

# **Sequence to Function**

**Enformer: Transformers in regulatory genomics**


Xinming, Ian

# Enformer

ARTICLES







<https://doi.org/10.1038/s41592-021-01252-x>

nature | methods

 Check for updates

OPEN

## Effective gene expression prediction from sequence by integrating long-range interactions

Žiga Avsec<sup>1</sup>  , Vikram Agarwal<sup>2,4</sup>, Daniel Visentin<sup>1,4</sup>, Joseph R. Ledsam<sup>1,3</sup>, Agnieszka Grabska-Barwinska<sup>1</sup>, Kyle R. Taylor<sup>1</sup>, Yannis Assael<sup>1</sup>, John Jumper<sup>1</sup>, Pushmeet Kohli<sup>1</sup>   and David R. Kelley<sup>2</sup>  



David Kelly  
PI at Calico Labs

# Preview

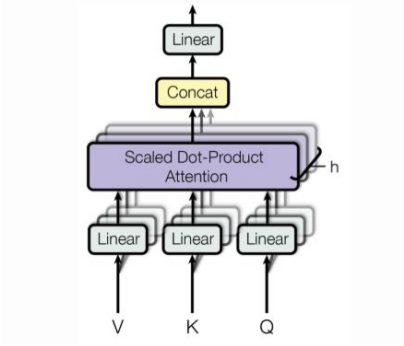


BIOINFORMATICS

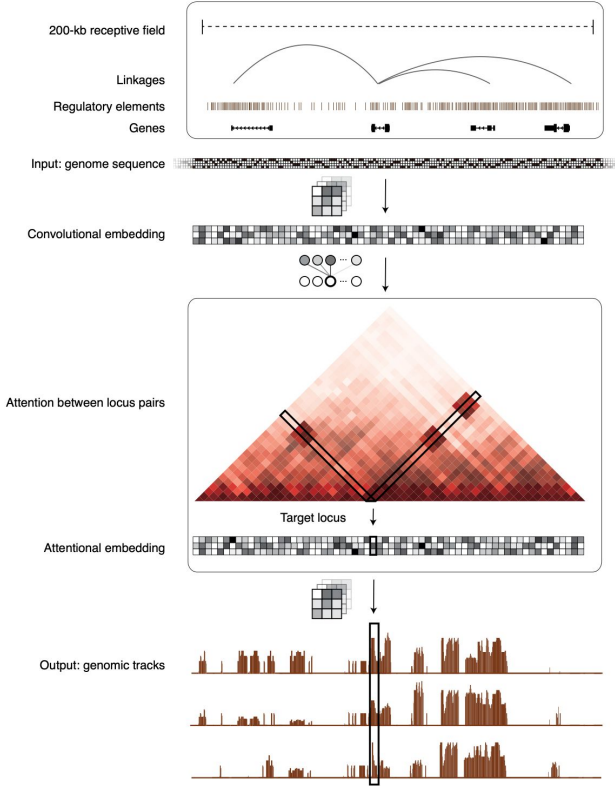
## A wider field of view to predict expression

A gene sequence-to-expression machine learning model achieves improved accuracy by incorporating information about potential long-range interactions.

Yang Young Lu and William Stafford Noble



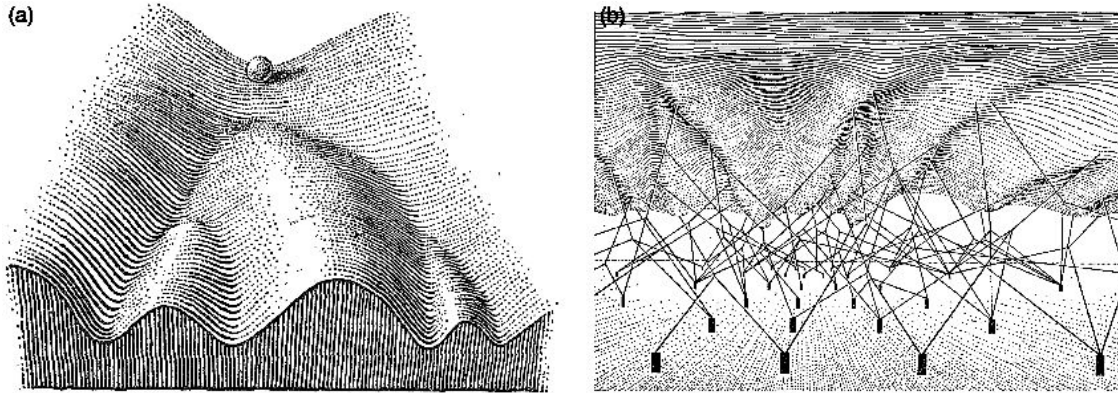
source: attention is all you need



# Outline

- Background
  - regulatory genomics
  - previous work
- Model: convolutions → self-attention
- Results
- Discussion

