

Deep Learning Theory and Applications

Sequence-based Bioinformatics applications

Yale

CPSC/AMTH 663

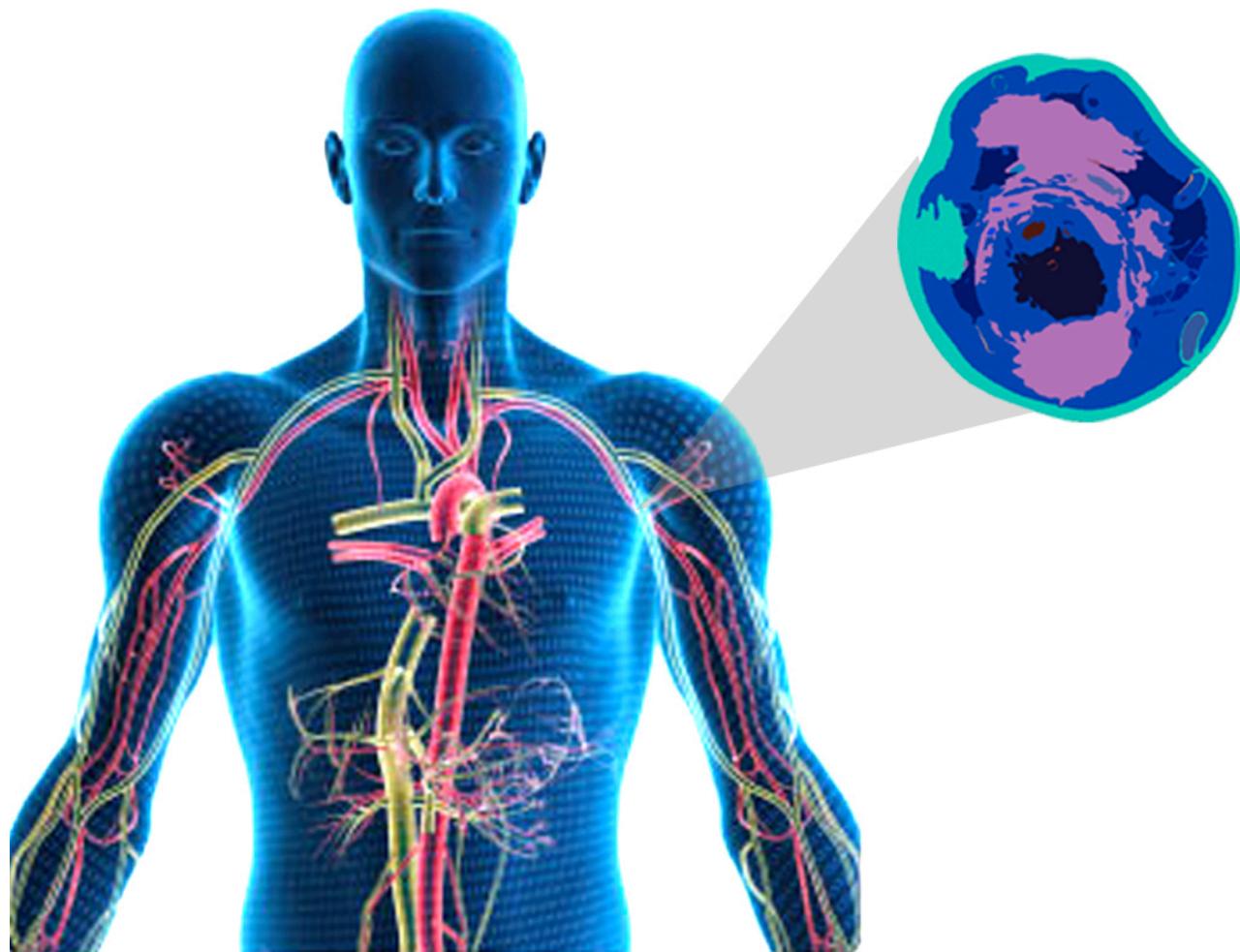




Outline

1. Autoencoders
2. Denoising Autoencoder
3. Contractive Autoencoder
4. Variational Autoencoder

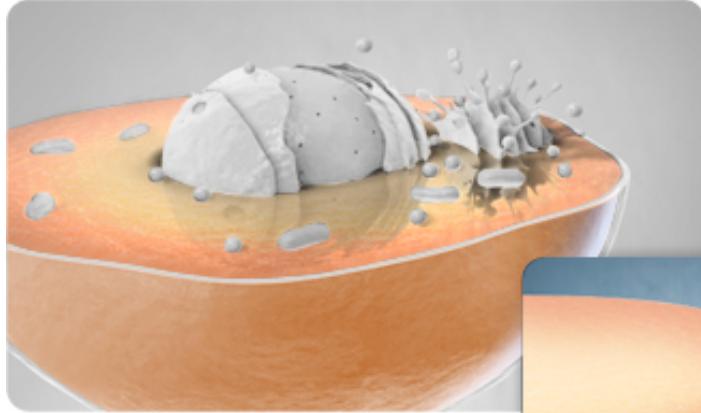
Biology Crash Course





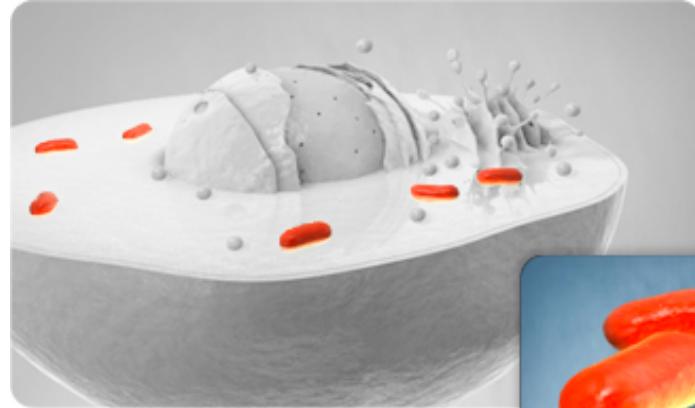
Parts of Cell

<https://ghr.nlm.nih.gov/primer/basics.pdf>



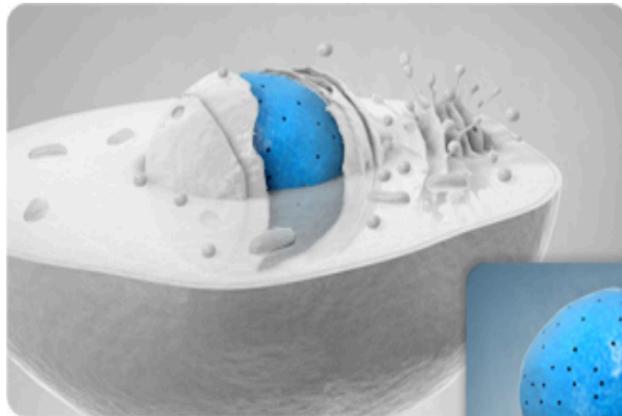
Cytoplasm

