



DEEP  
LEARNING  
INSTITUTE

# Deep Learning for Genomics Using DragoNN

Johnny Israeli

Biophysics PhD Candidate & SIGF Bio-X Fellow, Stanford University

Charles Killiam

Certified Instructor, NVIDIA Deep Learning Institute  
NVIDIA Corporation

May 9th 2017

A photograph of a woman with glasses, wearing a maroon long-sleeved shirt, speaking into a microphone. She is gesturing with her right hand, pointing upwards. The background shows a purple wall and other audience members.

# DEEP LEARNING INSTITUTE

## DLI Mission

Helping people solve challenging problems using AI and deep learning.

- Developers, data scientists and engineers
- Self-driving cars, healthcare and robotics
- Training, optimizing, and deploying deep neural networks

# TOPICS

- Lab Perspective
- Genomics
- DragoNN / SimDNA / Keras / DeepLIFT
- Lab
  - Discussion / Overview
  - Lab Environment
  - Lab Review
- Example Application
- State of the Field

# LAB PERSPECTIVE

# WHAT THIS LAB IS

- Introduction to large-scale data in genomics
- Applications of deep learning in genomics
- Guided, hands-on exercise using the DragoNN (Deep Regulatory Genomic Neural Network) toolkit for learning how to train and interpret simple deep learning models for genomics

# WHAT THIS LAB IS NOT

- Introduction to machine learning from first principles
- Detailed explanation of biology and genomics
- Rigorous mathematical formalism of neural networks
- Survey of all the features and options of the DragoNN package

# ASSUMPTIONS

- You are familiar with:
  - Convolutional neural networks (CNNs)
- Helpful to have:
  - Rudimentary coding background

# TAKE AWAYS

- Ability to setup your own DragoNN network for genomic research
- Know where to go for more info on DragoNN
- Familiarity with DragoNN workflow
- Downstream applications

# GENOMICS



# Decoding genome function

TGCCAAGCAGCAAAGTTTGCTGCTGTTATTTTAGCTCTTACTATATT  
CTACTTTACCATTGAAAATATTGAGGAAGTTATTATATTCTATTTTTAT  
ATATTATATATTATGTATTAAATTACTATTACACATAATTATTTTAT  
ATATATGAAGTACCAATGACTCCTTCAGAGCAATAATGAAATTTCAC  
AGTATGAAAATGGAAGAAATCAATAAAATTATACGTGACCTGTGGCGAAG  
TACCTATCGTGGACAAGGTGAGTACATGGTGTATCACAAATGCTCTTCC  
AAAGCCCTCTCCGCAGCTTCCCTATGACCTCTCATCATGCCAGCATT  
CCTCCCTGGACCCCTTCTAACGATGTCTTGAGATTCTAACGAATTCTTA  
TCTGGCAACATCTGTAGCAAGAAAATGTAAGTTCTGTCCAGAGCC  
TAACAGGACTTACATATTGACTGCAGTAGGCATTATATTAGCTGATGA  
ATAATAGGTTCTGTCATAGTAGATAGGGATAAGCCAAATGCAATAAG  
AAAAACCATCCAGAGGAAACTCTTTTTCTTTCTTTTTTTTCCA  
GATGGAGTCTCGACTTCTGTCAACCCTGGGCTGGAGCGCAGTGGTCAA  
TCTGGCTCACTGCAACCTCCACCTCCTGGGTCAGGTGATTCTCCCACCTC  
AGCCTCCCGAGTAGTAGCTGGAATTACAGGTGCGCGCTCCACACCTGGC  
TAATTTTTGTATTCTTAGTAGAGATGGGGTTCACCATGTTGCCAGGCT  
GGTCTCAAACCTCCTGCCCTCAGGTGATCTGCCACCTTGGCCTCCAGTGT  
TGGGTTACAGGCCTGAGCCACCGCGCTGGCCTGGAGGAAACTCTTAT  
AACTACCGAGTGGTGATGCTGAAGGGAGACACAGCCTGGATATGCGAG  
GACGATGCAGTGCTGGACAAAGGCAGGTATCTAAAAGCCTGGGGAGC  
CAACTCACCCAAGTAACTGAAAGAGAGAAACAAACATCAGTGCAGTGG  
AGCACCCAAGGCTACACCTGAATGGTGGGAAGCTTTGCTGCTATATAA  
AATGAATCAGGCTCAGCTACTATTATT .....

Function?

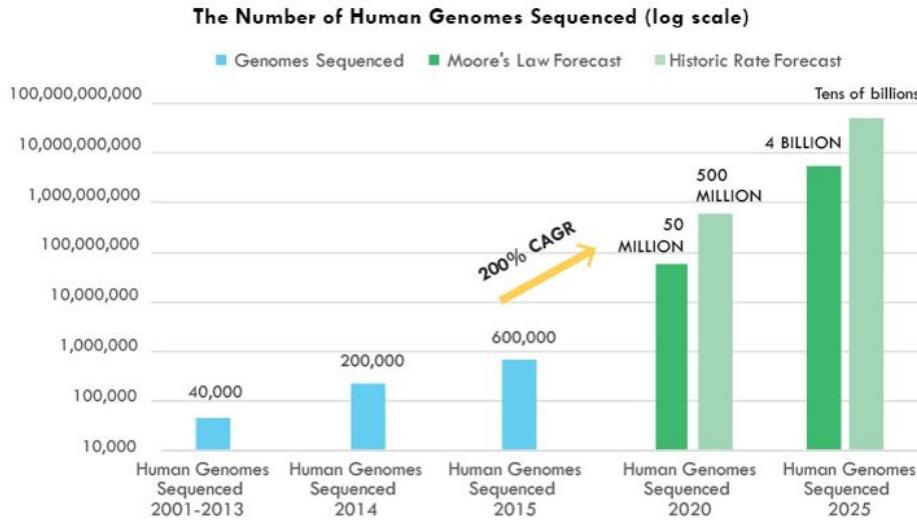
The Human Genome Project



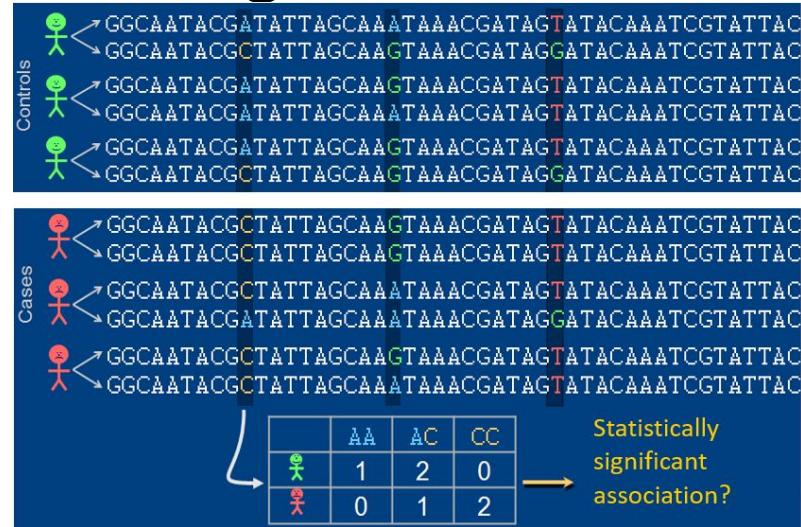
2003

~ 3 billion nucleotides

# Population sequencing identifies genetic variants



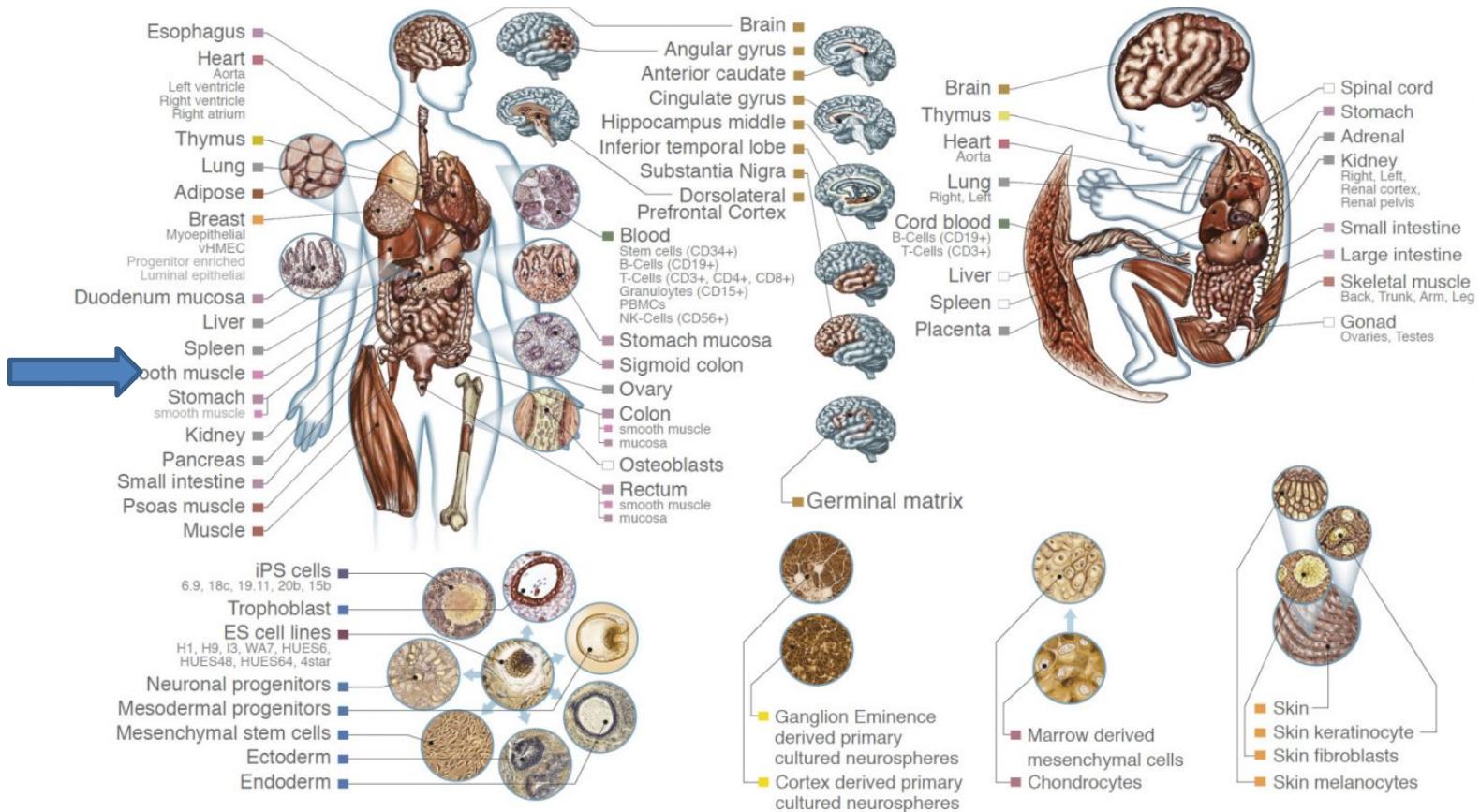
Source: National Human Genome Research Institute (NHGRI), ARK Investment Management LLC



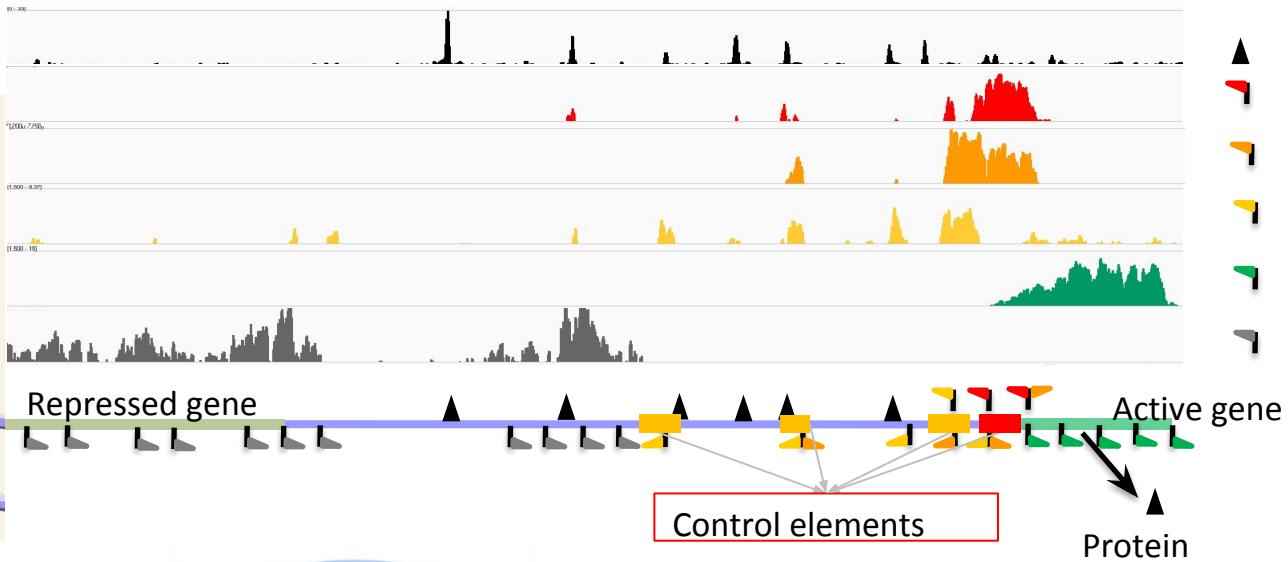
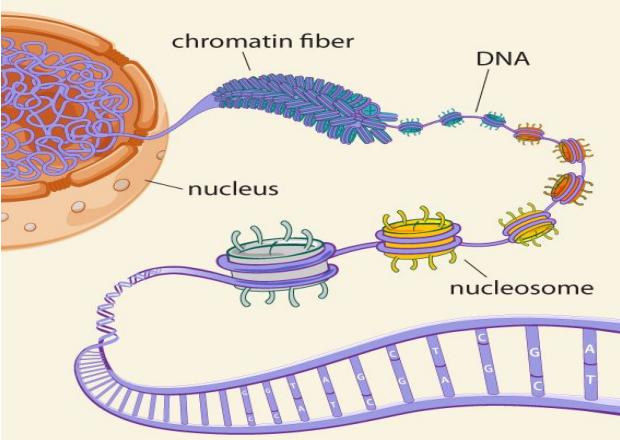
- 100,000s of personal genomes (population-scale sequencing)
- Millions of genetic variants (mutations) across individuals
- Which variants are benign and which ones are related to disease?
- What functional DNA words are these disease-associated variants disrupting?

# One genome ↔ many cell types

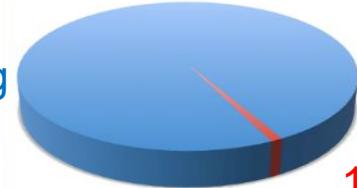
ACCAGTTACGACGG  
 TCAGGGTACTGATA  
 CCCCAAACCGTTGA  
 CCGCATTACAGAC  
 GGGGTTTGGGTTTT  
 GCCCCCACACAGGTA  
 CGTTAGCTACTGGT  
 TTAGCAATTACCG  
 TTACAACGTTACA  
 GGGTTACGGTTGGG  
 ATTTGAAAAAAAGT  
 TTGAGTTGGTTTT  
 TCACGGTAGAACGT  
 ACCTTACAAA.....



