



# PROTEIN EMBEDDING MODELS: LEARNING THE LANGUAGE OF PROTEINS

# Overview of this workshop

/HDMCYIRCAHTSIRDHDRDGGIN**INTRODUCTION**EYFFNKTTDMEFPWYAPTFTRLESTLPFLVAFWNP  
D  
V  
W

PNFISGDQWAVCGHFSPQPFFM**HISTORY**FTLGMGEEKIDDG  
M  
L  
P

DMW**STATE-OF-THE-ART**VWEHGPNFISGDQWYRSVHRGQDQHRHMIY  
P  
H

IRCAHTSIRDHDRDGGESIDNFTLPFLVAFWNPQTWLFKY**OUTLOOK**LHEIADCYWRVLDIWH  
W  
N

RYASEYFFNKTTD**PRACTICE**TFTRLESQDQWAVCGHFSPQPFFMIRIFHQWCFTLTNGTKFAFFVFFYVWEHGPM

# The rapid progress in natural language processing

Natural language processing: the field of computer science that is concerned with automated text and language analysis.



Thought experiment generation  
10 Thought experiments



Humor  
500+ Openers for Tinder wri...



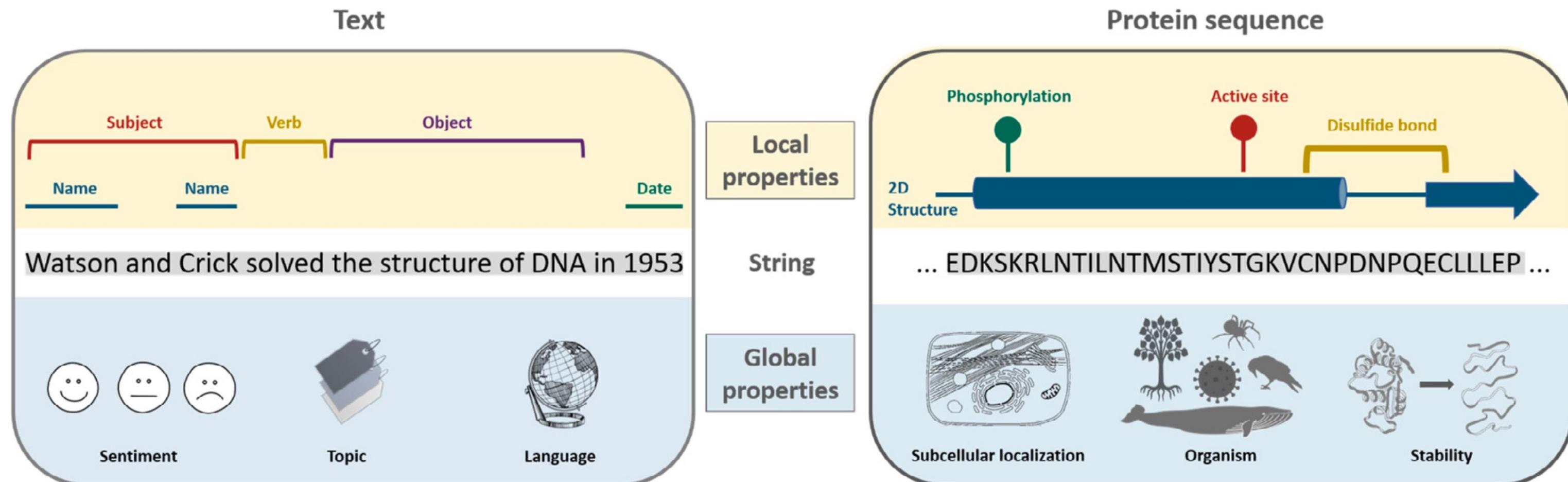
Deepfakes  
AI Eminem  
[View details](#)



Image Generation  
DALL·E by OpenAI

OpenAI's GPT-3 examples

# Proteins are similar to natural language...



## **... but not entirely similar**

- We cannot ‘read’ proteins
- I love you vs. I loved you
- Protein 3D structure -> distant interactions
- Sequencing bias

# Are proteins learnable?

## The linguistic hypothesis

The space of naturally occurring proteins occupies a learnable manifold. This manifold emerges from evolutionary pressures that heavily encourage the reuse of components at many scales: from short motifs of secondary structure, to entire globular domains.

*Mohammed AlQuraishi*

# Keep it simple: use a bag-of-words

|  | about | bird | heard | is | the | word | you |
|--|-------|------|-------|----|-----|------|-----|
| About the <b>bird</b> , the<br>bird, <b>bird bird bird</b> | 1     | 5    | 0     | 0  | 2   | 0    | 0   |
| You heard about<br>the <b>bird</b>                         | 1     | 1    | 1     | 0  | 1   | 0    | 1   |
| The <b>bird</b> is the<br>word                             | 0     | 1    | 0     | 1  | 2   | 1    | 0   |

# Keep it simple: use a bag-of-words

|                   | DMC | MCY | YIR | HRG | SVH | QDQ | CAH |
|-------------------|-----|-----|-----|-----|-----|-----|-----|
| DMCYIRCAHTSIRD... | 1   | 1   | 1   | 0   | 0   | 0   | 0   |
| WYRSVHRGQDQH...   | 0   | 0   | 0   | 1   | 1   | 1   | 1   |

- Different k-mer lengths
- Overlapping vs. non-overlapping
- Binary vs. count vs. frequency

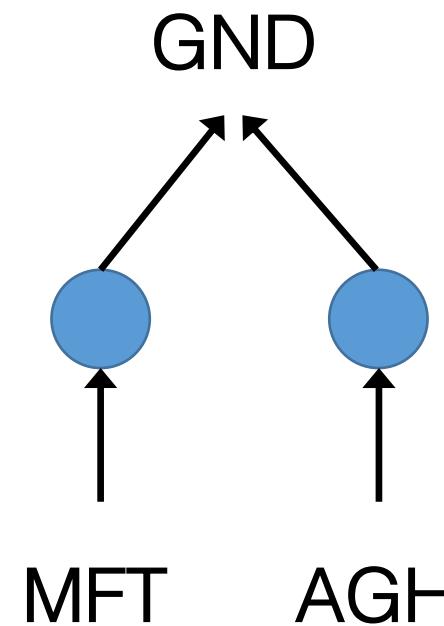
# Embedding methods for protein language

