

# LMM for Gene Repression Correlation

Data

- Rabani\_Expression\_A+

id	1h	2h	3h	4h	5h	6h	7h	8h	10h
S0_M_T1	1.1502	1.1256	1.4005	0.2332	0.73195	-0.47038	-0.57411	-0.25983	-0.76564

- Mean RBP Data

Motif	Mean_RBP
TGTGGAT	7.17826412

# Theory

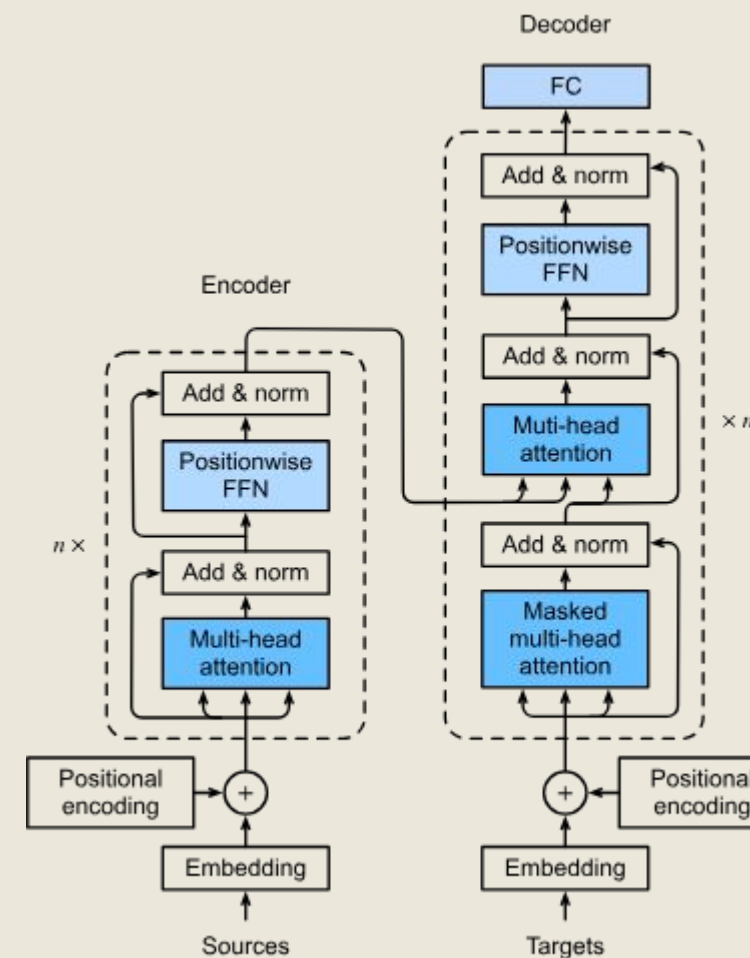
- Given the 'Motif Composition' of a sequence, we can predict its Repression Rate
- **Motif Composition**: Motifs that make up a sequence; motifs with higher Mean RBP have higher significance
- We can measure **correlation** using LLM

GGGTCCCCCTGACAG

# Methodology

GGGTCCCCCTGACAG → TOKENIZER

[GGG], [TCCCCCT], [GACAG]



Prediction  
(Scalar)