# Biochemical markers of functional elements

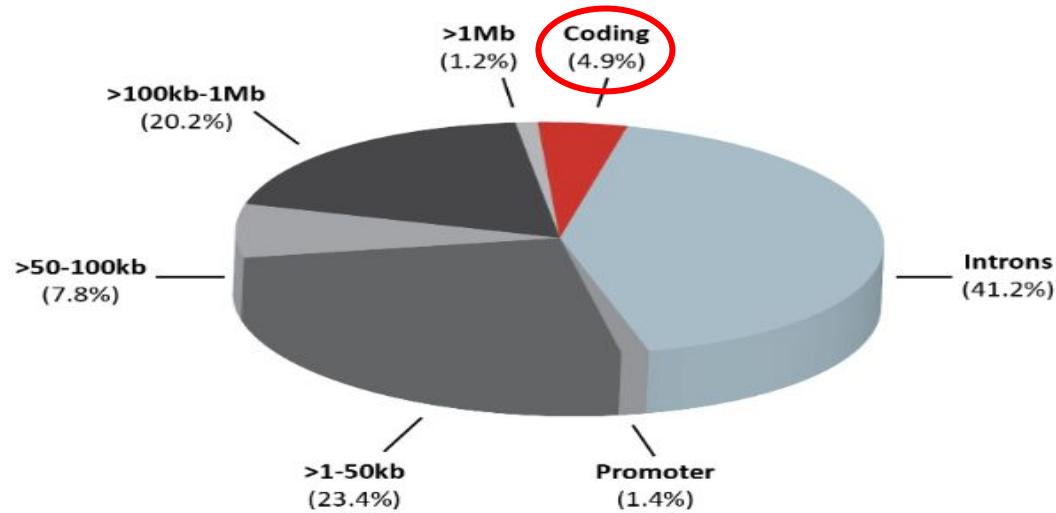
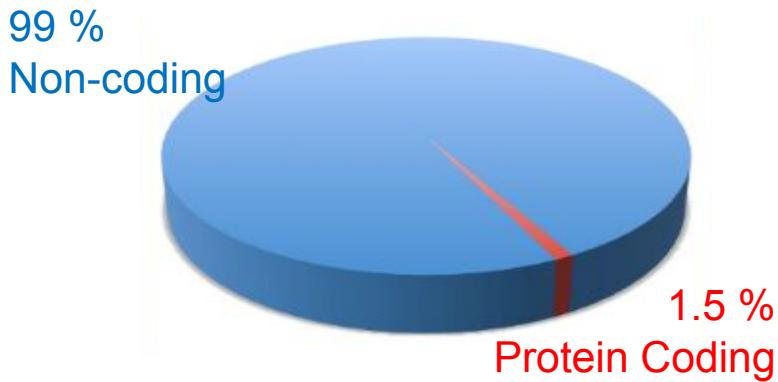


99 %  
Non-coding



1.5 % Protein Coding

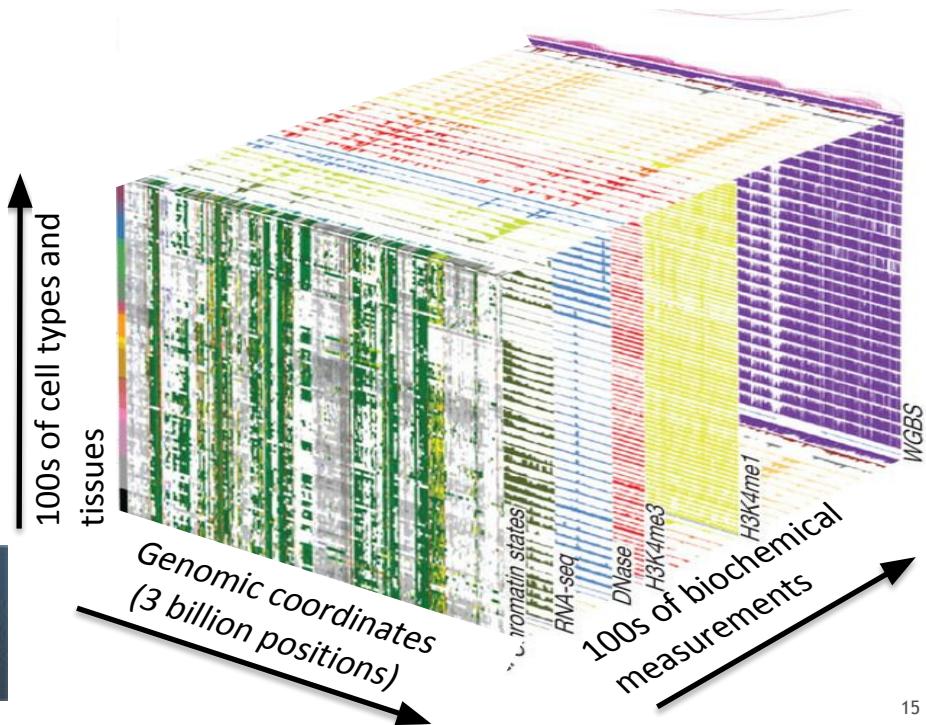
# Why are control elements relevant in disease?



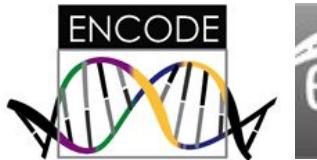
>90% of genetic variants associated with complex diseases do not disrupt protein-coding genes! Instead disrupting control elements!

# Large-scale functional genomic data

- **“Functional Genomics”**  
Application of sequencing technology for **profiling**  
**100s of biochemical markers of function** across the entire genome in 100s of different cell types and tissues

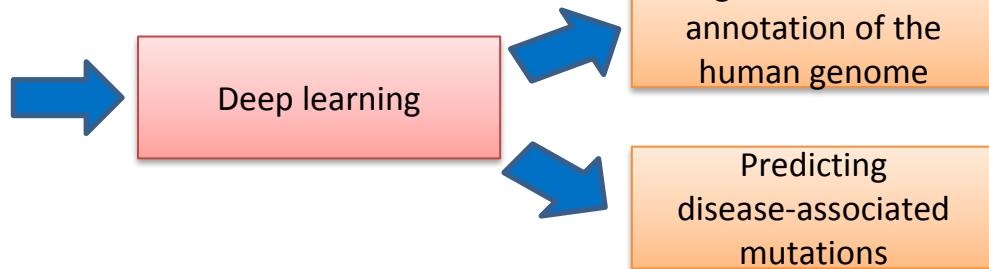
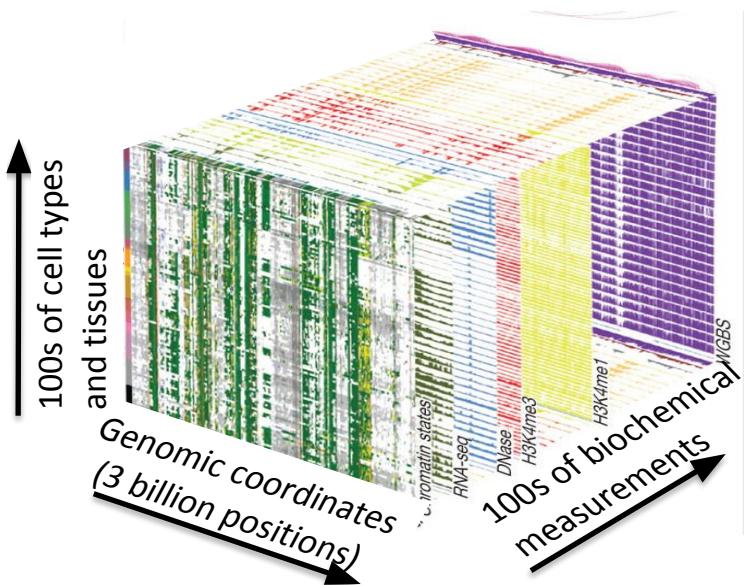


NIH funded collaborative consortia

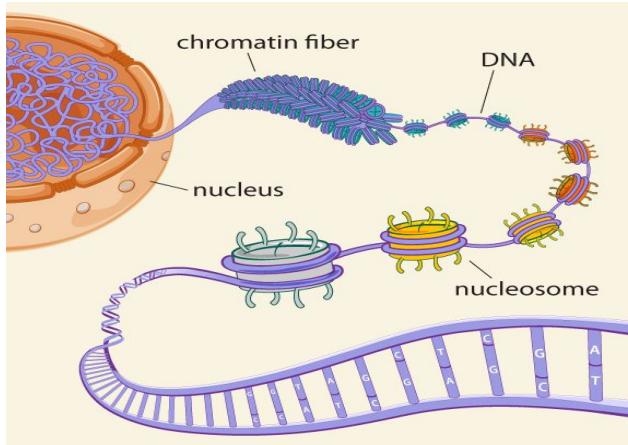


# Why deep learning and functional genomics?

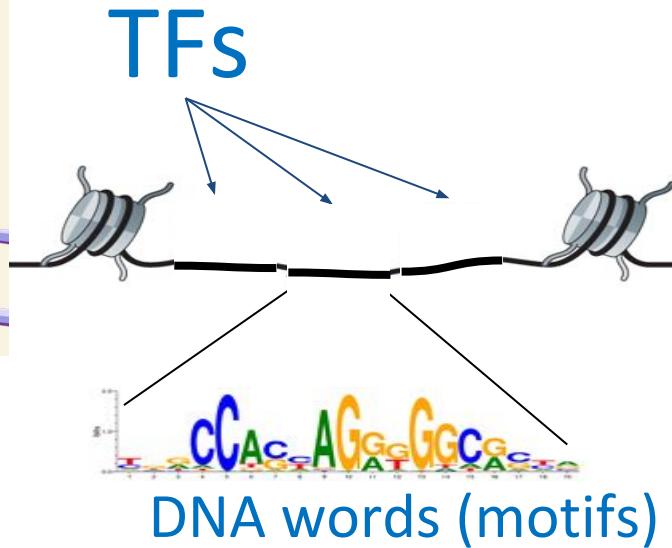
- **Terascale data cube** is rich and complex - deep learning is poised to integrate the heterogeneous data to decode genome function and its role in disease



# Proteins called transcription factors (TFs) bind functional DNA words to activate control elements

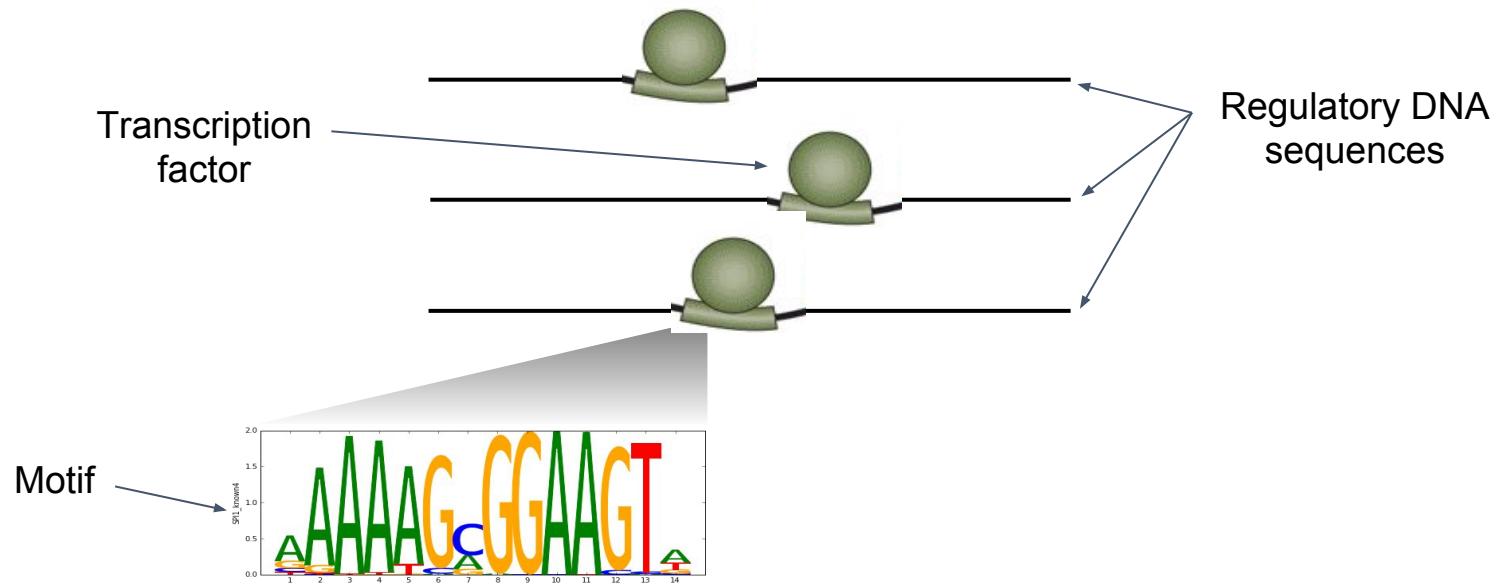


<https://www.broadinstitute.org/news/1504>



Adapted from Shlyueva et al. (2014) *Nature Reviews Genetics*.

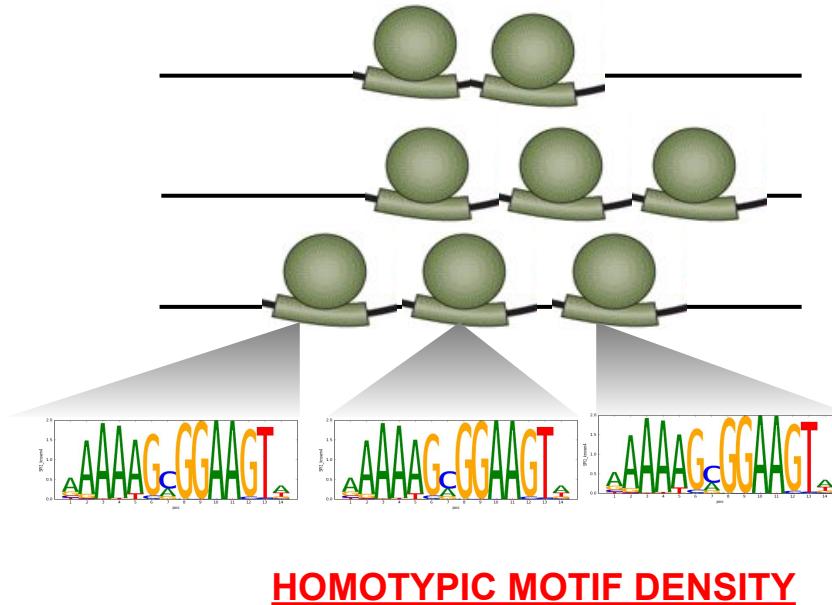
# Key properties of control element sequences



## TRANSCRIPTION FACTOR BINDING

Regulatory proteins called **transcription factors (TFs)** bind to high affinity sequence patterns (**motifs**) in regulatory DNA

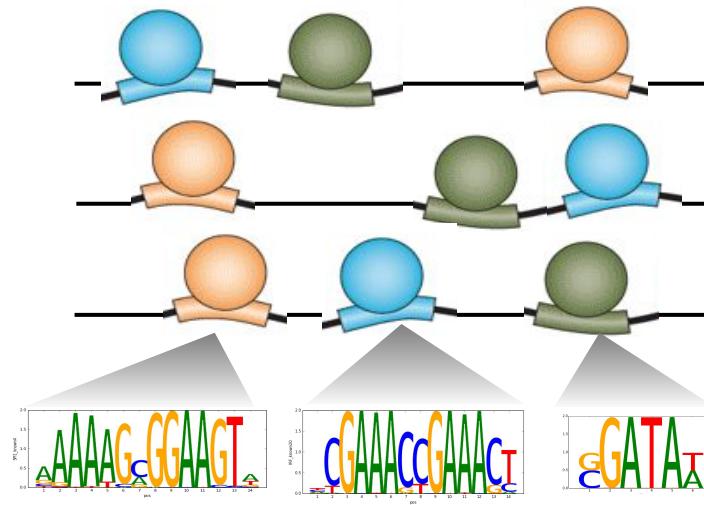
# Key properties of control element sequences



## HOMOTYPIC MOTIF DENSITY

Regulatory sequences often contain more than one binding instance of a TF resulting in homotypic clusters of motifs of the same TF

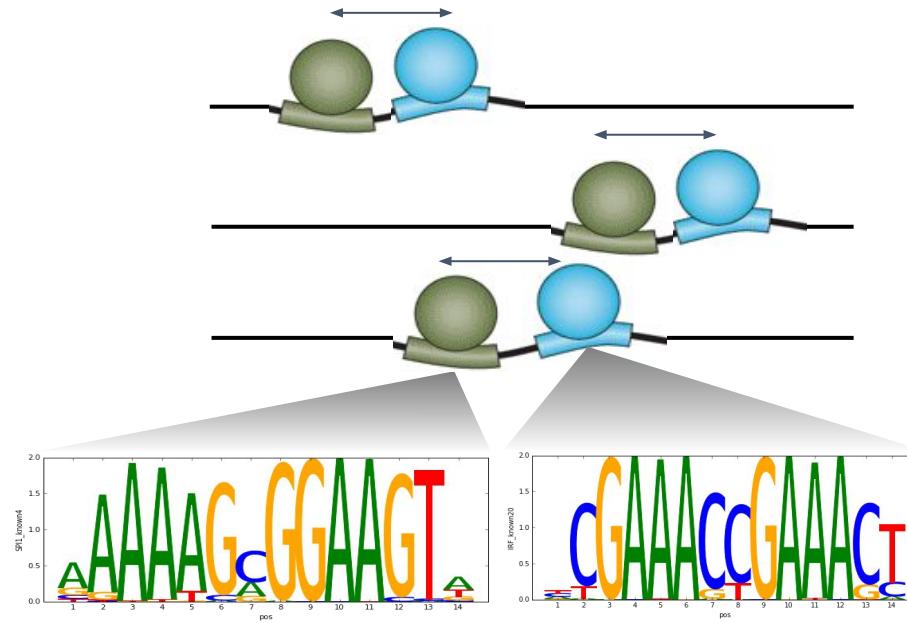
# Key properties of control element sequences



## HETEROTYPIC MOTIF COMBINATIONS

Regulatory sequences often bound by combinations of TFs resulting in heterotypic clusters of motifs of different TFs

# Key properties of control element sequences



## SPATIAL GRAMMARS OF HETEROGENEOUS MOTIF COMBINATIONS

Regulatory sequences are often bound by combinations of TFs with specific spatial and positional constraints resulting in distinct motif grammars

# DragoNN / SimDNA / KERAS / DeepLIFT

# DragoNN

- A pedagogical tool for deep learning for genomics.
- Systematic model design and interpretation. Runs on
  - SimDNA for simulations
  - Keras for model development
  - DeepLIFT for model interpretation
- IPython notebook tutorials.
- A command line interface for modeling and interpretation.
- Docker build + image for cloud usage.

