

Applying machine learning to biological promoters

Nathan Hu

What is the research problem

Project overview

- Largely based off two papers applying machine learning to promoter prediction

ARTICLE

<https://doi.org/10.1038/s41467-020-15977-4>

OPEN

Model-driven generation of artificial yeast promoters

Benjamin J. Kotopka¹ & Christina D. Smolke^{1,2}✉

Promoters play a central role in controlling gene regulation; however, a small set of promoters is used for most genetic construct design in the yeast *Saccharomyces cerevisiae*. Generating and utilizing models that accurately predict protein expression from promoter sequences would enable rapid generation of useful promoters and facilitate synthetic biology efforts in this model organism. We measure the gene expression activity of over 675,000 sequences in a constitutive promoter library and over 327,000 sequences in an inducible promoter library. Training an ensemble of convolutional neural networks jointly on the two data sets enables very high ($R^2 > 0.79$) predictive accuracies on multiple sequence-activity prediction tasks.

Genome analysis

DNABERT: pre-trained Bidirectional Encoder Representations from Transformers model for DNA-language in genome

Yanrong Ji^{1,†}, Zhihan Zhou^{2,†}, Han Liu^{2,*} and Ramana V. Davuluri^{3,*} 

¹Division of Health and Biomedical Informatics, Department of Preventive Medicine, Northwestern University Medicine, Chicago, IL 60611, USA, ²Department of Computer Science, Northwestern University, Evanston, IL 60201, USA, ³Department of Biomedical Informatics, Stony Brook University, Stony Brook, NY 11794, USA

*To whom correspondence should be addressed.

[†]The authors wish it to be known that, in their opinion, the first two authors should be regarded as Joint First Authors.

Associate Editor: Dr. Janet Kelso

Received on September 10, 2020; revised on December 31, 2020; editorial decision on January 25, 2021; accepted on February 1, 2021

Abstract

Motivation: Deciphering the language of non-coding DNA is one of the fundamental problems in genomics. Gene regulatory code is highly complex due to the existence of polysemy and distant semantic relationships. Previous informatics methods often fail to capture especially in data-scarce scenarios.

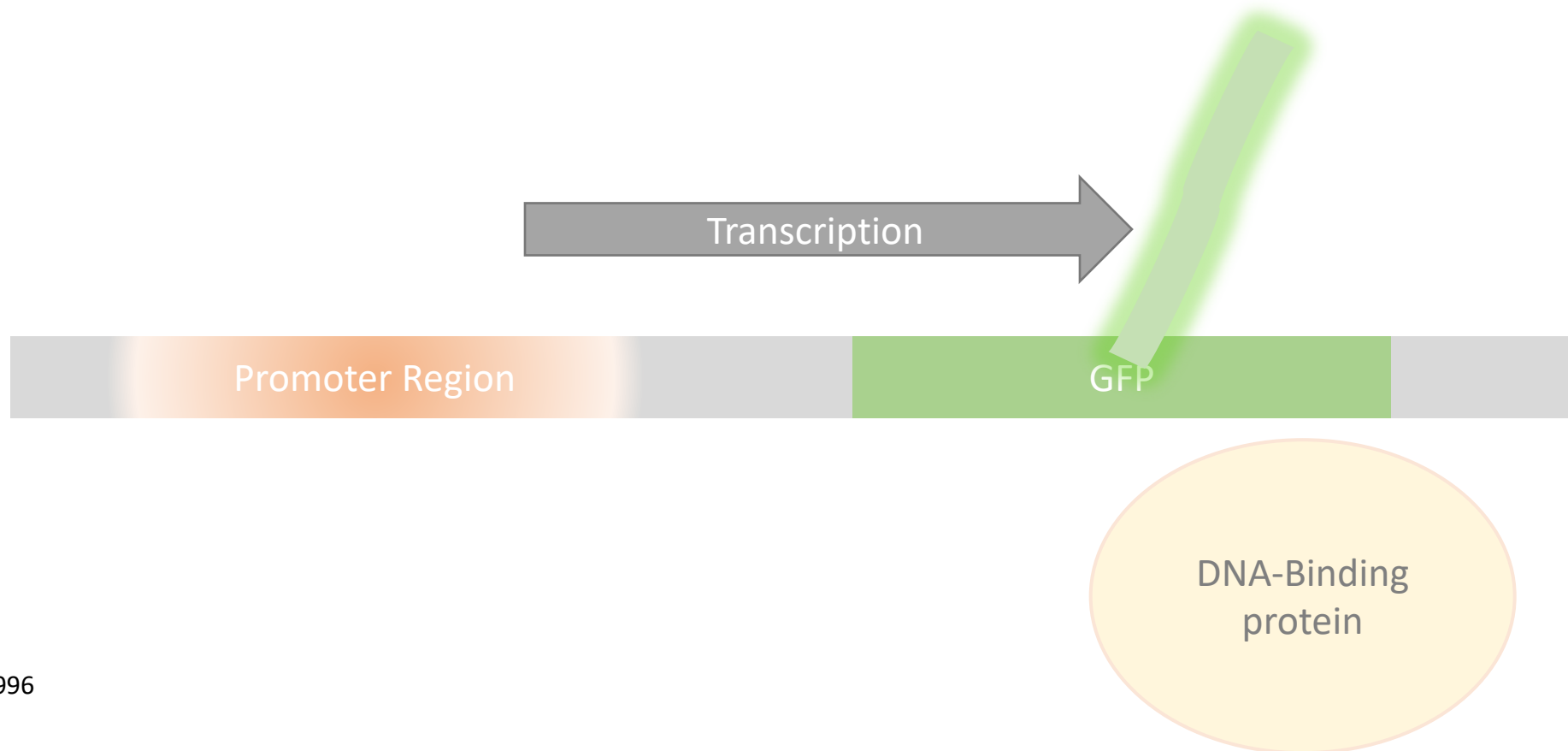
What is a promoter

- Promoters are the key non-coding regions of the genome²
- Found near the transcription start site of a gene²
- Important for initiation of transcription and for regulating gene expression¹ by serving as recognition sites for necessary proteins²

¹Ayoubi & Van De Yen 1996

²Oubounyt et al. 2020

Green Fluorescence Protein (an example)



Motivation

- Precise gene expression control is required in³:
 - Engineered metabolic pathways
 - Gene circuits
- Finding promoters or creating artificial promoters with useful properties may aid gene construct design³
- Studying DNA with machine learning may aid understanding of promoters and DNA regulation in general

Question:

How can we use machine learning to make predictions on promoters?

Research question

How can we use machine learning
to make predictions on promoters?



Regression

*How strong is this
promoter?*
(number value)

Example: Creating artificial promoters

ATGTAAGTGACASTARTGACAGGACATAGACATTACATAT

Unknown New
Sequence

Regression



Model says it's a strong promoter



Verify model prediction *in vivo*

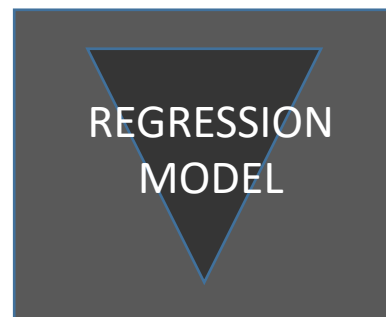


New promoter
sequence found!