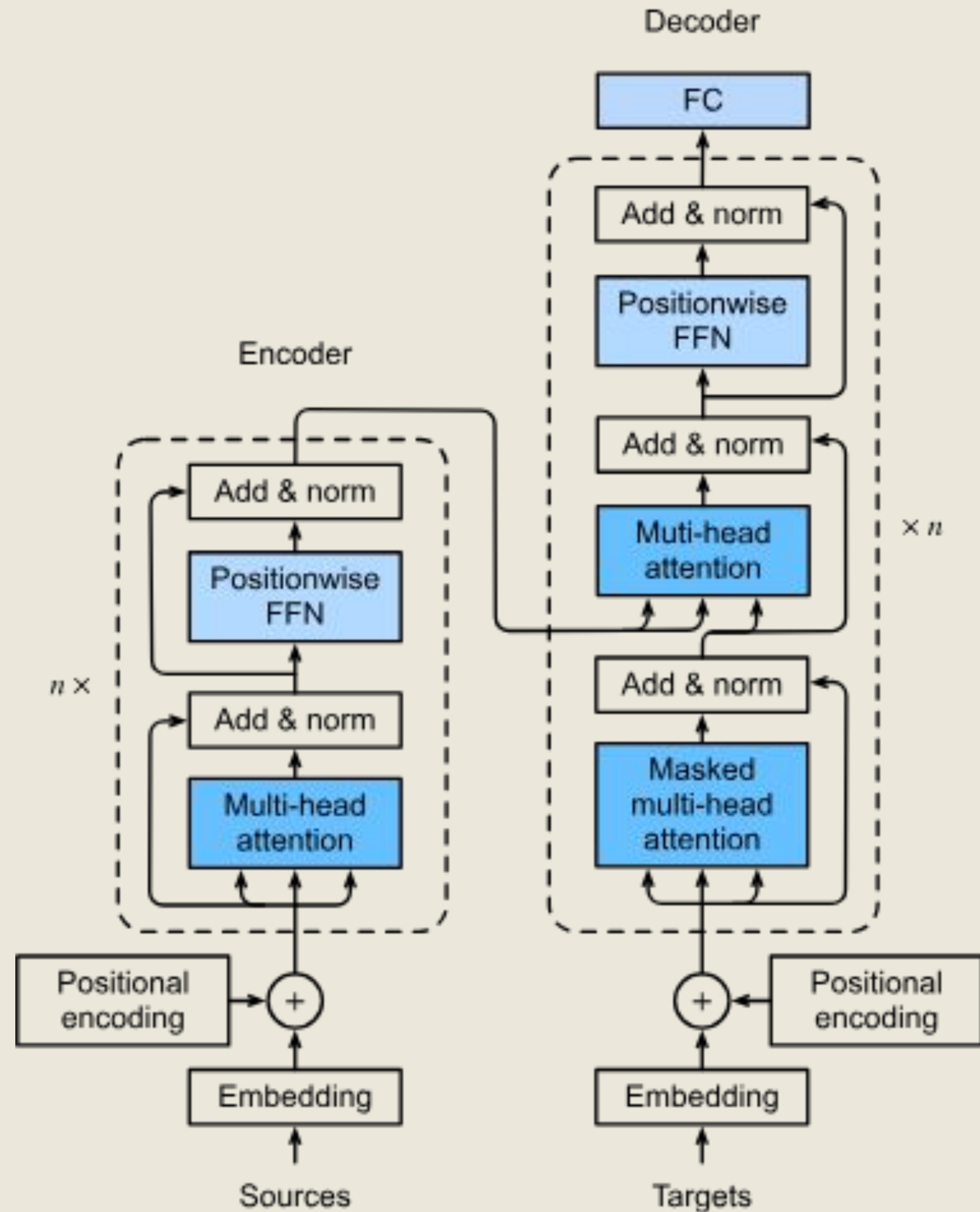
# Large Language Model

Originally made for translating languages.

Takes tokens as numbers and performs matrix operations to achieve output

Has no concept of 'Order' or 'Characters', only tokens

Fine Tuning: Utilizing pretrained model for quick learning



# Experimentation with Transformers

## DNA-MORPH

- BERT-Model from scratch
- Transformer from scratch
  - Ranked motifs by Mean RBP
  - Found ‘words’ by Motif Ranking according to Morphological Segmentation