Word2Vec

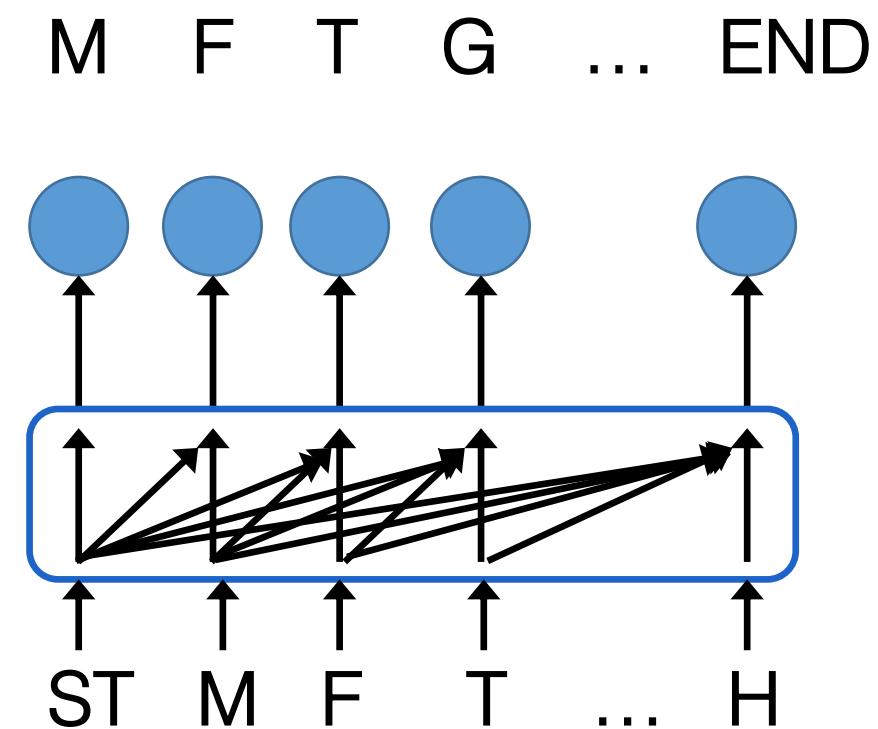


# Embedding methods for protein language

Word2Vec

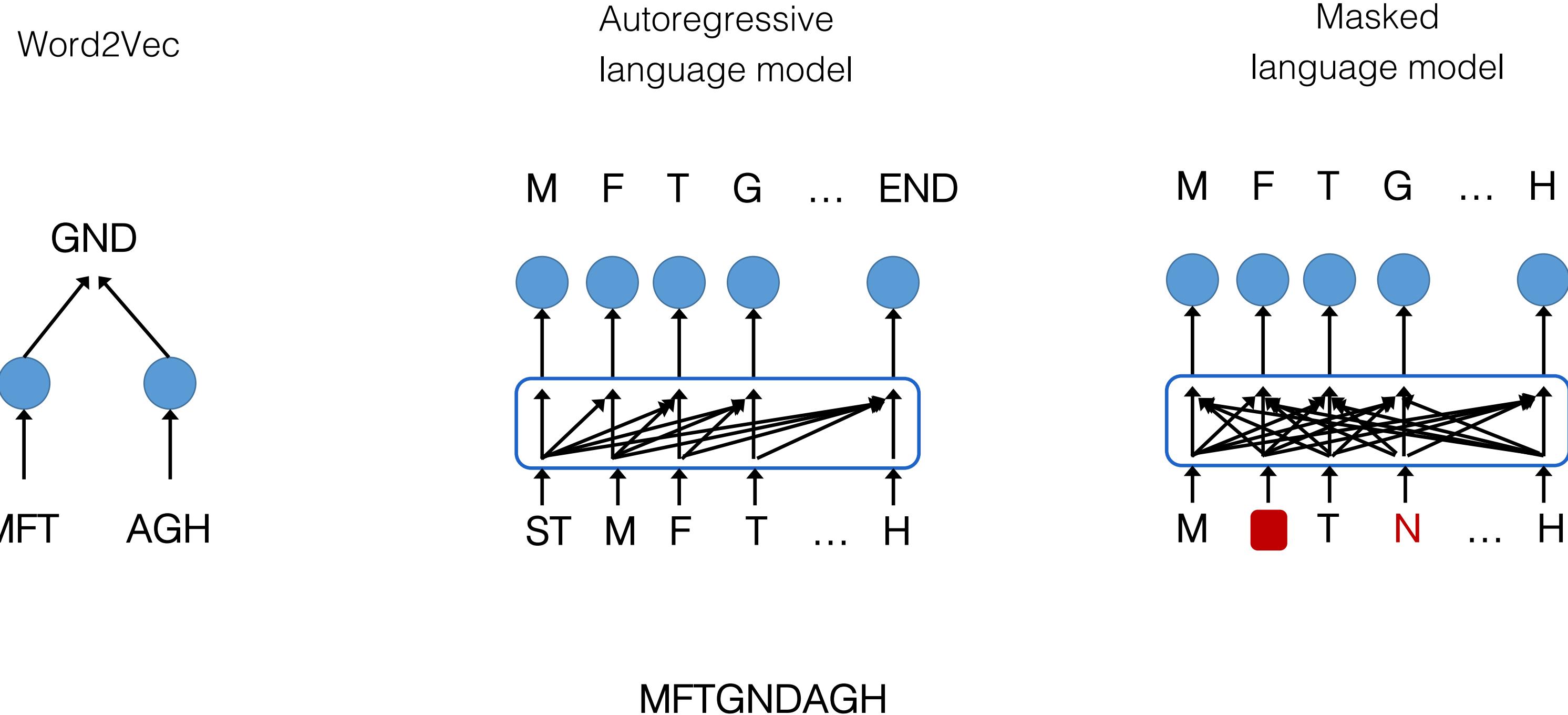


Autoregressive  
language model



MFTGNDAGH