# DragoNN

DragoNN provides a toolkit to learn how to model and interpret regulatory sequence data using deep learning.

 kundajelab / dragonn

 Unwatch ▾ 15

 Unstar 75

 Fork 26

 Code

 Issues 2

 Pull requests 0

 Projects 0

 Wiki

 Pulse

 Graphs

 Settings

A toolkit to learn how to model and interpret regulatory sequence data using deep learning.

Edit

<http://kundajelab.github.io/dragonnn/>

deep-learning genomics Manage topics

 107 commits

 6 branches

 3 releases

 5 contributors

 MIT

Systematic simulations

Tutorials

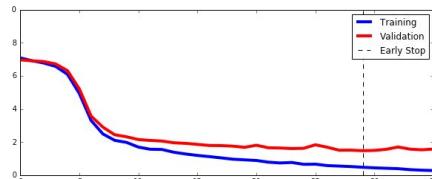
Command-line interface

Cloud-ready

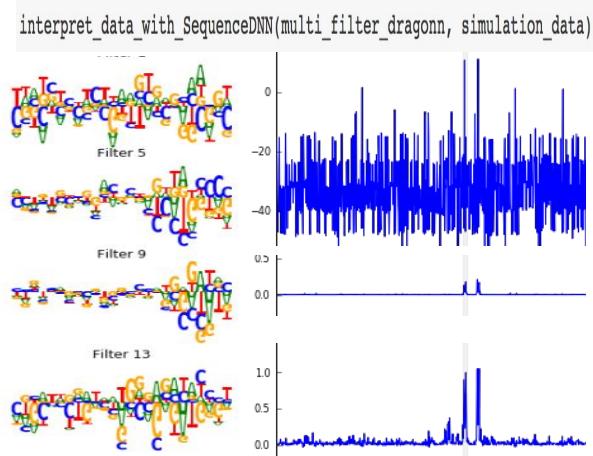
```
motif_density_localization_simulation_parameters = {  
    "motif_name": "TAL1_known4",  
    "seq_length": 1000,  
    "center_size": 150,  
    "min_motif_counts": 2,  
    "max_motif_counts": 4,  
    "num_pos": 3000,  
    "num_neg": 3000,  
    "GC_fraction": 0.4}
```

```
one_filter_dragonn_parameters = {  
    'seq_length': 1000,  
    'num_filters': [1],  
    'conv_width': [10],  
    'pool_width': 35}
```

SequenceDNN\_learning\_curve(one\_filter\_dragonn)



interpret\_SequenceDNN\_filters(multi\_layer\_dragonn, simulation\_data)



IPython Notebook  
Tutorials

Command Line  
Interface

DragoNN

SimDNA

Keras

DeepLIFT

TensorFlow

Theano

CPU

GPU

Locally or on the cloud

usage: dragonn [-h] {train,test,predict,interpret}

main script for DragoNN modeling of sequence data.

positional arguments:

{train,test,predict,interpret}

draggond command help

model training help

model testing help

model prediction help

model interpretation help

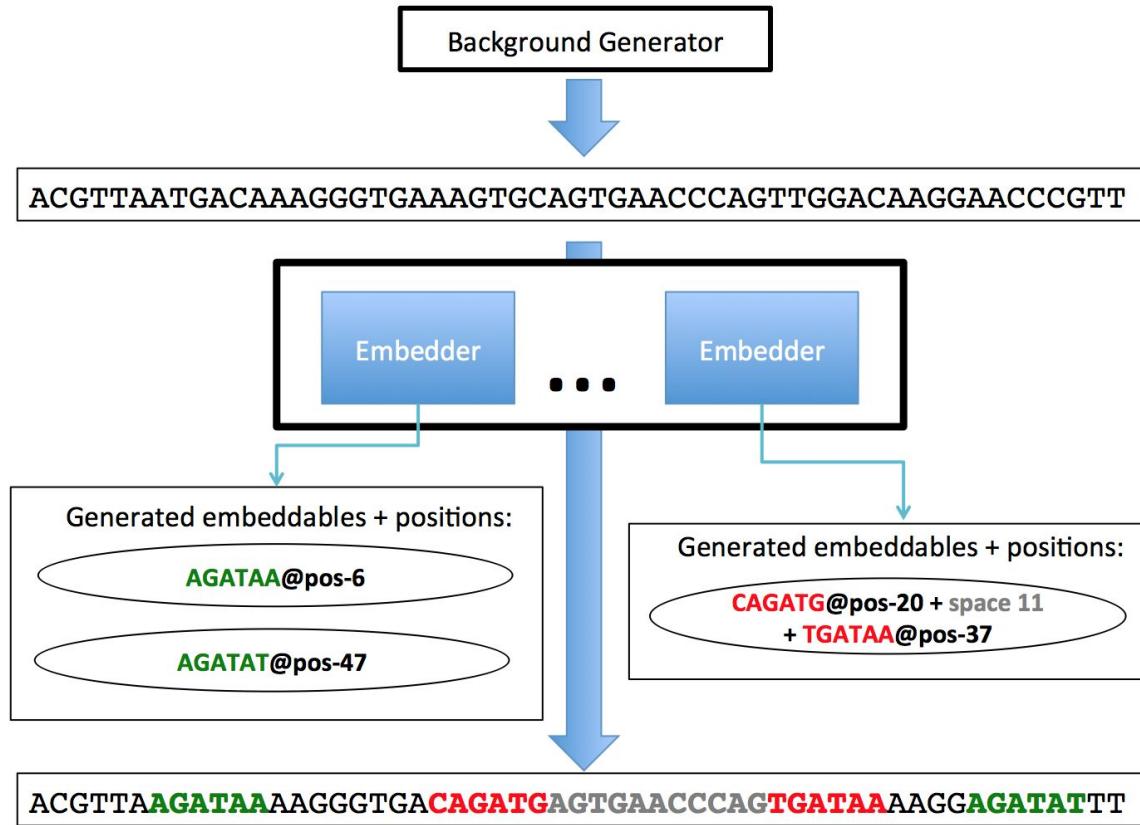
train

test

predict

interpret

# SimDNA: Simulations of DNA

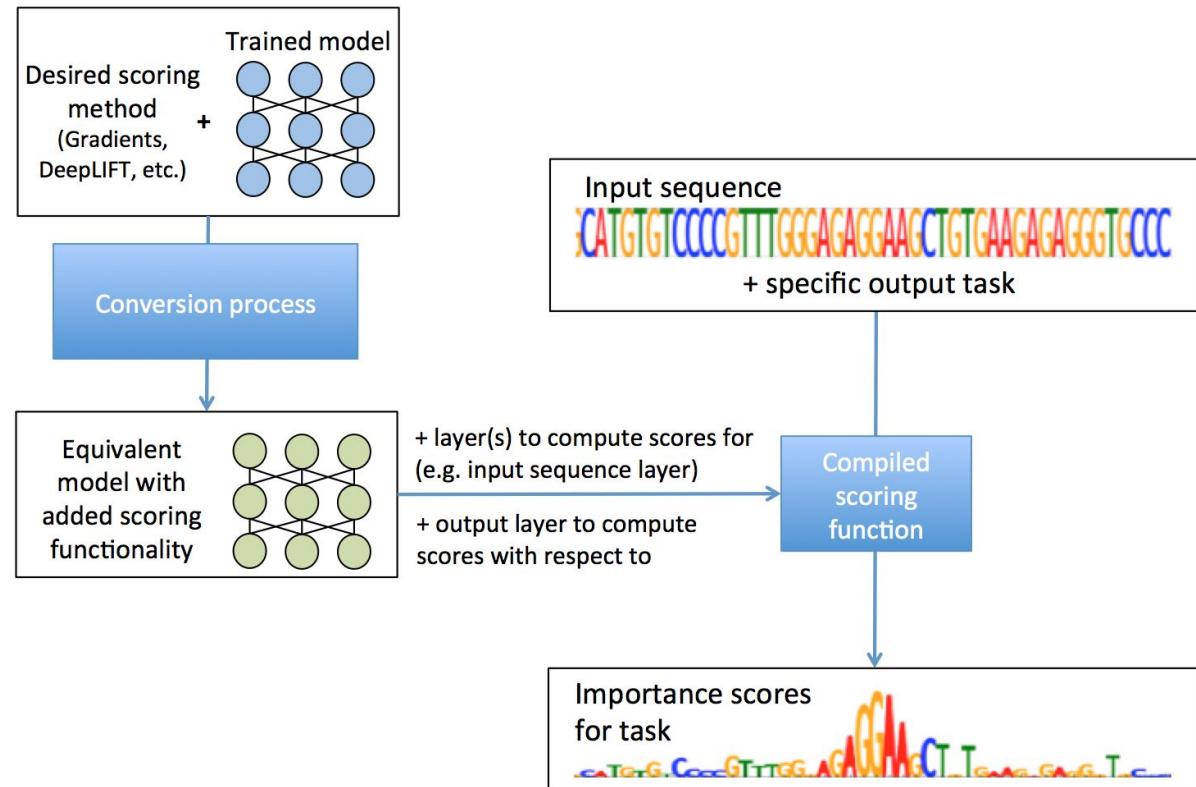


# KERAS: Neural Network API

- Modular neural network written in Python
- Runs on TensorFlow and Theano
  - Theano excels at RNNs / LSTMs
- Keras library allows for easy and fast prototyping
- Runs on GPUs and CPUs
- Compatible with Python 2.7 - 3.5

# DeepLIFT: Neural Network Interpretation API

- Modular feature importance API in Python 2.7 - 3.5
- Runs on TensorFlow and Theano
- Allows for easy interpretation of neural networks
- Runs on GPUs and CPUs



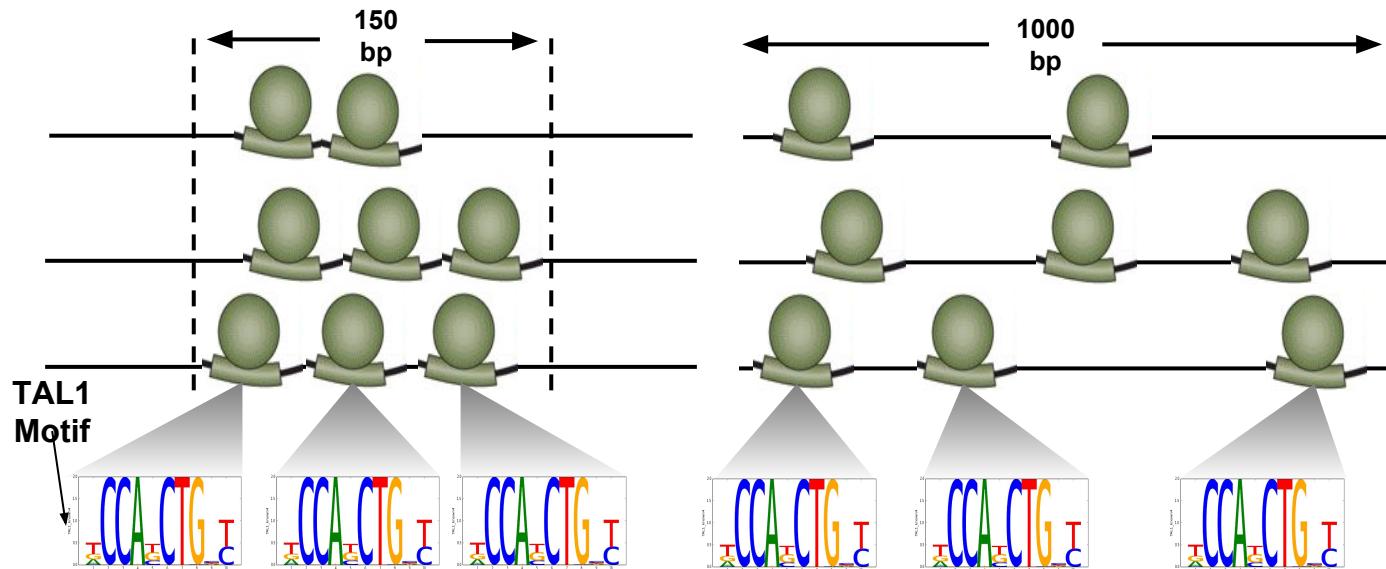
# LAB DISCUSSION / OVERVIEW

# LAB PROCESS

1. Simulate data for homotypic motif density localization
2. Architect DragonN model
3. Train your model
4. Visualize results / interpret data
5. Modify architecture to enhance predictions and interpretation

# Simulate homotypic motif density localization

```
simulation_parameters = {  
    "motif_name":  
        "TAL1_known4",  
    "seq_length": 1000,  
    "center_size": 150,  
    "min_motif_counts": 2,  
    "max_motif_counts": 4,  
    "num_pos": 3000,  
    "num_neg": 3000,  
    "GC_fraction": 0.4}
```



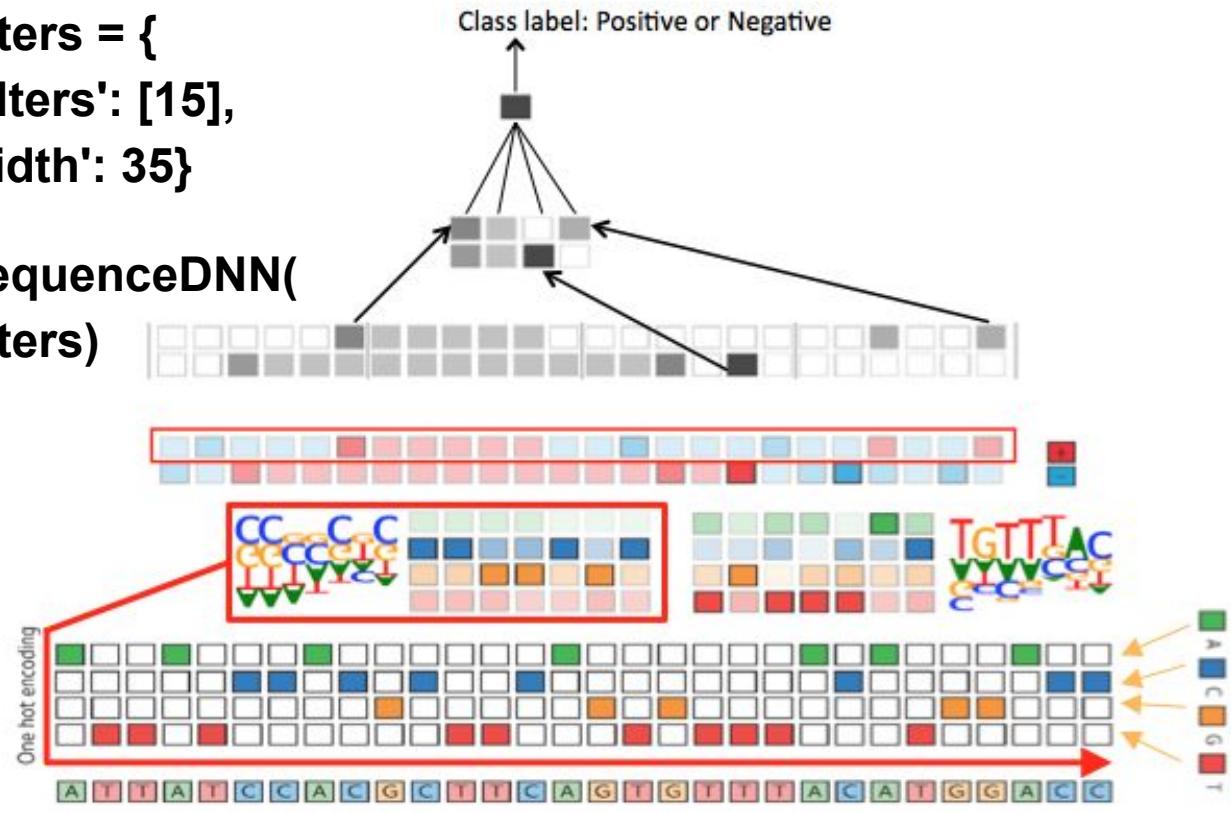
**Positive set:** 1000 base pairs (bp) sequences with 2-4 instances of TAL1 motif in central 150bp

