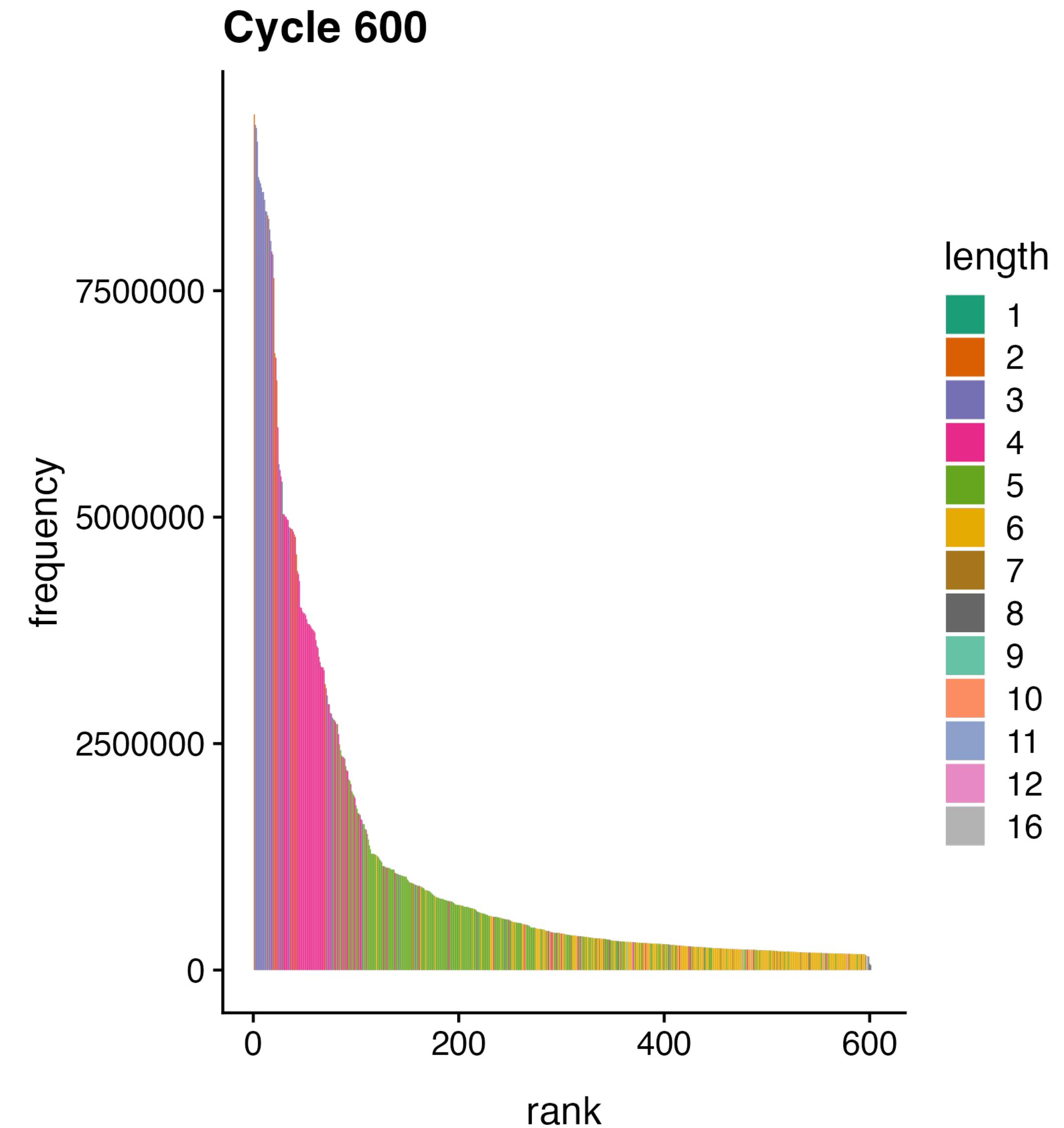
T | C -> TC

A | C-> AC

G | G-> GG

A | TT-> ATT

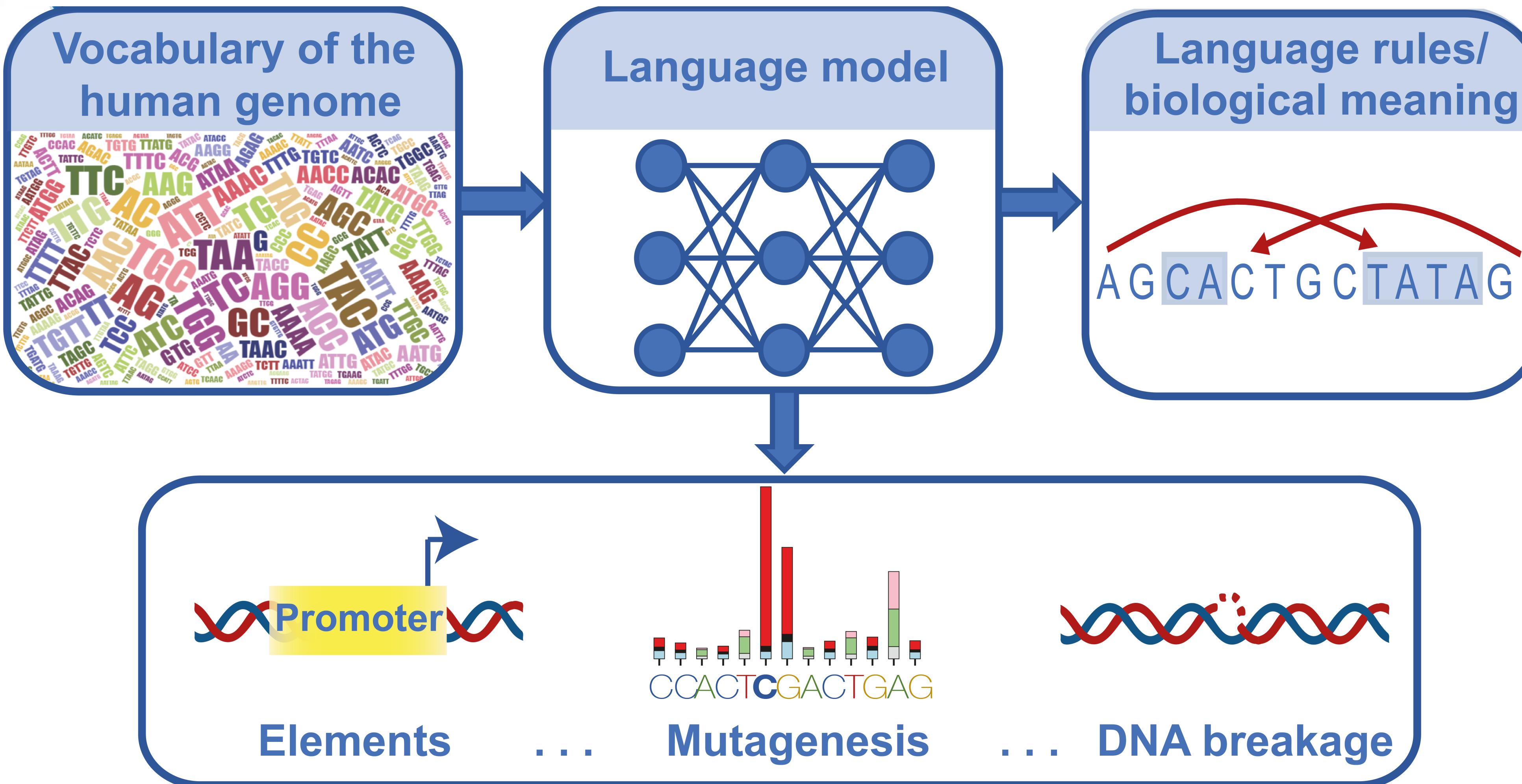
600 cycles builds the “best” vocabulary



What does the vocabulary look like?



What can we use GROVER for?



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





What is GROVER learning?

It learns to differentiate token characteristics

Sanabria *et al* BioRxiv



What is GROVER learning?

It learns to see tokens in context

Sanabria *et al* BioRxiv



What is GROVER learning?

It learns some context

Sanabria *et al* BioRxiv



Fine-tuning tasks to understand genome biology

Sanabria *et al* BioRxiv



What does GROVER learn?