# Embedding methods for protein language

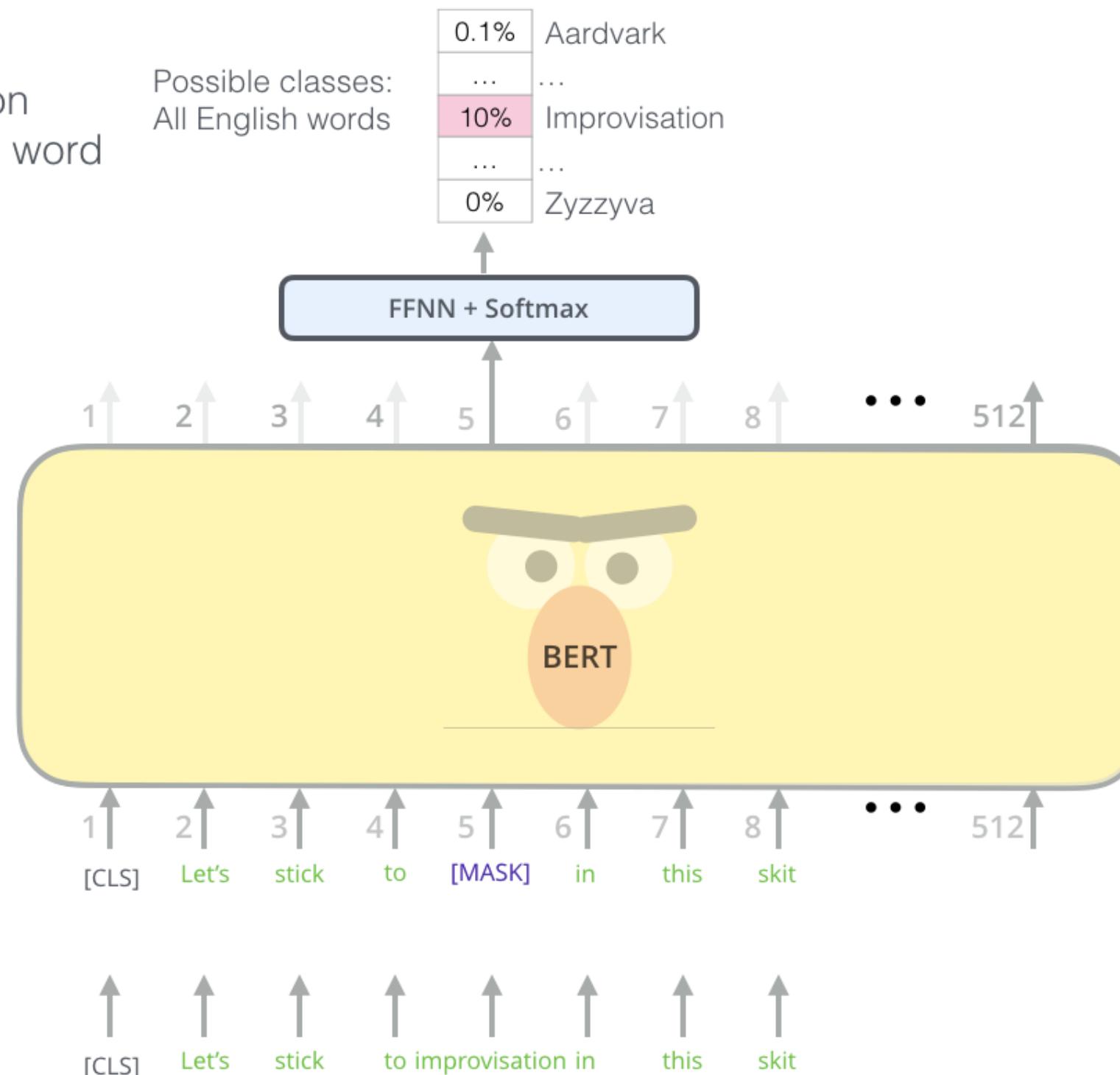


# BERT model training

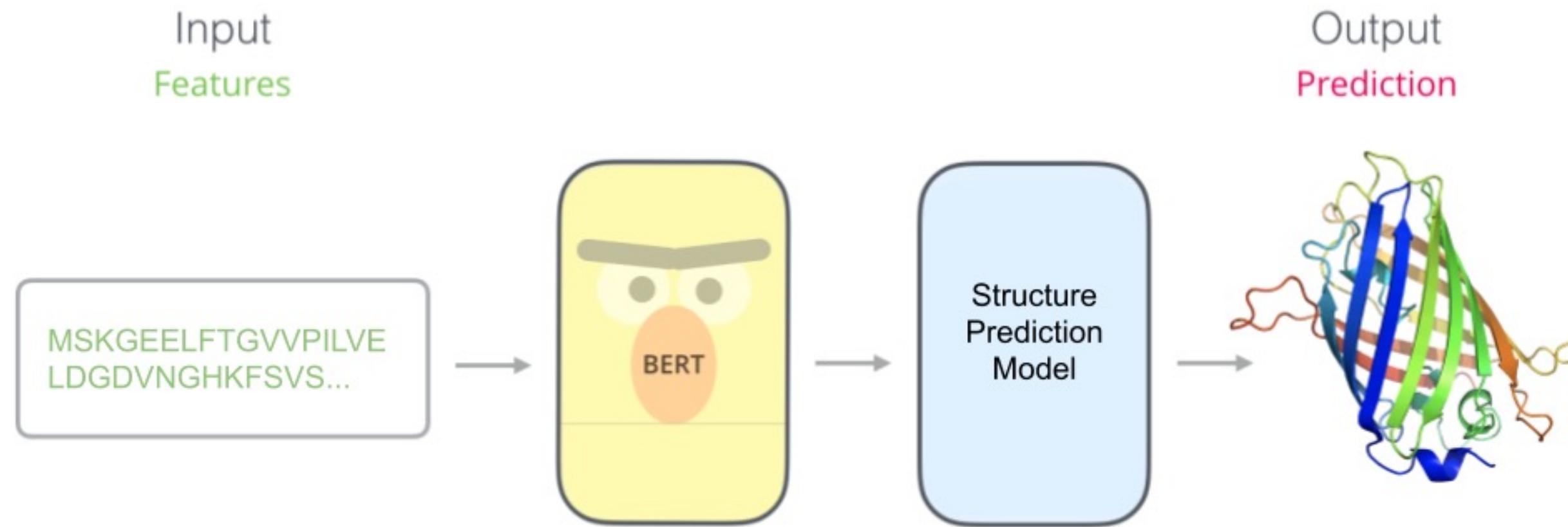
Use the output of the masked word's position to predict the masked word