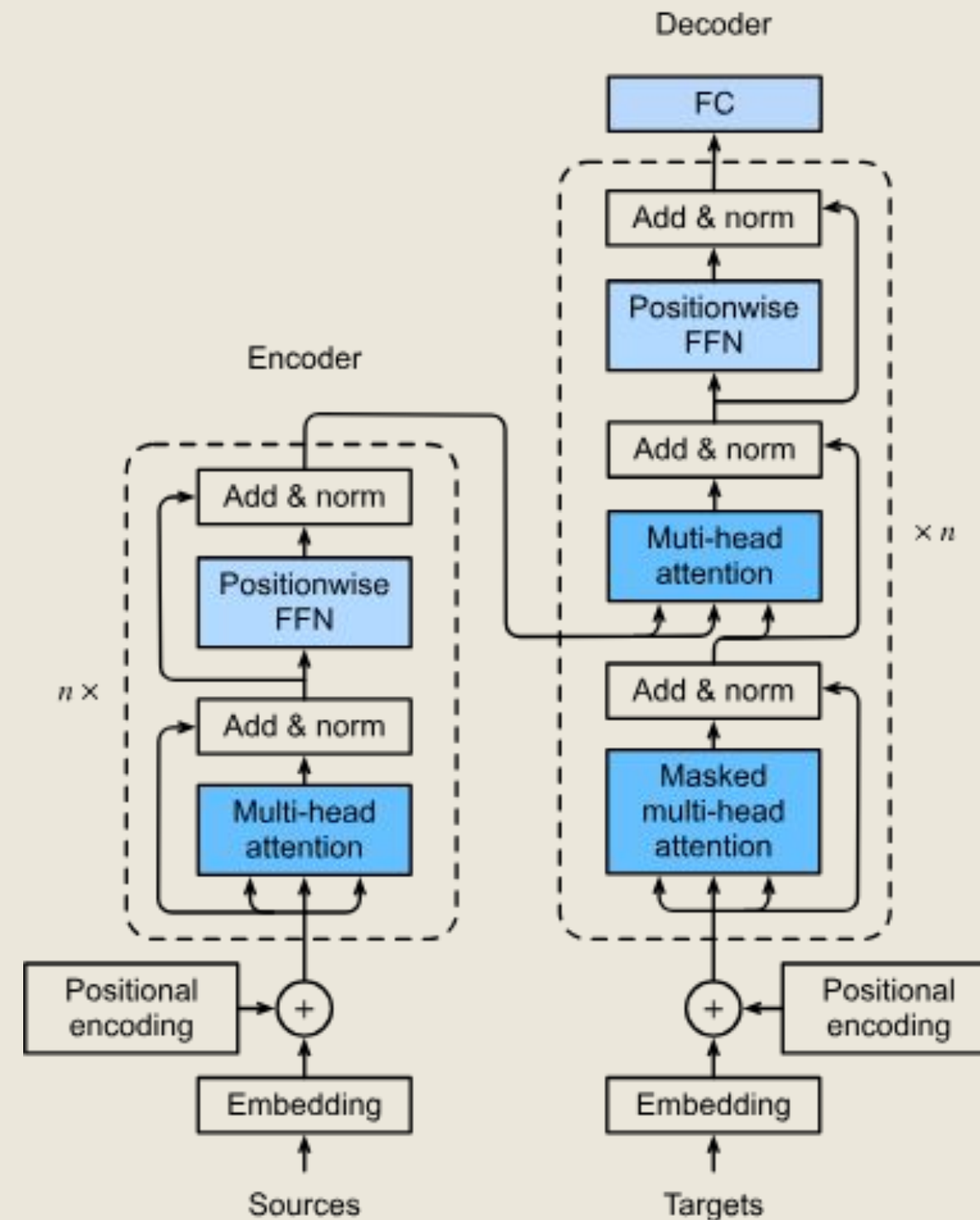
<h3>MORPH_Dmitri</h3>	<h3>MORPH_Owen</h3>
<ul style="list-style-type: none"><li>● Trained on Dmitri Data</li></ul>	<ul style="list-style-type: none"><li>● Trained on Owen Data</li></ul>

# Experimentation with Transformers

Do LLMs Work?

DNABERT2:

- $r^2 = .335$



# Experimentation with Transformers

DNAMORPH has 2 functions:

- Segment Sequences into Motifs:

TGTCCCC: 12.12

GTCCCCG : 1.01

TCCCCGG: 4.43

TGTCCCCGGGTCTT → [TGTCCCC,  
TCCCCGG,...]

- ‘Tokenize’ Motifs :

[TGTCCCC, TCCCCGG,...] → [123233, 13332, ...]

# Experimentation with Transformers

Seq = TGTCCCCGGGTCTTCCAACGGACTGG...

Sliding Window(Seq) = GGGTCTT CACTGGG ACTGGGG...