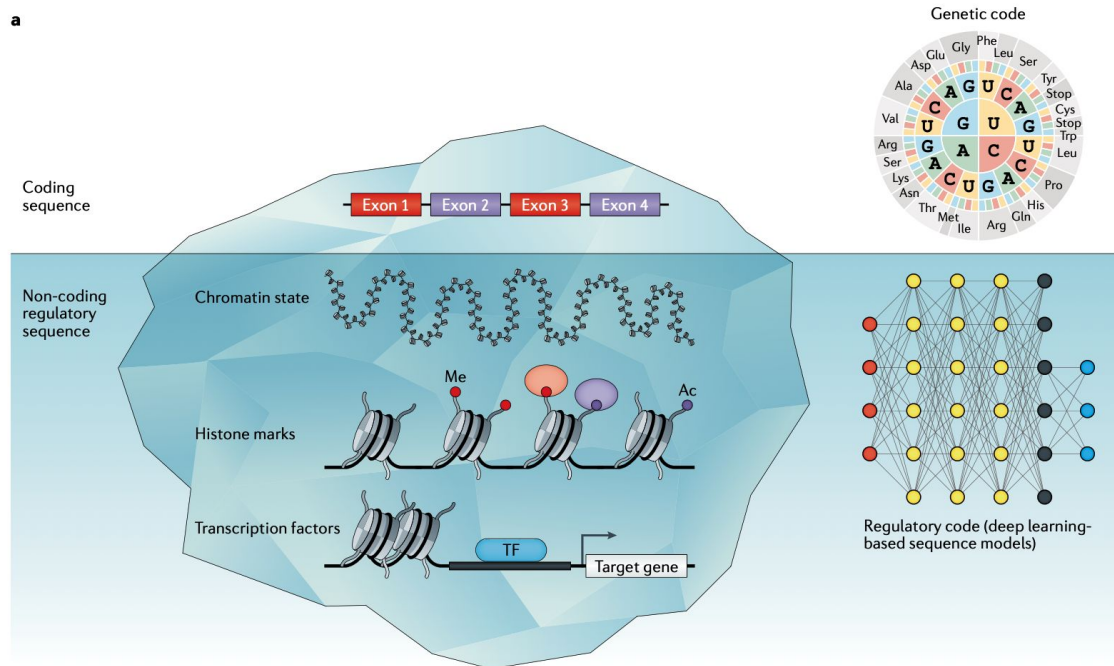
# Regulatory genomics



Epigenetics landscape, Waddington 1957

# Background: non-coding regions

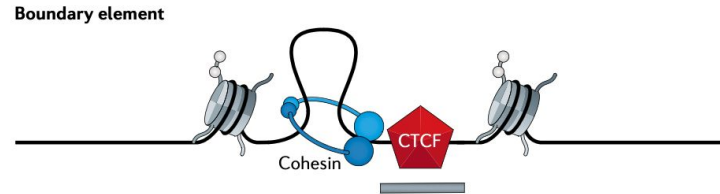
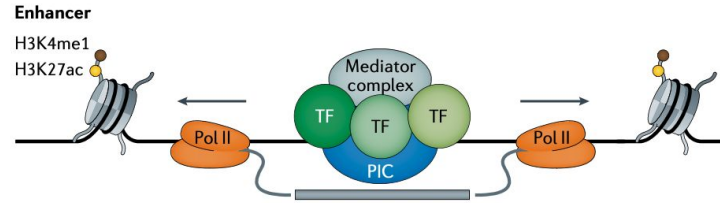
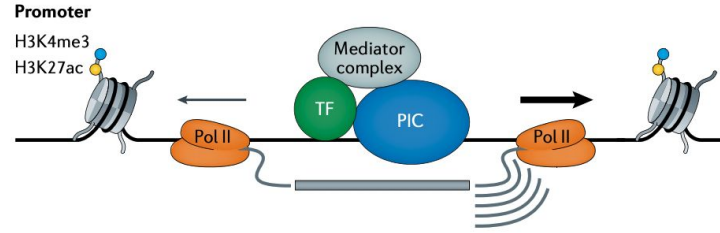
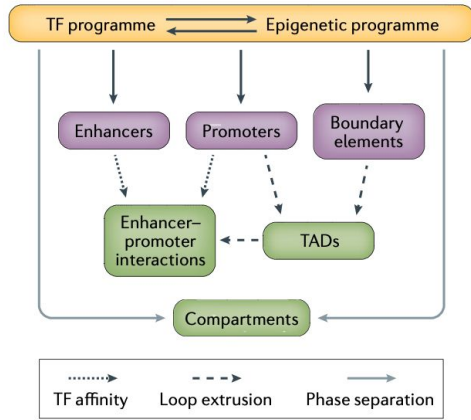
a



Wong, A. K., Sealfon, R. S. G., Theesfeld, C. L. & Troyanskaya, O. G. Decoding disease: from genomes to networks to phenotypes. *Nat Rev Genet* 1–17 (2021) doi:10.1038/s41576-021-00389-x.

# Background: regulatory elements

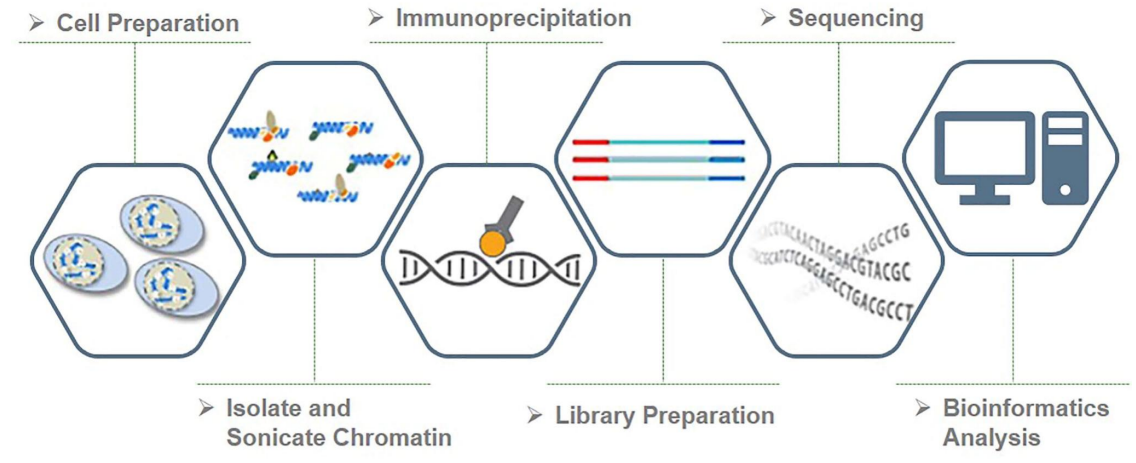
- Promoter
- Enhancer
- Boundary element
  - Insulator (CTCF)



Oudelaar, A. M. & Higgs, D. R. The relationship between genome structure and function. *Nat Rev Genet* **22**, 154–168 (2021).