U.S. National Library of Medicine



Mitochondria

U.S. National Library of Medicine



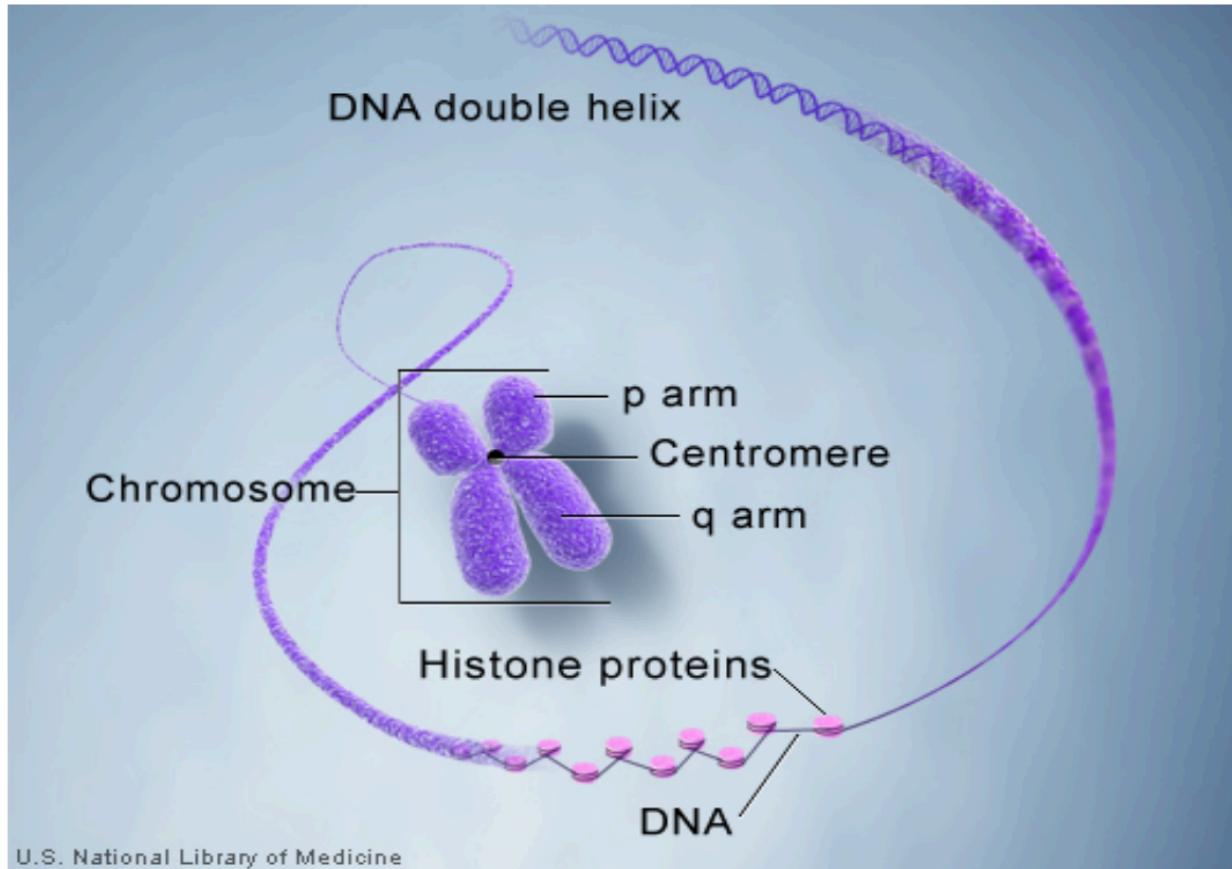
Nucleus



U.S. National Library of Medicine



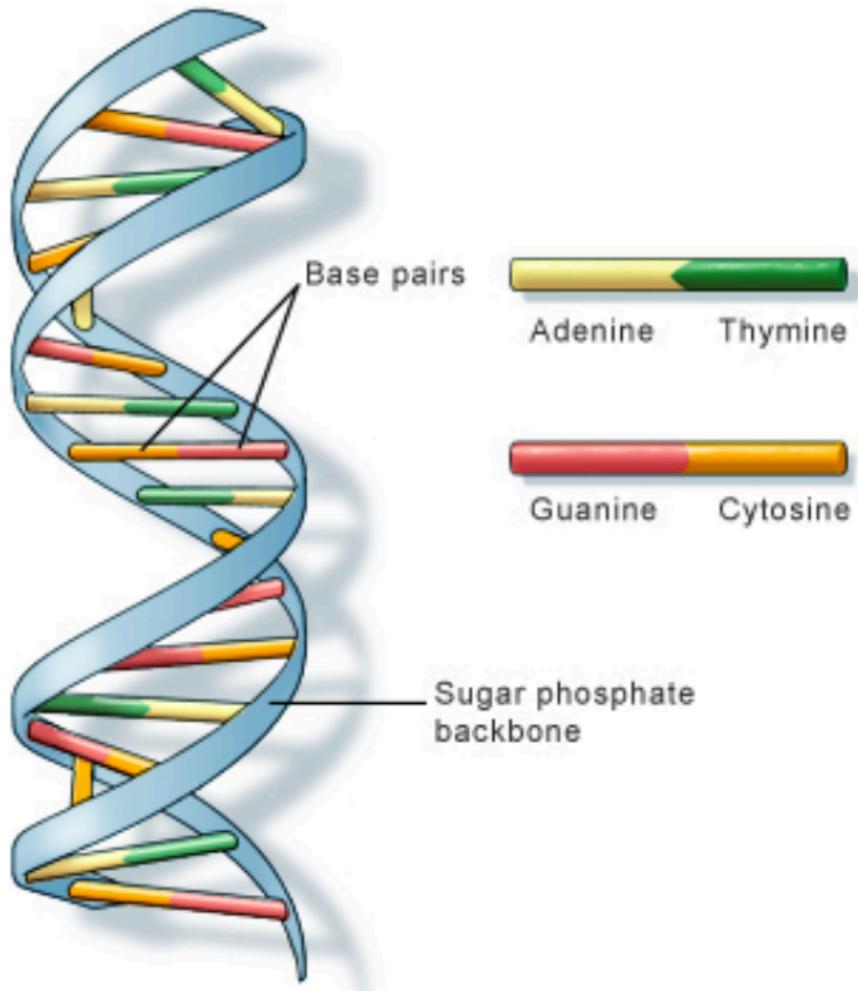
Chromosomes



DNA and histone proteins are packaged into structures called chromosomes.

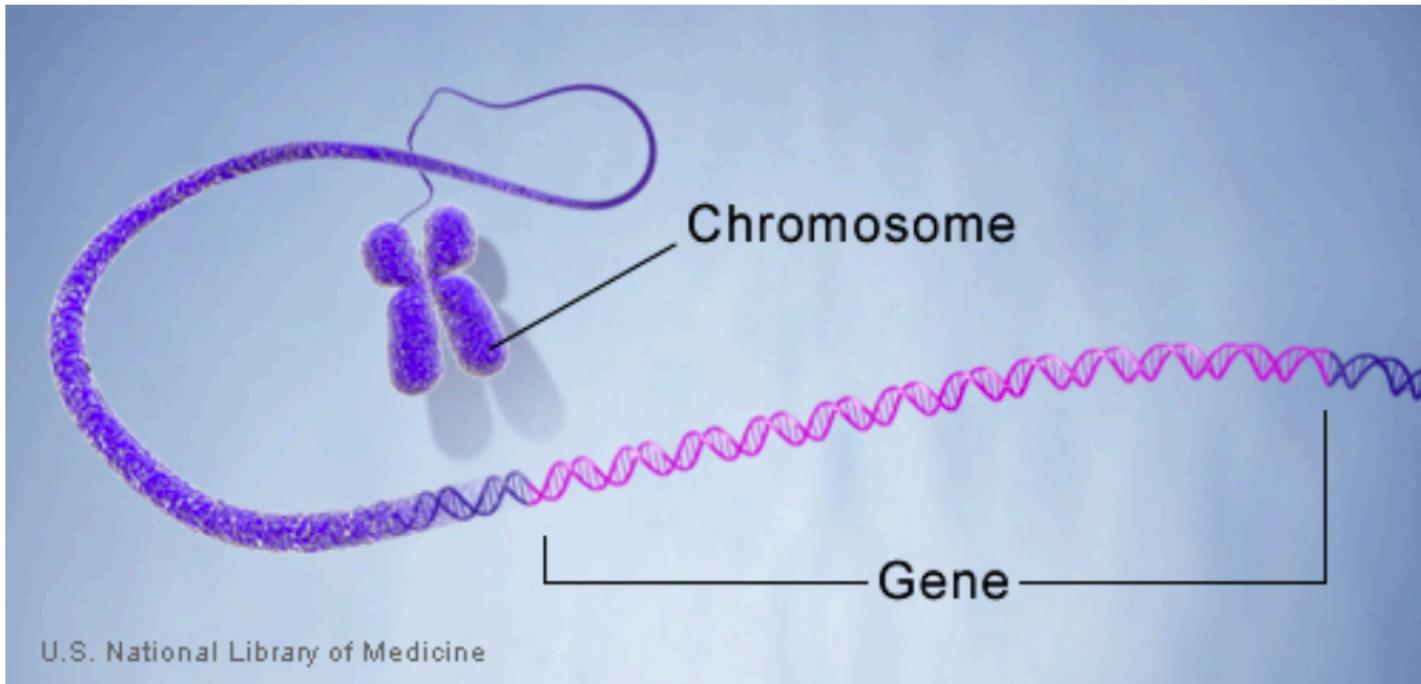


DNA





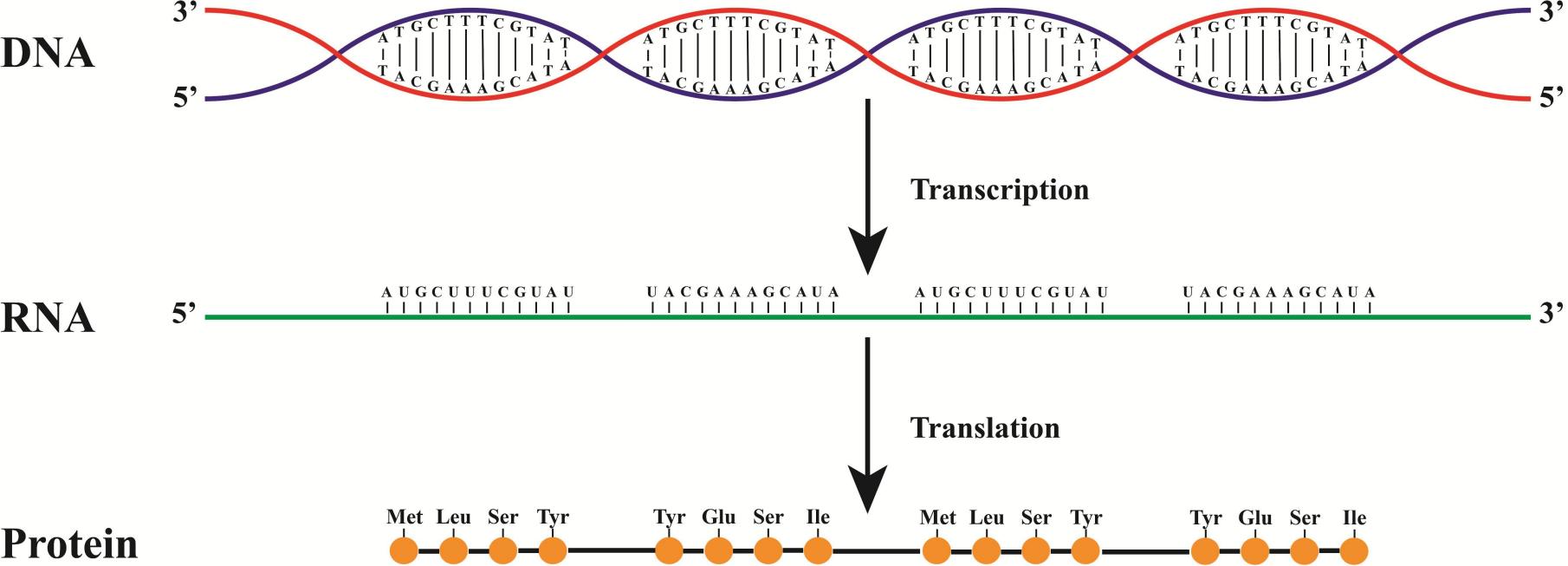
Genes



Genes are made up of DNA. Each chromosome contains many genes.