- MORPH\_O:

['gggtctt', 'cactggg', 'actgggg']

- MORPH\_D

- ['gggtctt', 'cactgg', '##g', 'actgg', '##gg']



# Experimentation

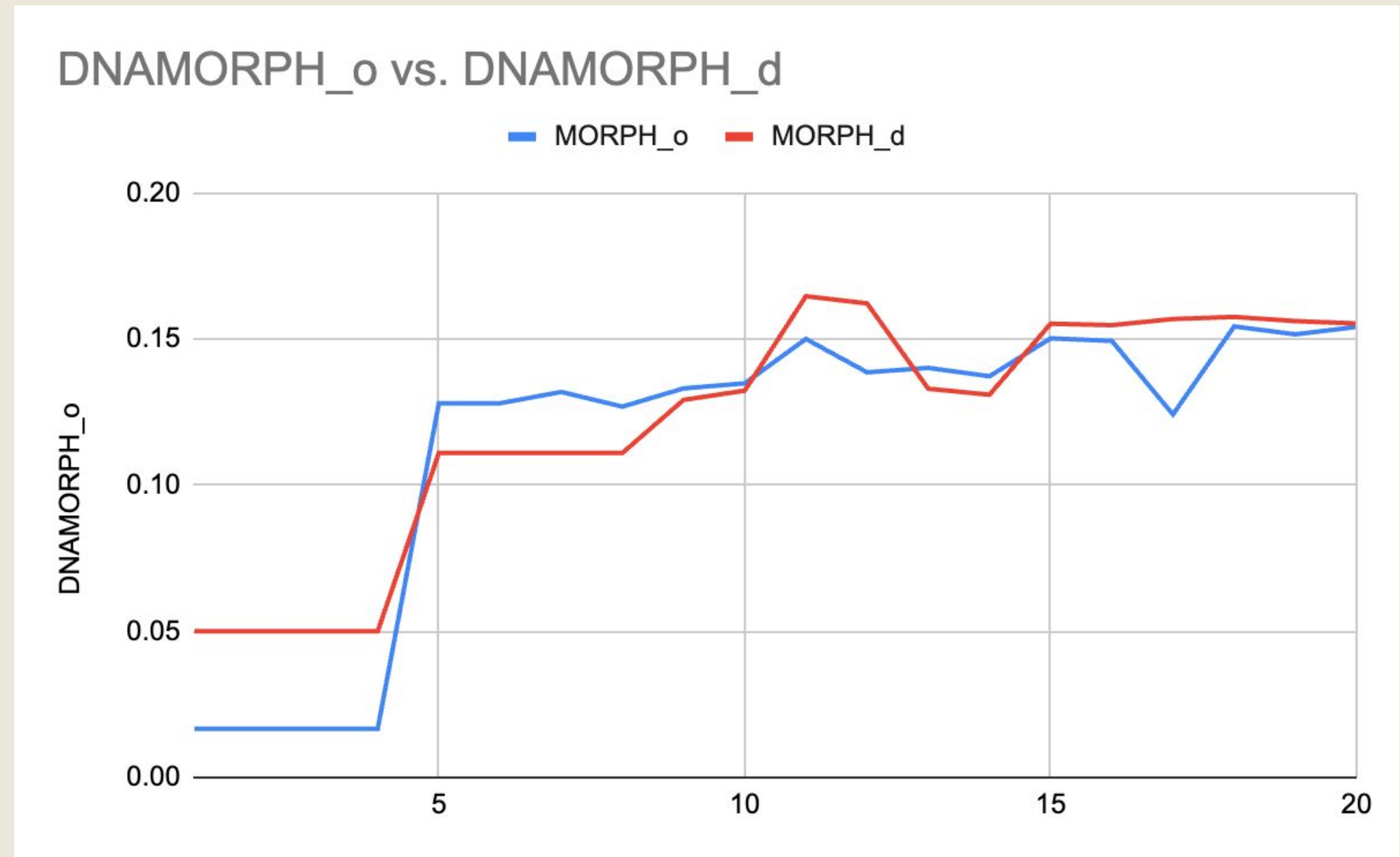
- MORPH\_O

- $R^2 = .178$

- MORPH\_D

- $R^2 = 0.174$

R<sup>2</sup>: Proportion of variance in the dependent variable that can be explained by the independent variable