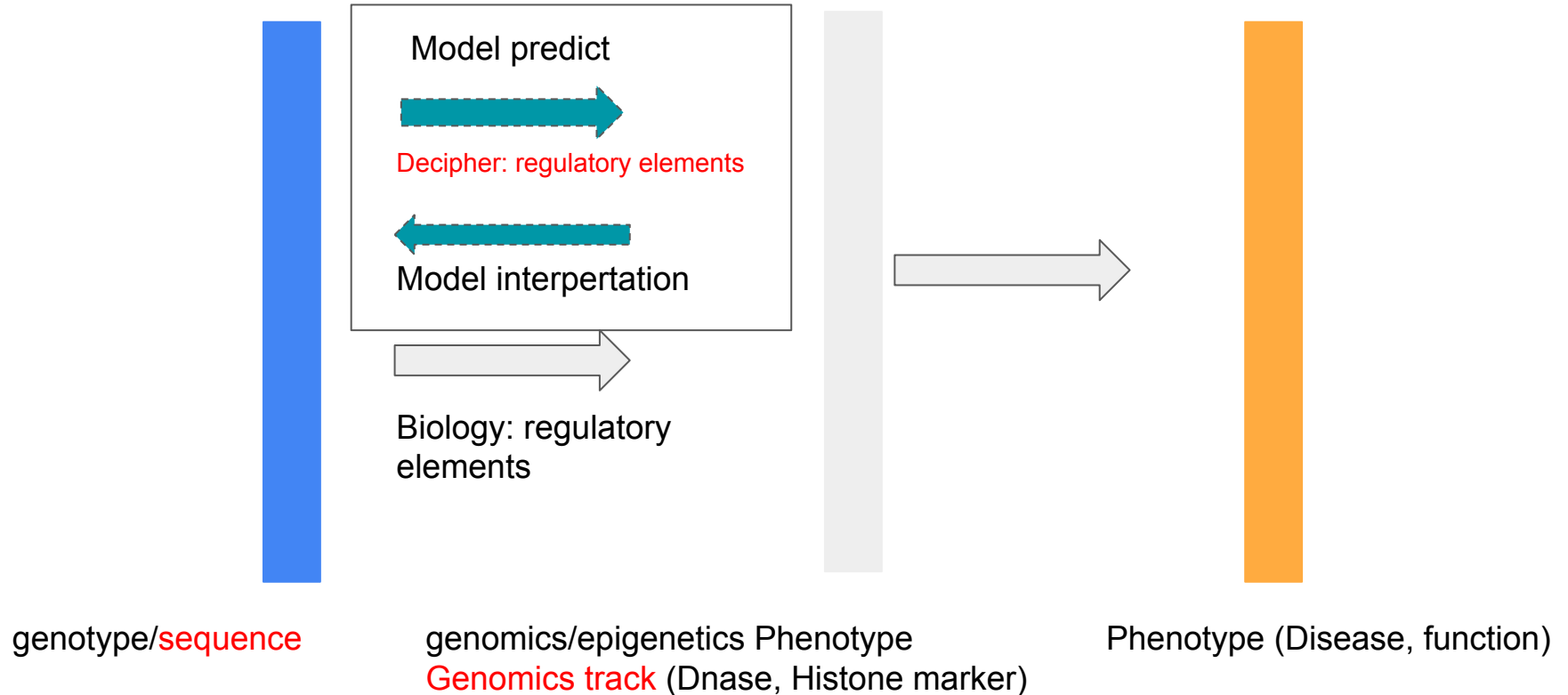
# Background - Genomics tracks

- DNase/ATAC
  - Chromatin accessibility
- CAGE
  - gene expression
- Chip Histone
  - Histone marker
- Chip TF
  - TF binding





# Computational regulatory genomics

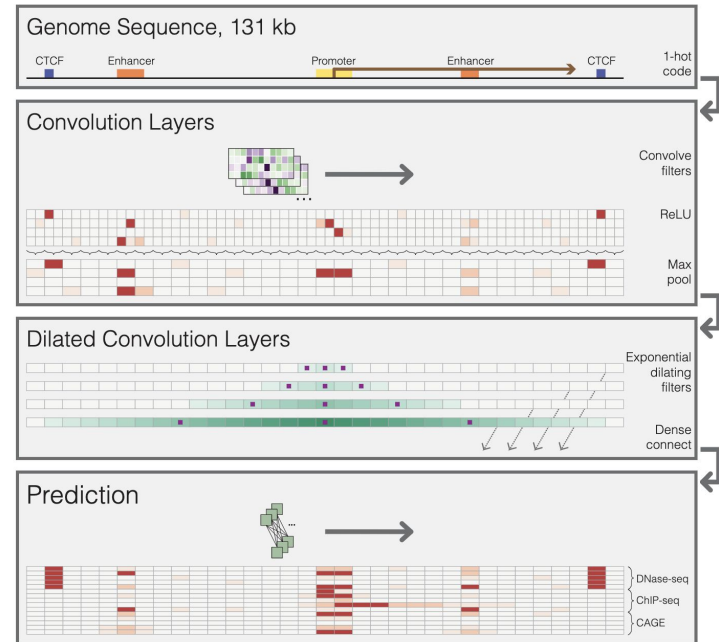


# Previous work in computational regulatory genomics

- The state of the art: CNN
  - DeepBind, DeepSEA, Basenji, etc

Table 1 | Methods for transcriptional/biochemical impact

Model overview			Input data
Model	Method	Prediction task	
DeepSEA <sup>25,26</sup>	Deep learning, CNN	Chromatin, TF binding	TF, HM, DHS
gkm-SVM/delta-SVM <sup>35</sup>	SVM	Chromatin, TF binding	TF, HM, DHS
DanQ <sup>32</sup>	Deep learning, CNN, BLSTM	Chromatin, TF binding	TF, DHS, HM
Basset <sup>37</sup>	Deep learning, CNN	Chromatin accessibility	DHS
DeepCpG <sup>39</sup>	Deep learning, CNN, GRU	CpG state	Bisulfite sequencing
ExPecto <sup>95</sup>	Deep learning, CNN, linear regression	Expression prediction	TF, HM, DHS, RNA-seq
Basenji <sup>35</sup>	Deep learning, CNN	Expression prediction	TF, DHS, HM, CAGE peaks
BPNet, DeepLIFT, TFModisco <sup>36</sup>	Deep learning, CNN	TF binding	TF
ChromDragoNN <sup>35</sup>	Deep learning, CNN, ResNet	Chromatin accessibility	DHS, RNA-seq
Xpresso <sup>33,94</sup>	Deep learning, CNN	Expression prediction	CAGE peaks, gene annotations
AMBER <sup>40</sup>	Auto machine learning, RNN, deep learning, CNN	Chromatin, TF binding	TF, HM, DHS



Kelley, D. R. et al. Sequential regulatory activity prediction across chromosomes with convolutional neural networks. *Genome Res* 28, 739–750 (2018).