Central Dogma of Biology

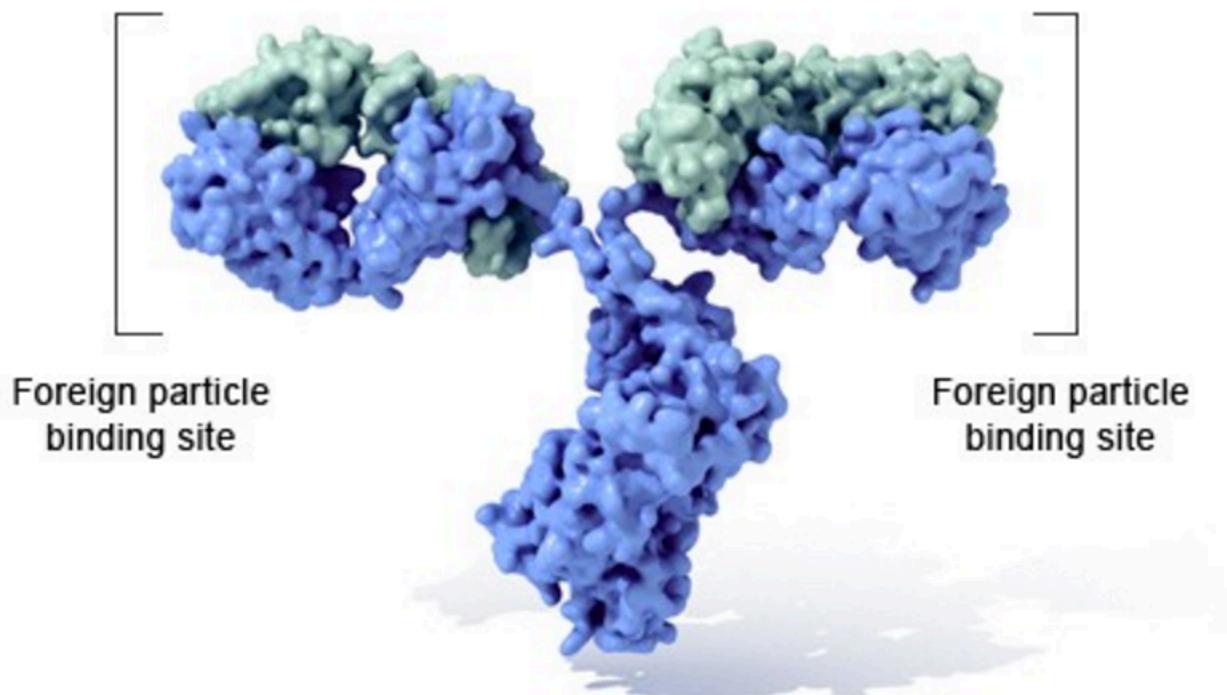


mRNA reveals what proteins the cell is making



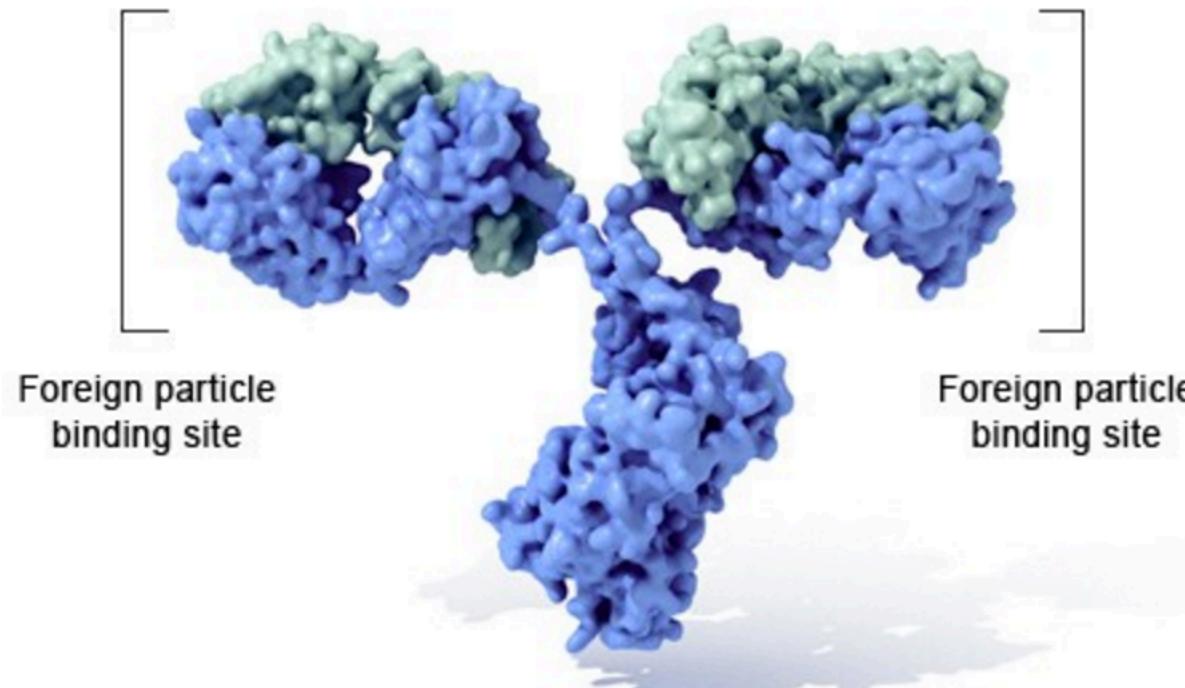
What do proteins do?

Binds to surface receptors, help cells recognize invaders, viruses, chemical signals





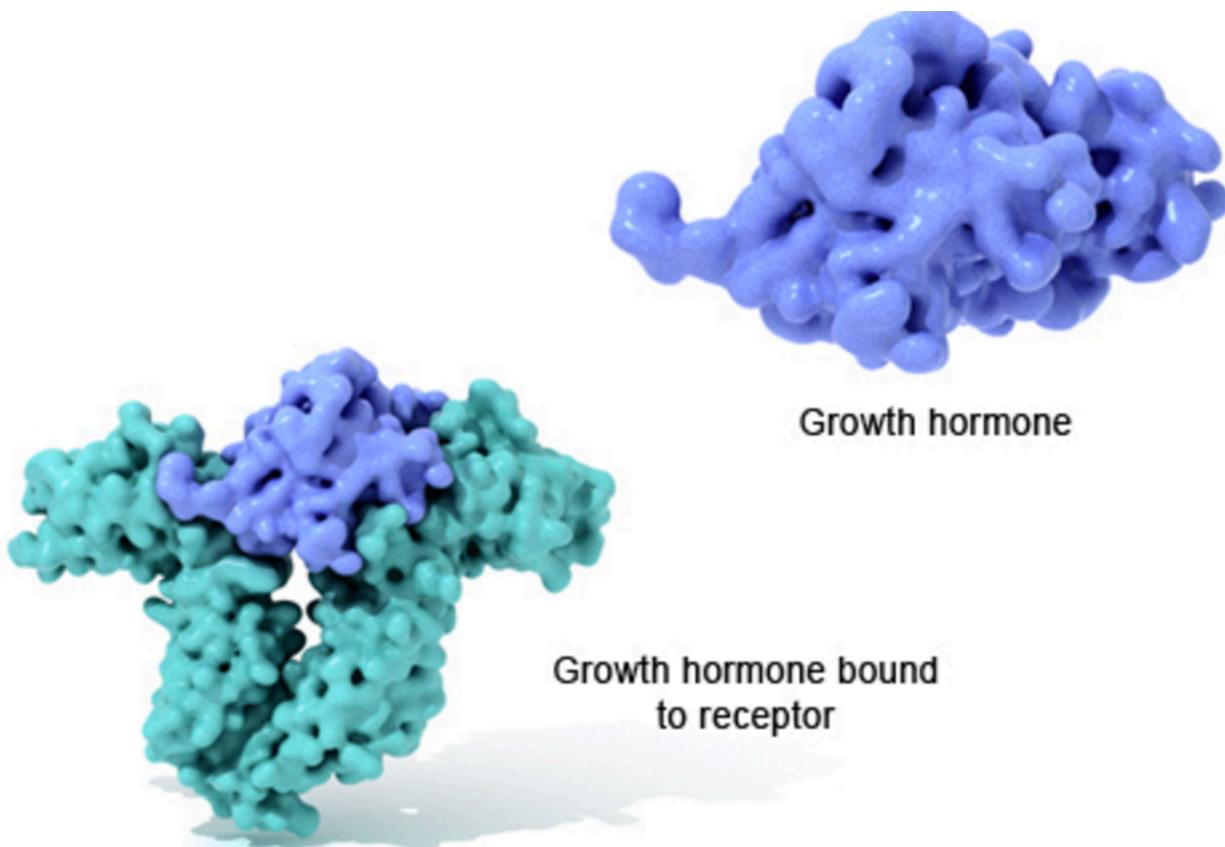
What do proteins do?