**Negative set:** 1000 base pairs (bp) sequences with 2-4 instances of TAL1 motif positioned anywhere in the 1000bp

# Architecting a DragoNN model

```
multi_filter_dragonn_parameters = {  
    'seq_length': 1000, 'num_filters': [15],  
    'conv_width': [10], 'pool_width': 35}
```

```
multi_filter_dragonn = get_SequenceDNN(  
    multi_filter_dragonn_parameters)
```



# TRAIN A DRAGONN MODEL

```
multi_filter_dragonn = get_SequenceDNN(multi_filter_dragonn_parameters)  
train_SequenceDNN(multi_filter_dragonn, simulation_data)
```

Epoch 1:

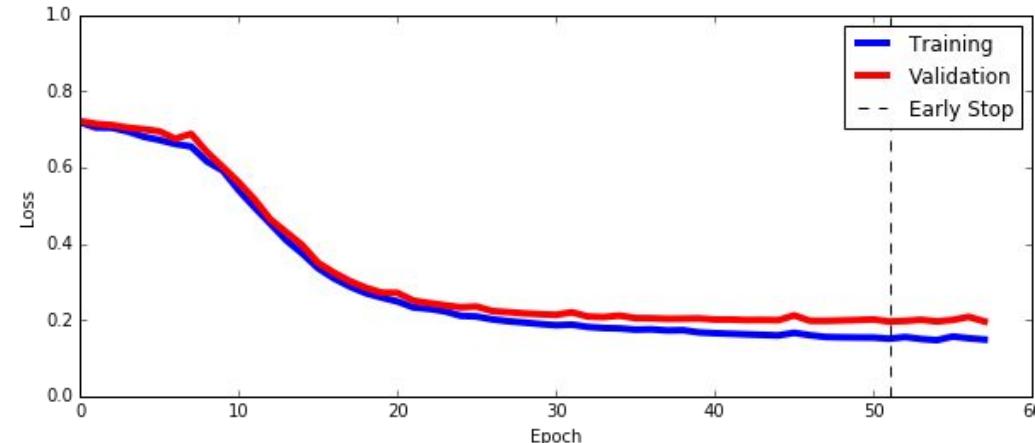
Train Loss: 0.7184	Balanced Accuracy: 50.27%	auROC: 0.502	auPRC: 0.492	Recall at 5% 10% 20% FDR: 0.1% 0.1% 0.1%	Num Positives: 1957	Num Negatives: 2043
Valid Loss: 0.7212	Balanced Accuracy: 48.60%	auROC: 0.470	auPRC: 0.498	Recall at 5% 10% 20% FDR: 0.0% 0.0% 0.0%	Num Positives: 528	Num Negatives: 472 *

.....

Epoch 58:

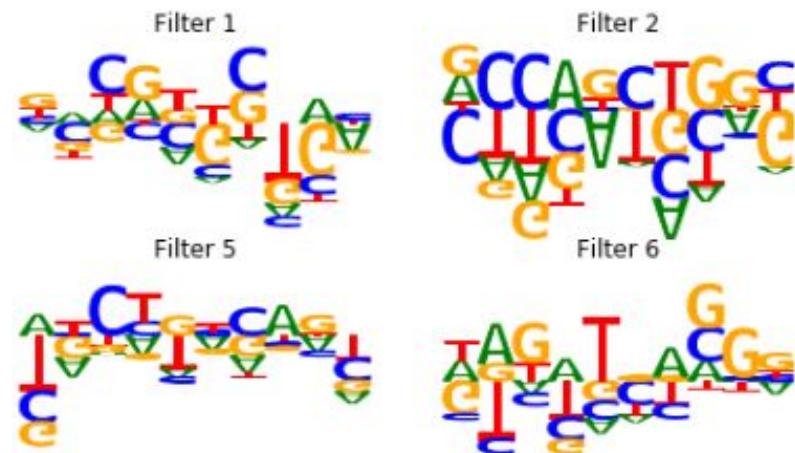
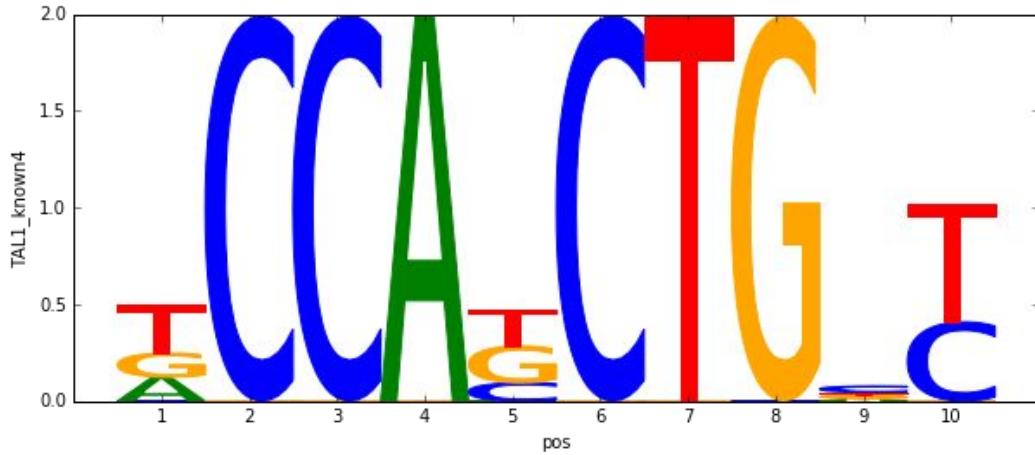
Train Loss: 0.1484	Balanced Accuracy: 94.40%	auROC: 0.988	auPRC: 0.987	Recall at 5% 10% 20% FDR: 93.4% 97.8% 99.7%	Num Positives: 1957	Num Negatives: 2043
Valid Loss: 0.1958	Balanced Accuracy: 91.97%	auROC: 0.976	auPRC: 0.977	Recall at 5% 10% 20% FDR: 88.6% 95.3% 99.2%	Num Positives: 528	Num Negatives: 472

Finished training after 58 epochs.



# Evaluate Learned Filters

*interpret\_SequenceDNN\_filters(multi\_filter\_dragonn, simulation\_data)*

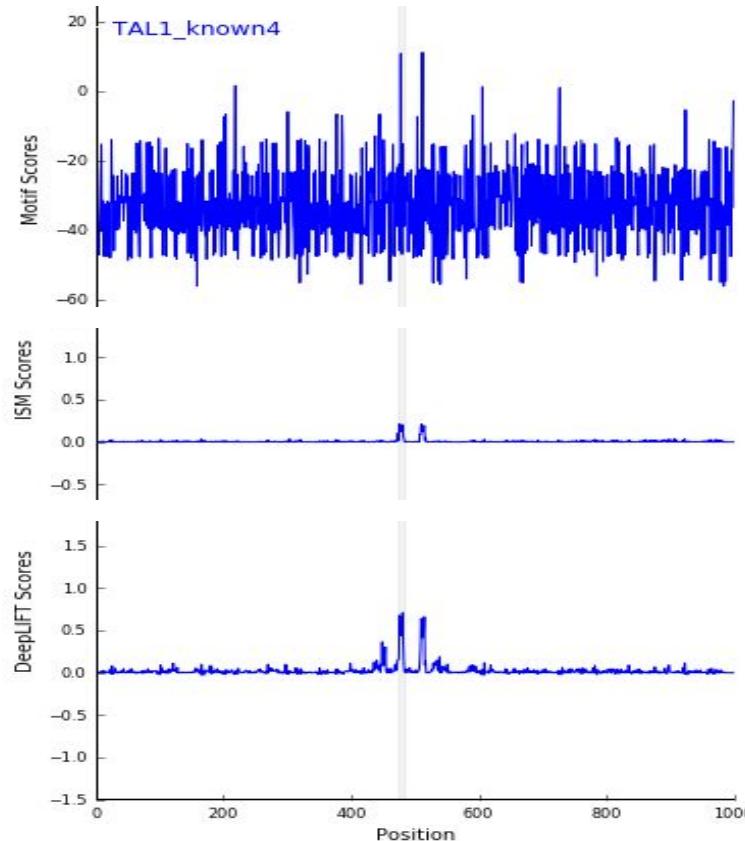


- Simulated motif is the “ground truth” feature
- Learned filters evaluated wrt simulated motif

# Evaluate Feature Importances

```
interpret_data_with_SequenceDNN(  
    multi_filter_dragonn,  
    simulation_data)
```

- Motif scores are the “ground truth” feature importances
- Inferred feature importances inferred evaluated wrt motif scores



# LAB ENVIRONMENT

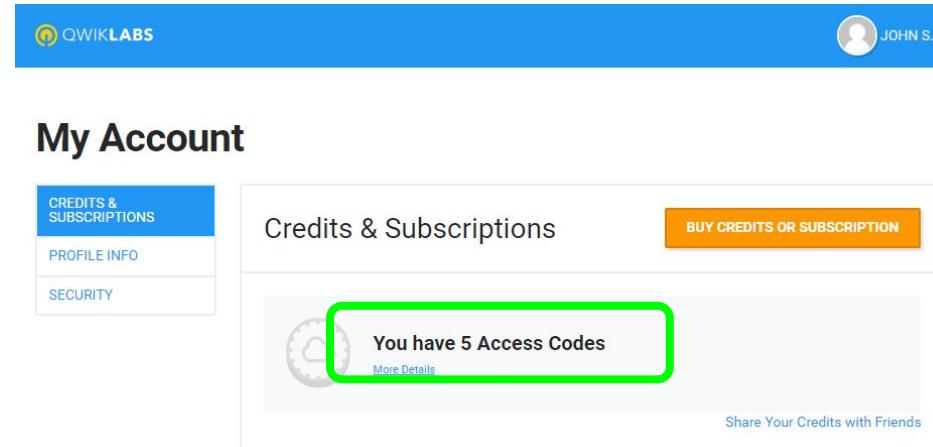
# LAB CONNECTION INSTRUCTIONS - Part 1

Go to [nvlabs.qwiklab.com](https://nvlabs.qwiklab.com)

Sign in or create an account

Check for Access Codes (each day):

- Click **My Account**
- Click **Credits & Subscriptions**



The screenshot shows the 'My Account' page of the QwikLabs website. At the top, there's a navigation bar with the QwikLabs logo and a user profile for 'JOHN S.'. Below the header, the main content area has a title 'My Account'. On the left, there's a sidebar with three options: 'CREDITS & SUBSCRIPTIONS' (which is highlighted in blue), 'PROFILE INFO', and 'SECURITY'. The main content area is titled 'Credits & Subscriptions'. It features a large button 'BUY CREDITS OR SUBSCRIPTION' in orange. Below it, there's a section with a circular icon and the text 'You have 5 Access Codes' (which is enclosed in a green rectangular box). At the bottom right of this section, there's a link 'More Details'. At the very bottom of the page, there's a link 'Share Your Credits with Friends'.

If no Access Codes, ask for paper one from TA.

Please tear in half once used

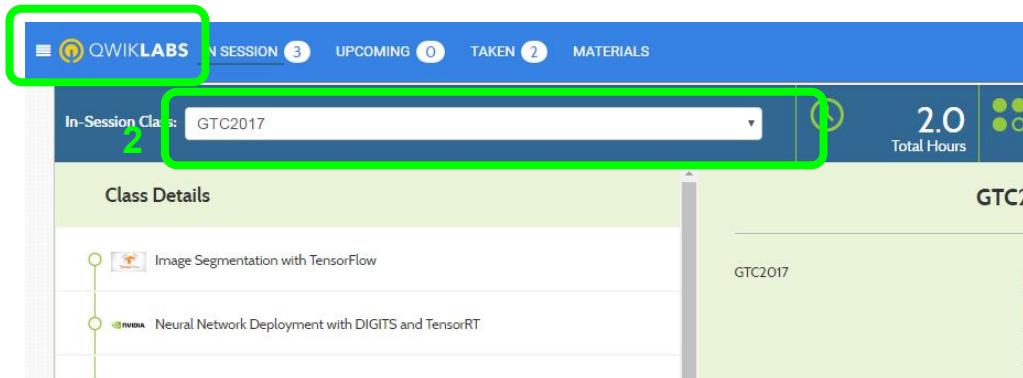
An Access Code is needed to start the lab

WIFI SSID: **GTC\_Hands\_On**

Password: **HandsOnGpu**

# LAB CONNECTION INSTRUCTIONS - Part 2

1. Click **Qwiklabs** in upper-left **1**
2. Select GTC2017 Class
3. Find lab and click on it
4. Click on Select
5. Click Start Lab



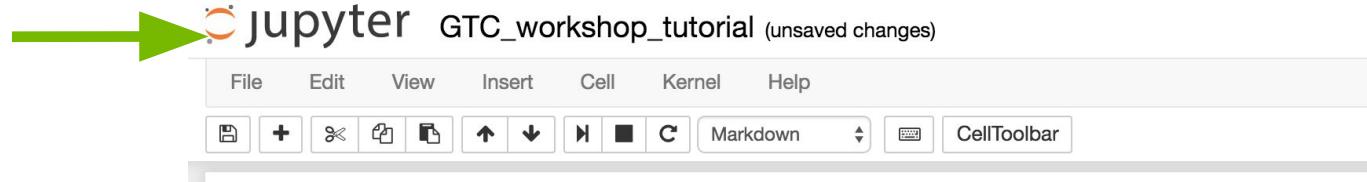
This screenshot shows the "Class Details" page for the "GTC2017" class. The "In-Session Class" dropdown is still set to "GTC2017". The "Total Hours" are listed as 2.0. On the right, there are stats for "Completed Labs" (2) and "Classes Taken" (2). The "Class Details" section lists four labs: "Image Segmentation with TensorFlow", "Neural Network Deployment with DIGITS and TensorRT", "Image Classification with DIGITS", and "Object Detection with DIGITS". At the bottom of the list, the "Photo Editing with Generative Adversarial Networks in Tensorflow and DIGITS" lab is highlighted with a green box and a green number "3" to its left. To the right of this highlighted lab is a blue "Select" button, which is also highlighted with a green box and a green number "4" above it.

WIFI SSID:  
**GTC\_Hands\_On**

Password:  
**HandsOnGpu**

# ACCESSING LAB INSTRUCTIONS

Should see  
Jupyter  
notebook



## How to train your DragoNN tutorial

**Tutorial length:** 25-30 minutes with a CPU.

### Outline

- \* How to use this tutorial
- \* Review of patterns in transcription factor binding sites
- \* Learning to localize homotypic motif density
- \* Sequence model definition
- \* Training and interpretation of
  - single layer, single filter DragoNN
  - single layer, multiple filters DragoNN
  - Multi-layer DragoNN
  - Regularized multi-layer DragoNN
- \* Critical questions in this tutorial:
  - What is the "right" way to get insight from a DragoNN model?
  - What are the limitations of different interpretation methods?
  - Do those limitations depend on the model and the target pattern?
- \* Suggestions for further exploration

# ACCESSING LAB INSTRUCTIONS

Place cursor  
in code block  
and click  
execute  
button

The screenshot shows a Jupyter Notebook interface with the title "jupyter GTC\_workshop\_tutorial (autosaved)". The toolbar includes File, Edit, View, Insert, Cell, Kernel, Help, and CellToolbar buttons. A green arrow points from the text "Place cursor in code block and click execute button" to the "Cell" button in the toolbar. Another green arrow points from the text "We start by loading dragonn's tutorial utilities and reviewing properties of regulatory sequence that transcription factors bind." to the code block In [1].

In [2]: `print_available_simulations??`

1. type function name
2. add “??”
3. click the play button

Signature: `print_available_simulations()`  
Source:  
`def print_available_simulations():  
 for function_name in get_available_simulations():  
 print function_name`  
File: `~/src/dmn4genomics_review/dragonn/tutorial_utils.py`  
Type: `function`

We start by loading dragonn's tutorial utilities and reviewing properties of regulatory sequence that transcription factors bind.