Wong, A. K., Sealfon, R. S. G., Theesfeld, C. L. & Troyanskaya, O. G. Decoding disease: from genomes to networks to phenotypes. *Nat Rev Genet* 1–17 (2021)  
doi:10.1038/s41576-021-00389-x.

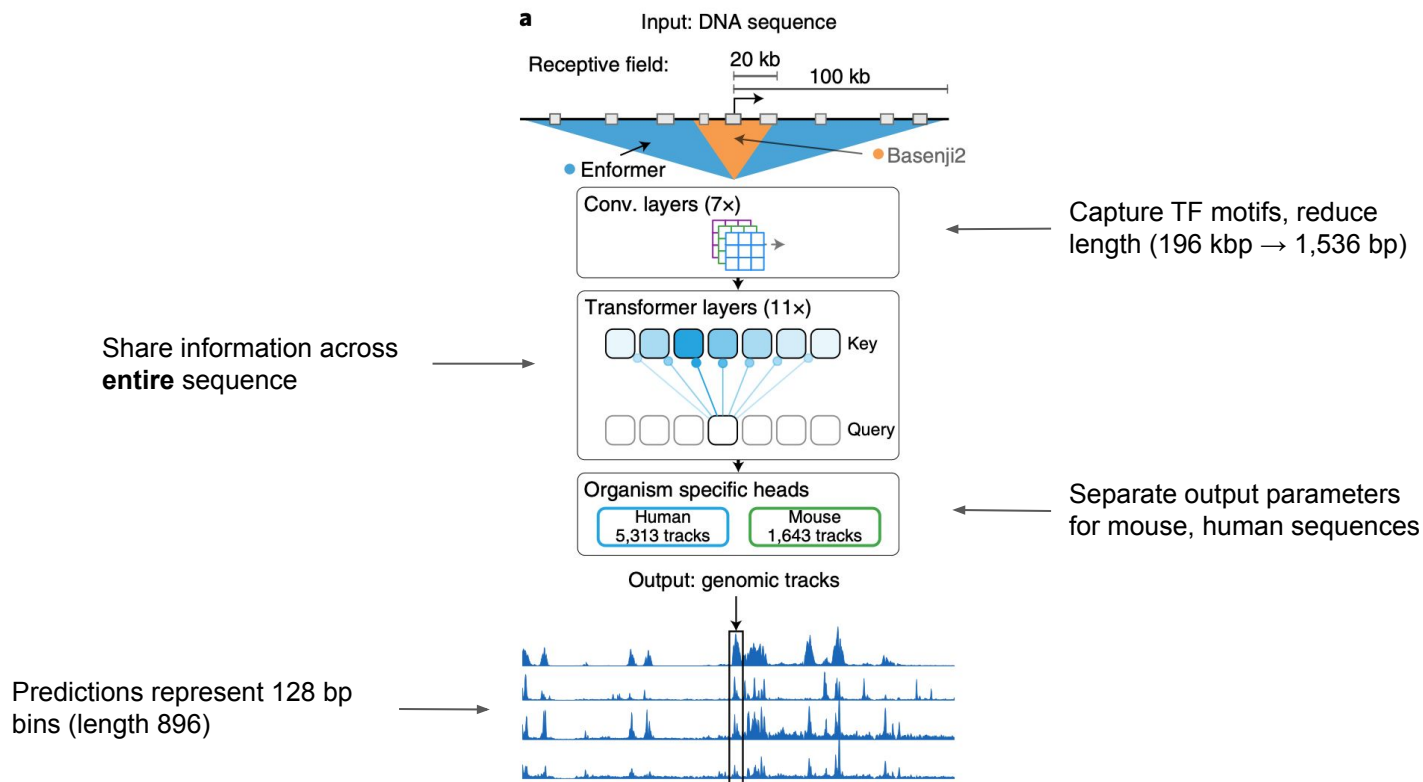
# Model

- A self-attention based sequence-to-sequence model (transformer)
- Enformer = enhancer + transformer

# Architecture



# Architecture



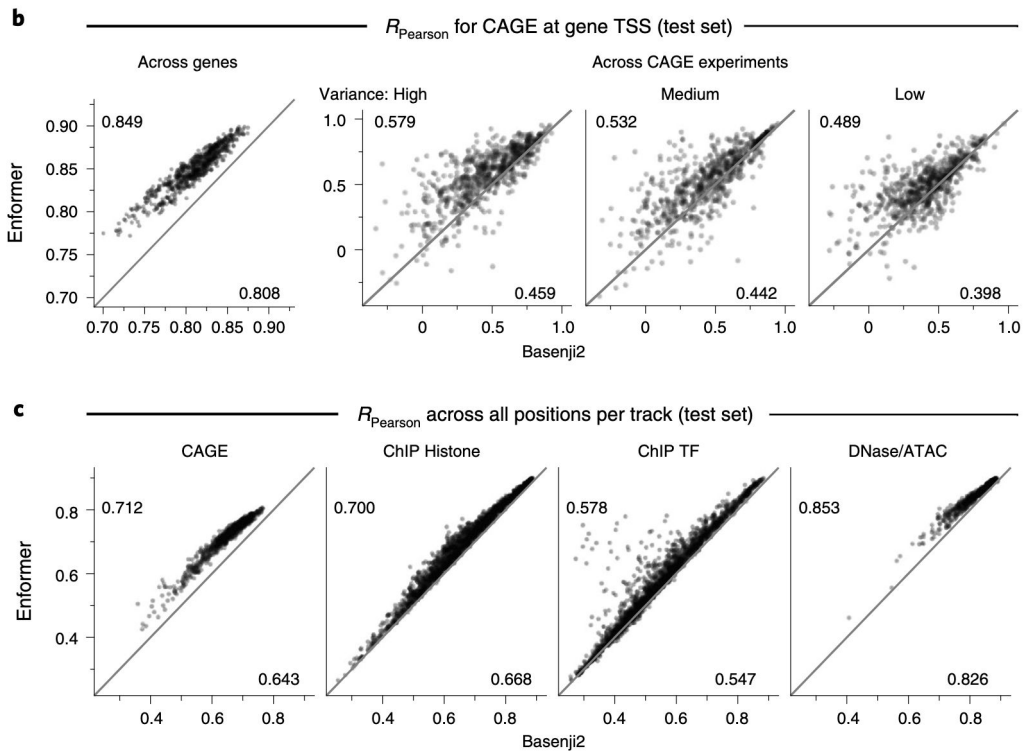
# Transformer layers

- Standard self-attention operations: keys, queries, values
- Attention based on keys/queries, and positional embeddings
- Modified positional embeddings
  - Standard NLP approach did not scale to large distances
  - Additive with key/query dot product (as in TransformerXL)
  - Combined several distance functions (exponential, central mask, Gamma)