Randomly mask 15% of tokens

Input



# Stacking supervised prediction models

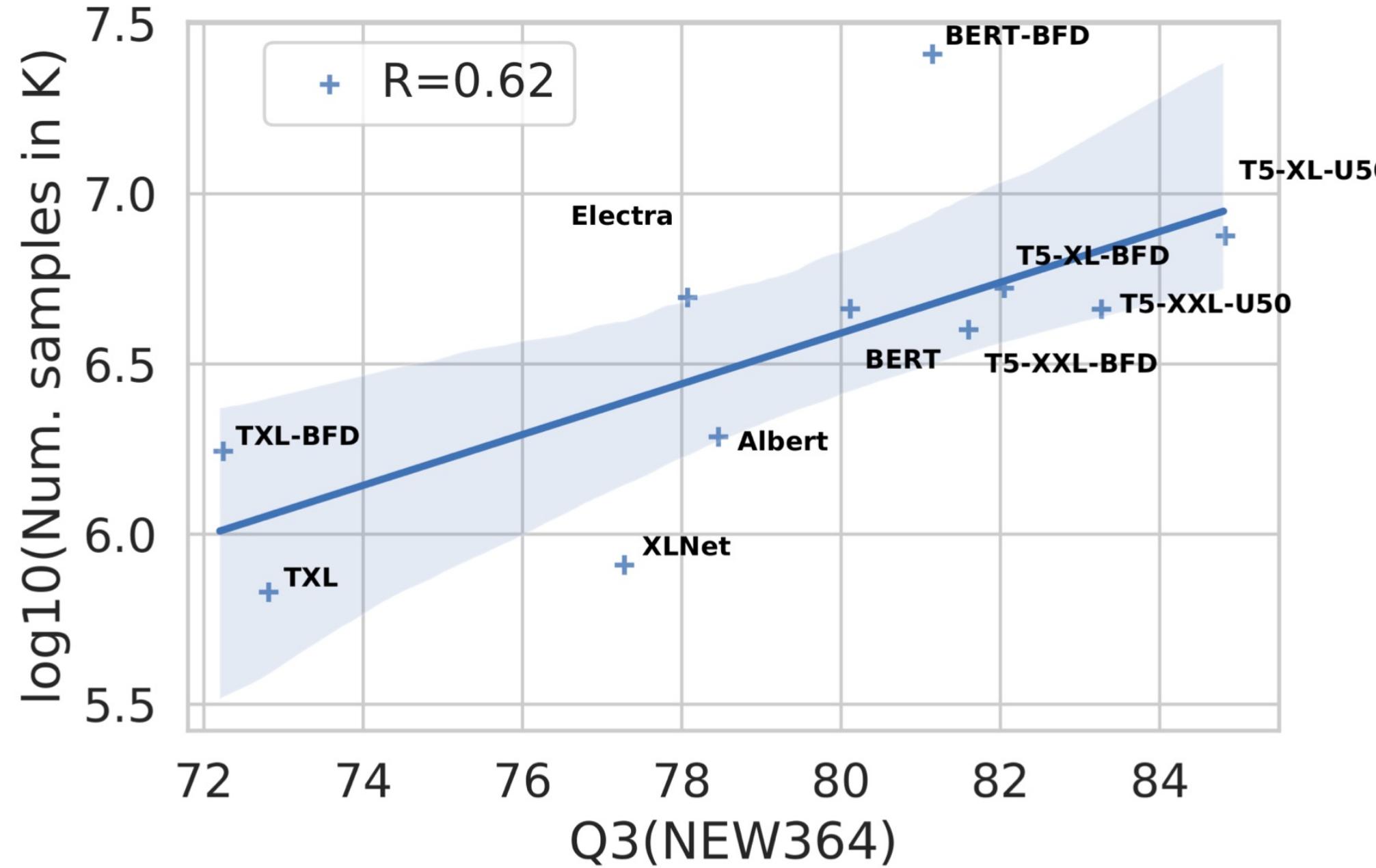


# Language models are data hungry

| <i>Data LM</i>                      | <i>UniRef50</i> | <i>UniRef100</i> | <i>BFD</i> |
|-------------------------------------|-----------------|------------------|------------|
| <i>Number proteins [in m]</i>       | 45              | 216              | 2,122      |
| <i>Number of amino acids [in b]</i> | 14              | 88               | 393        |
| <i>Disk space [in GB]</i>           | 26              | 150              | 572        |

TABLE 1: Data Protein LM - UniRef50 and UniRef100 cluster the UniProt database at 50% and 100% pairwise sequence identity (100% implying that duplicates are removed) [41]; BFD combines UniProt with metagenomic data keeping only one copy for duplicates [24], [42]. Units: number of proteins in millions (m), of amino acids in billions (b), and of disk space in GB (uncompressed storage as text).