In [1]: `%reload_ext autoreload  
%autoreload 2  
from dragonn.tutorial_utils import *\br/>%matplotlib inline`

# LAB REVIEW

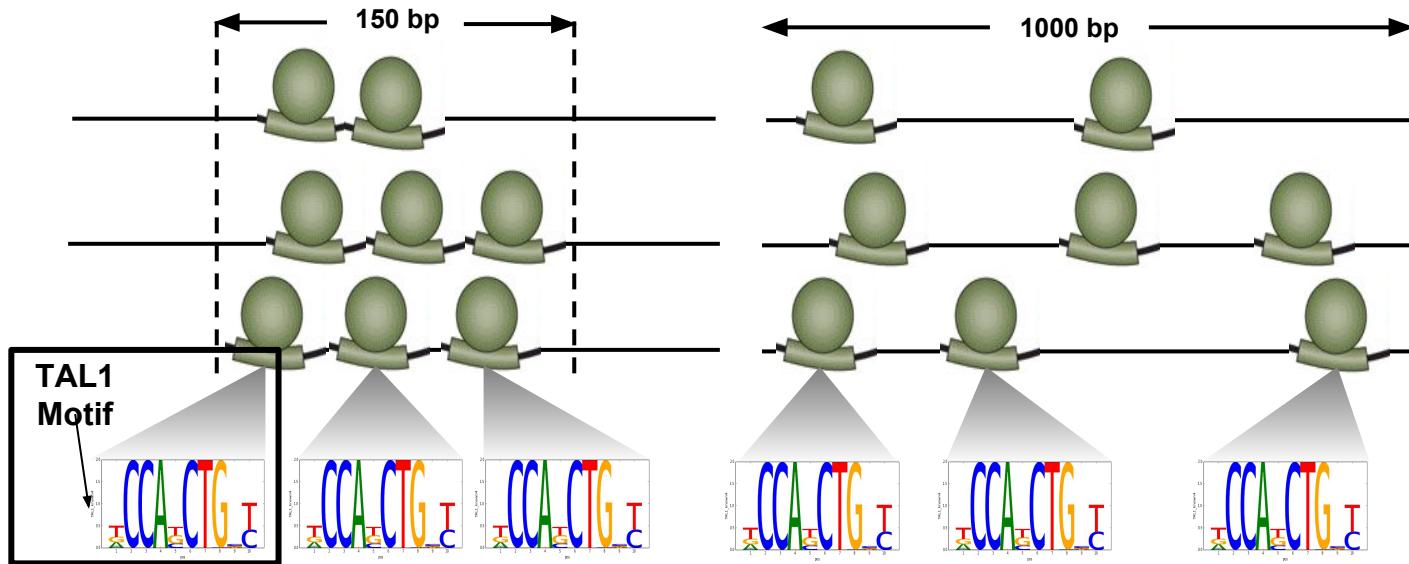
# LAB REVIEW

1. Generate simulation data
2. Architect and train single filter DragoNN model
3. Architect and train multi-filter model
4. Visualize results
5. Interpret model and data
6. Introduce enhancements

# Generating simulation data

# Homotypic motif density localization

```
simulation_parameters = {  
    "motif_name":  
        "TAL1_known4",  
    "seq_length": 1000,  
    "center_size": 150,  
    "min_motif_counts": 2,  
    "max_motif_counts": 4,  
    "num_pos": 3000,  
    "num_neg": 3000,  
    "GC_fraction": 0.4}
```

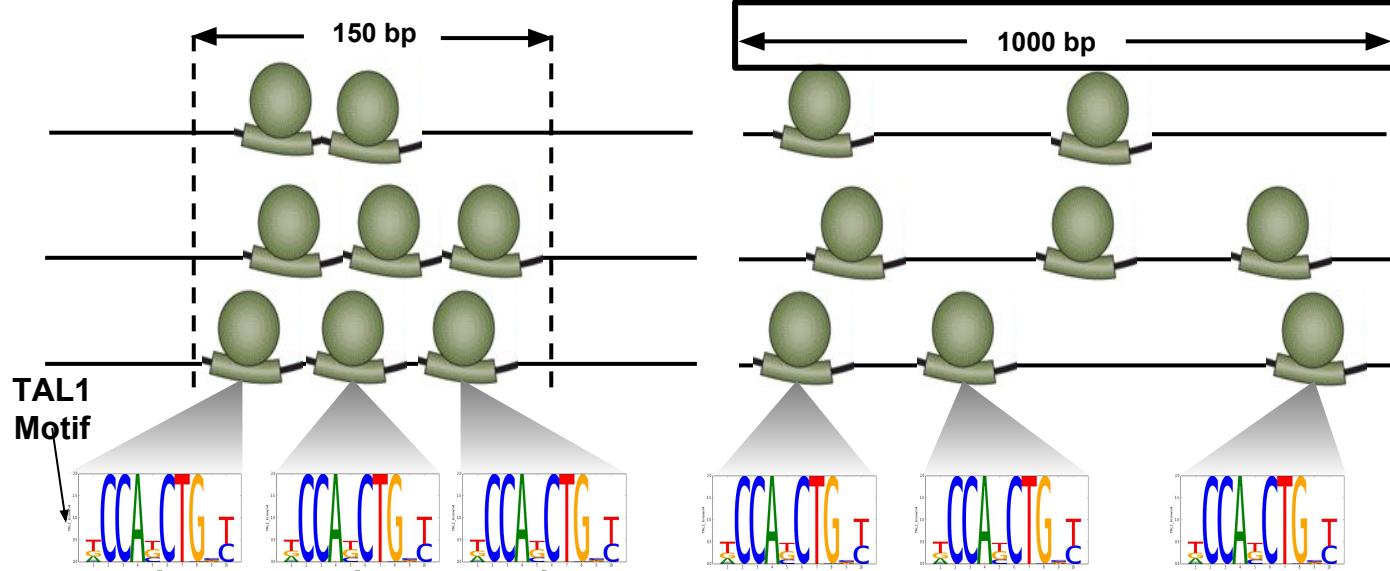


**Positive set:** 1000 base pairs (bp)  
sequences with 2-4 instances of  
TAL1 motif in central 150bp

**Negative set:** 1000 base pairs (bp)  
sequences with 2-4 instances of  
TAL1 motif positioned anywhere in  
the 1000bp

# Homotypic motif density localization

```
simulation_parameters = {  
    "motif_name":  
        "TAL1_known4",  
    "seq_length": 1000,  
    "center_size": 150,  
    "min_motif_counts": 2,  
    "max_motif_counts": 4,  
    "num_pos": 3000,  
    "num_neg": 3000,  
    "GC_fraction": 0.4}
```

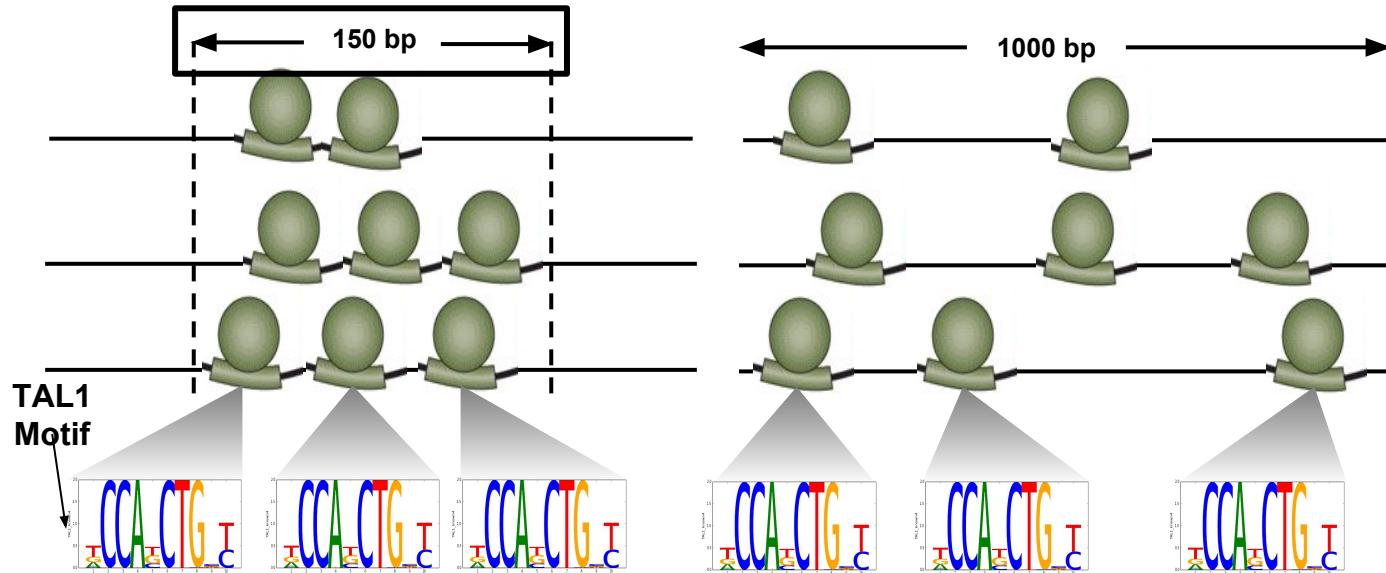


**Positive set:** 1000 base pairs (bp) sequences with 2-4 instances of TAL1 motif in central 150bp

**Negative set:** 1000 base pairs (bp) sequences with 2-4 instances of TAL1 motif positioned anywhere in the 1000bp

# Homotypic motif density localization

```
simulation_parameters = {  
    "motif_name":  
        "TAL1_known4",  
    "seq_length": 1000,  
    "center_size": 150,  
    "min_motif_counts": 2,  
    "max_motif_counts": 4,  
    "num_pos": 3000,  
    "num_neg": 3000,  
    "GC_fraction": 0.4}
```

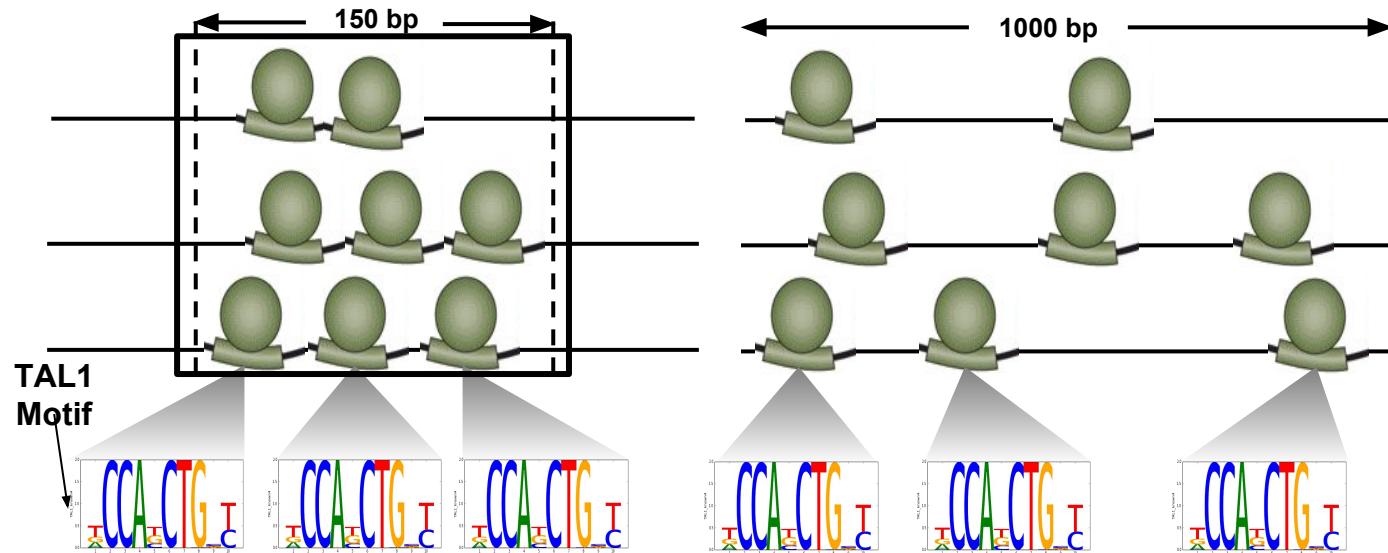


**Positive set:** 1000 base pairs (bp) sequences with 2-4 instances of TAL1 motif in central 150bp

**Negative set:** 1000 base pairs (bp) sequences with 2-4 instances of TAL1 motif positioned anywhere in the 1000bp

# Homotypic motif density localization

```
simulation_parameters = {  
    "motif_name":  
        "TAL1_known4",  
    "seq_length": 1000,  
    "center_size": 150,  
    "min_motif_counts": 2,  
    "max_motif_counts": 4,  
    "num_pos": 3000,  
    "num_neg": 3000,  
    "GC_fraction": 0.4}
```



**Positive set:** 1000 base pairs (bp)  
sequences with 2-4 instances of  
TAL1 motif in central 150bp

**Negative set:** 1000 base pairs (bp)  
sequences with 2-4 instances of  
TAL1 motif positioned anywhere in  
the 1000bp

# Architecting a DragoNN model

# ARCHITECT one\_filter\_dragonn MODEL

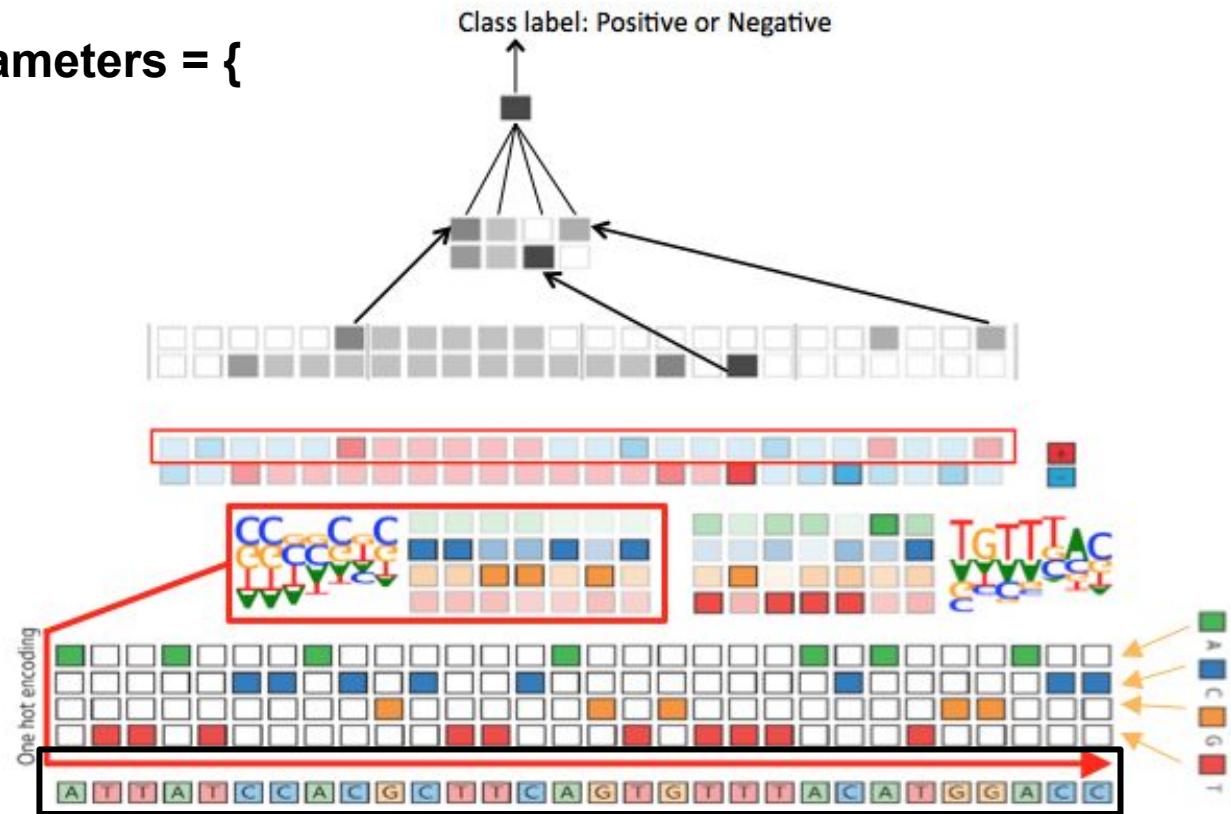
```
one_filter_dragonn_parameters = {
```

```
    'seq_length': 1000,
```

```
    'num_filters': [1],
```