Goal

AGGACATAGACATTACATAT

Input DNA sequence



Blackbox model



0.67

Model output
(*promoter strength*)

DNA vector representation

DNA is a sequence of 4 base pairs

How do we represent a sequence such as **ATTGCT** in a computer readable way?

One-hot encoding

Each base represented by a vector:

$$\textcolor{red}{A} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \quad \textcolor{green}{G} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}, \quad \textcolor{blue}{C} = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}, \quad \textcolor{orange}{T} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$$

A sequence is represented by a 2D matrix:

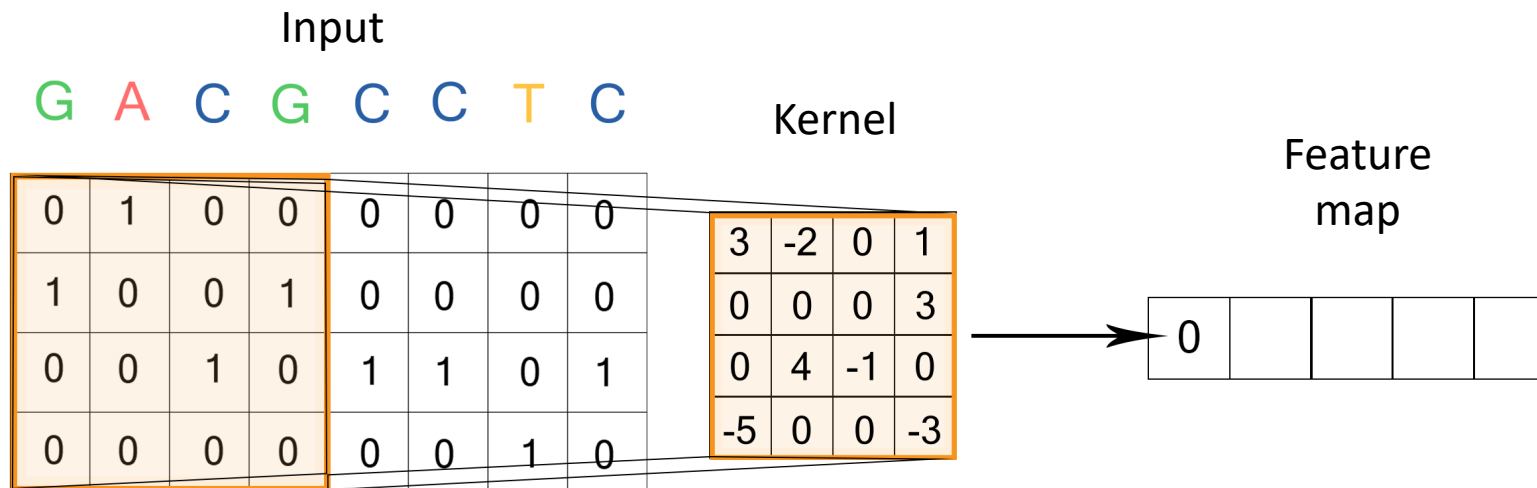
ATTGCT becomes

$$\begin{matrix} \textcolor{red}{A} & \textcolor{orange}{T} & \textcolor{orange}{T} & \textcolor{green}{G} & \textcolor{blue}{C} & \textcolor{blue}{T} \\ \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 \end{bmatrix} \end{matrix}$$

What architectures can we
use?

Convolutional neural network (CNN)

- Usually applied to images and video
- Based off the concept of a convolution
- Focuses on the input 1 section at a time



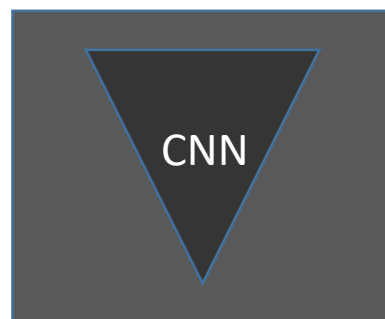
Pipeline 1

AGGACATAGACATTACATAT

Input DNA sequence

$$\begin{bmatrix} 1 & \dots & 0 \\ 0 & \ddots & 0 \\ 0 & \dots & 0 \\ 0 & \dots & 1 \end{bmatrix}$$

One-hot encoding



Blackbox model

0.67

Model output
(*promoter strength*)

How else can we represent a sequence

We have one-hot encoding

Another approach is to learn a representation for DNA

DNABERT: pre-trained Bidirectional Encoder Representations from Transformers model for DNA-language in genome

5

DNABERT⁵

Based off BERT⁷, a state-of-the-art natural language model

Able to take in sentences as input and make different predictions

⁷Devlin et al. 2018

BERT⁷ Language Example

Sentiment analysis⁶

"I hate you"



BERT

I think this
sentence is
NEGATIVE

⁶Hugging Face 2020

⁷Devlin et al. 2018

Sentence representations

BERT can 'embed' or represent each sentence as a vector

"I hate you"



