



# ProtGPT2

Dominik Grabarczyk

# Parallels between Natural language and Proteins

## **NATURAL LANGUAGES**

Letters  
↓  
Words  
↓  
Sentences  
↓  
Meaning

## **PROTEINS**

Amino Acids  
↓  
Domains  
↓  
Proteins  
↓  
Function

# Parallels between Natural language and Proteins

## NATURAL LANGUAGES

A, B, C, D, E



I, To Be, Student



I am a student



*Description of occupation*

## PROTEINS

g, a, v, l, l, t, s



lqvgqvelgg, agslqp



lqvgqvelgggpgagslqpl



*Regulation of carbohydrate metabolism*



Question



CAN NLP DO THIS?

# What is ProtGPT2

- Protein LLM
- Transformer decoder
- Using many common practices from NLP
  - Decoding strategies
  - Tokenisation

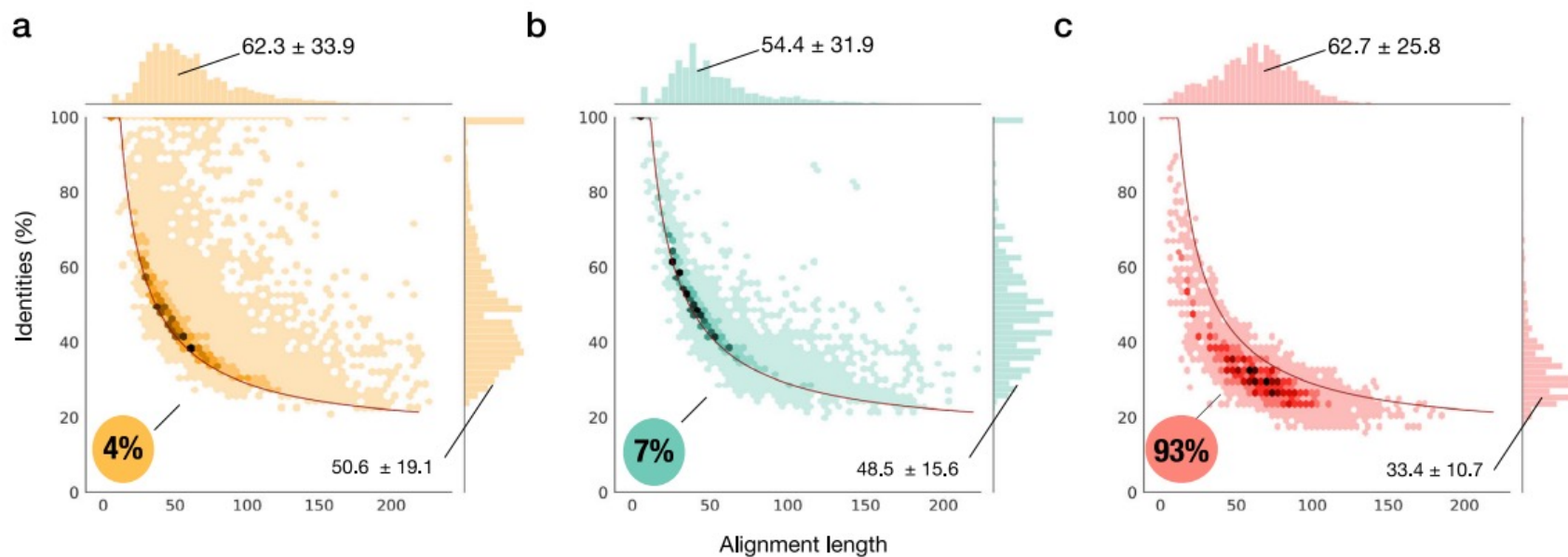
# Did they succeed

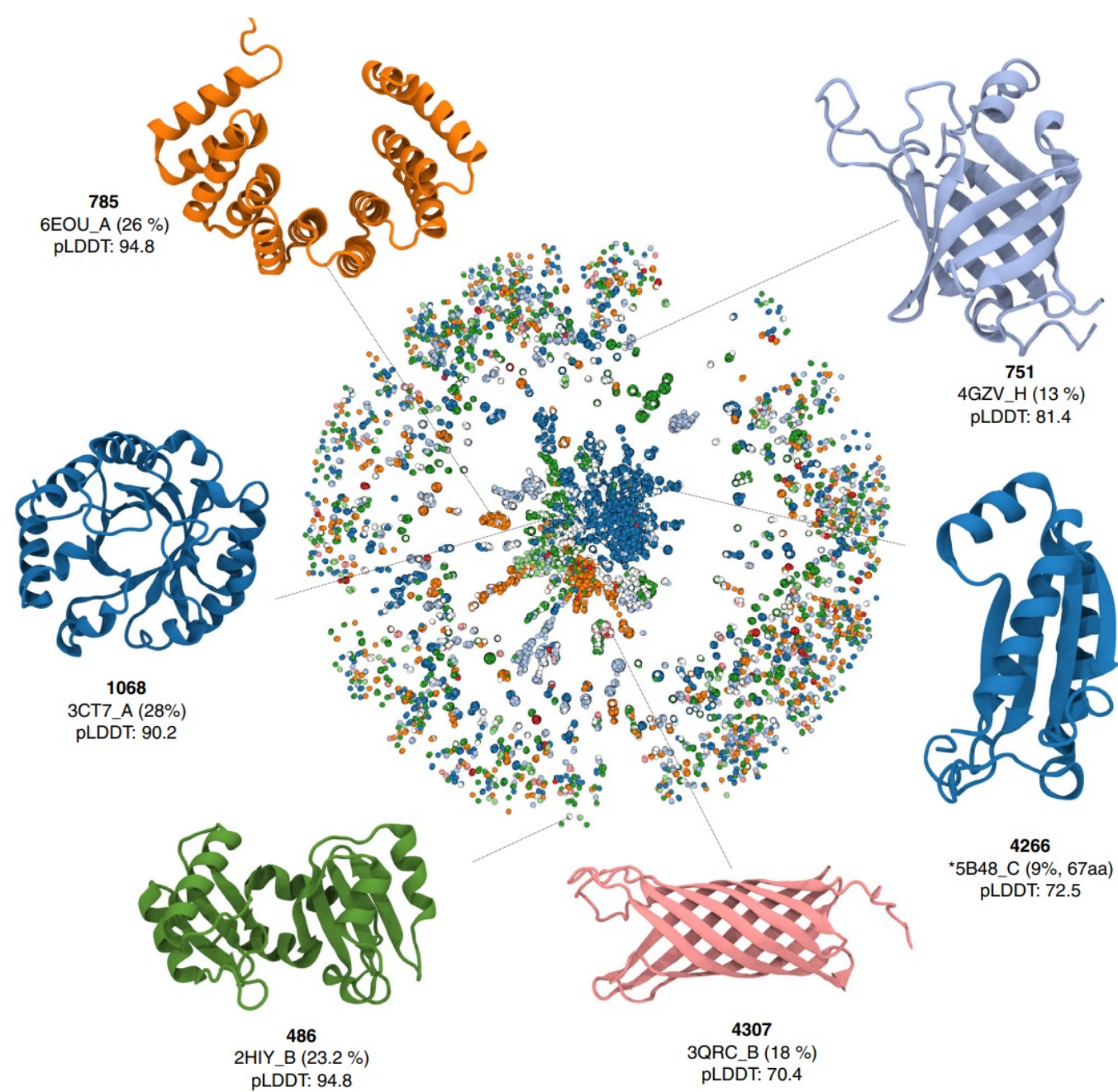
- Depends on your criteria?
  - No “unified” benchmark
  - No downstream tasks
  - BUT results are interesting