Antibodies: Binds foreign bodies, viruses, alerts immune system



Chemical Messenger

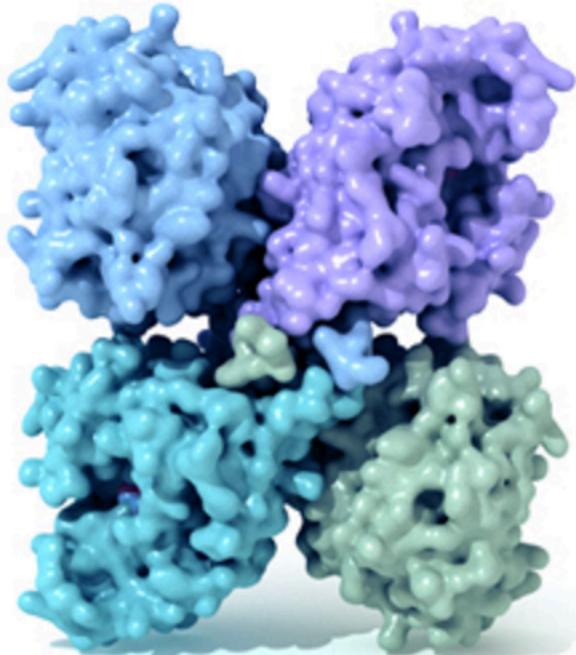


Hormones, growth factors and other extracellular signals



Enzyme

Carries out chemical reactions



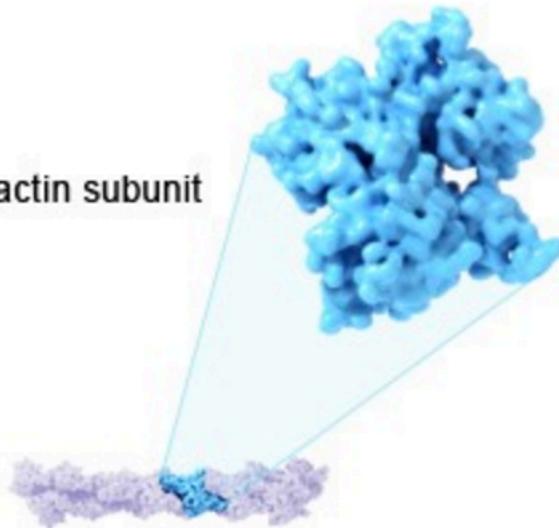
Phenylalanine hydroxylase
protein consisting of 4 subunits



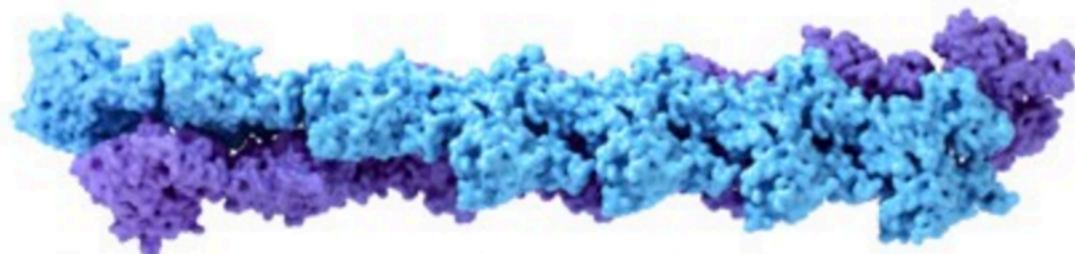
Structural Components

Actin

Single actin subunit



Actin filament consisting
of multiple subunits

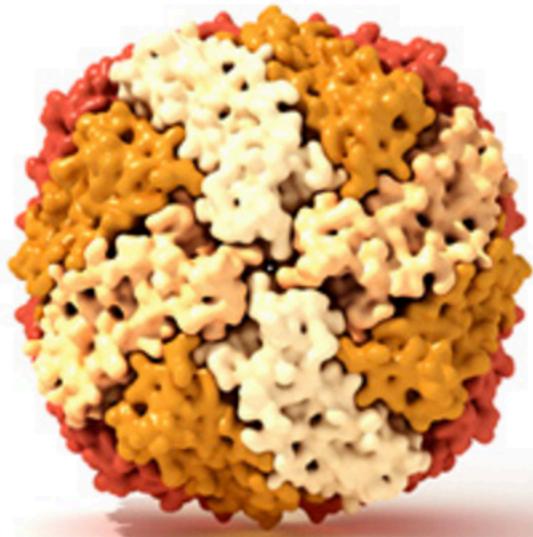




Transport/Storage

Ferritin

Ferritin protein consisting
of 24 subunits



Single ferritin subunit

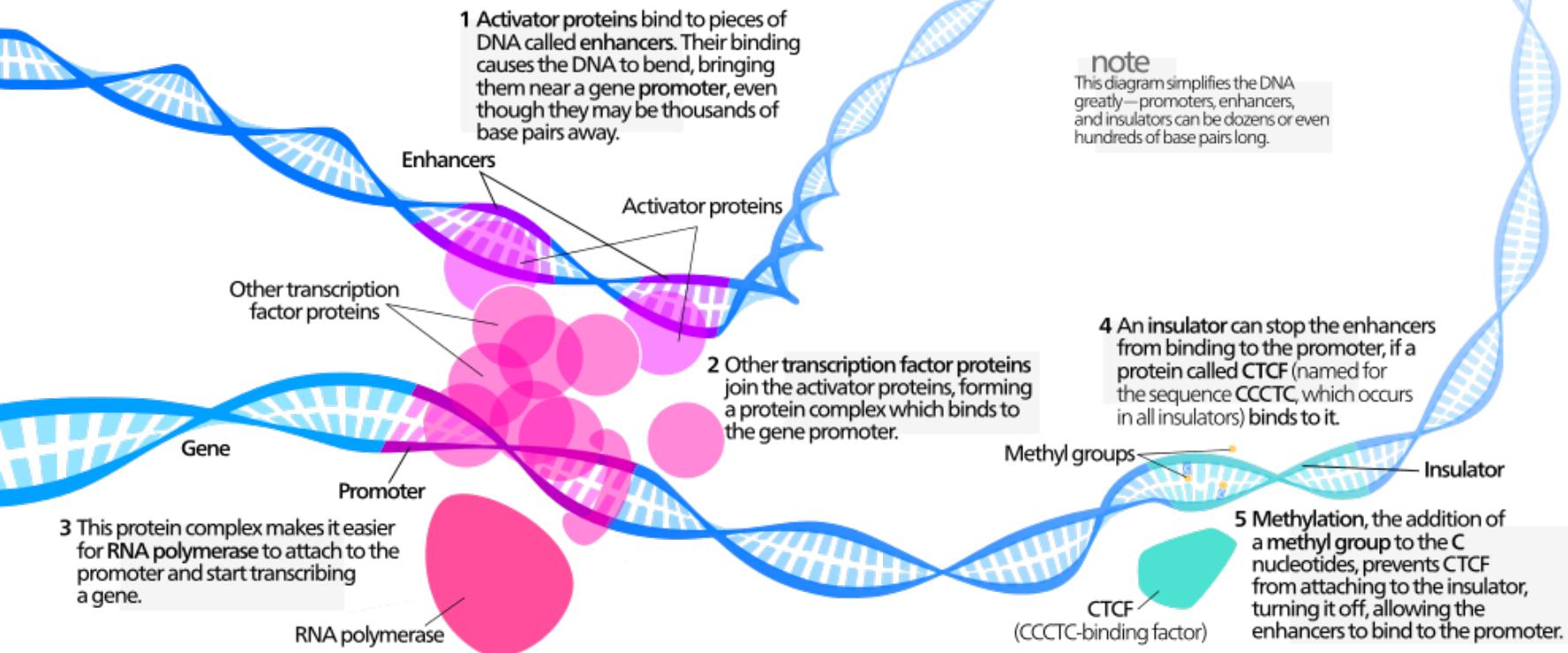


Cross section



Proteins binding to DNA

transcription factors of eukaryotic cells





DNA sequence understanding

- DNA and RNA-binding proteins can be "approximately sequence specific"
- Knowing the sequence specificity can allow us to develop a model of the regulatory network in a biological system
- Mutations can disrupt the regulatory system due to this sequence specificity
- Can also identify disease-causing variants



Deep Bind

Analysis | Published: 27 July 2015

Predicting the sequence specificities of DNA- and RNA-binding proteins by deep learning