# Key differences

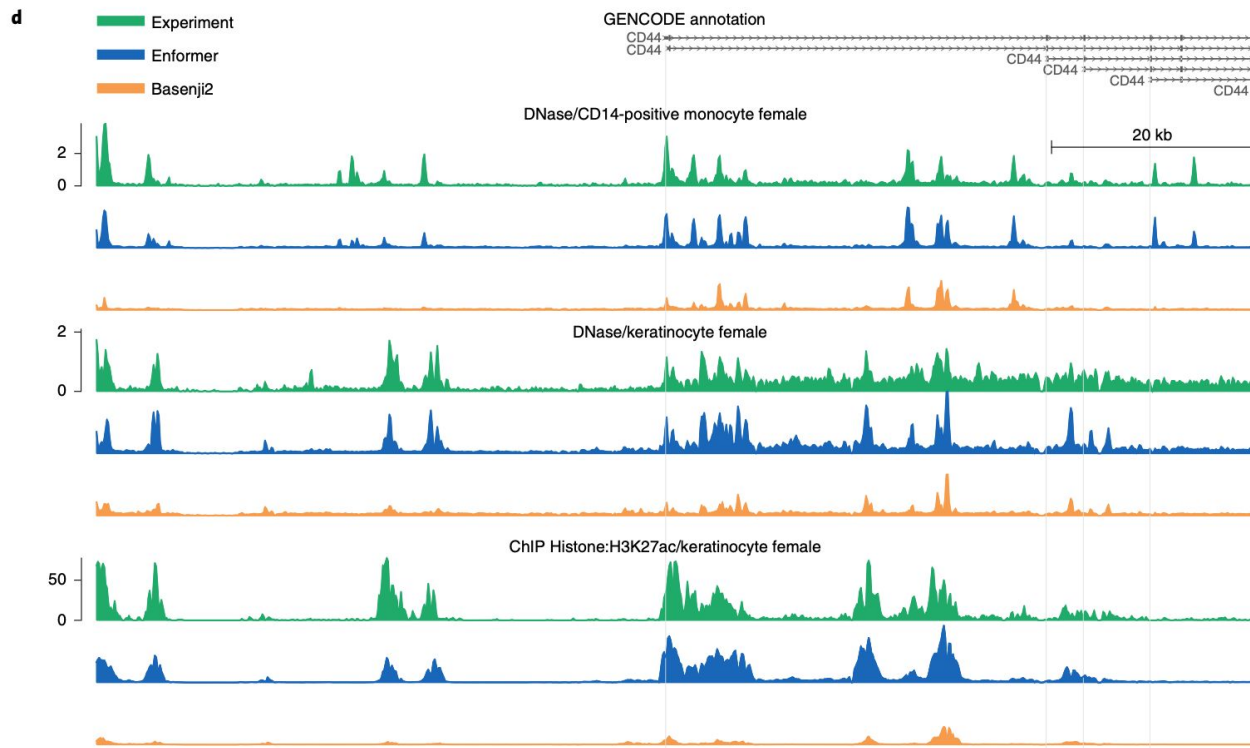
- Relative to DNA sequence models: large receptive field via self-attention, more channels, longer input (vs. Basenji)
- Relative to NLP: positional embeddings designed for long-range interactions

# Results





# Results (qualitative)



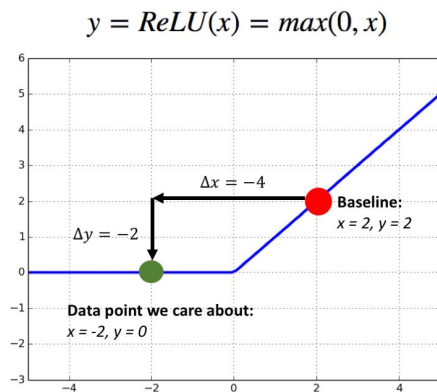
# Ablations

How to attribute performance improvement?

- Enformer outperforms Basenji (dilated convolutions)
- Consistent across various model sizes, # layers, # data points
- Reducing receptive field size hurts performance
- Standard positional embeddings hurt performance