**Table 1 | Disorder and secondary structure predictions of the natural and ProtGPT2 dataset**

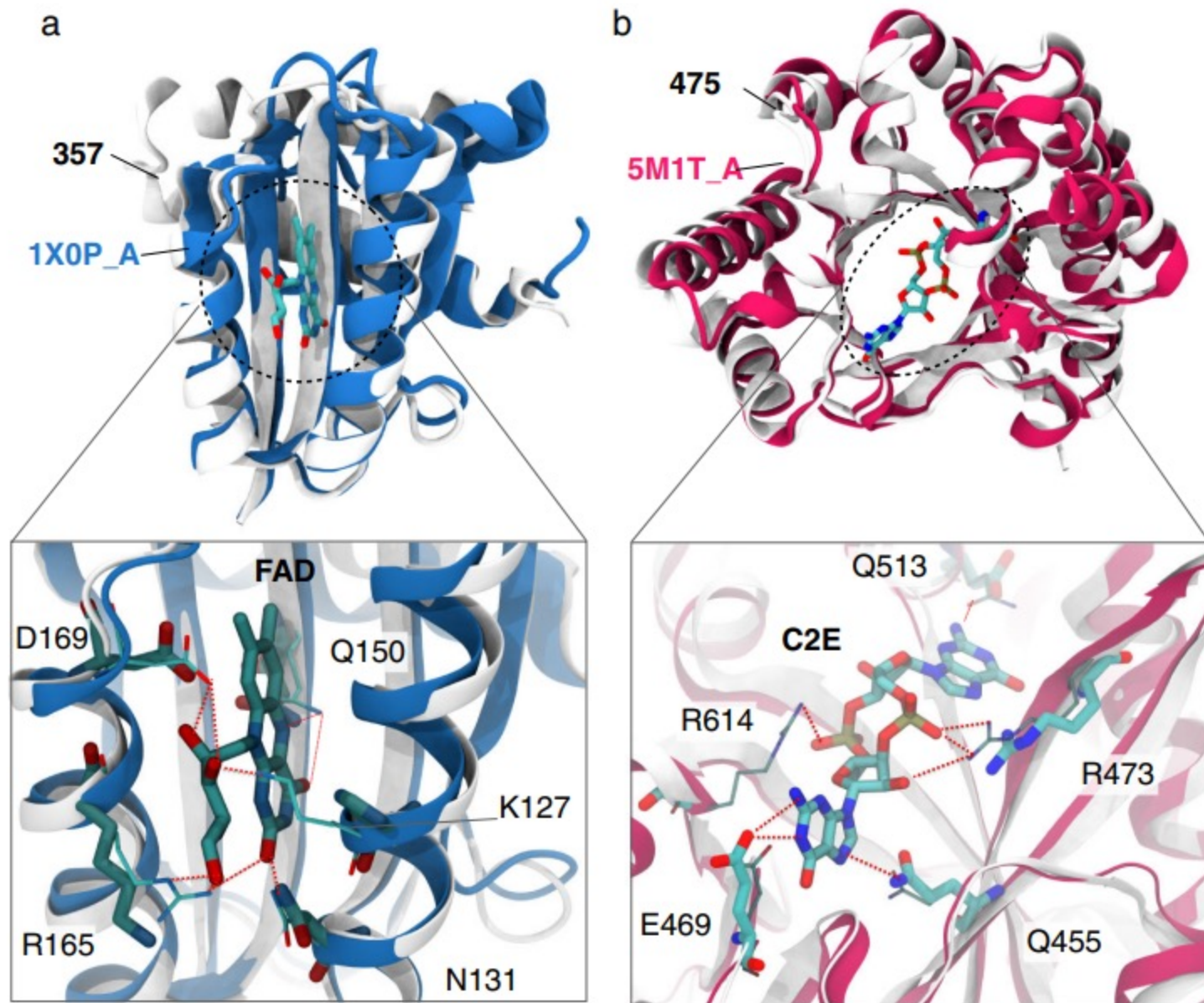
	Natural dataset	ProtGPT2 dataset
<b>IUPred3 (globular domains)</b>	88.40%	87.59%
<b>Ordered content</b>	79.71%	82.59%
<b>Alpha-helical content</b>	45.19%	48.64%
<b>Beta-sheet content</b>	41.87%	39.70%
<b>Coil content</b>	12.93%	11.66%

( $n = 10,000$  independent sequences/dataset).