Babak Alipanahi, Andrew Delong, Matthew T Weirauch & Brendan J Frey ✉

Nature Biotechnology **33**, 831–838 (2015) | Download Citation ↓



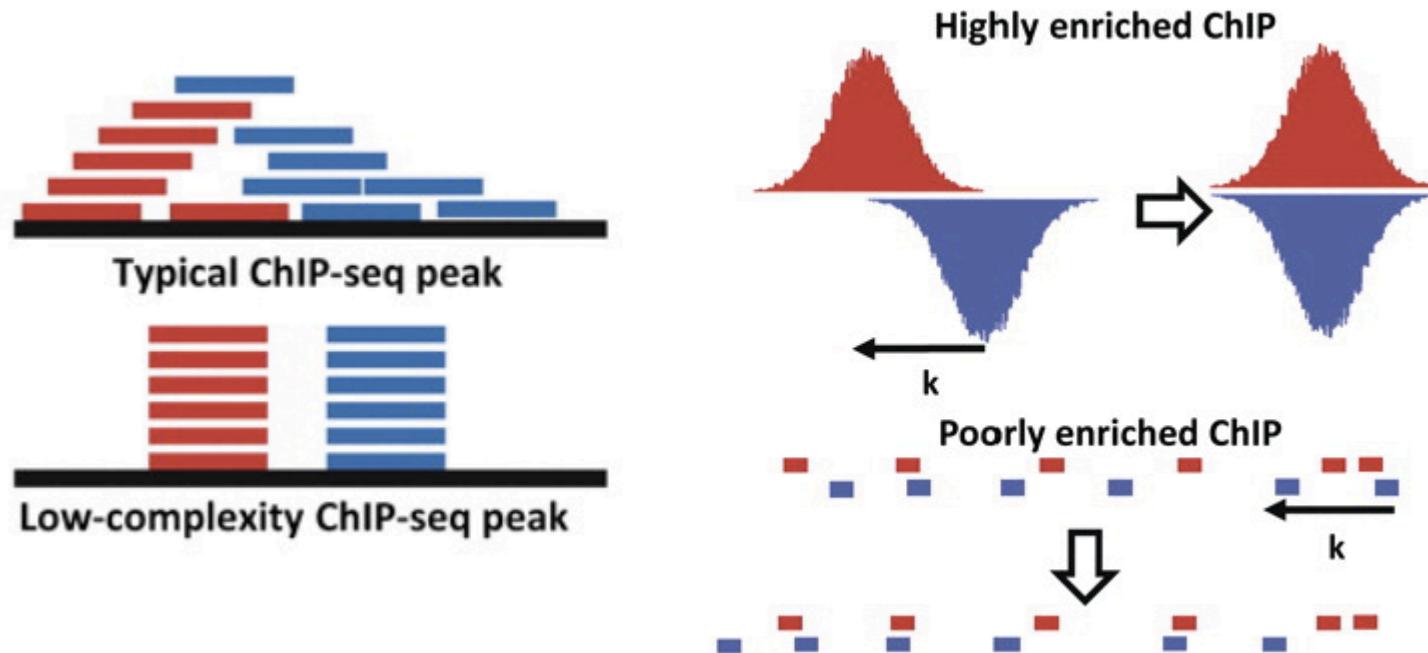
Where does binding data come from?

- Protein binding microarrays: contains a specificity coefficient for each probe sequence
- Chromatin immunoprecipitation (ChIP)-seq: provides a ranked list of putatively bound sequences of varying length
- HT-SELEX: generates a set of very high affinity sequences.
- LOTS of sequences:
 - A typical high-throughput experiment measures between 10,000 and 100,000 sequences



Artifacts, biases

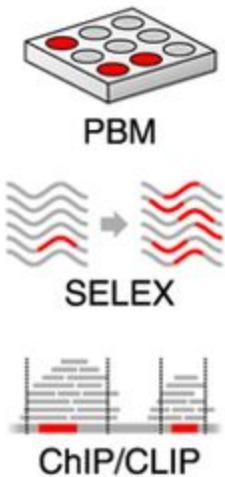
- Each data acquisition technology has its own artifacts, biases and limitations
- ChIP-seq reads often localize to “hyper-ChIPable” regions of the genome near highly expressed genes



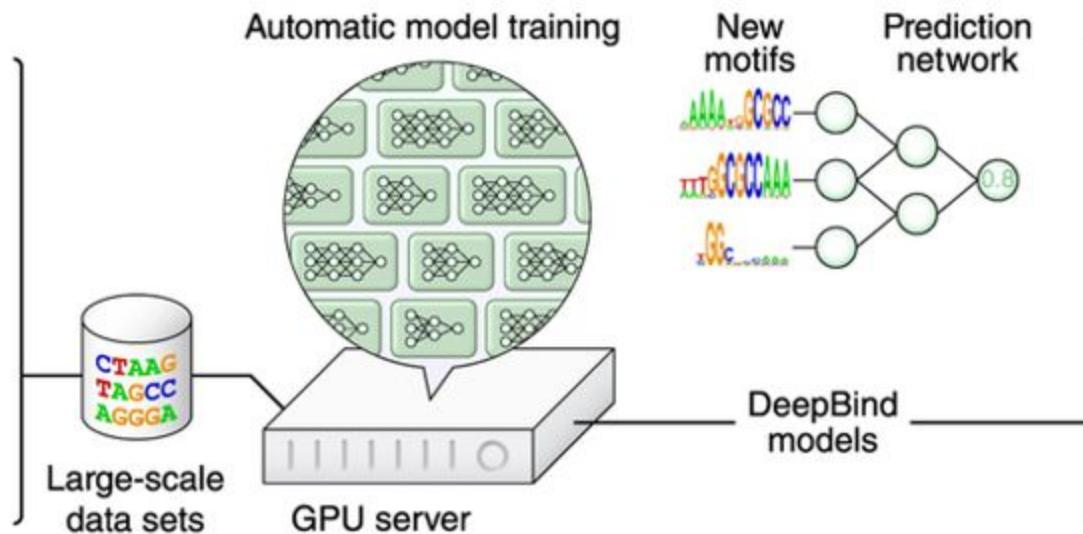


DeepBind

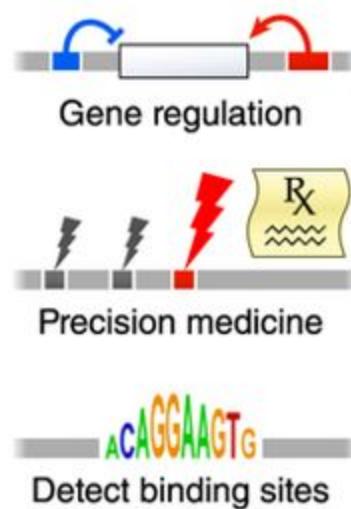
1. High-throughput experiments

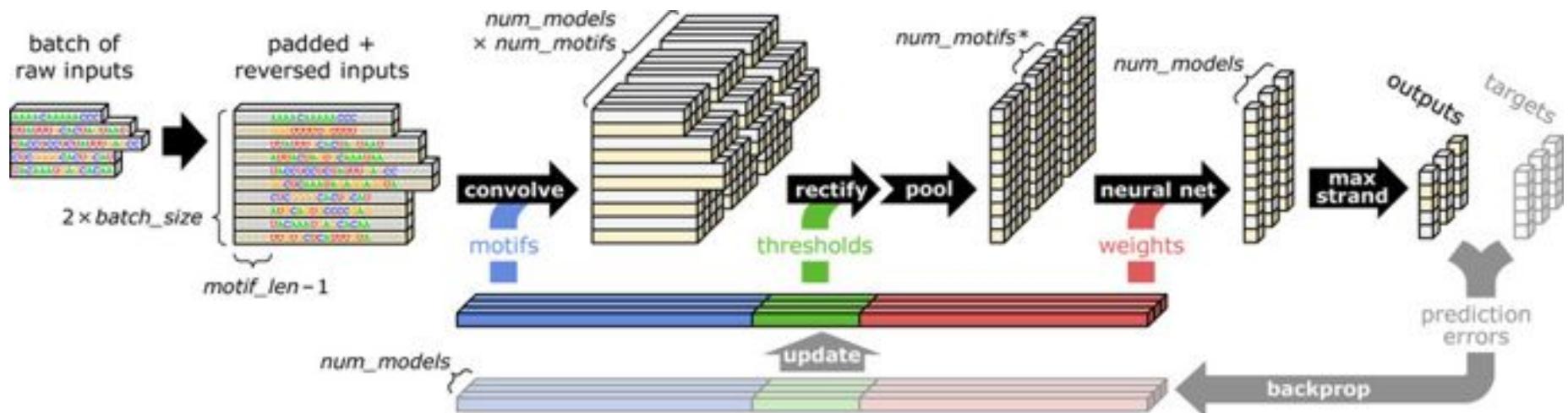


2. Massively parallel deep learning



3. Community needs







Training

- Uses a set of sequences with experimentally determined “binding scores”
- Sequences are of varying length 14-101 nucleotides
- Binding scores can be real-valued or binary measurements