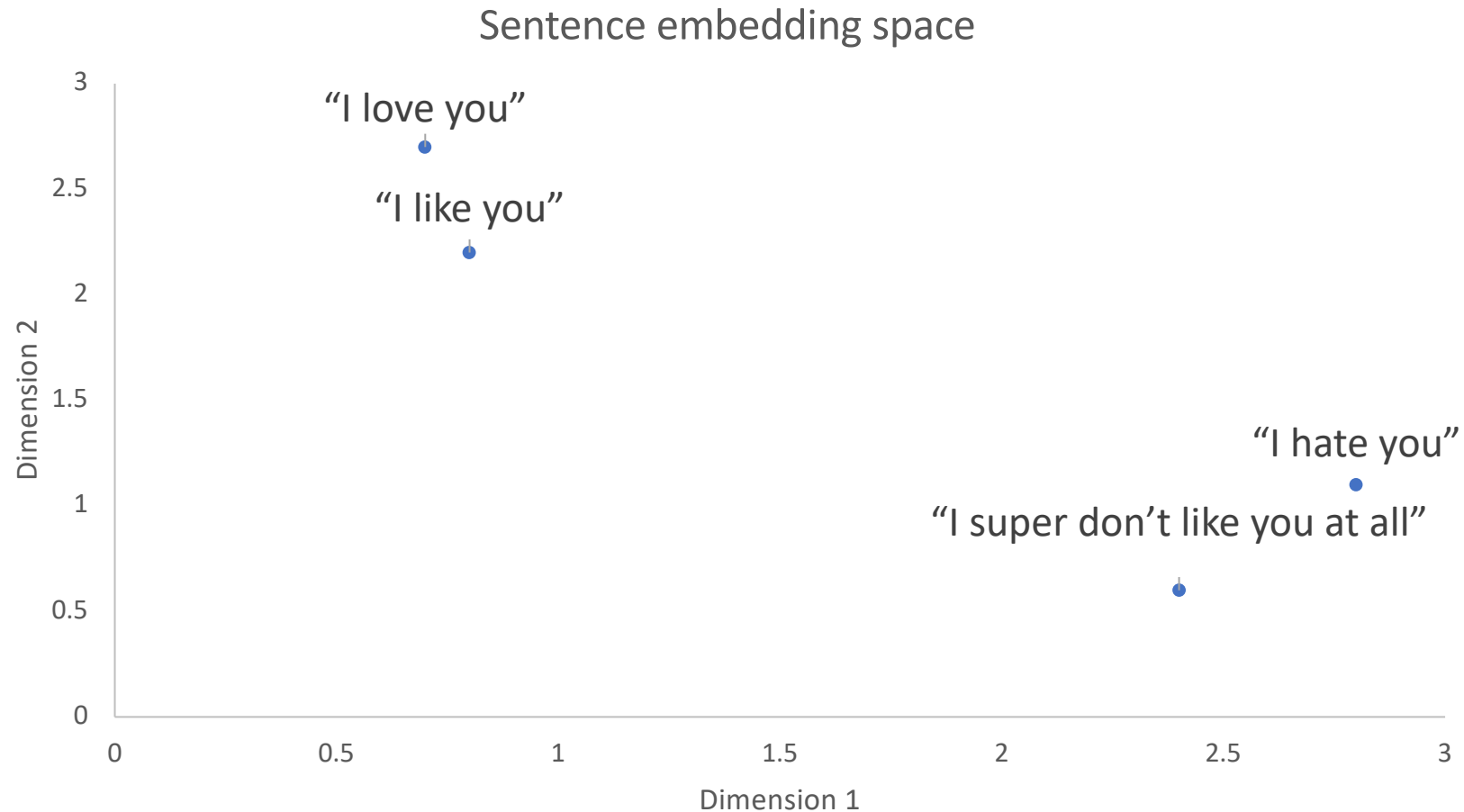
$[1.3221, -11.823, 0.332 \dots 2.398]$

"I love you"



$[-3.928, 0.991, 6.122 \dots 7.082]$

Sentence representations



Attention mechanism

BERT can separate different sentences by understanding context

Understands context by using the “attention mechanism” on each word

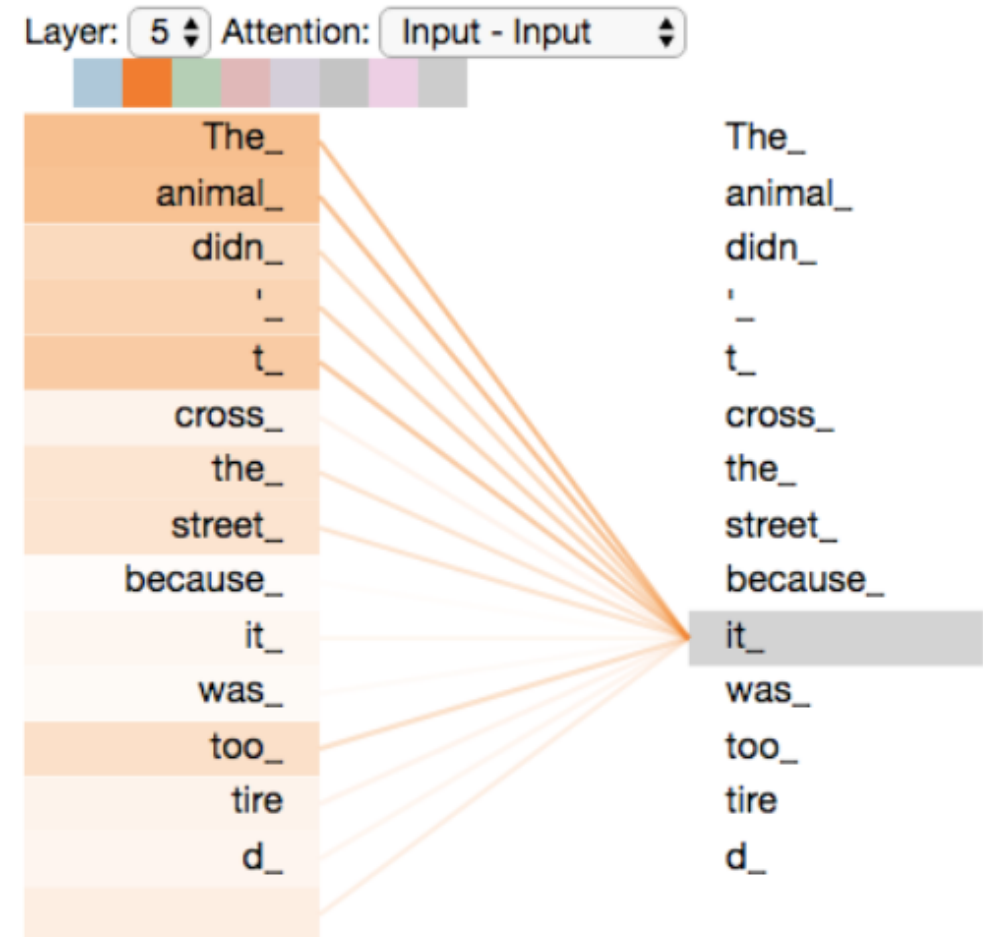
“I love you”



“I”

“love”

“you”



Link to DNA

A sequence of DNA is analogous to a sentence

GACAGGACATAGACATTACAT is a 'sentence'

DNABERT

DNABERT⁵ embeddings are first trained using the human genome for 25 days on 8 GPUS

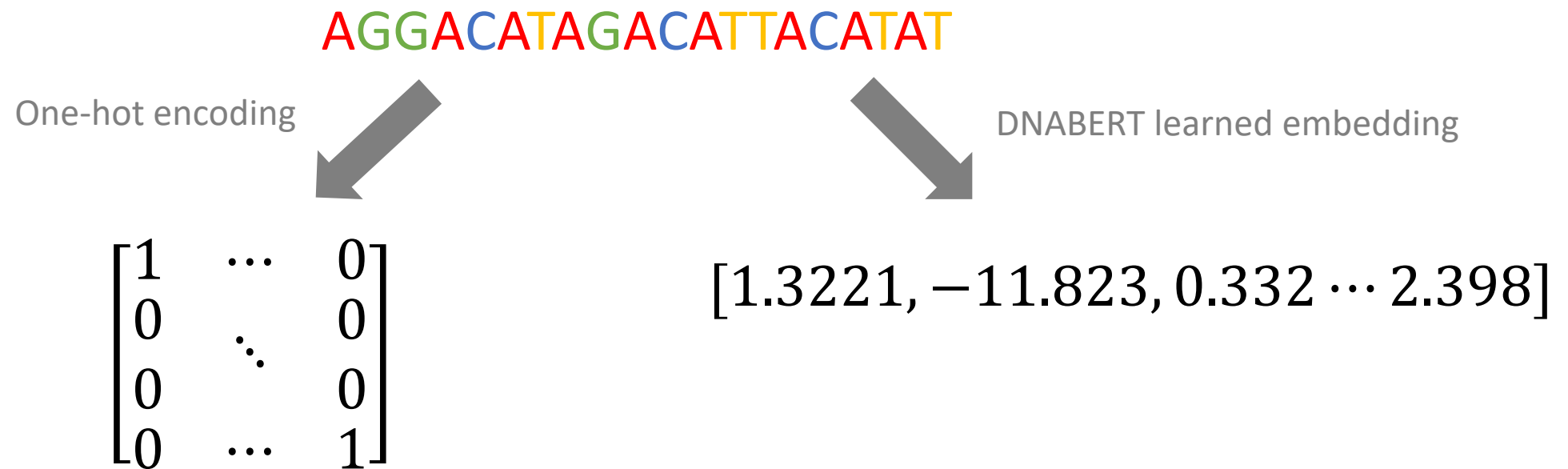
Provides it with a general and transferrable understanding of DNA⁵