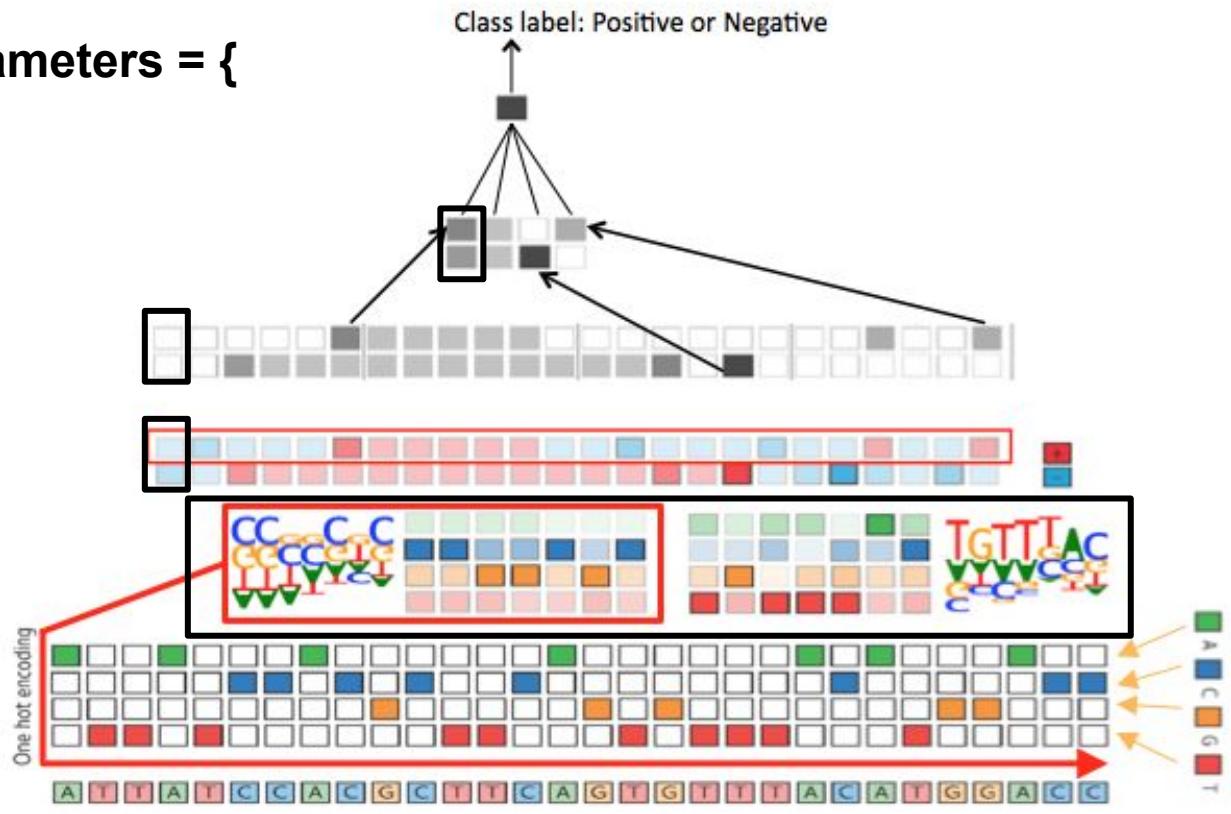
```
    'conv_width': [10],
```

```
    'pool_width': 35}
```



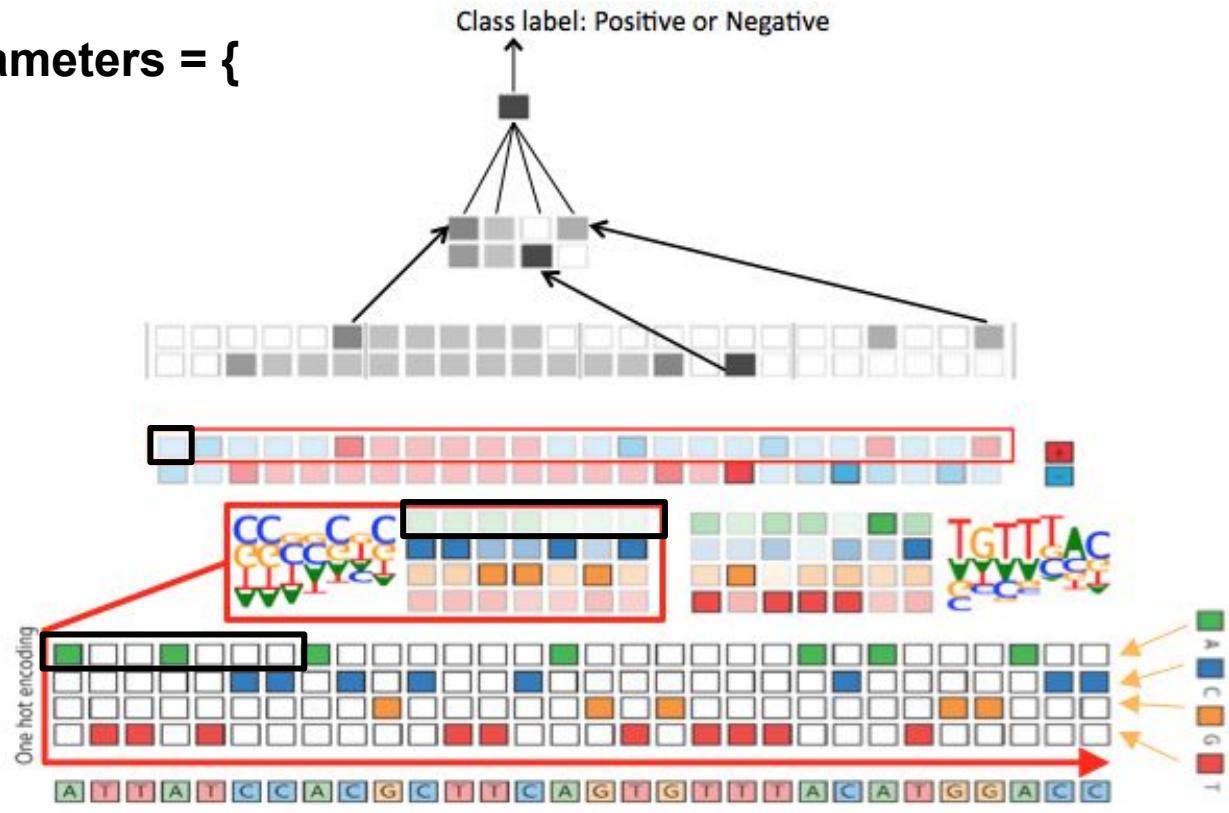
# ARCHITECT one\_filter\_dragonn MODEL

```
one_filter_dragonn_parameters = {  
    'seq_length': 1000,  
    'num_filters': [1],  
    'conv_width': [10],  
    'pool_width': 35}
```



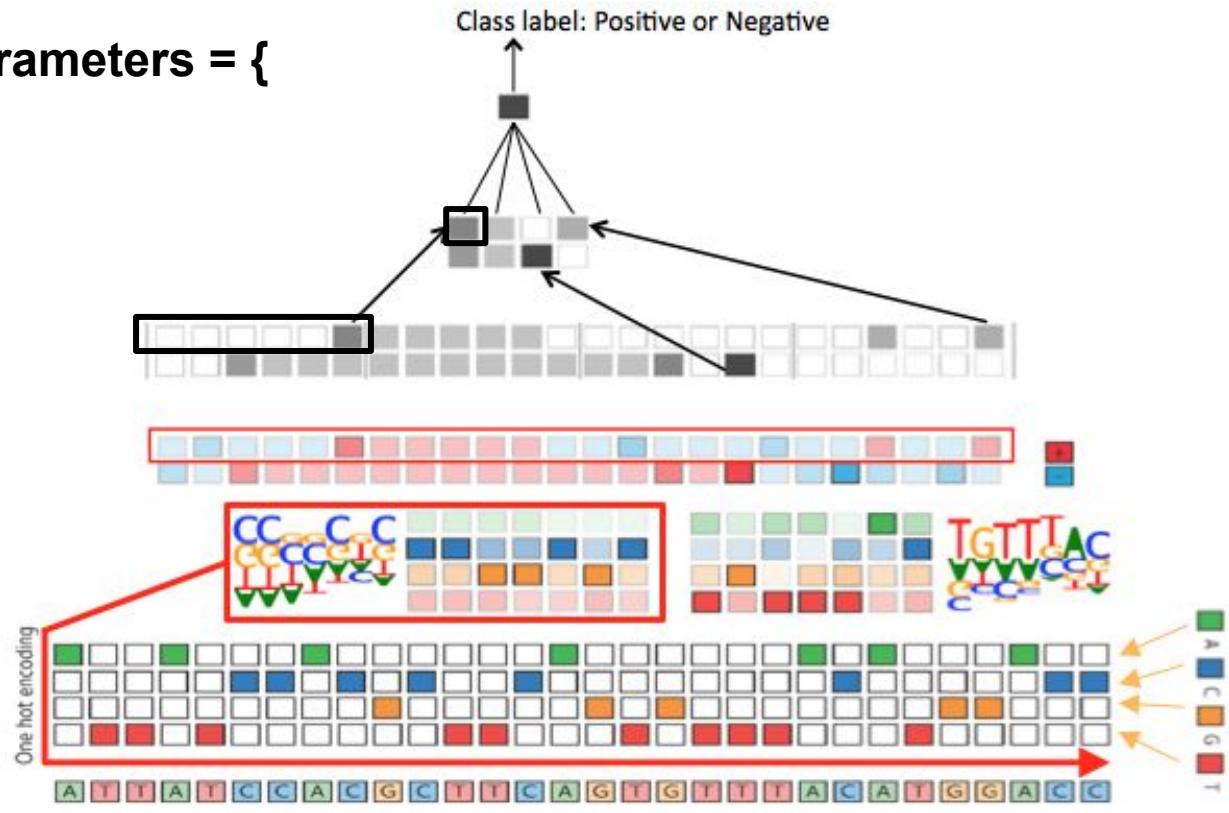
# ARCHITECT one\_filter\_dragonn MODEL

```
one_filter_dragonn_parameters = {  
    'seq_length': 1000,  
    'num_filters': [1],  
    'conv_width': [10],  
    'pool_width': 35}
```



# ARCHITECT one\_filter\_dragonn MODEL

```
multi_filter_dragonn_parameters = {  
    'seq_length': 1000,  
    'num_filters': [1],  
    'conv_width': [10],  
    'pool_width': 35}
```



# Training the Model

# TRAIN one\_filter\_dragonn MODEL

```
one_filter_dragonn = get_SequenceDNN(one_filter_dragonn_parameters)
train_SequenceDNN(one_filter_dragonn, simulation_data)
```

Epoch 1:

Train Loss: 0.7134      Balanced Accuracy: 50.86%    auROC: 0.509    auPRC: 0.502    Recall at 5%|10%|20%  
FDR: 0.0%|0.0%|0.0%    Num Positives: 1957    Num Negatives: 2043  
Valid Loss: 0.7077      Balanced Accuracy: 53.47%    auROC: 0.524    auPRC: 0.552    Recall at 5%|10%|20%  
FDR: 0.6%|0.6%|0.6%    Num Positives: 528    Num Negatives: 472 \*

.....

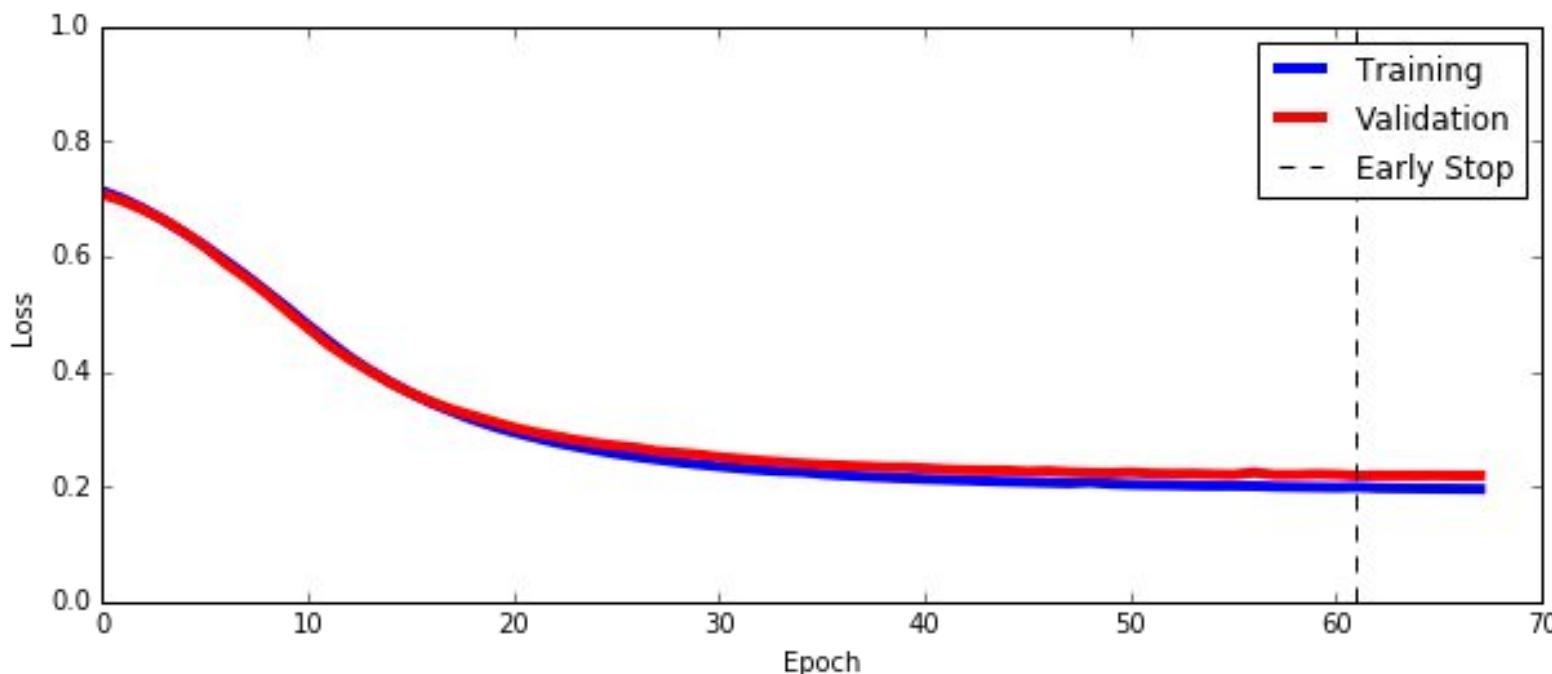
Epoch 68:

Train Loss: 0.1956      Balanced Accuracy: 91.63%    auROC: 0.977    auPRC: 0.975    Recall at 5%|10%|20%  
FDR: 86.1%|93.5%|98.6%    Num Positives: 1957    Num Negatives: 2043  
Valid Loss: 0.2189      Balanced Accuracy: 91.40%    auROC: 0.970    auPRC: 0.972    Recall at 5%|10%|20%  
FDR: 78.4%|94.9%|98.9%    Num Positives: 528    Num Negatives: 472

Finished training after 68 epochs.