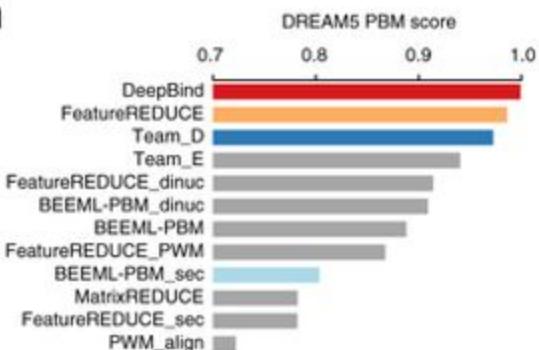
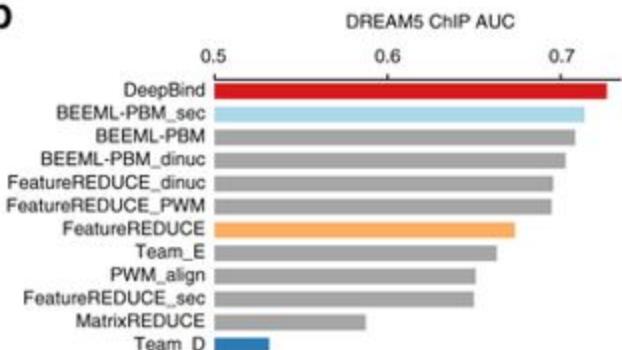
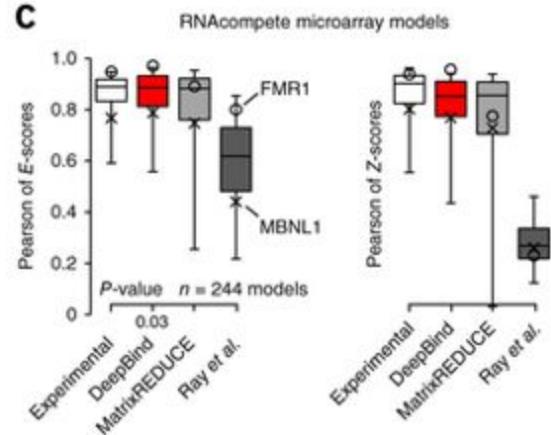
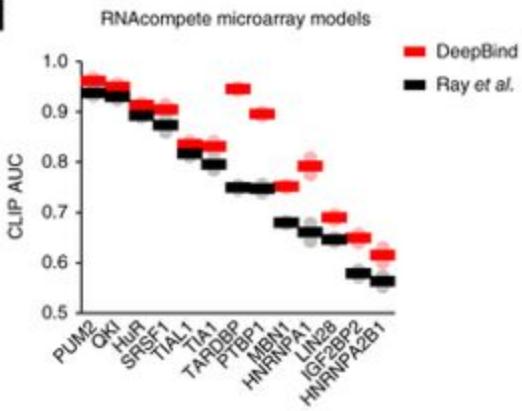
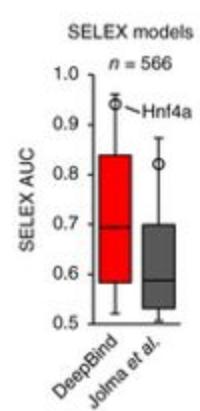
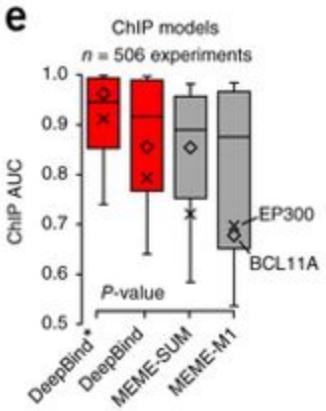
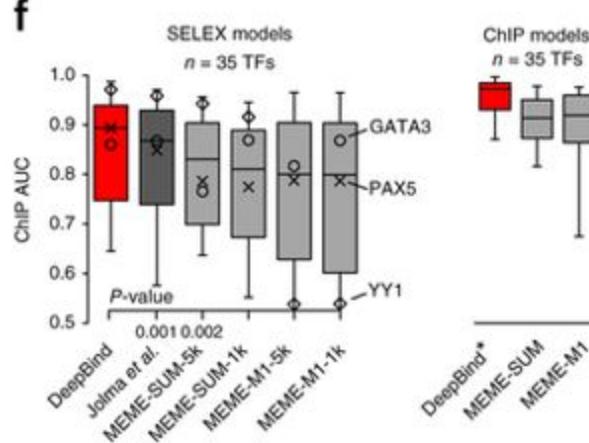
Architecture

- Convolutional stage: scans a set of ***motif detectors*** across the sequence space
 - Each motif detector is a 4xm matrix
- Rectification stage isolates positions with good pattern match by shifting the response detector of M_k by b_k and clamping negative values to 0
- Pooling stage computes maximum and average of each motif detector's rectified response across the sequence
 - Maximizing helps identify the presence of longer motifs
 - Averaging helps identify cumulative effect of shot motifs
- Values are fed into non-linear neural network which combines responses to produce score



Repository

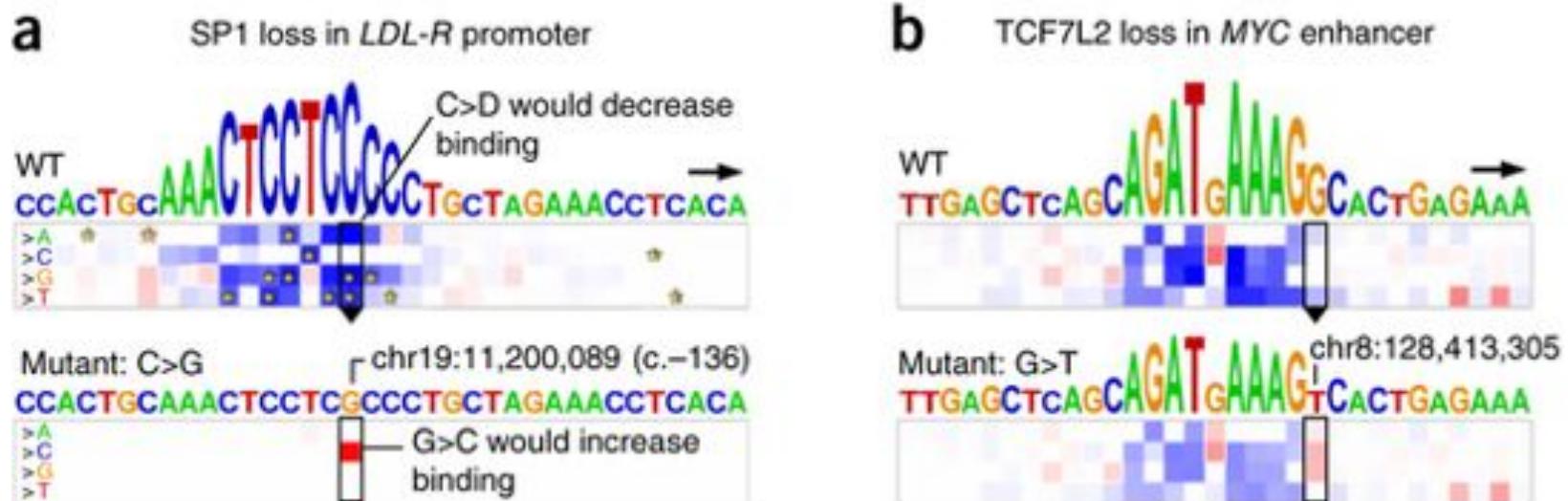
- (<http://tools.genes.toronto.edu/deepbind/>)
- 927 DeepBind models representing 538 distinct transcription factors and 194 distinct RBPs

**a****b****c****d****e****f**



Identifying damaging genetic variants

- ‘mutation map’ illustrates the effect that every possible point mutation in a sequence may have on binding affinity



Visualization shows how important each base is for the analysis by the height of the Base letter

Heat map indicates how much mutation will increase or decrease binding score



Coda

Denoising genome-wide histone ChIP-seq with convolutional neural networks

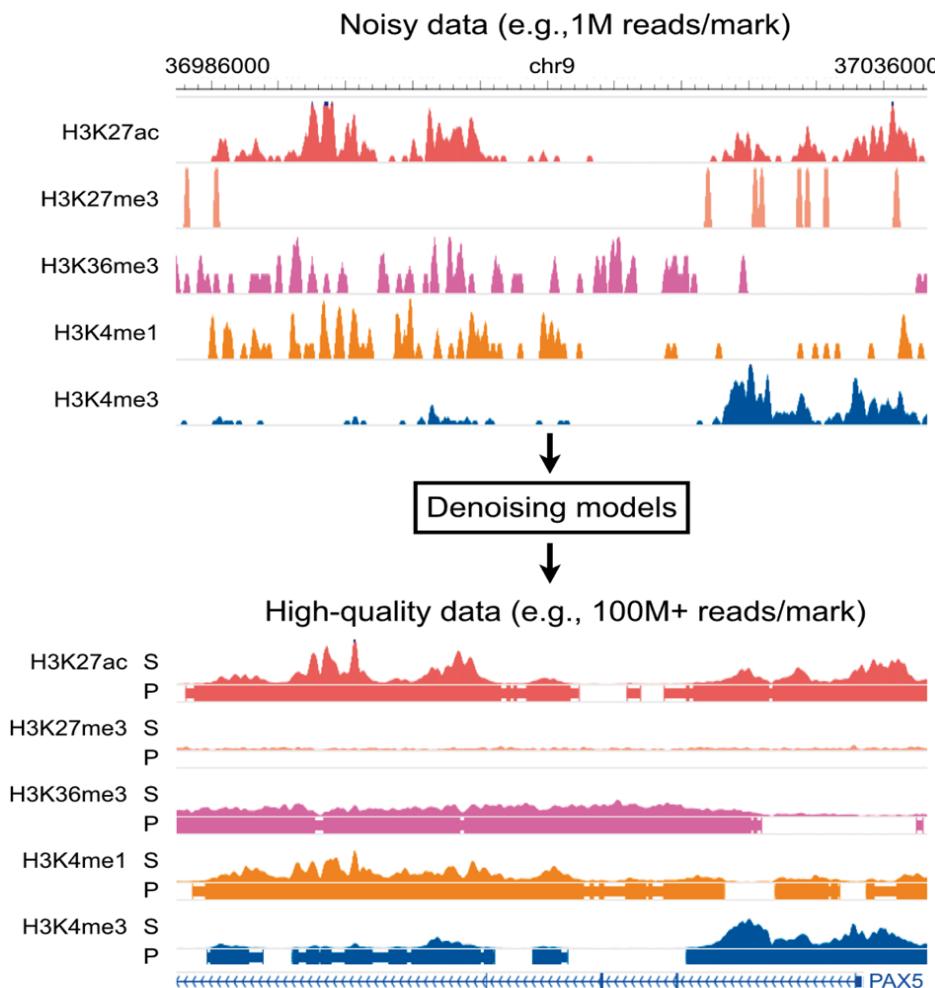
Pang Wei Koh, Emma Pierson, Anshul Kundaje  Author Notes