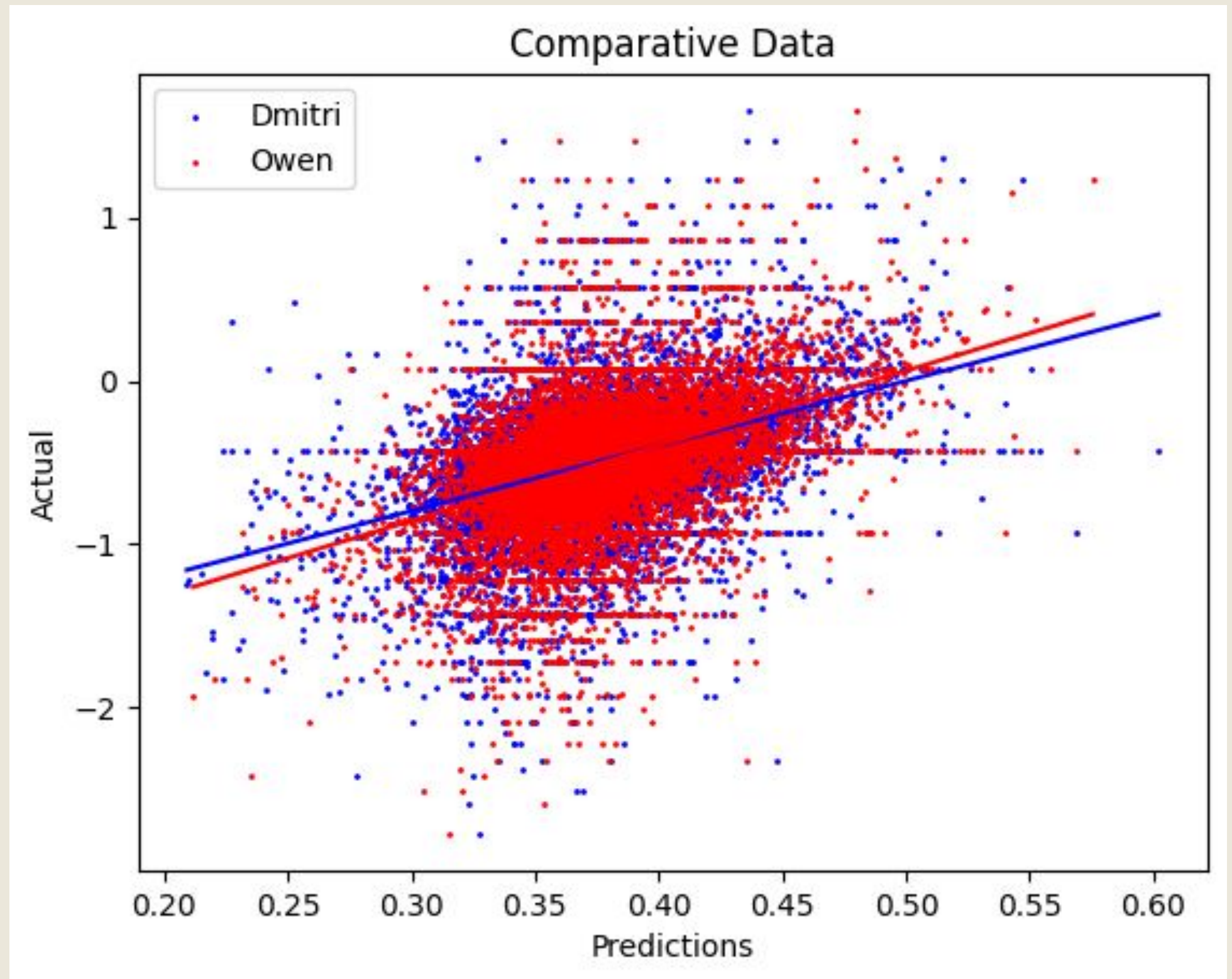


# Measuring Correlation with Models

	Model_d	Model_o
Mean Squared Error	.144	.141
Mean absolute error	.265	.261
Root Mean Squared Error	.379	.375

# Measuring Correlation with Models

Compare this with canonical data to see which data gives stronger predictions, thus has more correlation



# Measuring Correlation with Models

	Model_d	Model_o	Model_DNA
Mean Squared Error	.144	.141	0.092
Mean absolute error	.265	.261	.216
Root Mean Squared Error	.379	.375	.304
$R^2$	.174	.178	.457



# Measuring Correlation with Models

As compared to  
DNABERT2