# VISUALIZE one\_filter\_dragonn RESULTS

SequenceDNN\_learning\_curve(one\_filter\_dragonn)



# ARCHITECT MULTI-FILTER DRAGONN MODEL

```
multi_filter_dragonnn_parameters = {  
    'seq_length': 1000,  
    'num_filters': [15], ## notice the change from 1 filter to 15 filters  
    'conv_width': [10],  
    'pool_width': 35}  
  
multi_filter_dragonnn = get_SequenceDNN(multi_filter_dragonnn_parameters)  
train_SequenceDNN(multi_filter_dragonnn, simulation_data)  
SequenceDNN_learning_curve(multi_filter_dragonnn)
```

# TRAIN MULTI-FILTER DRAGONN MODEL

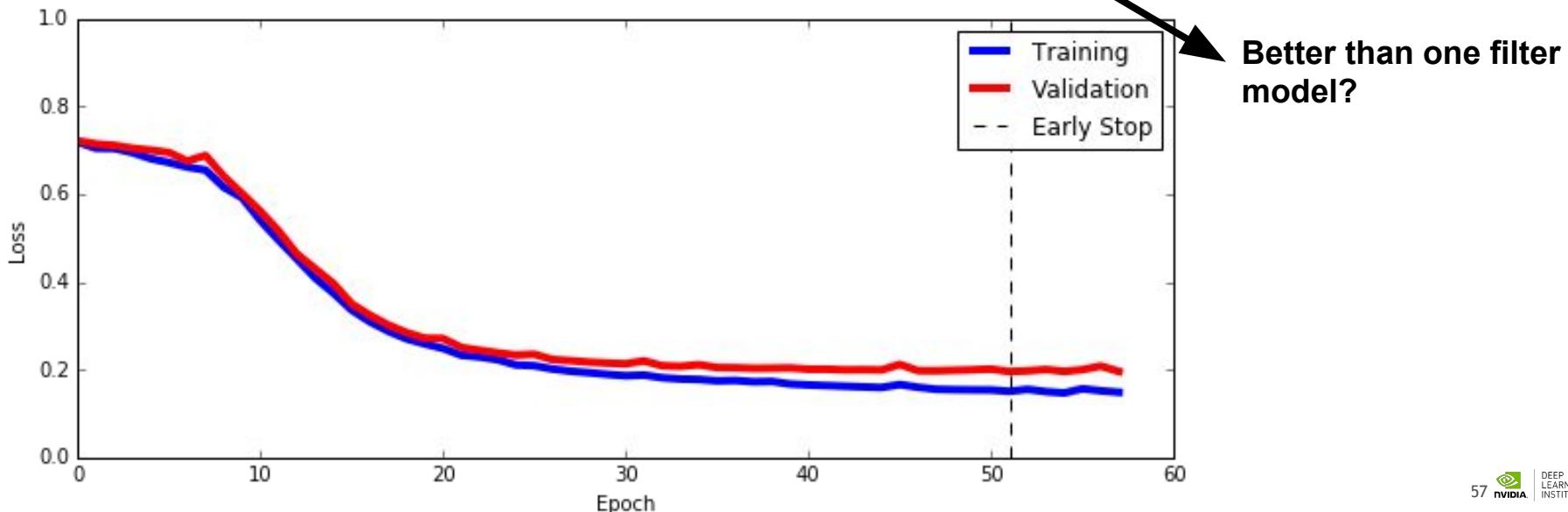
Epoch 58:

Train Loss: 0.1484      Balanced Accuracy: 94.40%    auROC: 0.988    auPRC: 0.987

Recall at 5%|10%|20% FDR: 93.4%|97.8%|99.7%    Num Positives: 1957    Num Negatives: 2043

Valid Loss: 0.1958      Balanced Accuracy: 91.97%    **auROC: 0.976**    auPRC: 0.977

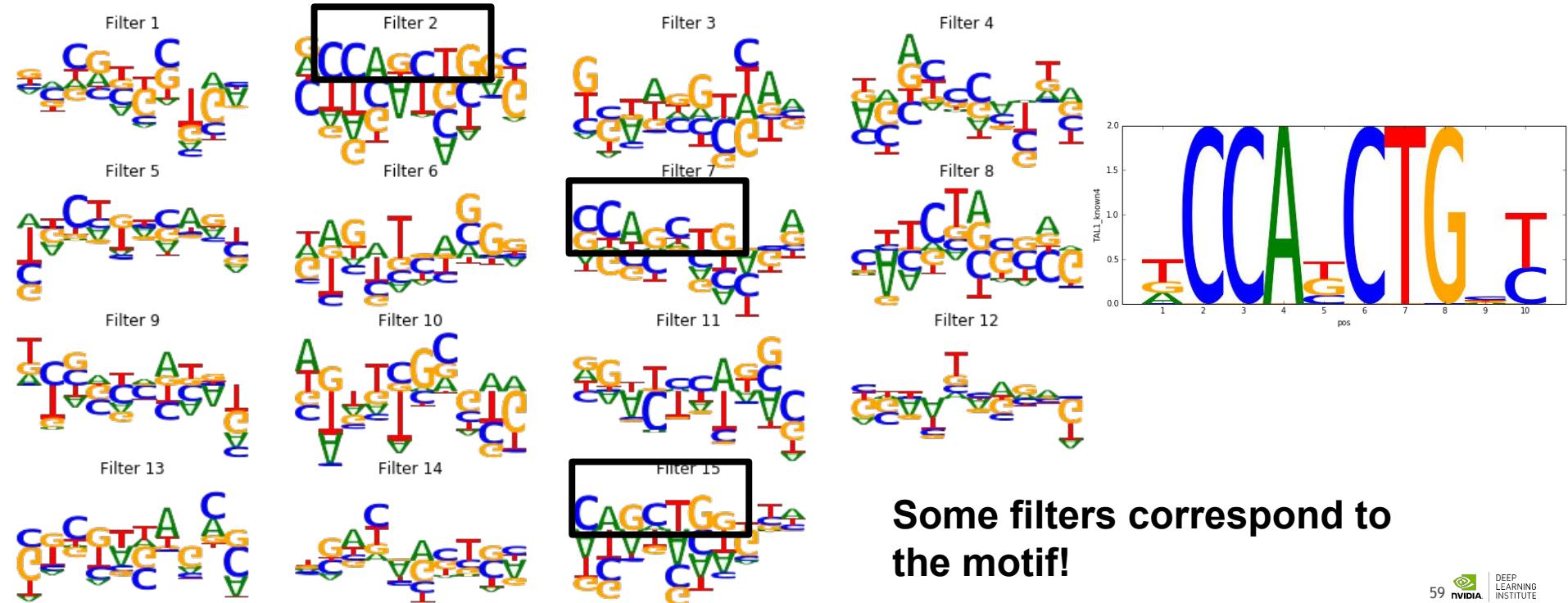
Recall at 5%|10%|20% FDR: 88.6%|95.3%|99.2%    Num Positives: 528    Num Negatives: 47



# Interpreting the Model & Data

# Compare learned filters to simulated motif

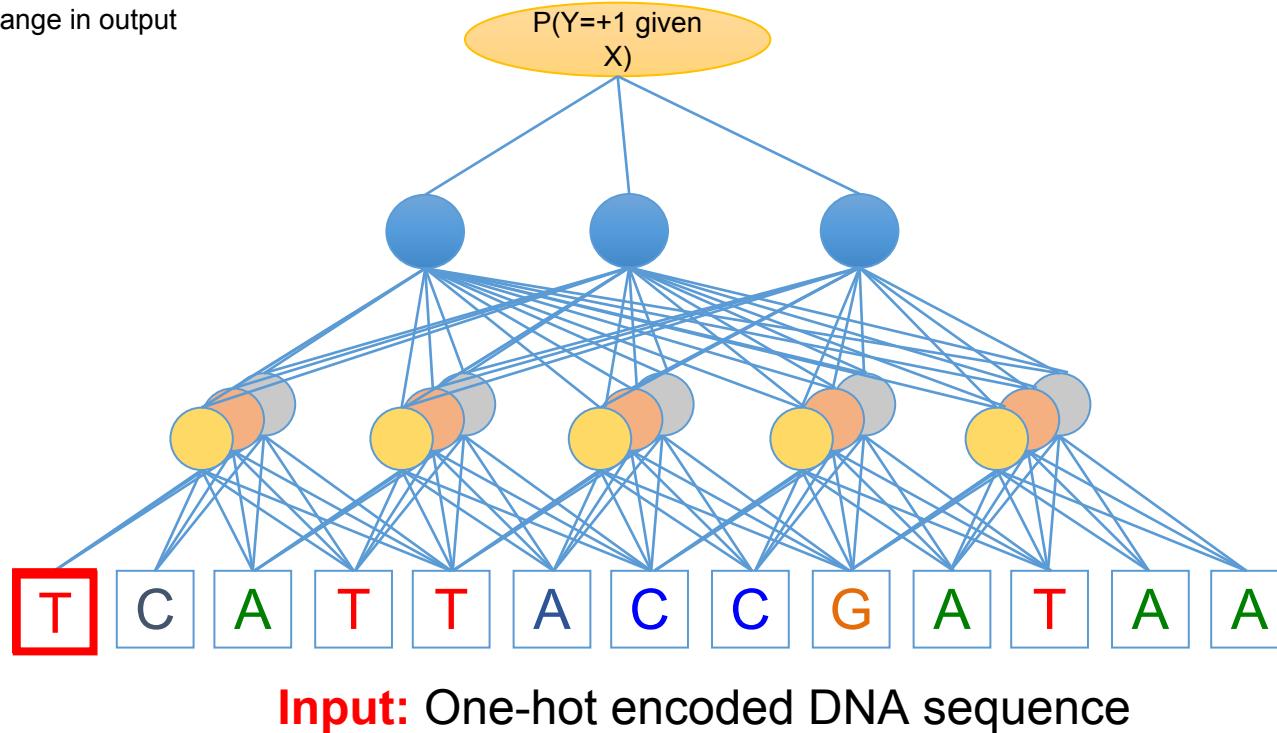
```
interpret_SequenceDNN_filters(multi_filter_dragonn, simulation_data)
```



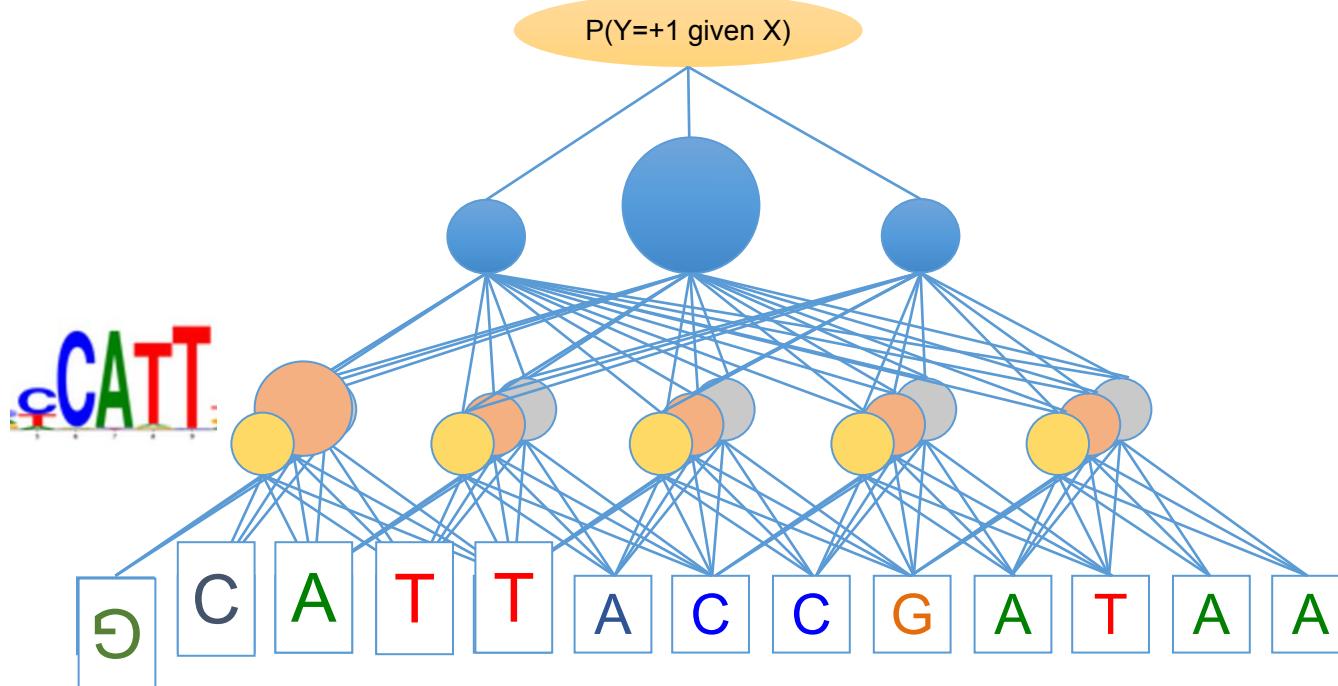
# Interpretation: In-silico mutagenesis

**Output:** Bound (+1) vs. not bound (0)

Assess change in output



# Interpretation: DeepLIFT (Deep learning feature importance)



Avanti Shrikumar

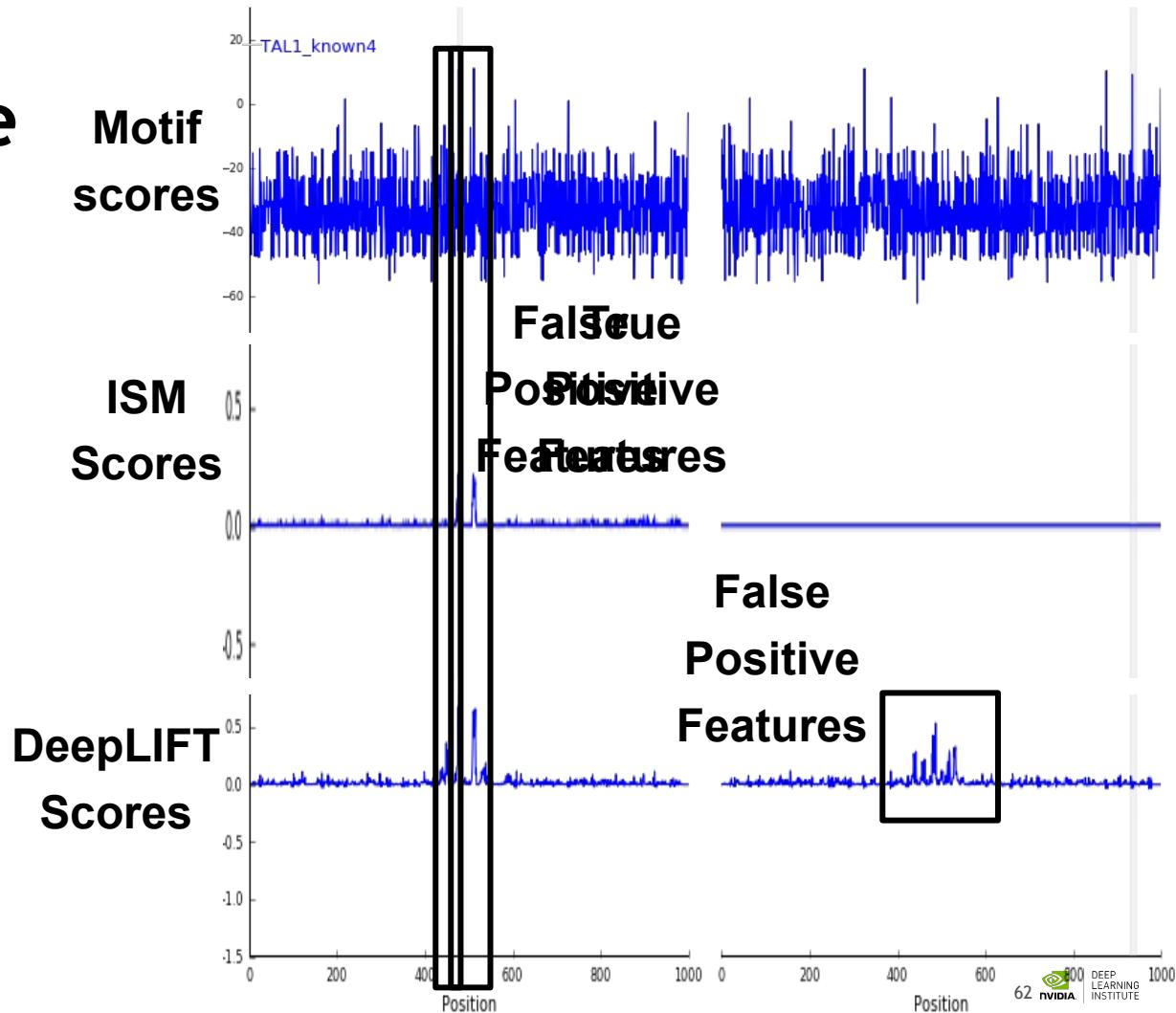
# Compare feature importances to motif scores

```
interpret_data_with_SequenceDNN(  
    multi_filter_dragonn,  
    simulation_data)
```

ISM



DeepLIFT



# Enhancement #1: multiple conv layers

# ARCHITECT MULTI-LAYER DRAGONN MODEL

```
multi_layer_dragonnn_parameters = {  
    'seq_length': 1000,  
    'num_filters': [15, 15, 15], ## notice the change to multiple filter values, one for each layer  
    'conv_width': [10, 10, 10],  
    'pool_width': 35}  
multi_layer_dragonnn = get_SequenceDNN(multi_layer_dragonnn_parameters)  
train_SequenceDNN(multi_layer_dragonnn, simulation_data)  
SequenceDNN_learning_curve(multi_layer_dragonnn)
```