Bioinformatics, Volume 33, Issue 14, 15 July 2017, Pages i225–i233,

<https://doi.org/10.1093/bioinformatics/btx243>

Published: 12 July 2017

Denoising Chip-seq reads



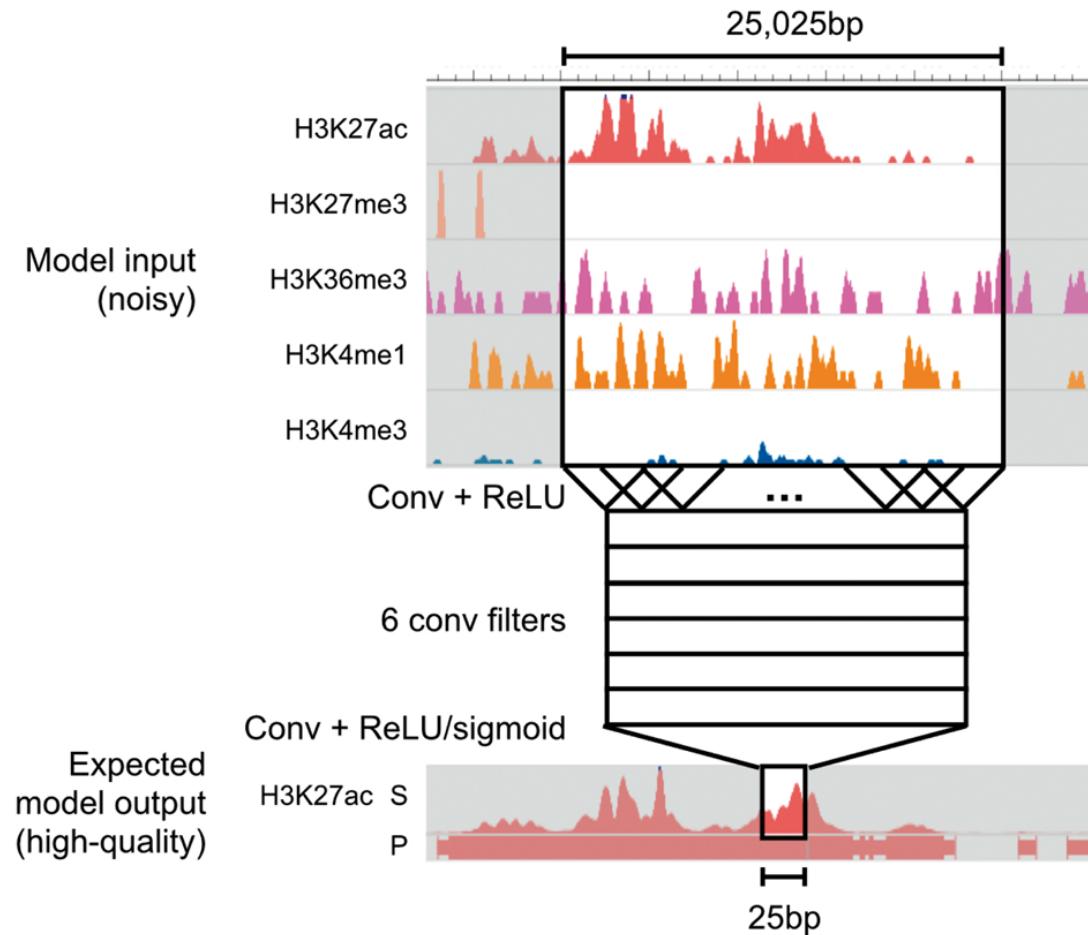


Training data

- Pairs of matching ChIP-seq datasets of the same histone modifications in the same cell type
 - one high-quality a
 - the other noisy
- The noisy dataset used in training can be derived:
 - computationally (e.g. by subsampling the high-quality data)
 - experimentally (e.g. by conducting the same ChIP-seq experiment with fewer input cells).



Architecture

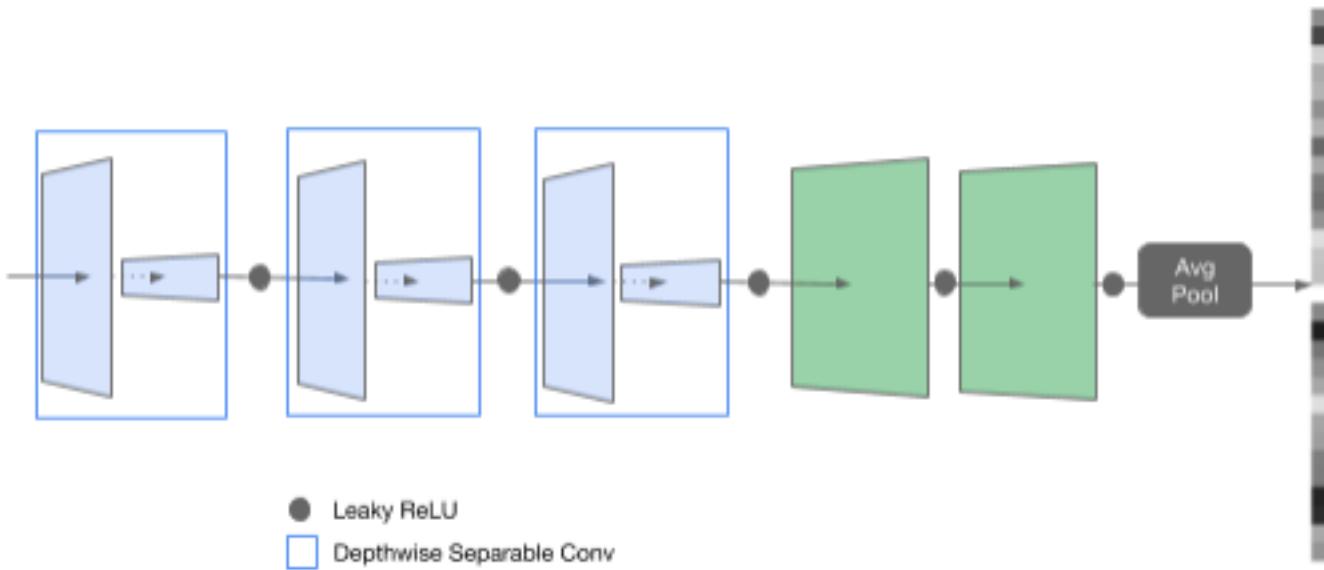


Alternative architectures for this task?



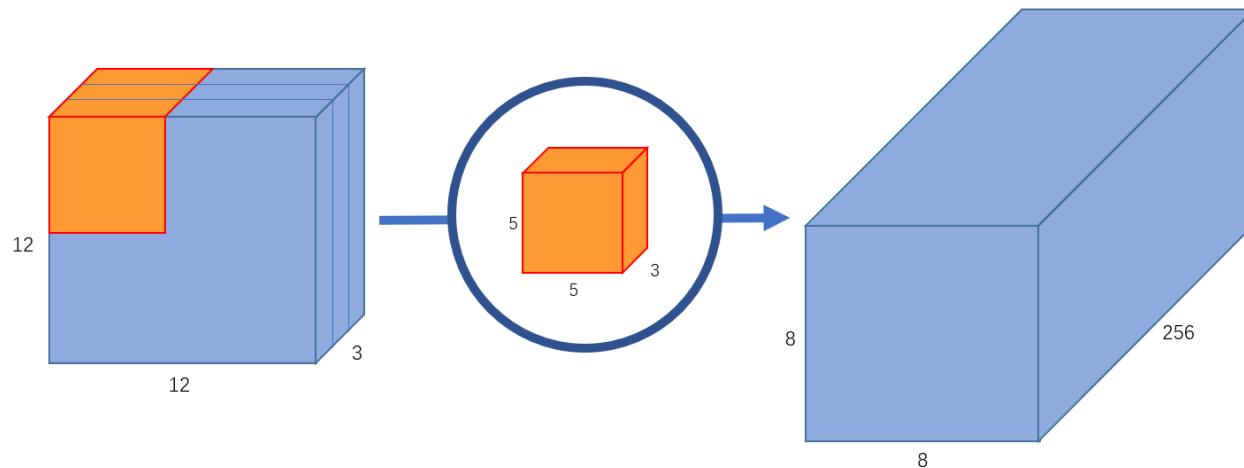
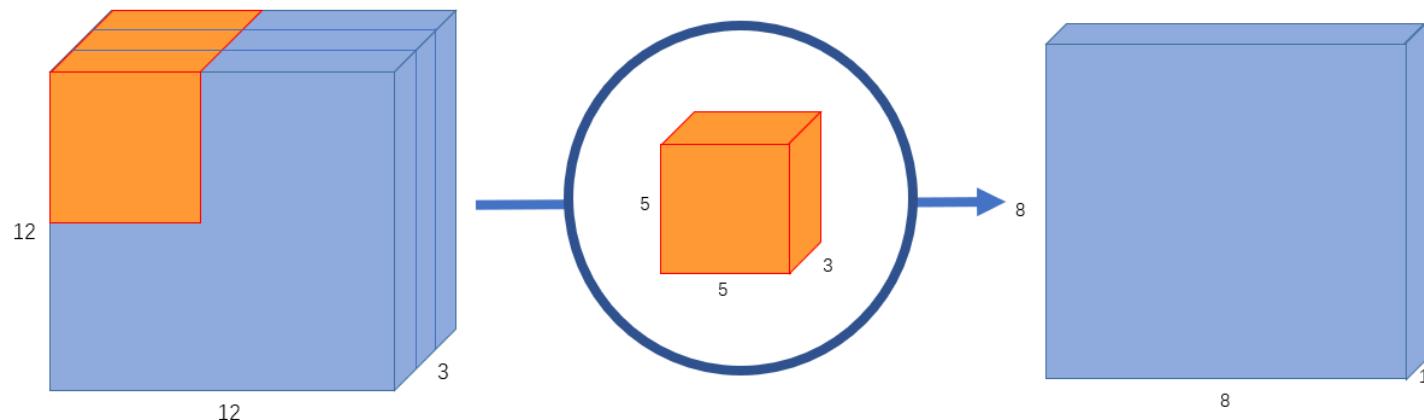