So can be transferred to different downstream tasks and organisms



Yeast promoter strength prediction

One-hot vs DNABERT



Pipeline 2

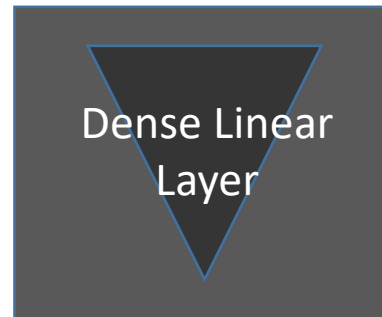
AGGACATAGACATTACATAT

Input DNA sequence



[1.3221, -11.823, 0.332 ... 2.398]

DNABERT embedding



Blackbox model



0.67

Model output
(*promoter strength*)

Pipeline 3

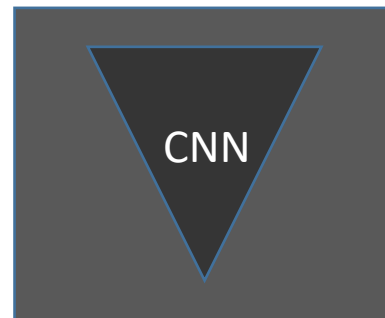
AGGACATAGACATTACATAT

Input DNA sequence



[1.3221, -11.823, 0.332 ... 2.398]

DNABERT embedding



Blackbox model



0.67

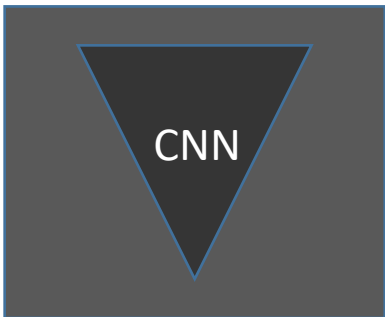
Model output
(*promoter strength*)

Summary of architectures

1. One-hot + CNN

AGGACATAGACATTACATAT

$\begin{bmatrix} 1 & \dots & 0 \\ 0 & \ddots & 0 \\ 0 & & 0 \\ 0 & \dots & 1 \end{bmatrix}$

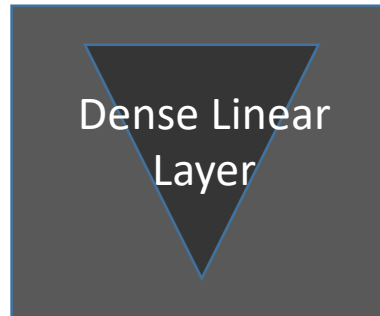


0.67

2. DNABERT + Dense

AGGACATAGACATTACATAT

$[1.3221, -11.823, 0.332 \dots 2.398]$

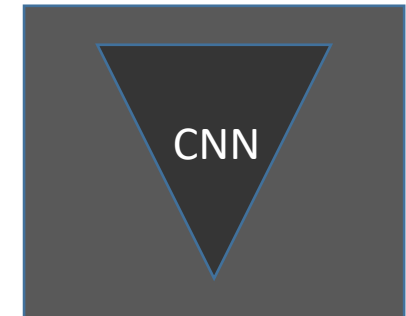


0.67

3. DNABERT + CNN

AGGACATAGACATTACATAT

$[1.3221, -11.823, 0.332 \dots 2.398]$



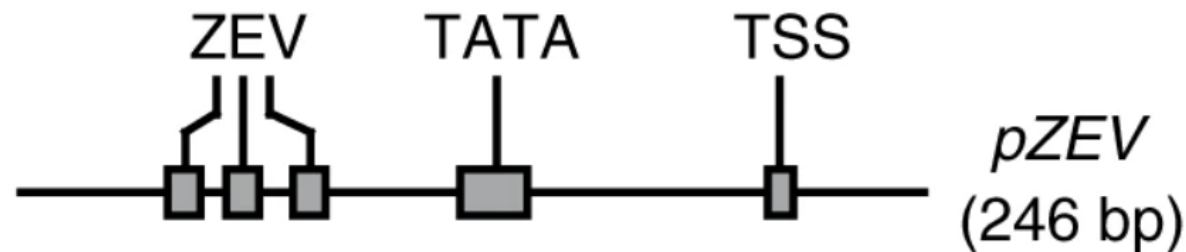
0.67

Methods & Results

The dataset³

327,000 yeast promoters containing a binding site for artificial transcription factor ZEV

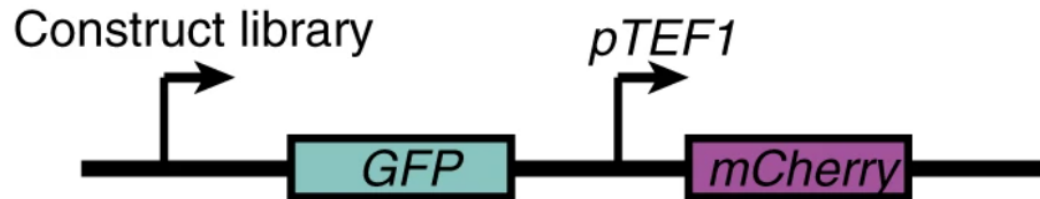
Each sequence is 246 base pairs and has a label for measured promoter strength