# interpretation pf the model

- Gradient based
- Attention based

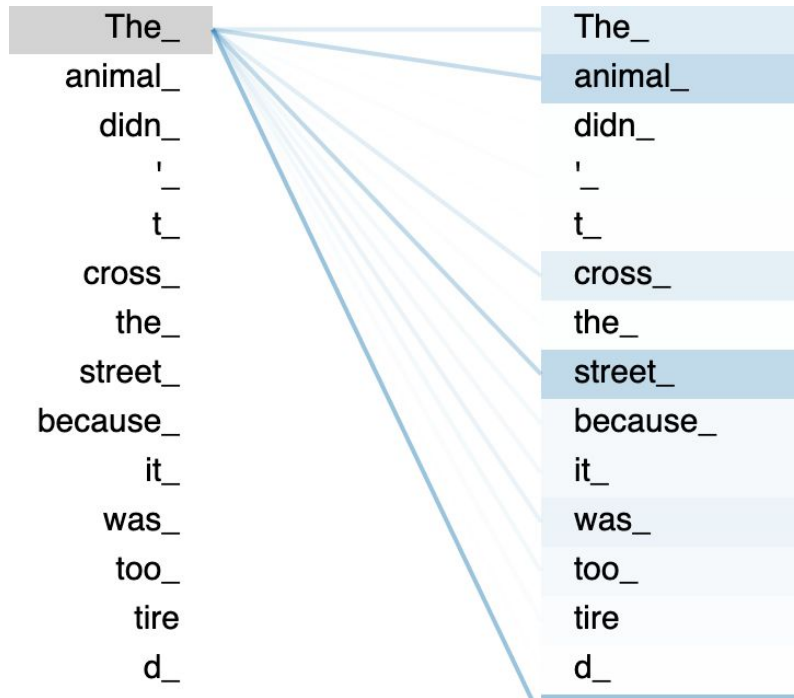


1. Calculating the slope

$$\frac{\Delta y}{\Delta x} = \frac{-2}{-4} = 0.5$$

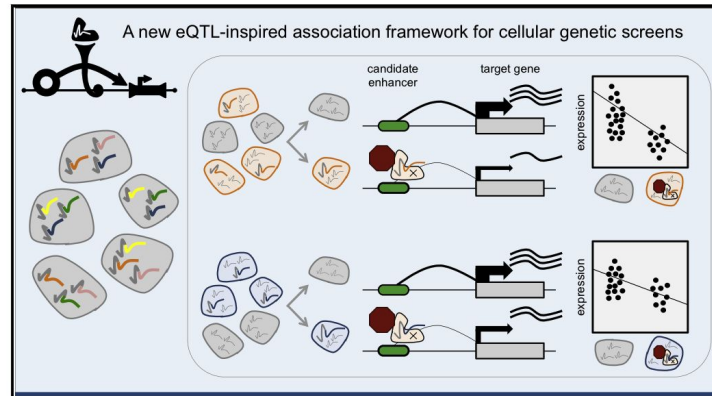
2. Finding the feature importance

$$\Delta x \times \frac{\Delta y}{\Delta x} = -4 \times 0.5 = -2$$



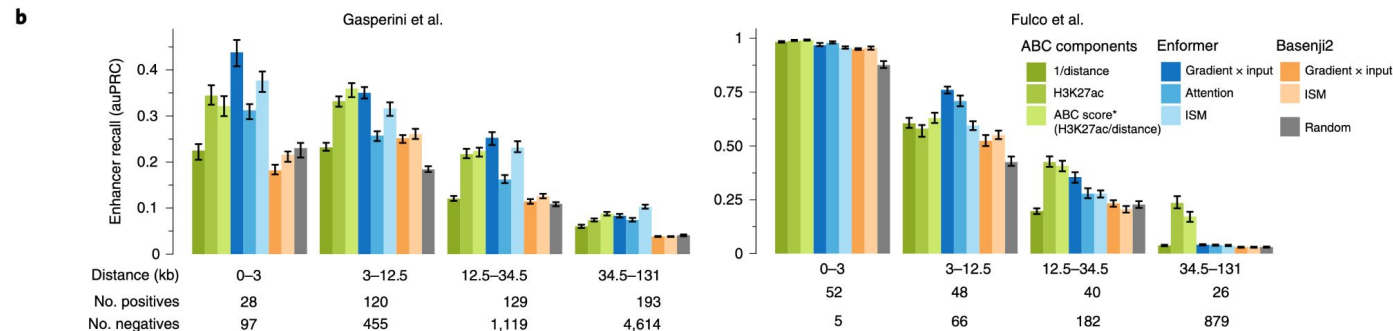
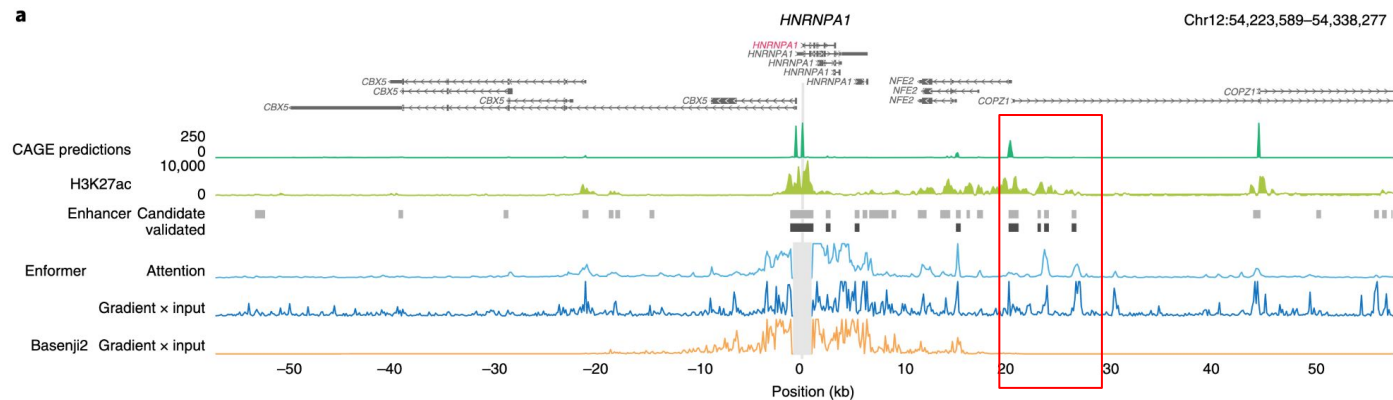
# Enhancer prioritization

- CRISPRi assay perturbing the enhancer of interest while measuring the expression change of the gene



Gasparini, M. et al. A Genome-wide Framework for Mapping Gene Regulation via Cellular Genetic Screens. *Cell* 176, 377-390.e19 (2019).

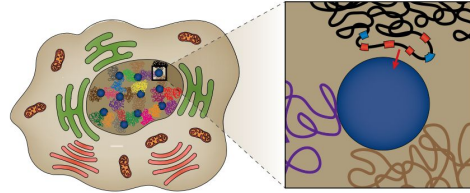
# Enformer attends to cell-type-specific enhancers