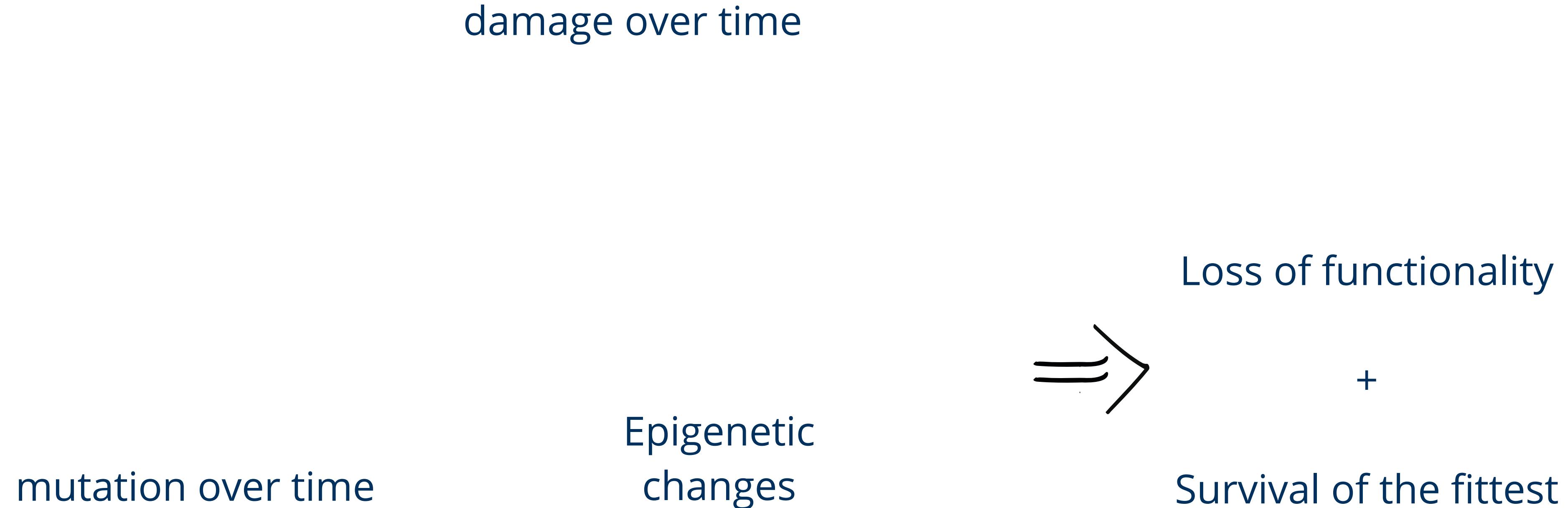
- Self-supervised context for text-completion
- token sequence
- token contexts
- biological functionality of sequence

How can we use it to discover biology?

Sanabria *et al* BioRxiv

Is Genome Instability is encoded in the DNA sequence?