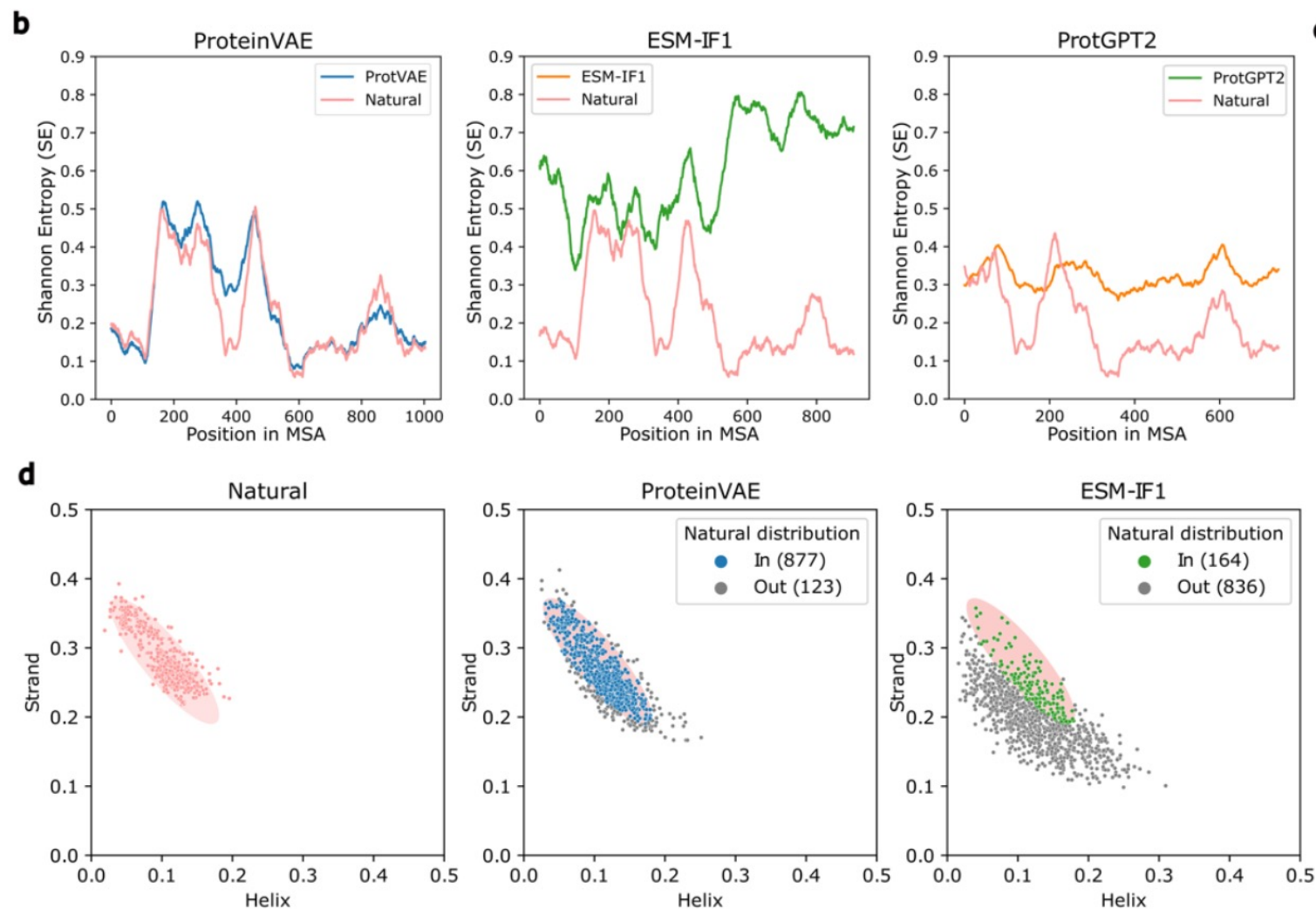


# What's the point

- Therapeutics development
  - Vaccines
  - Personalised medicine
- Structure and function prediction
- Genetic modification

# What's next?

- Different architectures
- Robust and standardised validation methods
- Development of real-world applications



Lyu et al. 2023



Thank you