- Many important problems in bioinformatics can be framed as determining a mapping from short biological query sequences to salient categorical or numerical labels
- Examples include:
 - Taxonomic classification or binning;
 - prediction of protein function,
 - gene properties, or pathogenicity;
 - read filtering for contaminants;
 - RNA-seq quantification





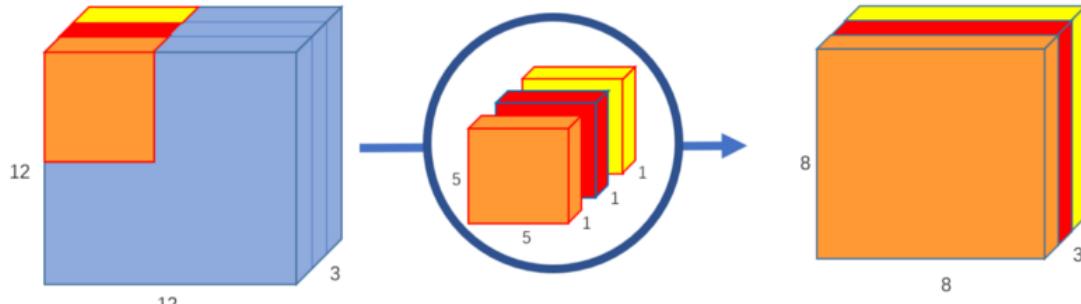
Normal Convolution



256 5x5x3 kernels
that move 8x8 times
1228800 multiplications



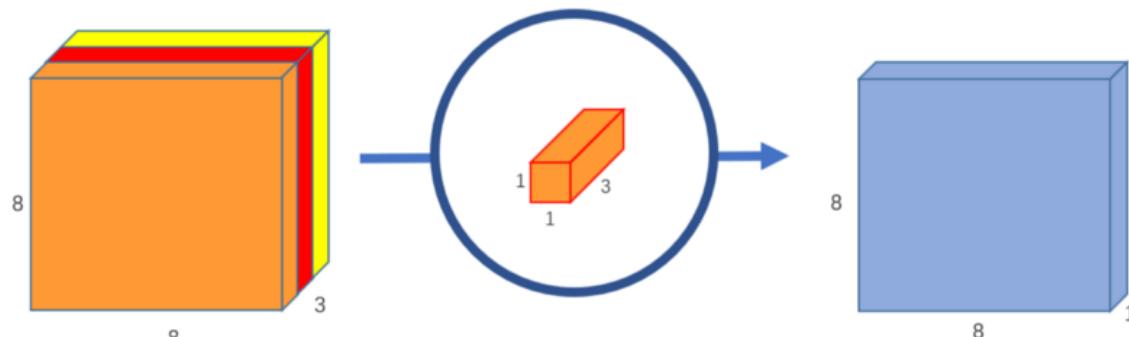
Depthwise Convolution



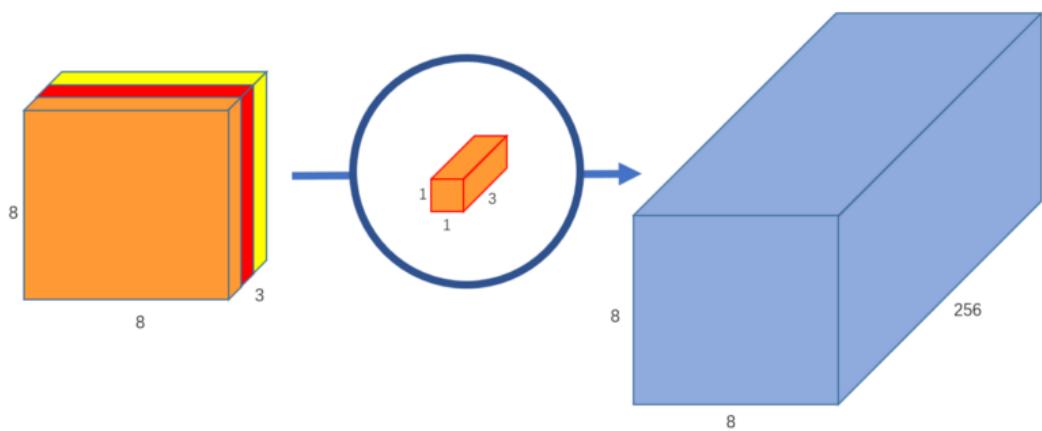
[Howard 2017]

Depthwise

3 5x5X1 kernels that move
8x8 times = 4800



Pointwise



256 1x1x3 kernels moving
8x8 times = 49152

Total of 53952 multiplications

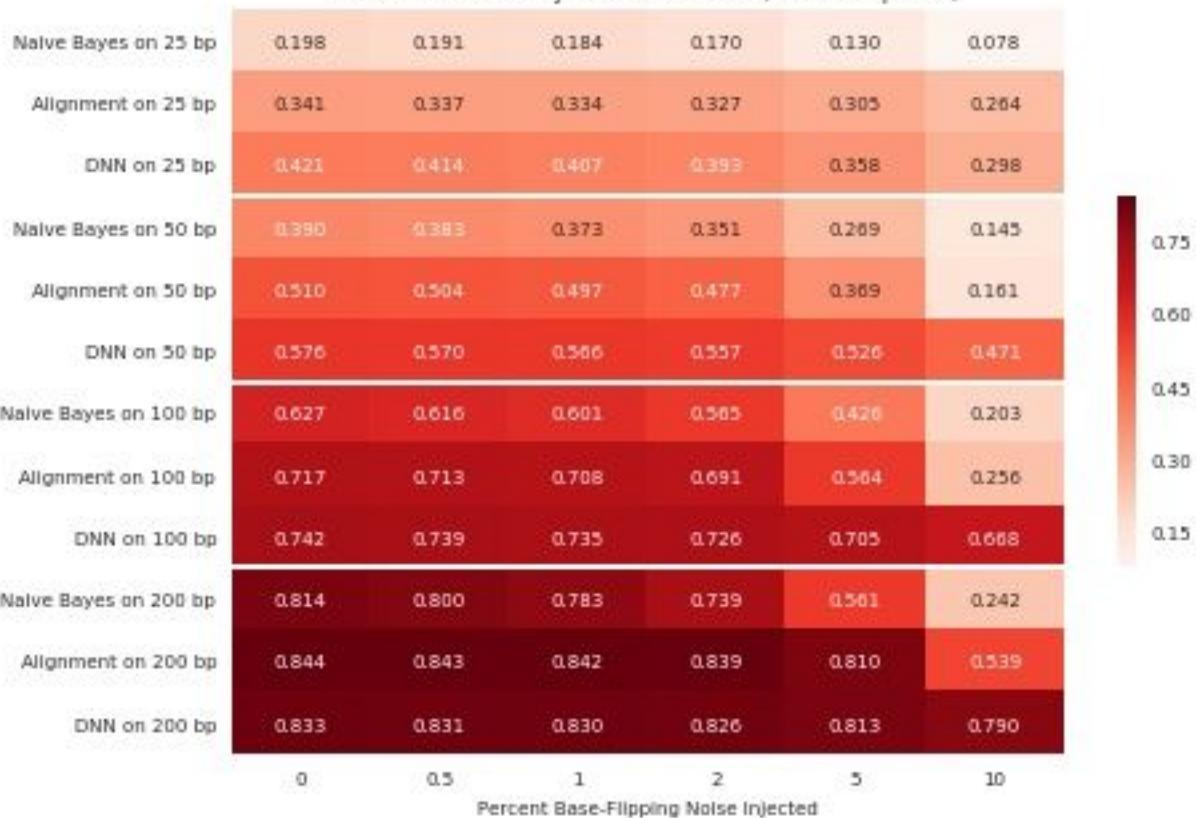


Ribosomal RNA

- 16S sequencing data
- Predicts species level taxonomy of microbiota
- Human Microbiome Project contains more than 14 terabytes of data and thousands of taxonomically characterized communities
- Database Project (RDP) SeqMatch and QIIME typically rely on explicit sequence matching against identified genomic sequences via a k-mer
 - Basic Local Alignment Search Tool



Genus-level Accuracy on Validation Set (Held-out Species)





Further reading

- CODA
<https://academic.oup.com/bioinformatics/article/33/14/i225/3953958>
- DeepBind <https://www.nature.com/articles/nbt.3300>
- MobileNet <https://arxiv.org/abs/1704.04861>
- Separable convolutions:
<https://towardsdatascience.com/a-basic-introduction-to-separable-convolutions-b99ec3102728>