How genomes age



Mario Aguilar



Fine-tuning tasks to predict the frequency of DNA double strand breaks



Fine-tuning tasks to predict the frequency of DNA double strand breaks

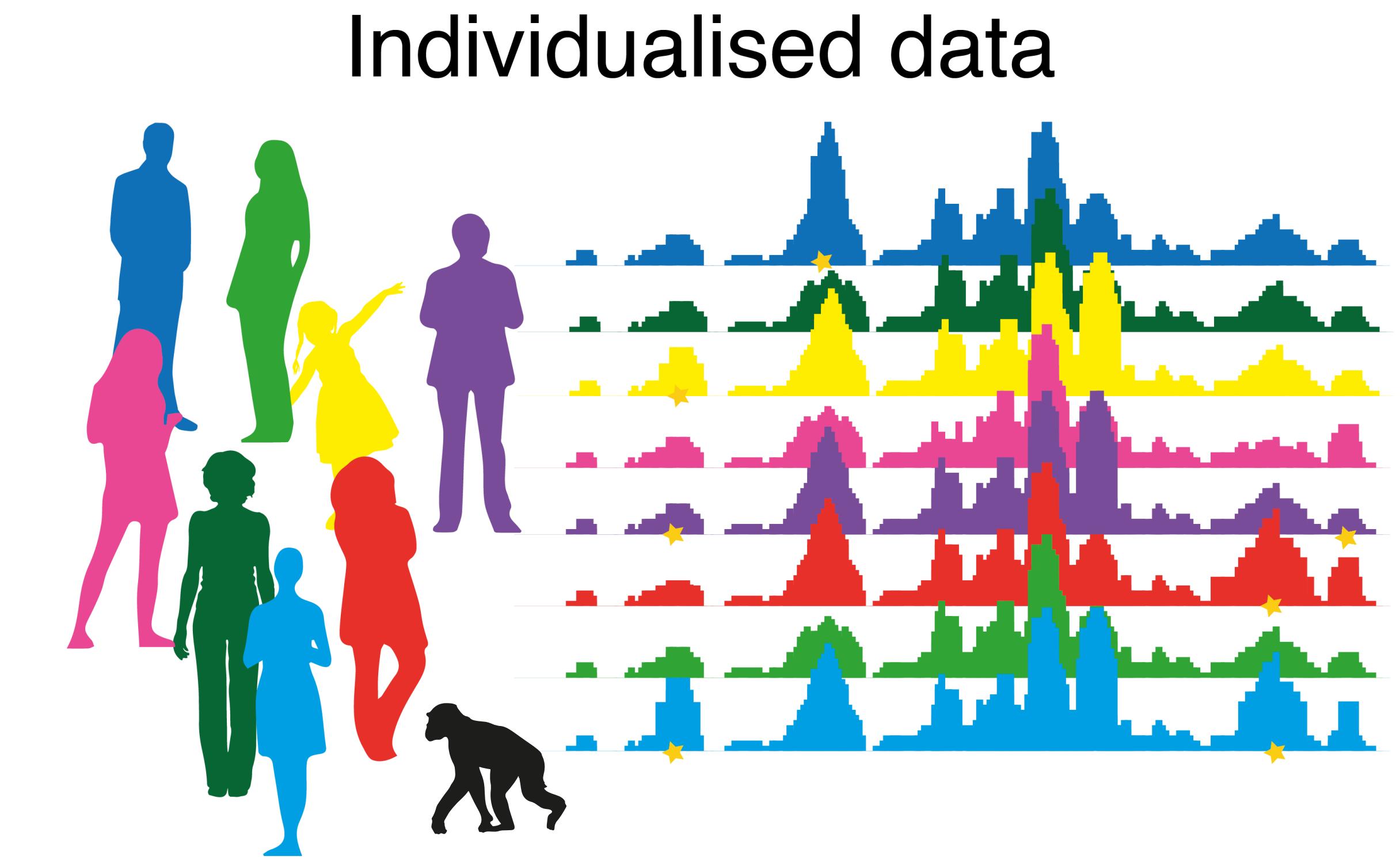
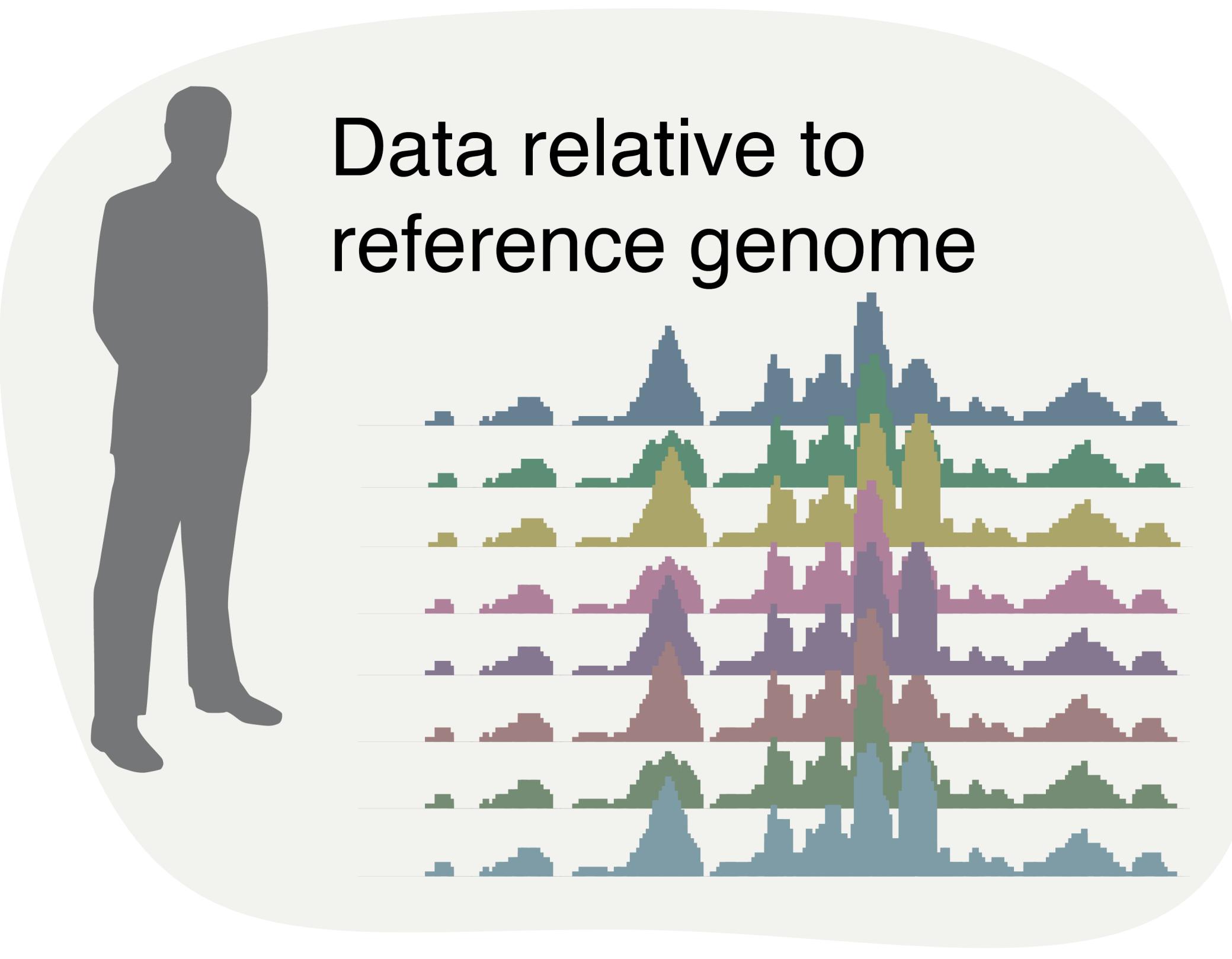