# Language models are data hungry



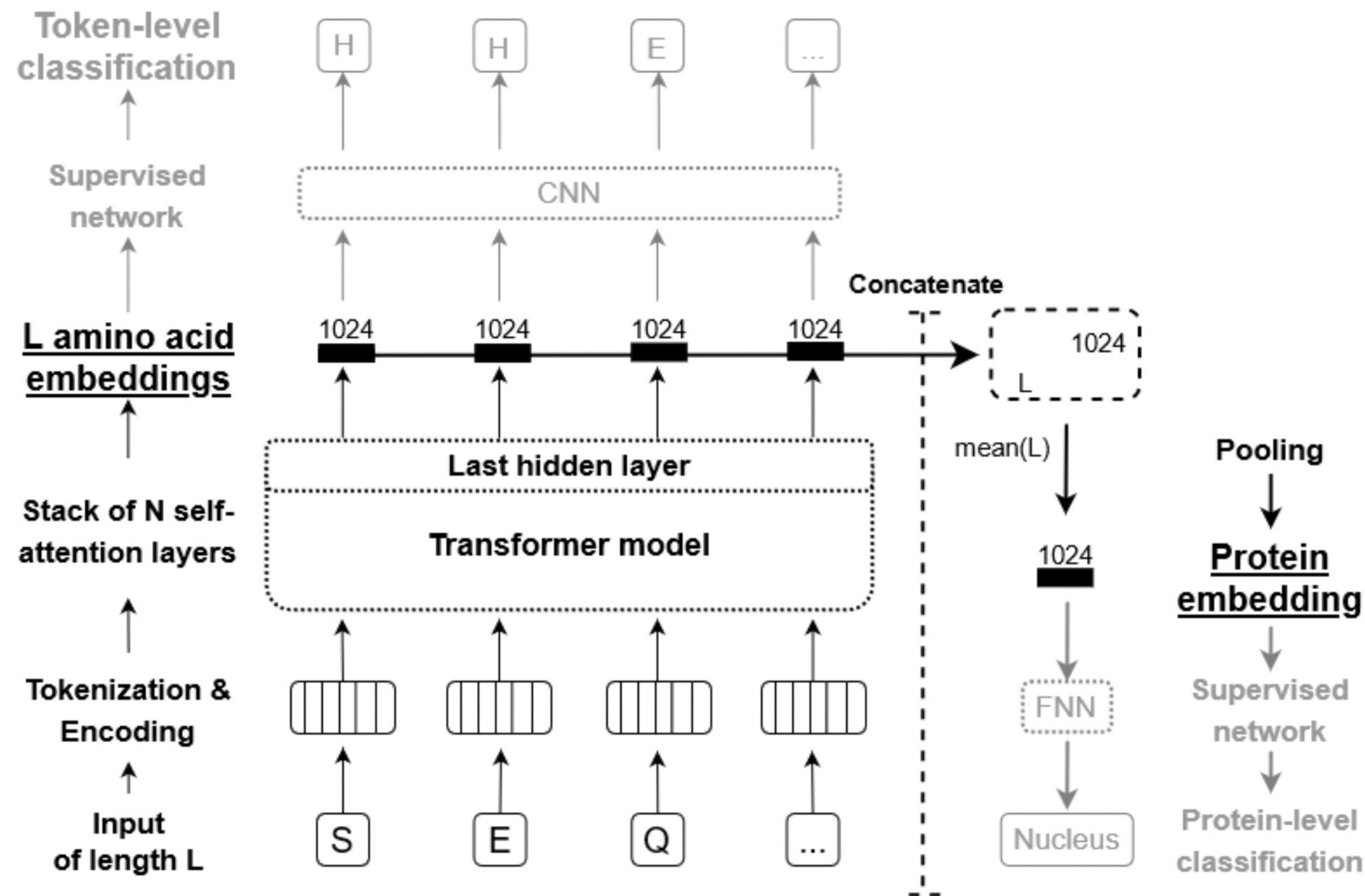
More data during pre-training seems better for downstream performance.

# Redundancy and noise

It's all about abundant and high-quality data!

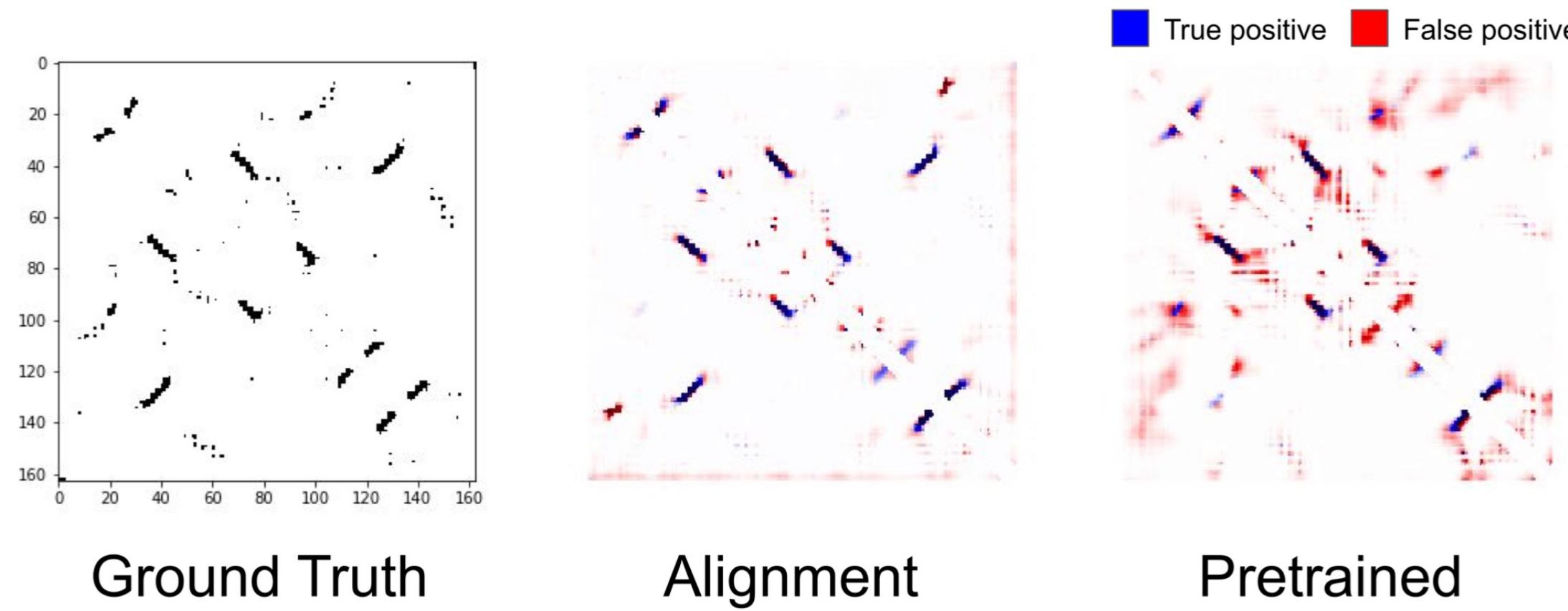
“Less noisy and less redundant corpora (e.g. UniRef50) improved over larger but more noisy and redundant corpora (e.g. BFD).” (Elnaggar *et al.*, 2021)