# Insulator

- TAD: Topological Association Domain

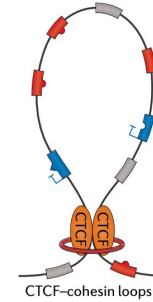
**b** Facilitating



- Relocation to transcription factories
- Phase-separated condensates



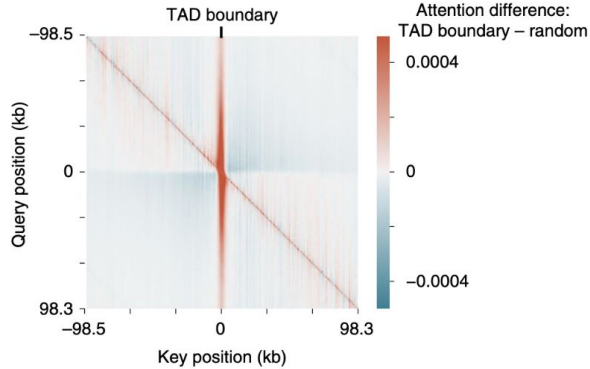
TADs and compartments



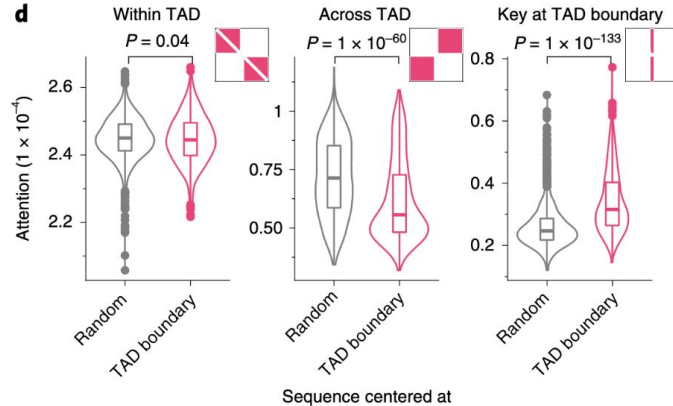
CTCF-cohesin loops

Schoenfelder, S. & Fraser, P. Long-range enhancer-promoter contacts in gene expression control. *Nat Rev Genet* 20, 437–455 (2019).

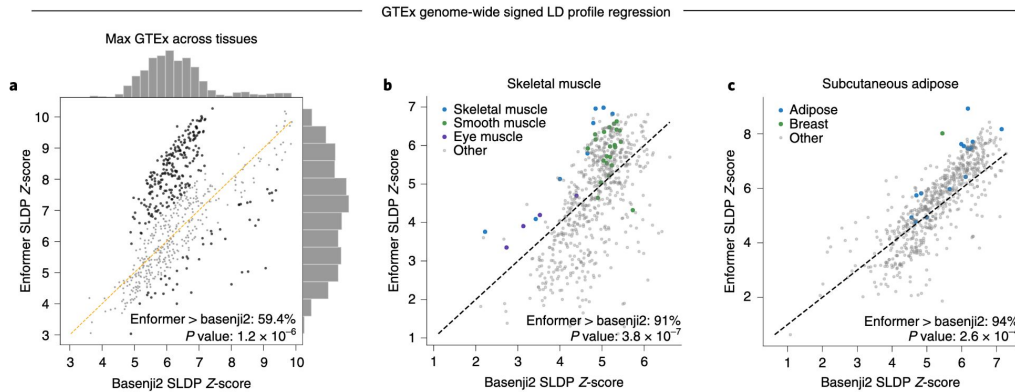
**c**



**d**

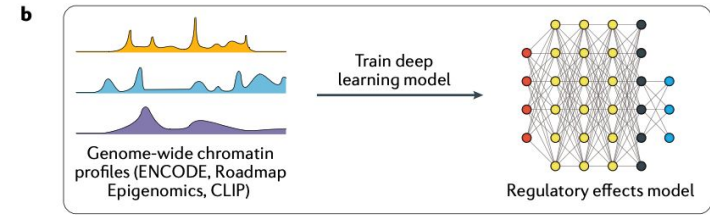


# Variant effect prediction on eQTL data.

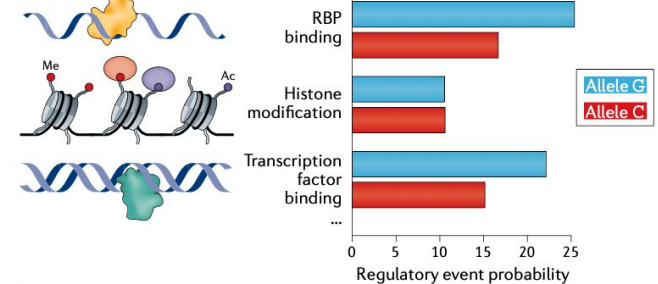


**a** Input: genomic sequences chr6:137918071 G > C

Allele G ...TGTTGTCAATGCTACGGAGC...  
Allele C ...TGTTGTCAATCGCTACGGAGC...

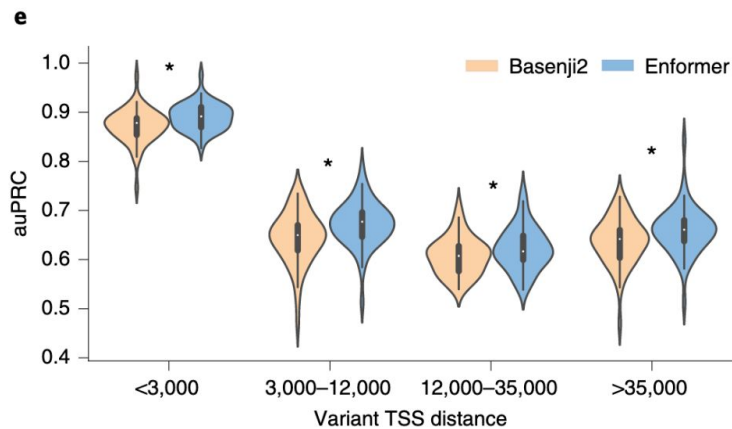
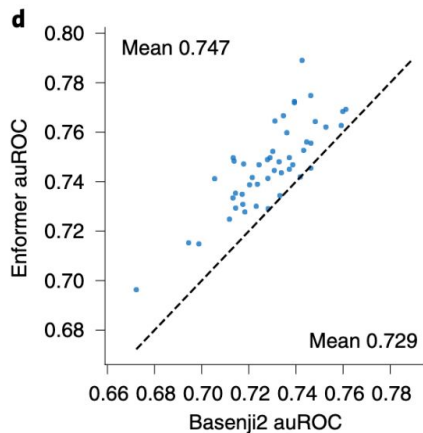