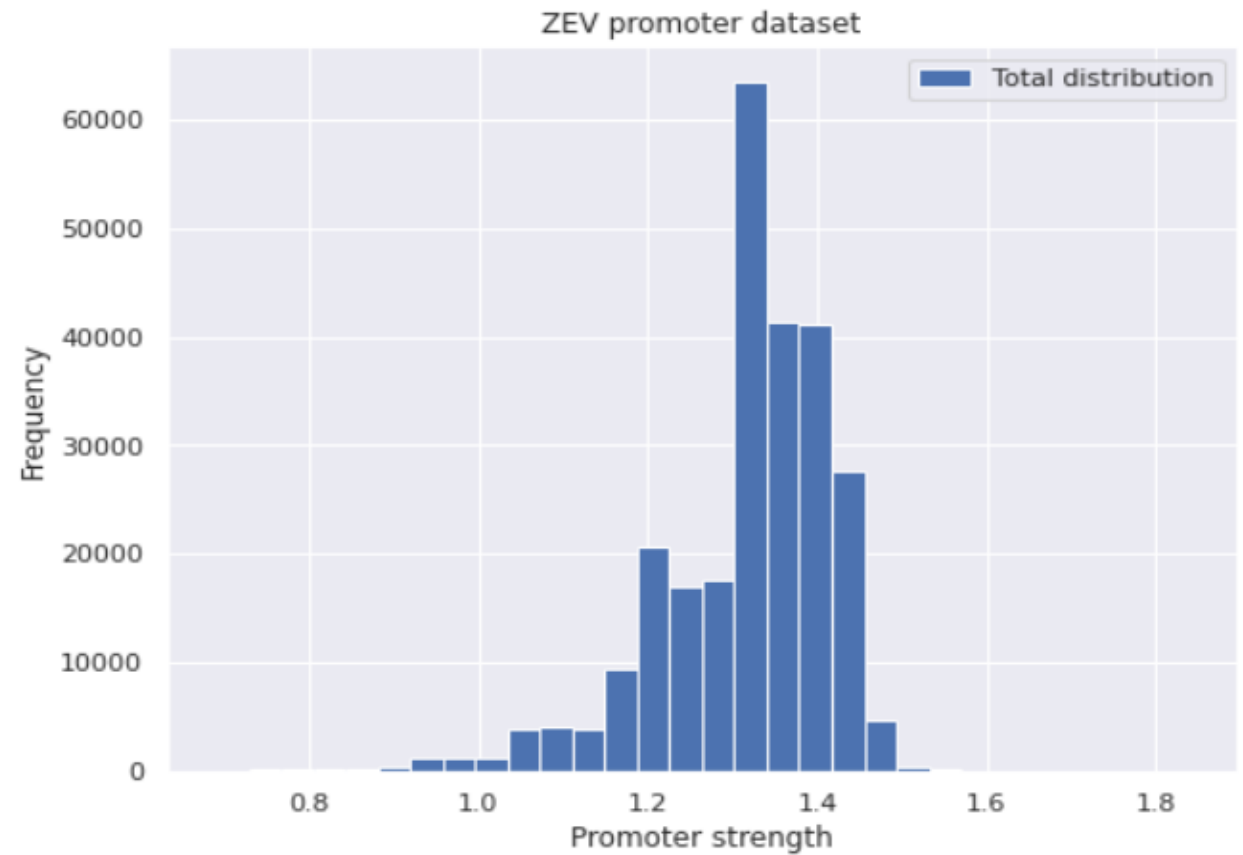
³Kotopka & Smolke 2020

The dataset

How was promoter strength measured?



³Kotopka & Smolke 2020



Training

Each model needs to be trained to make accurate predictions

We split the 327,000 promoters into

- 261,584 **training** examples
- 32,708 **validation** examples
- 32,708 **test** examples

One-hot+CNN Training

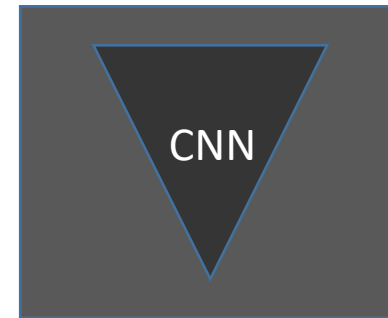
One-hot + CNN

AGGACATAGACATTACATAT

Not learned

$$\begin{bmatrix} 1 & \dots & 0 \\ 0 & & 0 \\ 0 & \ddots & 0 \\ 0 & \dots & 1 \end{bmatrix}$$

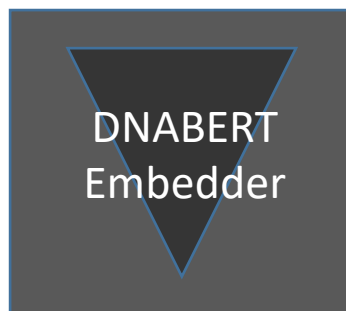
Learned



0.67

DNABERT Training

Learned

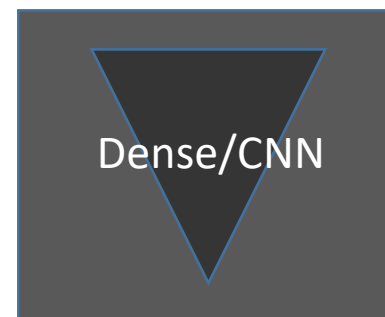


Learned

DNABERT + Dense/CNN

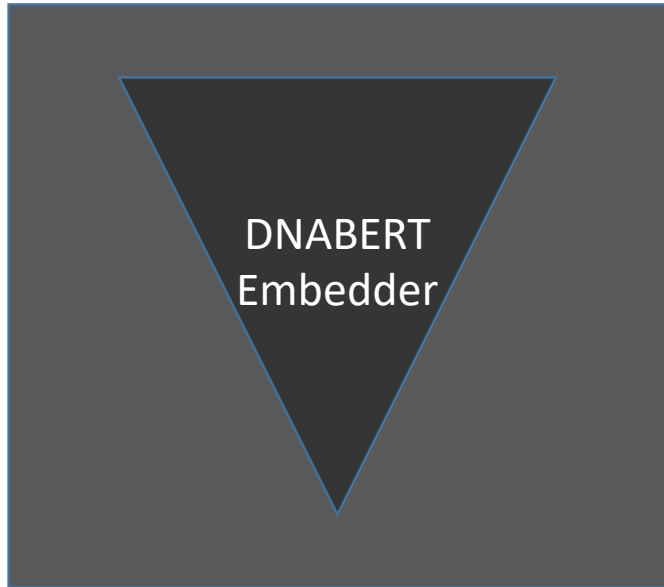
AGGACATAGACATTACATAT

[1.3221, -11.823, 0.332 ... 2.398]



0.67

DNABERT Training



Very big!!

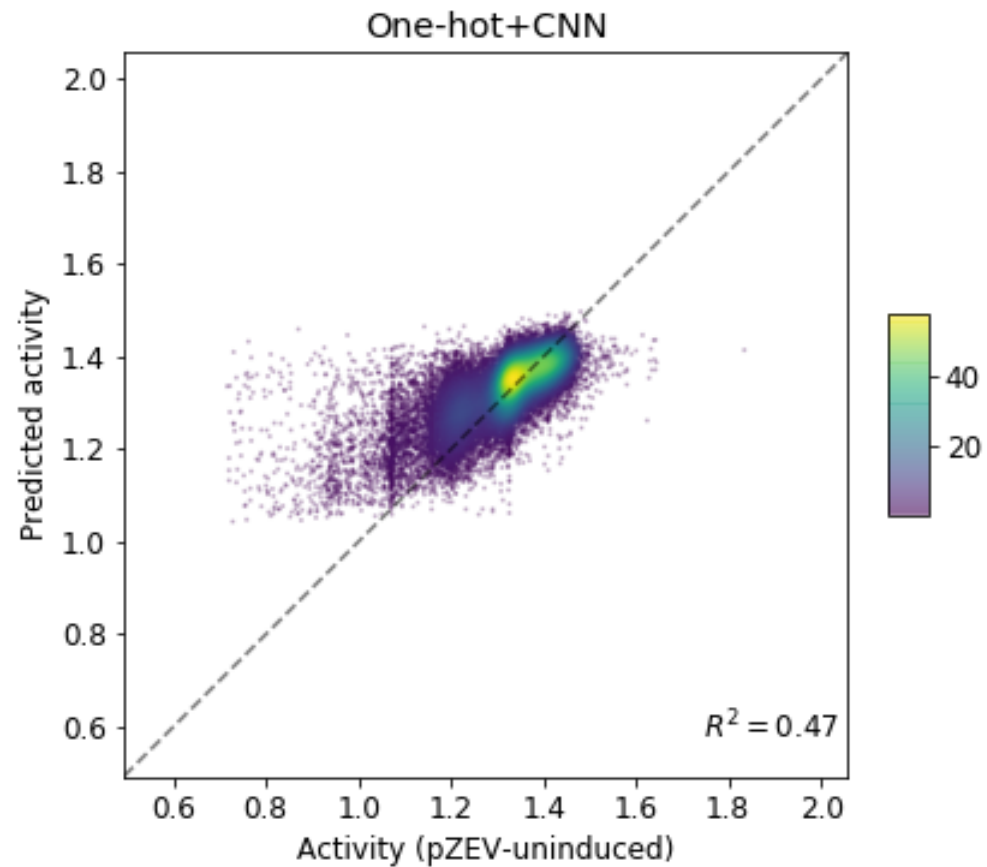
DNABERT has lots of layers

We download a pre-trained version that is trained using the entire **human** genome