# Recent applications



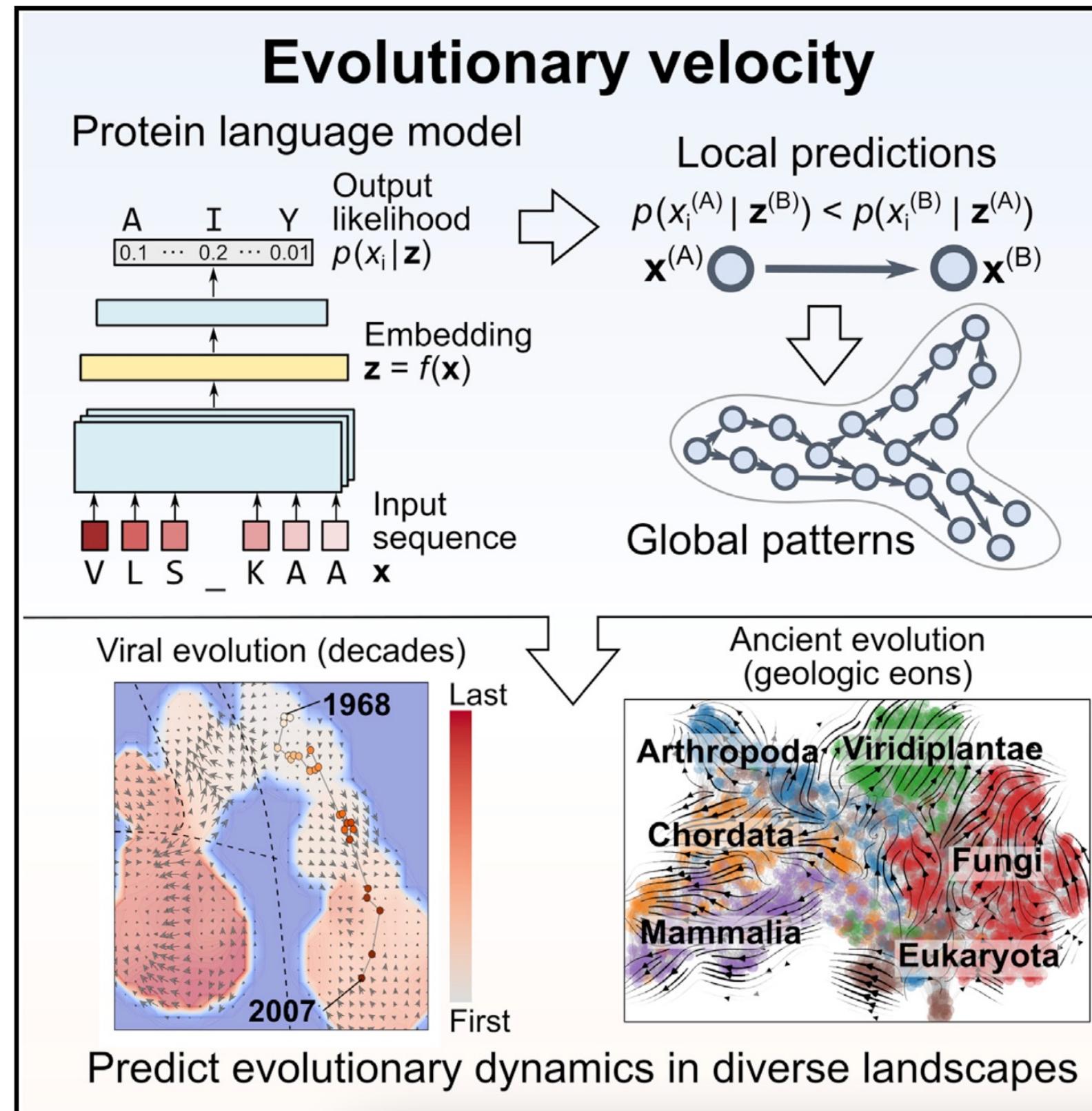
# Recent applications

Simpler methods using evolutionary information can (still) outperform DL methods.



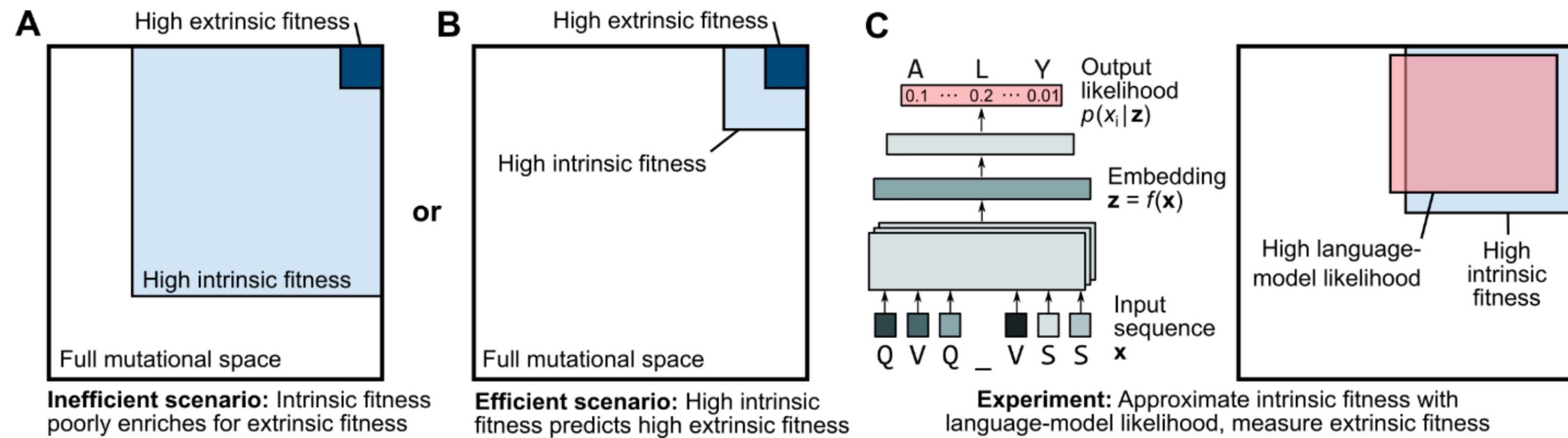
*The true contact map for the same protein as above (left), predicted contacts from a model with alignment-based features (center), predictions from a **pretrained** LSTM (right).*

# Recent applications



Protein language models can predict mutational effects across diverse proteins.