**c** Output: allele-specific regulatory predictions



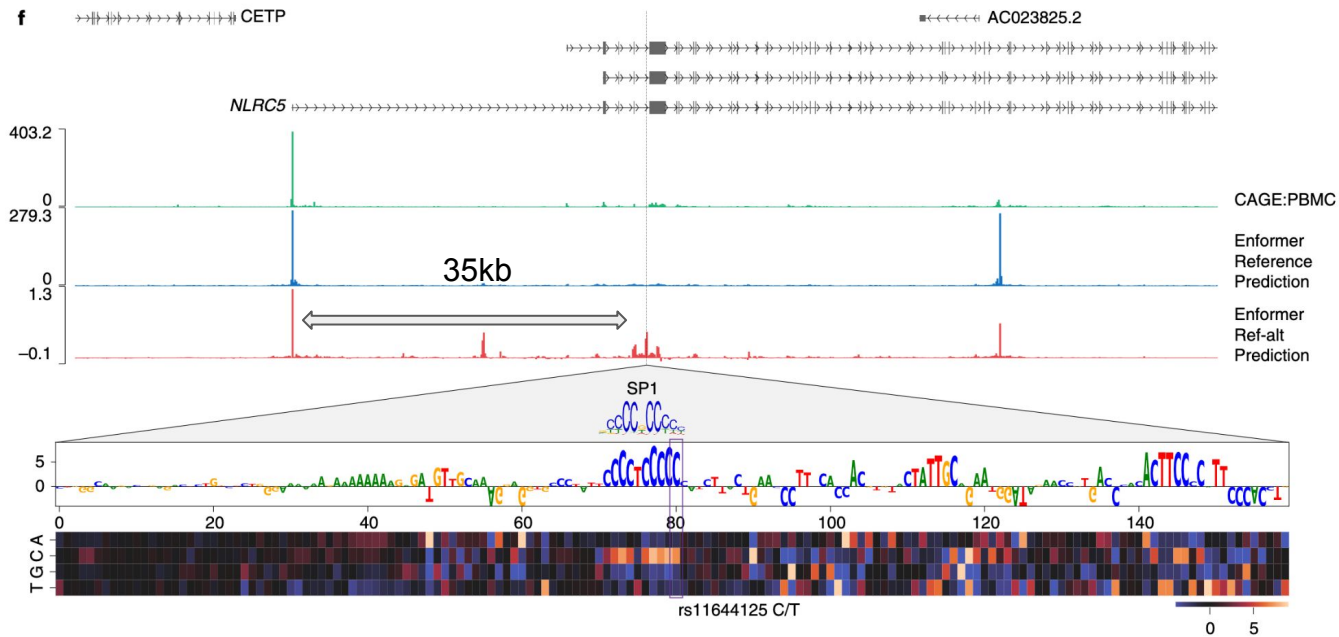
# Variant effect prediction on eQTL data.

GTEx SuSiE fine-mapped variant classification

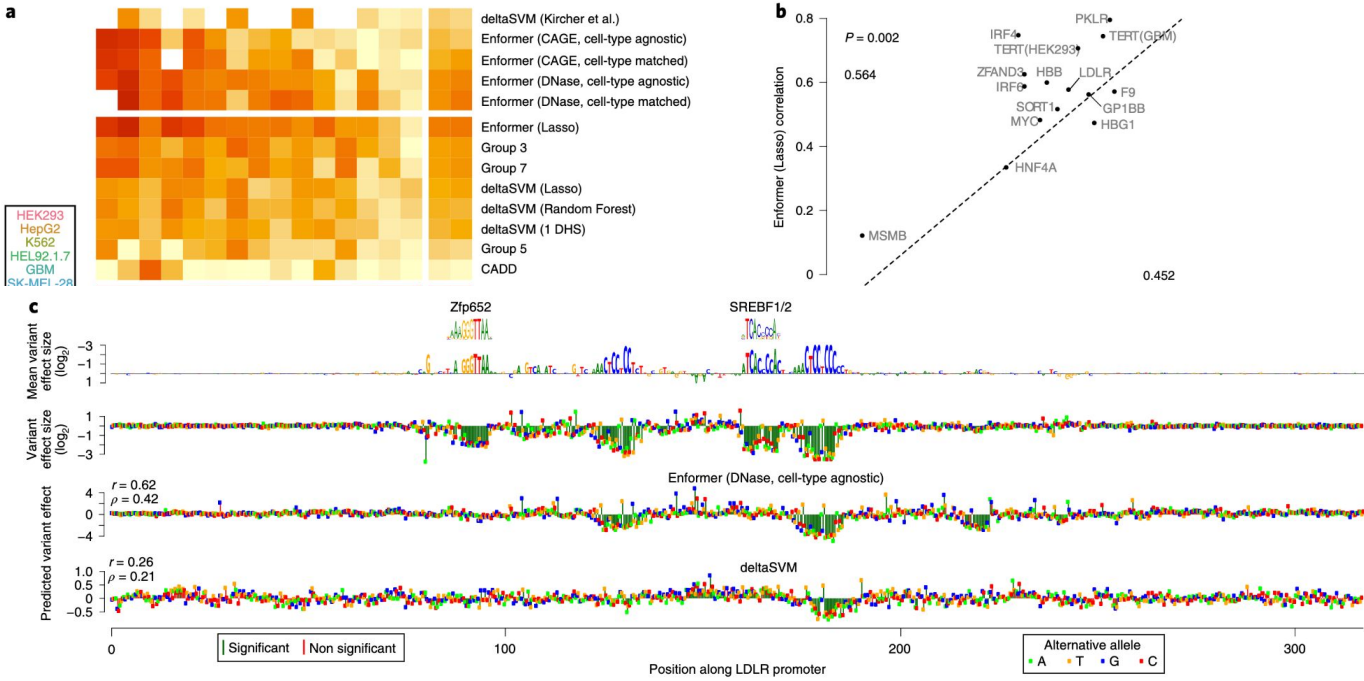




# Variant effect prediction on eQTL data.



# Enformer improves MPRA mutation effect prediction



# Take-home message

- Self-attention can extend the receptive field
- Self-attention can capture long-distance interaction
- Enformer can help us build a more comprehensive regulatory map than before

# Discussion

- Cis-regulatory, how about trans-regulatory?
- Covariates matrix - like gene expression(Expecto)?
- Opportunities for large pre-trained models in genomics?