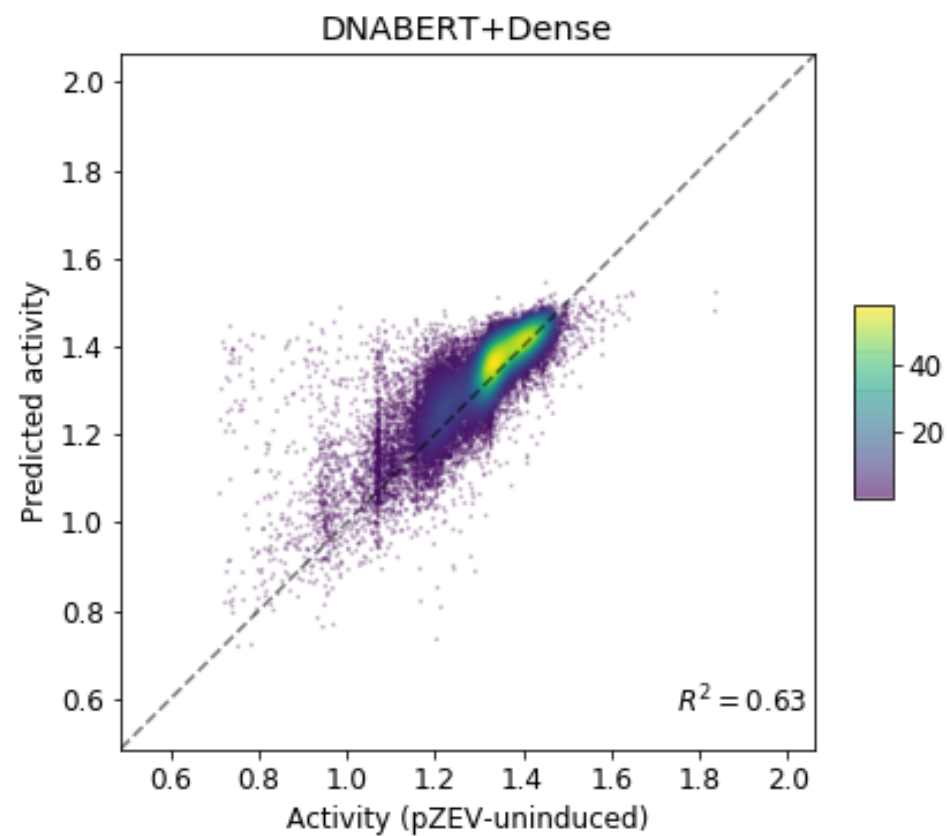
We then make fine adjustments using the **yeast** promoter dataset