Fine-tuning tasks to predict the frequency of DNA double strand breaks

Summary

- DNA language models learn sequence content and context
- They can be implemented with two strategies:
 - self-supervised pretraining and fine tuning
 - direct learning of the question at hand
- Combining models can help to differentiate information content between sequence only and epigenetic information
- GROVER can help with a variety of questions in genomics

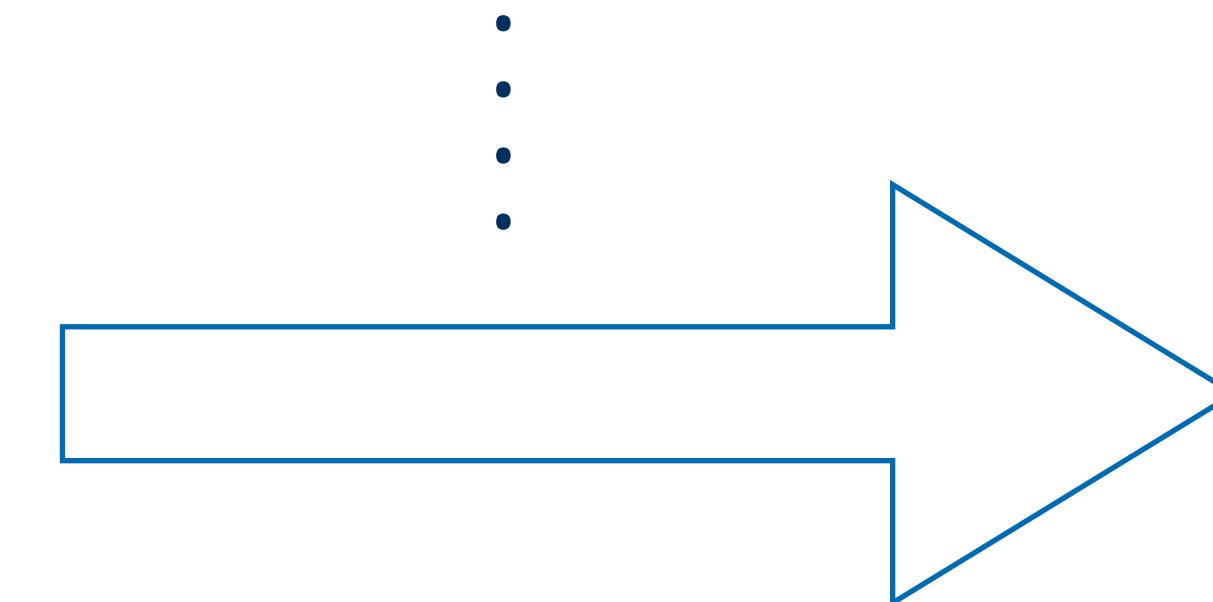
Translation - individualisation of omics data