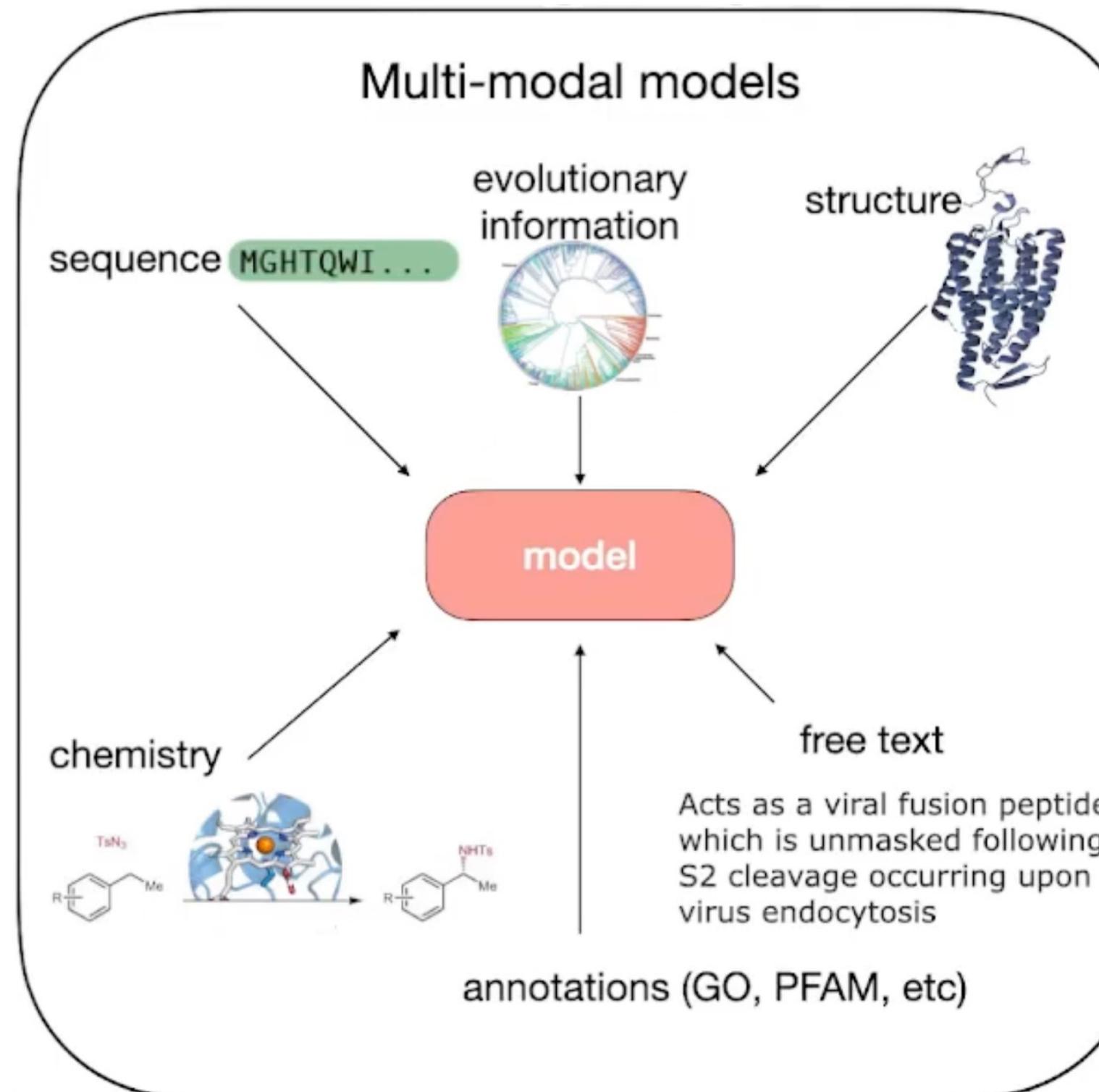
# Recent applications



# Outlook: trends and challenges

Big data needs large models need loads of computational resources.

# Outlook: trends and challenges



Kevin Yang 楊凱笙  
@KevinKaichuang

Senior Researcher in computational biology @MSFTResearch (@MSRNE). He/him.

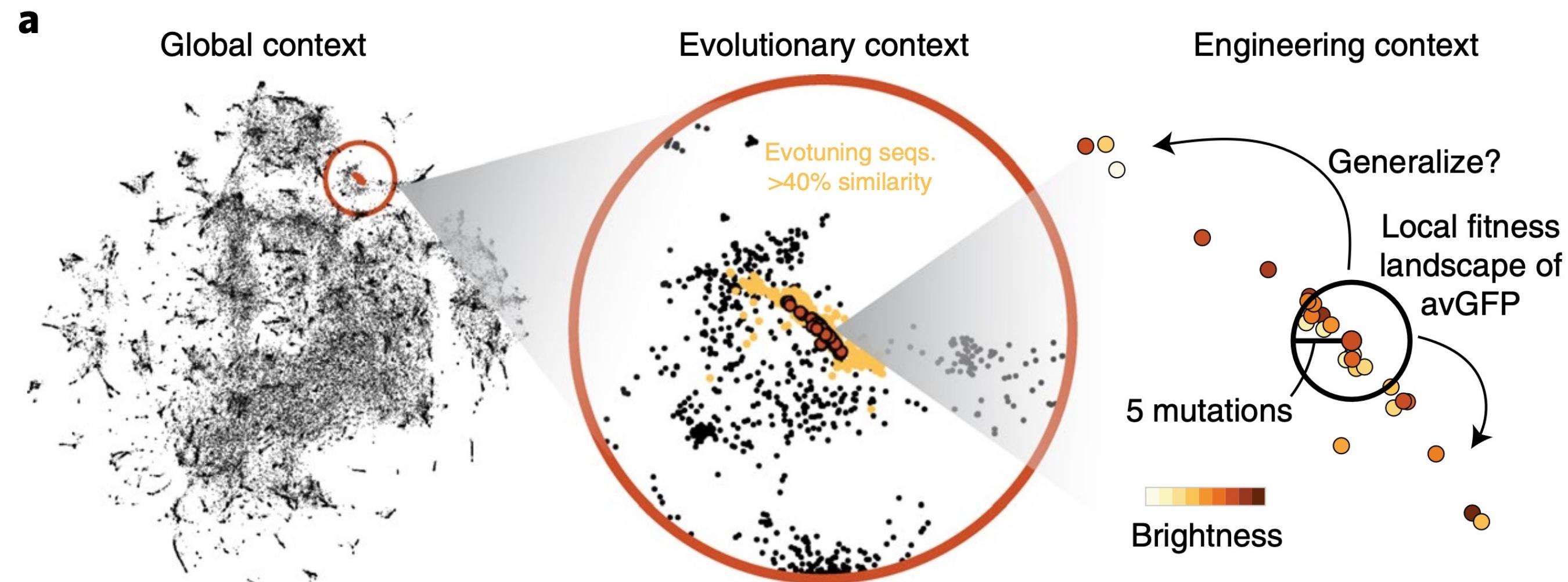
📍 Massachusetts, USA ↗ yangkky.github.io