1. One-hot + CNN

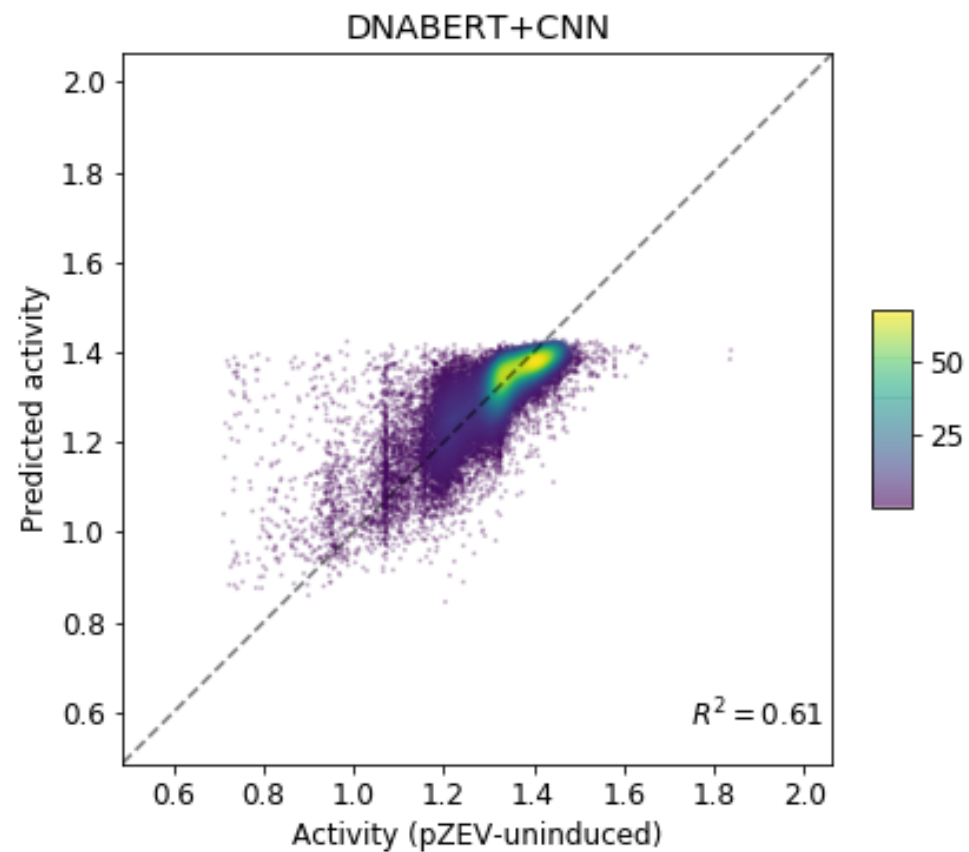
1. One-hot + CNN



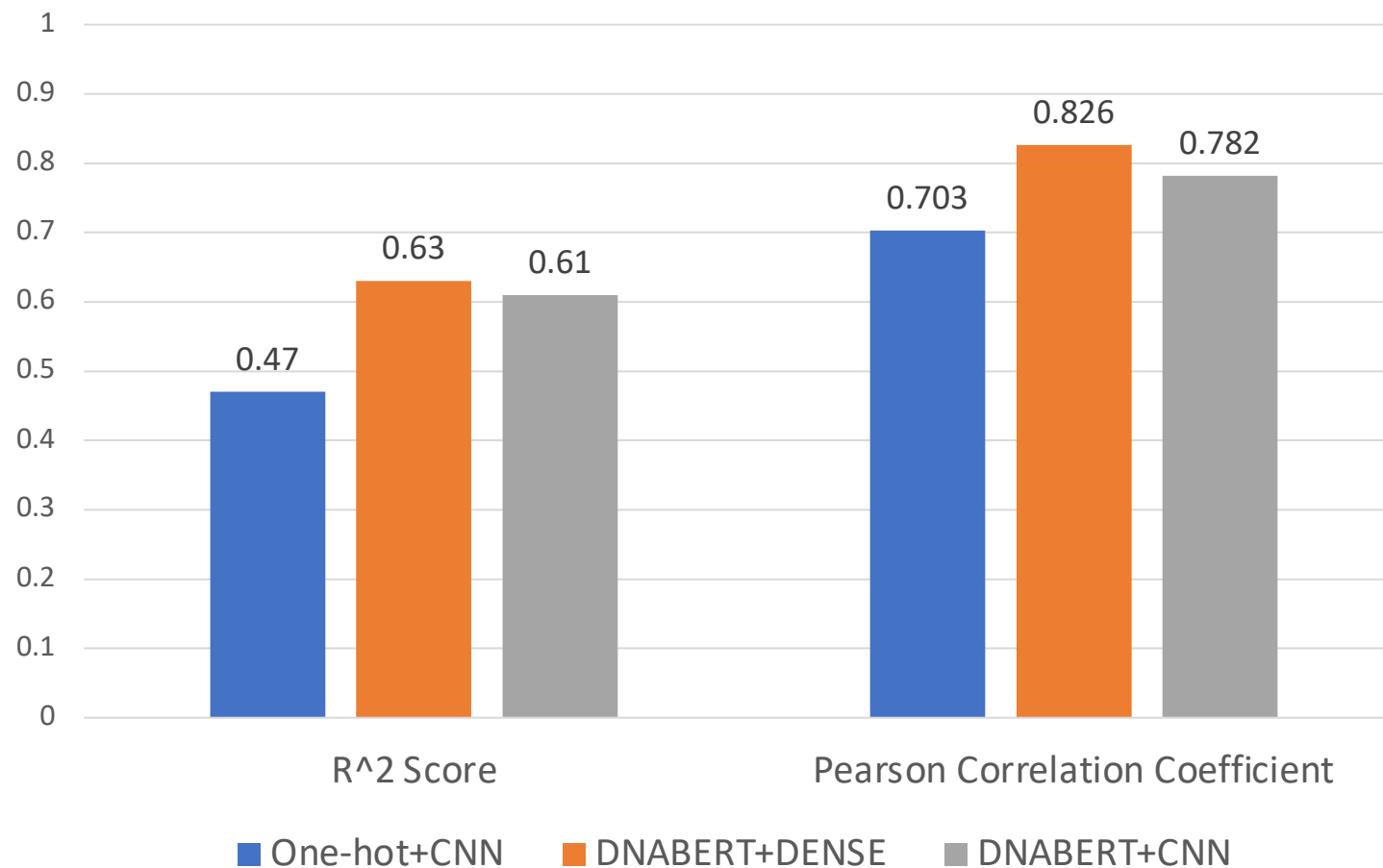
2. DNABERT + Dense



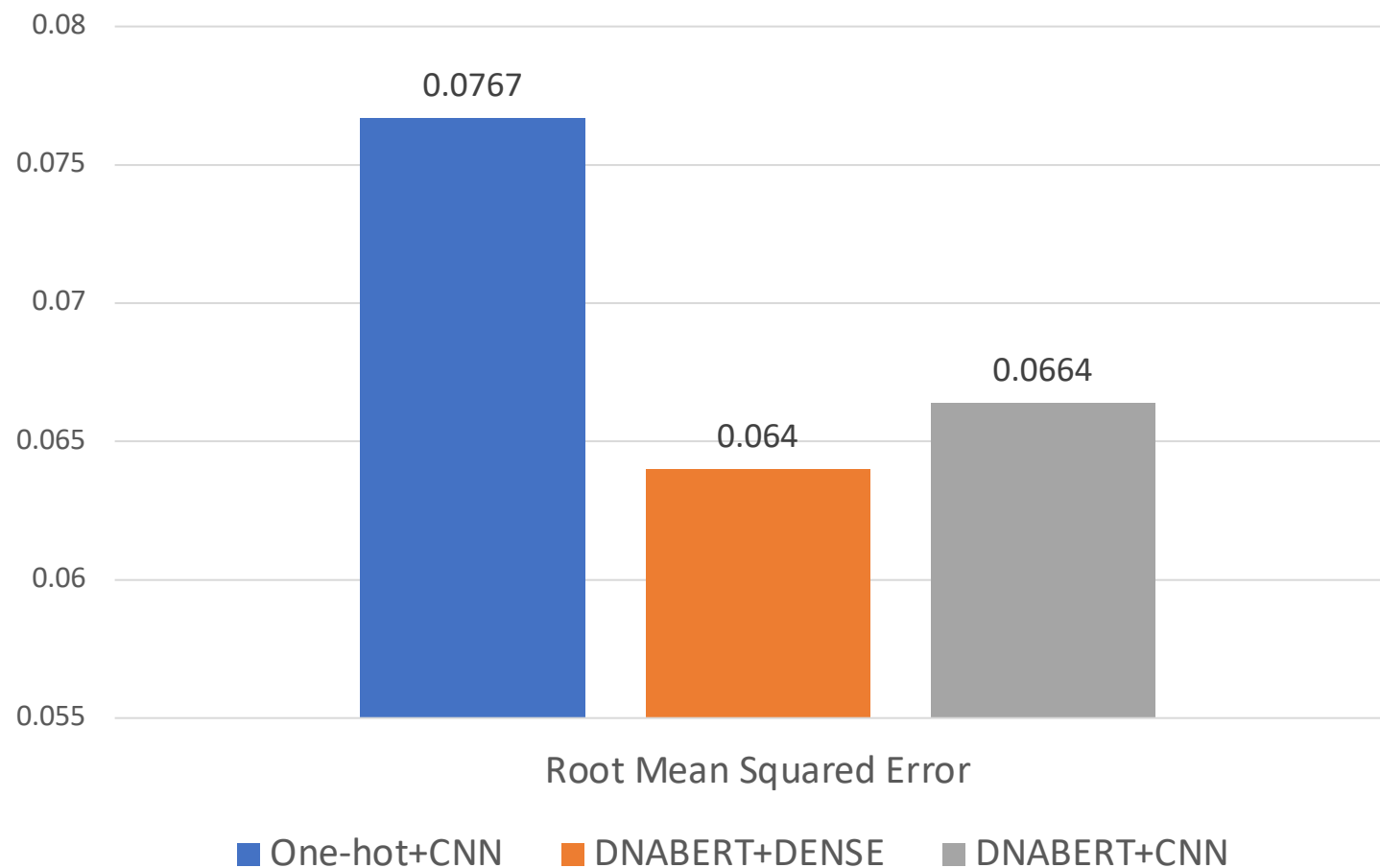
3. DNABERT + CNN



Summary of results



Summary of results



Interpretation

DNABERT generalises well to different tasks and organisms

DNABERT Performs better than one-hot+CNN

CNN is restricted by receptive field, unable to easily take in global contextual information of sequence

Results could imply this is important

Not enough time to spend hyperparameter tuning, these scores could be improved slightly

Future work

Pre-train DNABERT on yeast genome and compare performance

Evaluate the three pipelines on a classification task instead of a regression task