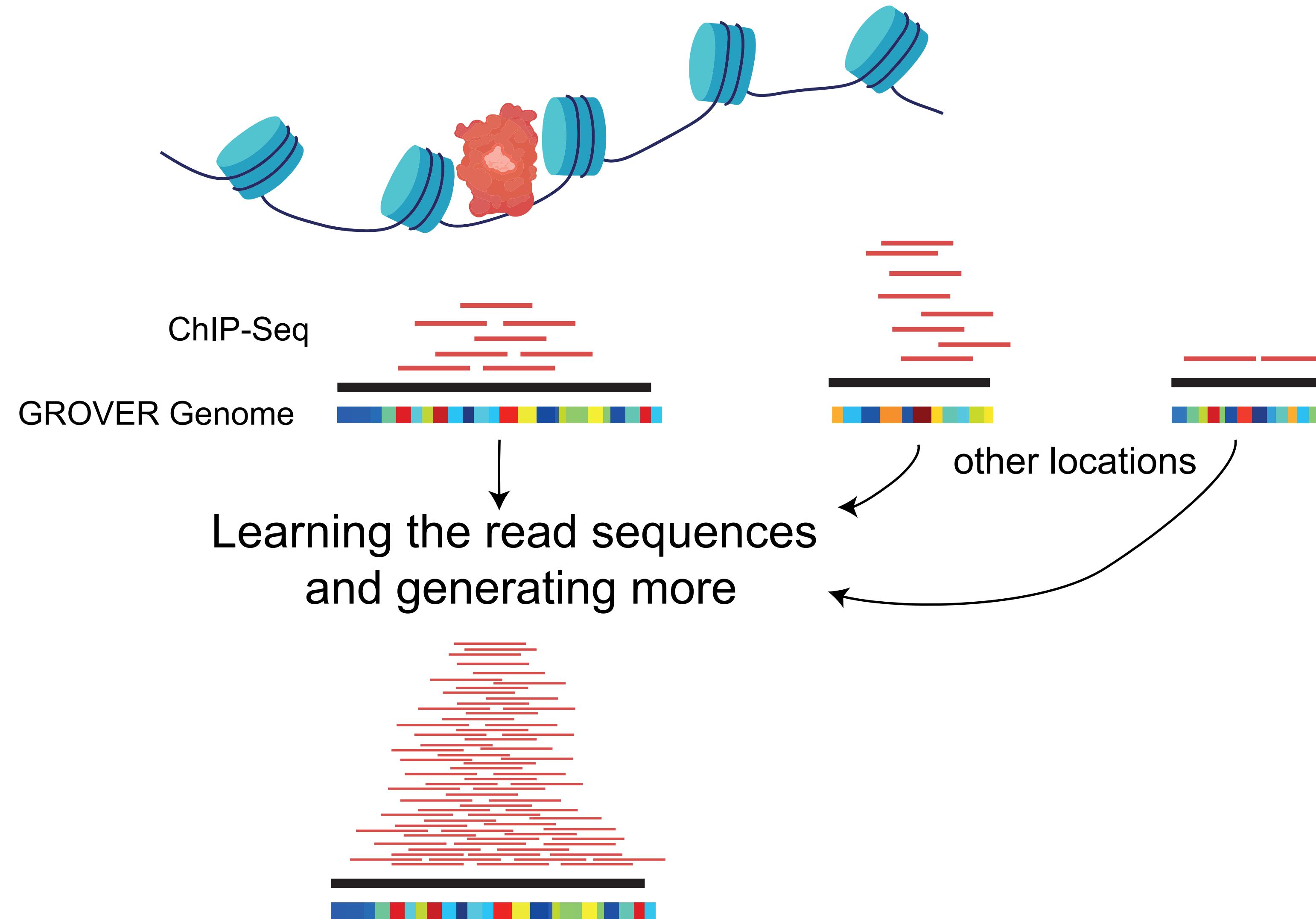
Translation AI - translation of genomics data across species