👤 Lid geworden in november 2016

[https://www.youtube.com/watch?v=0I2G9\\_LrY\\_c](https://www.youtube.com/watch?v=0I2G9_LrY_c)

# Outlook: trends and challenges

Actual impact in biological applications

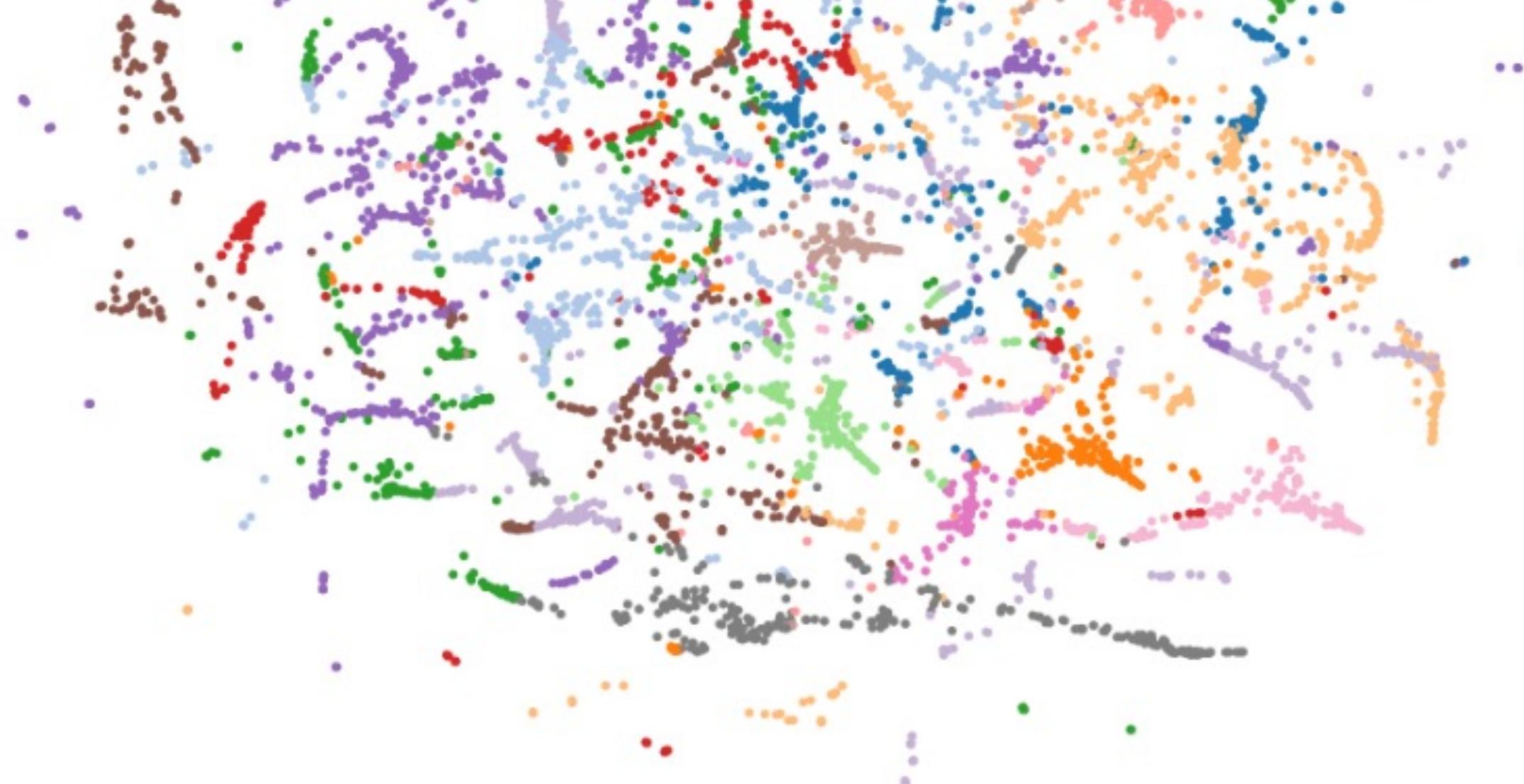


# Open questions

Foundation models for biology?

“Bio” DeepMind or OpenAI company?

Deep fake content, deep fake life forms?



# Questions?



@dimiboeckaerts



Dimi Boeckaerts

# Embeddings in practice!

The screenshot shows a GitHub repository page for 'dimiboeckaerts / ProteinLanguageWorkshop'. The repository is private. The main navigation bar includes links for Pull requests, Issues, Marketplace, and Explore. On the right, there are buttons for Unwatch (1), Fork (0), and Star (0). Below the navigation, there are tabs for Code, Issues, Pull requests, Actions, Projects, Security, Insights, and Settings. The 'Code' tab is selected.

Key statistics shown: main branch, 1 branch, 0 tags. There are 5 commits in the main branch. The commits listed are:

- dimiboeckaerts Update README.md (a21d9d0, 4 days ago)
- .gitignore (Initial commit, 12 days ago)
- LICENSE (Initial commit, 12 days ago)
- README.md (Update README.md, 4 days ago)

A large preview window displays the contents of the README.md file, which includes:

## A workshop on protein language

Building on top of the successes of word embeddings and transformer models for language, increasingly more of these architectures are now being used to learn the 'language of proteins'. This workshop introduces you, both theoretically and practically, to this latest trend in protein sequence analysis and feature engineering.

### Theory slides

The slides cover the following topics:

**About**

An introductory workshop to protein language models

**Code**

- Readme
- MIT license
- 0 stars
- 1 watching
- 0 forks

**Releases**

No releases published  
Create a new release

**Packages**

No packages published  
Publish your first package

# GPT-3 thought experiments

## The Omelet Parable

Imagine that you are on a train. You have always wanted to make the perfect omelet. A genie appears and offers you two options: you can have a delicious omelet, or you can have a device that will let you make unlimited omelets for the rest of your life.



Would you rather have a delicious omelet?  
Or the device?

## The Red Pill/Blue Pill

Imagine that there is a drug that will cure all diseases. It also has a side effect of making people feel like they are living in a computer simulation. They are in a simulation, and after they take the pill they can see the pixels and their world is filled with triangles.



Would you take it?