- Need to create a non-promoter dataset

Visualise what DNABERT is focusing its 'attention' on

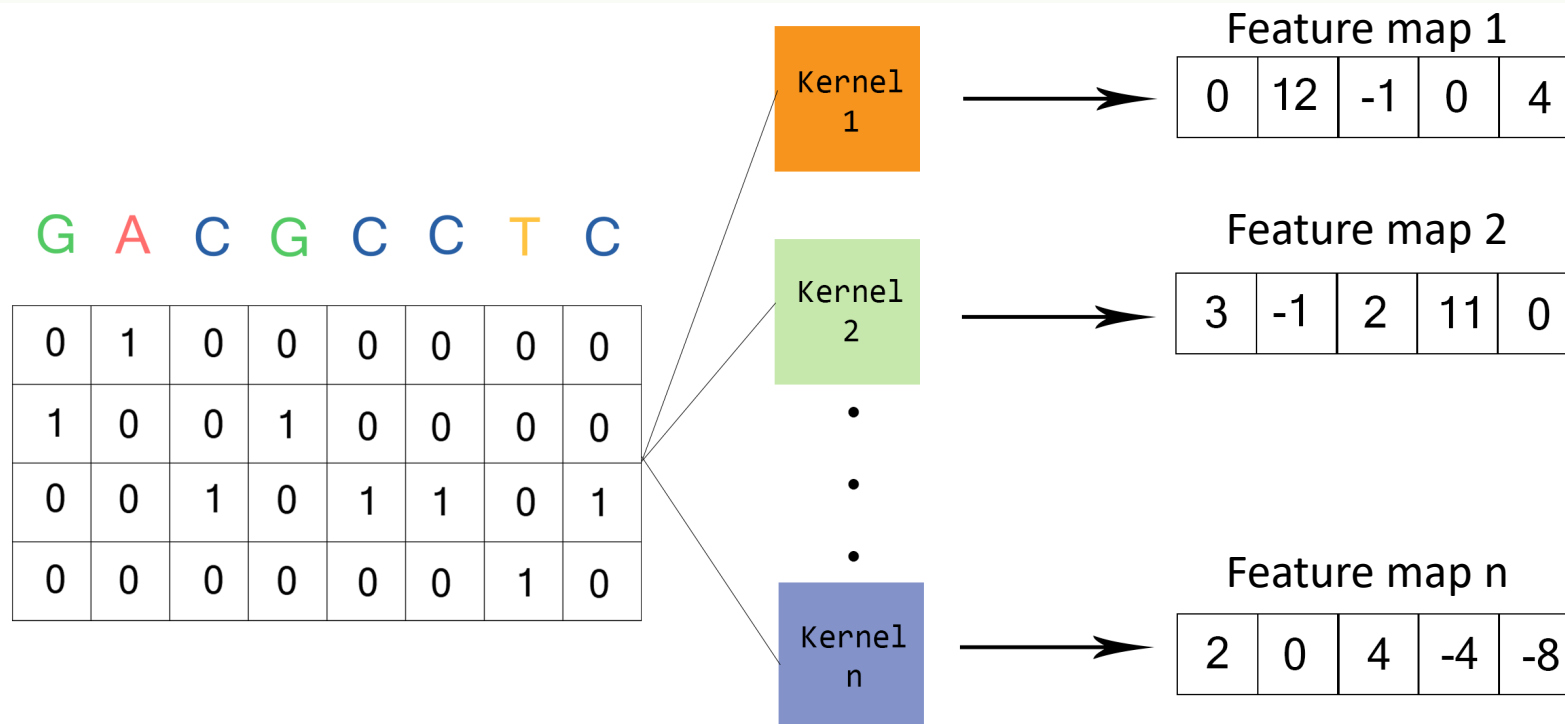
- Could elucidate important features of promoters that make it strong

References

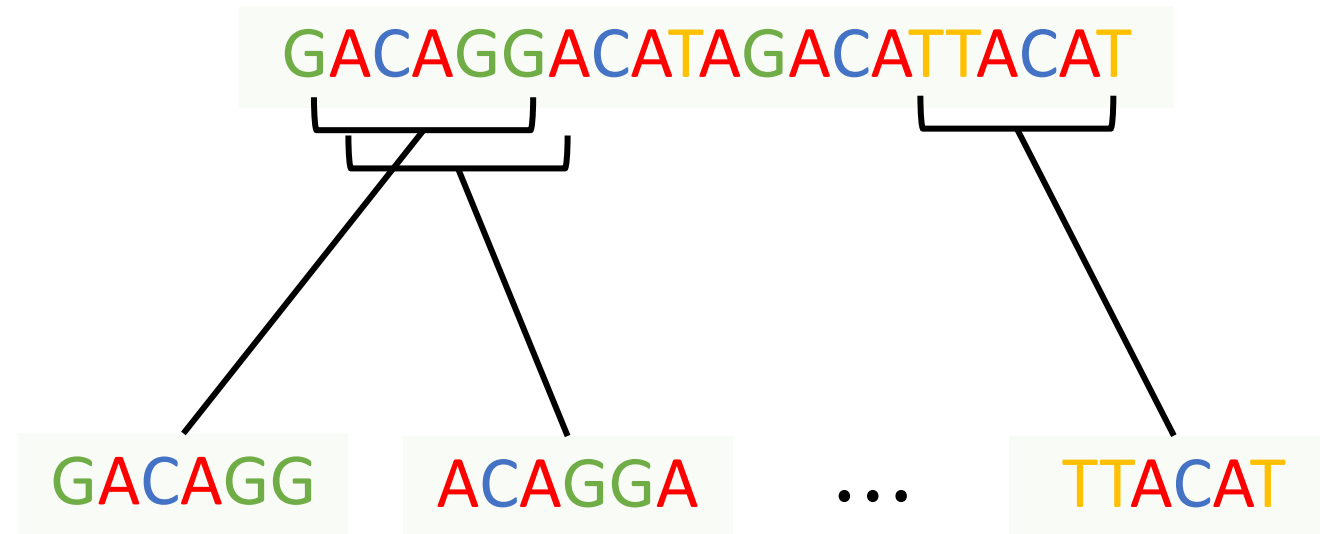
- ¹Ayoubi, T.A.Y. and Van De Yen, W.J.M. (1996), Regulation of gene expression by alternative promoters. The FASEB Journal, 10: 453-460. <https://doi.org/10.1096/fasebj.10.4.8647344>
- ²Oubounyt, M., Louadi, Z., Tayara, H., & Chong, K. T. (2019). DeePromoter: Robust Promoter Predictor Using Deep Learning. Frontiers in genetics, 10, 286. <https://doi.org/10.3389/fgene.2019.00286>
- ³Kotopka, B.J., Smolke, C.D. Model-driven generation of artificial yeast promoters. Nat Commun 11, 2113 (2020). <https://doi.org/10.1038/s41467-020-15977-4>
- ⁴Alammar, J. (2018) The Illustrated Transformer.
<http://jalammar.github.io/illustrated-transformer/>
- ⁵Yanrong Ji, Zhihan Zhou, Han Liu, Ramana V Davuluri, DNABERT: pre-trained Bidirectional Encoder Representations from Transformers model for DNA-language in genome, *Bioinformatics*, 2021, btab083, <https://doi.org/10.1093/bioinformatics/btab083>
- ⁶Hugging Face (2020). Summary of the tasks
https://huggingface.co/transformers/task_summary.html
- ⁷Devlin J. et al. (2018) Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv : 1810.04805*.

Convolutional neural network (CNN)

- Usually create more than one feature map from the same input using different kernels



k -mers



6-mers

4096 possible
6-mers