annotation
regulatory genomics data



Generative AI - Boosting of genomics data



TextToImage - GenotypeToPhenotype

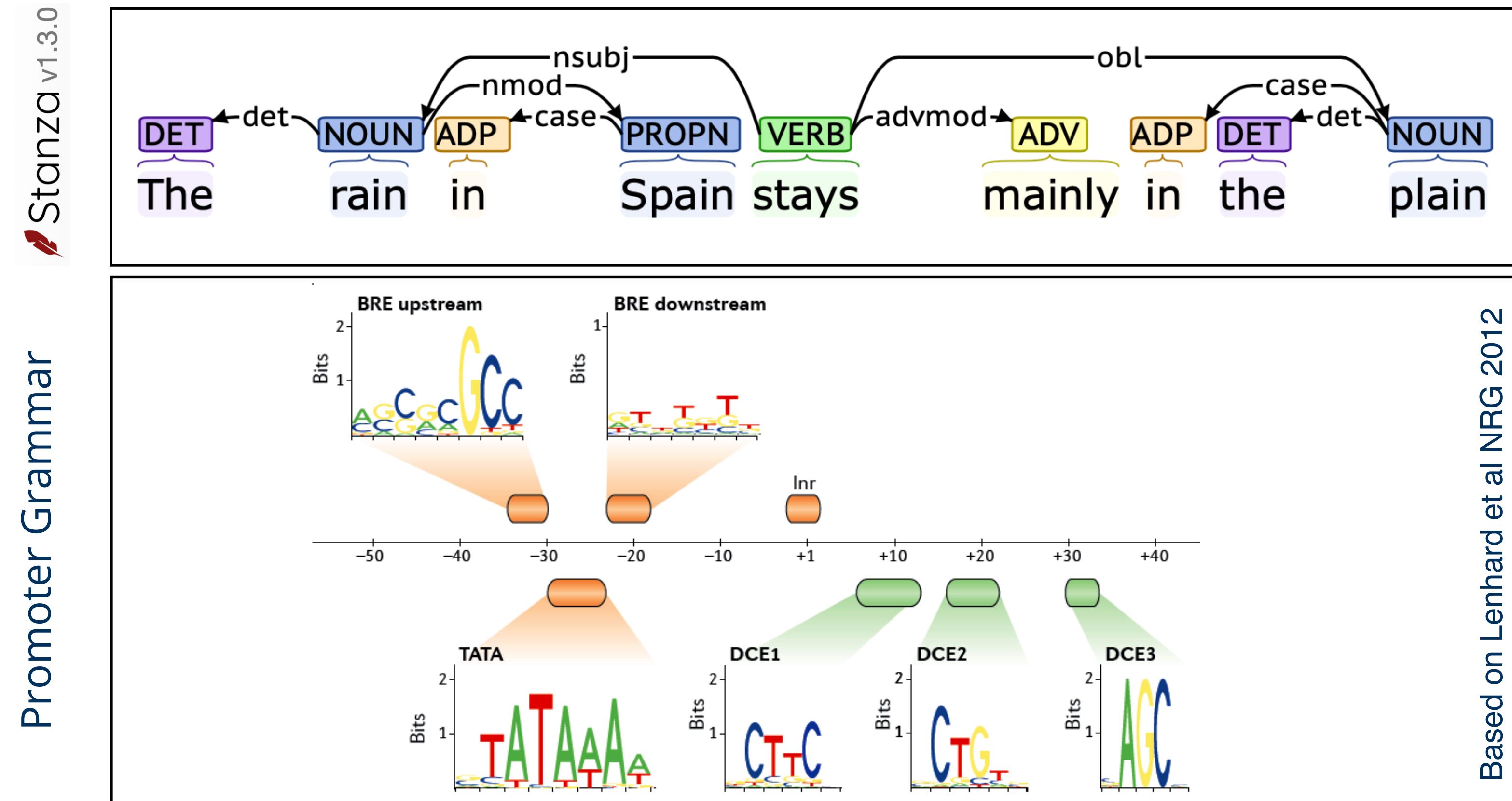
“DallE2, please give me an illustration of damaged DNA in comic style”:



“...ACGCGTAAAATCGATTAGCGATTGCAA...”:

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

Computational linguistics - Genome linguistics



GROVER tutorial

- predicting the use of CTCF binding sites