# TRAIN MULTI-LAYER DRAGONN MODEL

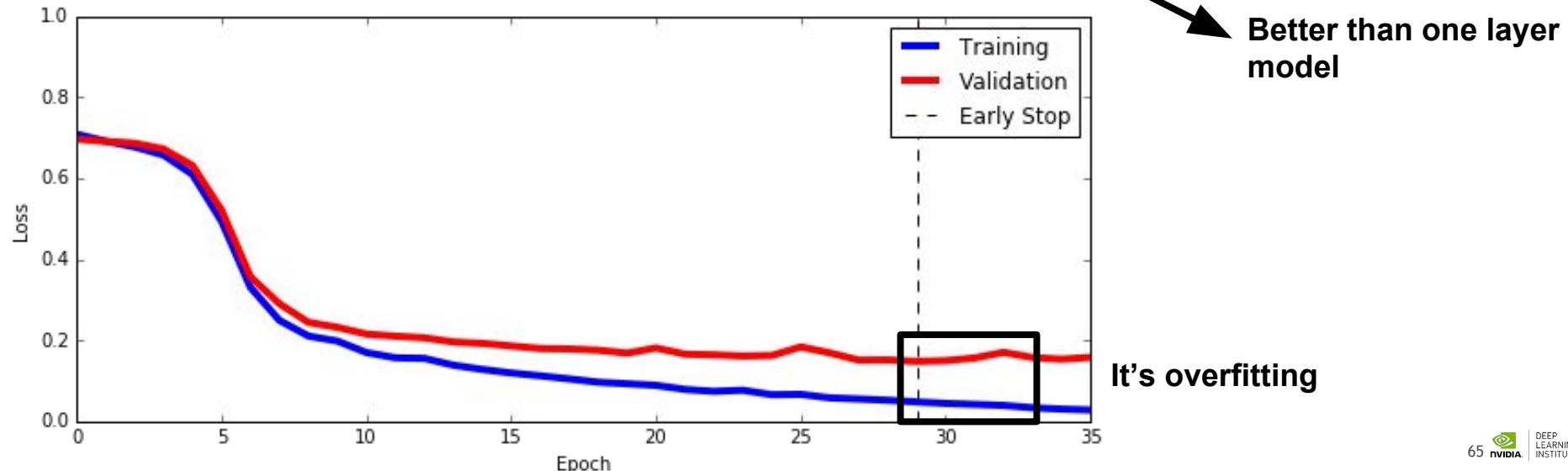
Epoch 36:

Train Loss: 0.0278      Balanced Accuracy: 99.73%    auROC: 1.000    auPRC: 1.000

Recall at 5%|10%|20% FDR: 100.0%|100.0%|100.0%      Num Positives: 1957    Num Negatives: 2043

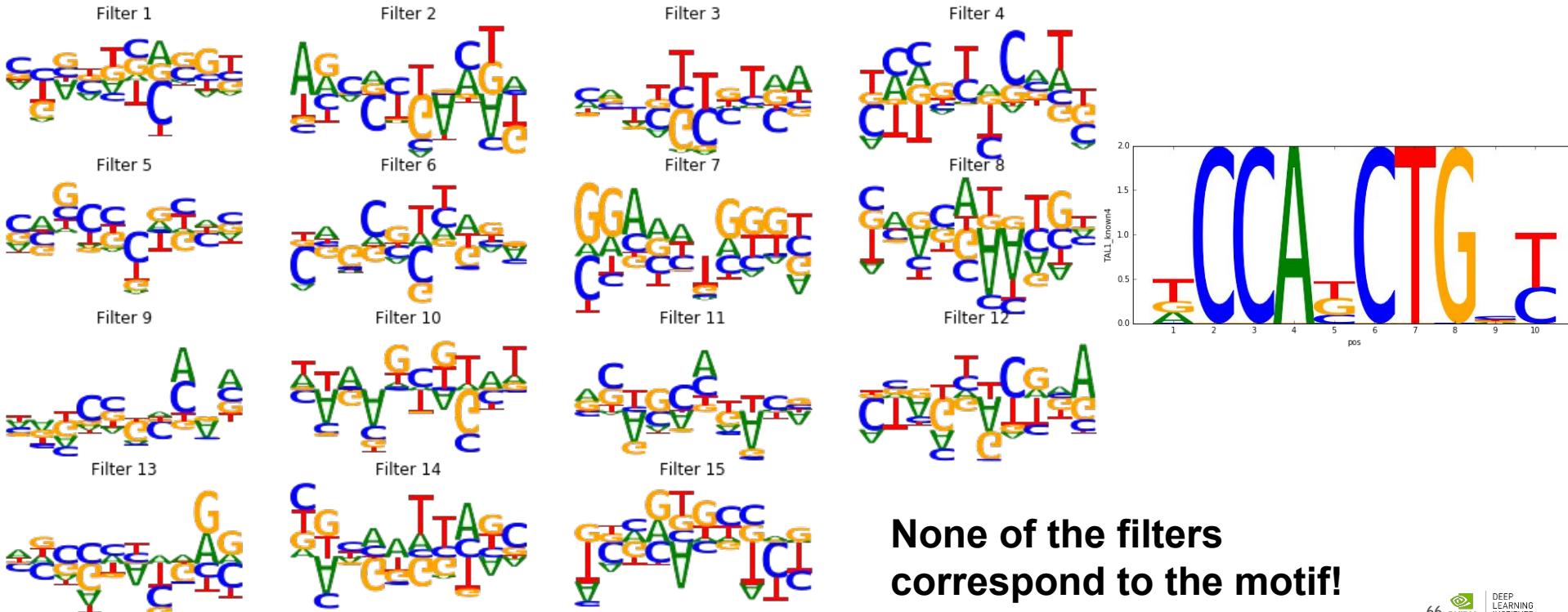
Valid Loss: 0.1577      Balanced Accuracy: 94.34%    **auROC: 0.987**    auPRC: 0.985

Recall at 5%|10%|20% FDR: 95.6%|99.2%|99.8%      Num Positives: 528    Num Negatives: 47



# Compare learned filters to simulated motif

```
interpret_SequenceDNN_filters(multi_layer_dragonn, simulation_data)
```



# Compare feature importances to motif scores

```
interpret_data_with_SequenceDNN(  
    multi_layer_dragonn,  
    simulation_data)
```

ISM

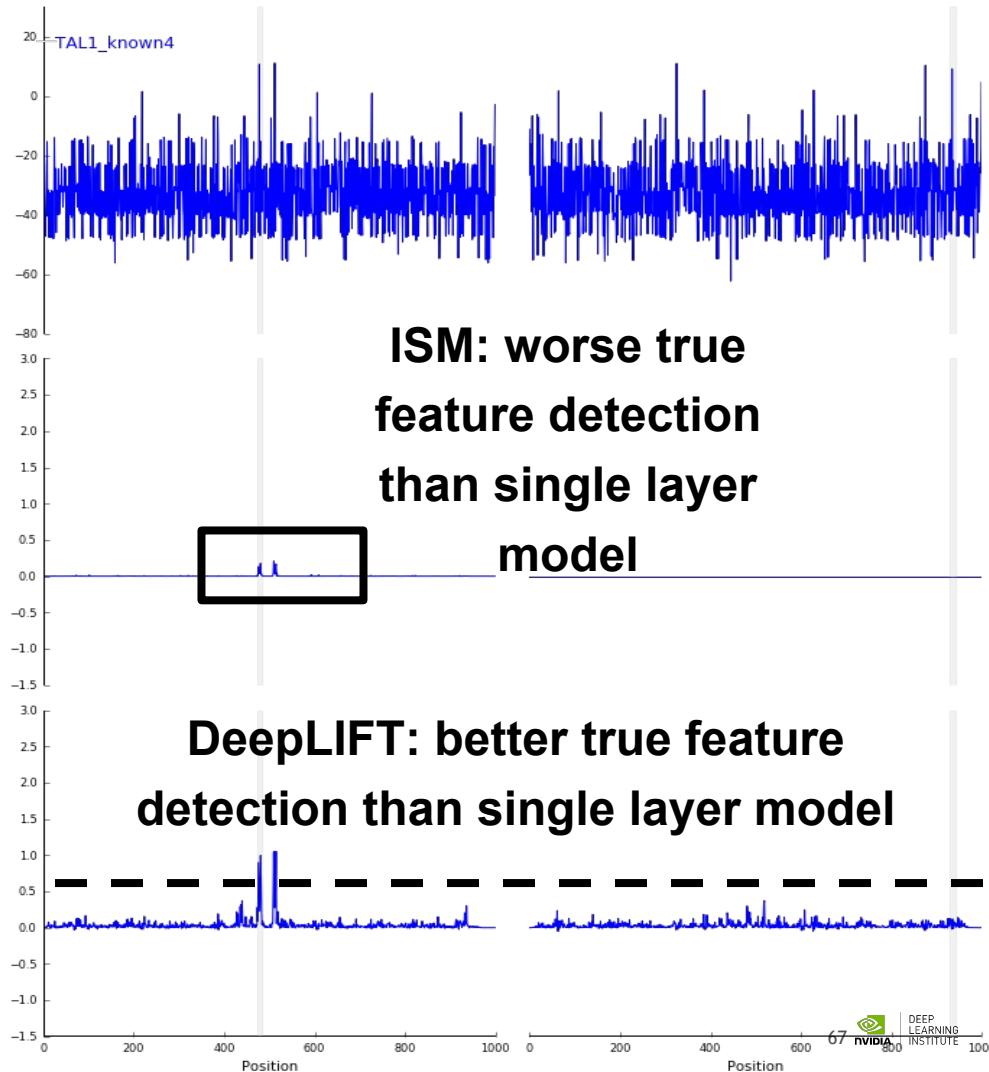


DeepLIFT

Motif  
scores

ISM  
Scores

DeepLIFT  
Scores



## Enhancement #2: dropout regularization

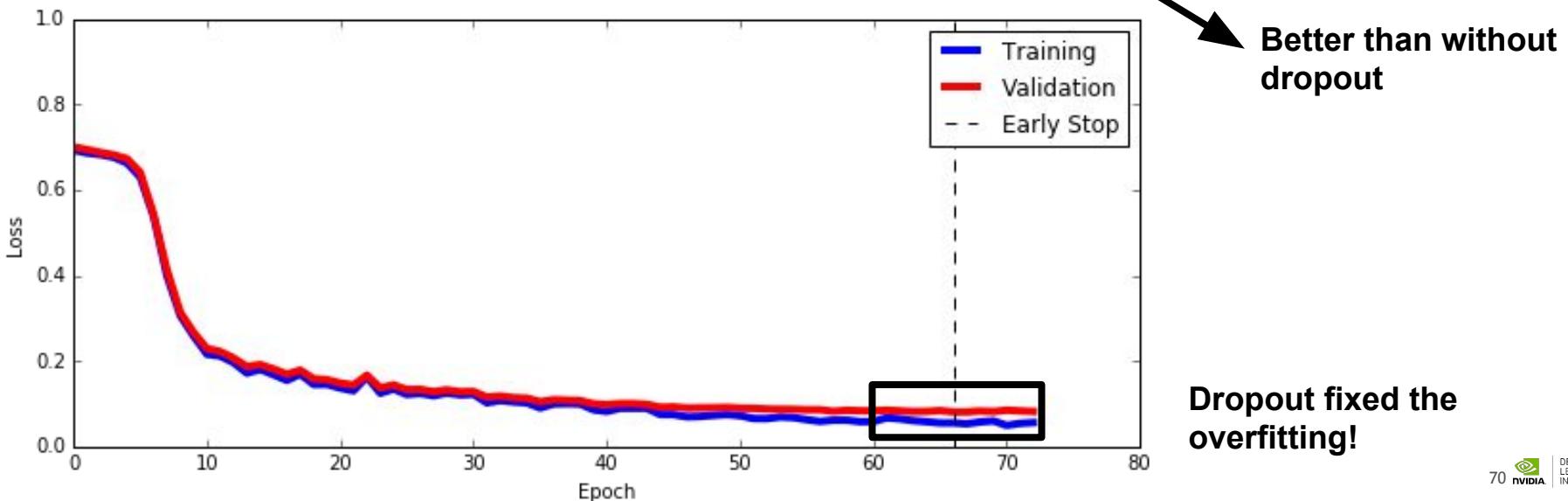
# ARCHITECT REGULARIZED DRAGONN MODEL

```
regularized_multi_layer_dragonn_parameters = {  
    'seq_length': 1000,  
    'num_filters': [15, 15, 15],  
    'conv_width': [10, 10, 10],  
    'pool_width': 35,  
    'dropout': 0.2} ## we introduce dropout of 0.2 on every convolutional layer for regularization  
regularized_multi_layer_dragonn = get_SequenceDNN(  
    regularized_multi_layer_dragonn_parameters)  
train_SequenceDNN(regularized_multi_layer_dragonn, simulation_data)  
SequenceDNN_learning_curve(regularized_multi_layer_dragonn)
```

# TRAIN REGULARIZED DRAGONN MODEL

Epoch 73:

Train Loss: 0.0558      Balanced Accuracy: 98.36%    auROC: 0.998    auPRC: 0.998  
Recall at 5%|10%|20% FDR: 99.8%|100.0%|100.0%      Num Positives: 1957    Num Negatives: 2043  
Valid Loss: 0.0817      Balanced Accuracy: 97.23%    **auROC: 0.994**    auPRC: 0.992  
Recall at 5%|10%|20% FDR: 98.5%|100.0%|100.0%      Num Positives: 528    Num Negatives: 472



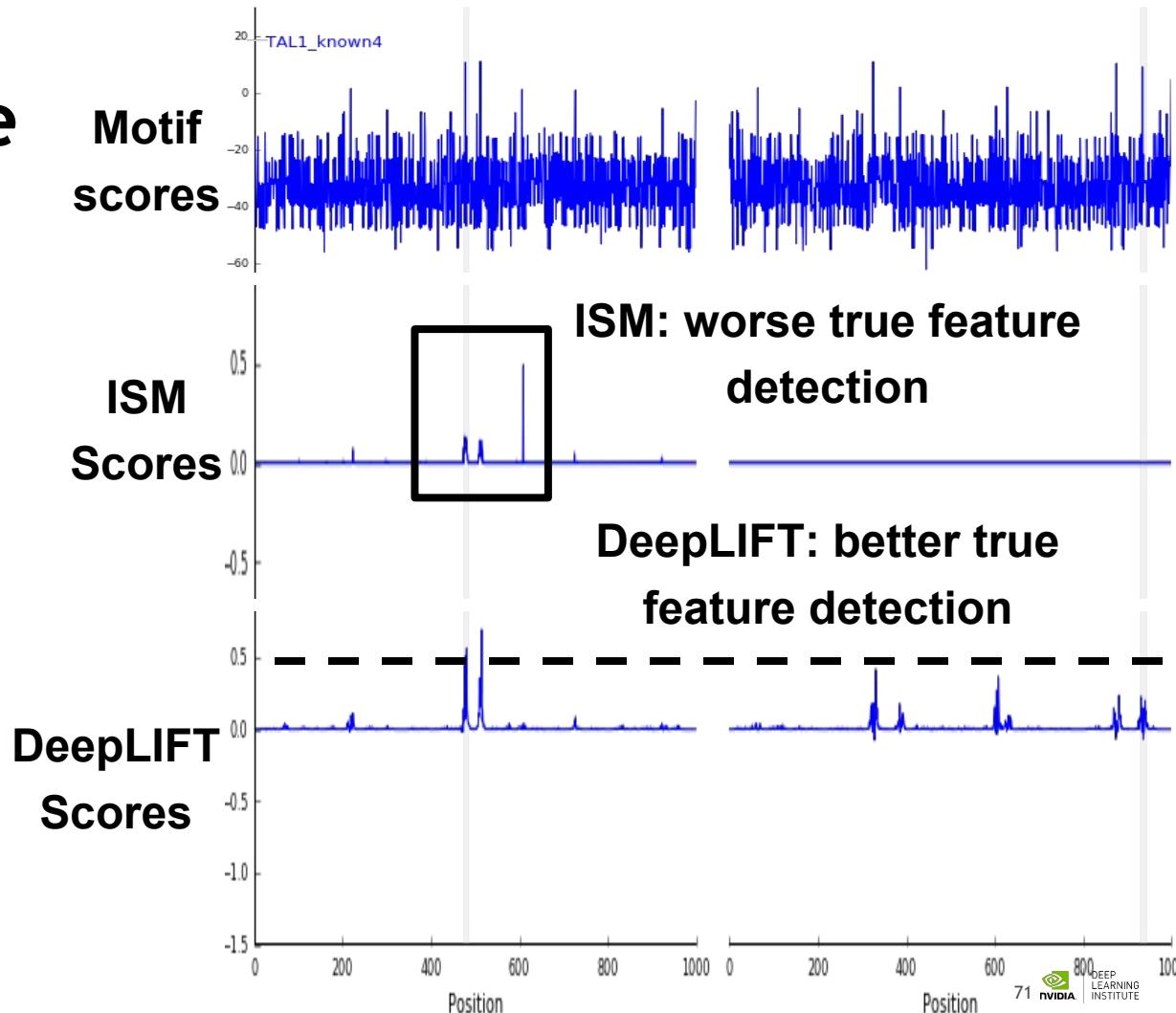
# Compare feature importances to motif scores

```
interpret_data_with_SequenceDNN(  
    regularized_multi_layer_dragonn,  
    simulation_data)
```

DeepLIFT!



ISM!



# Discussion

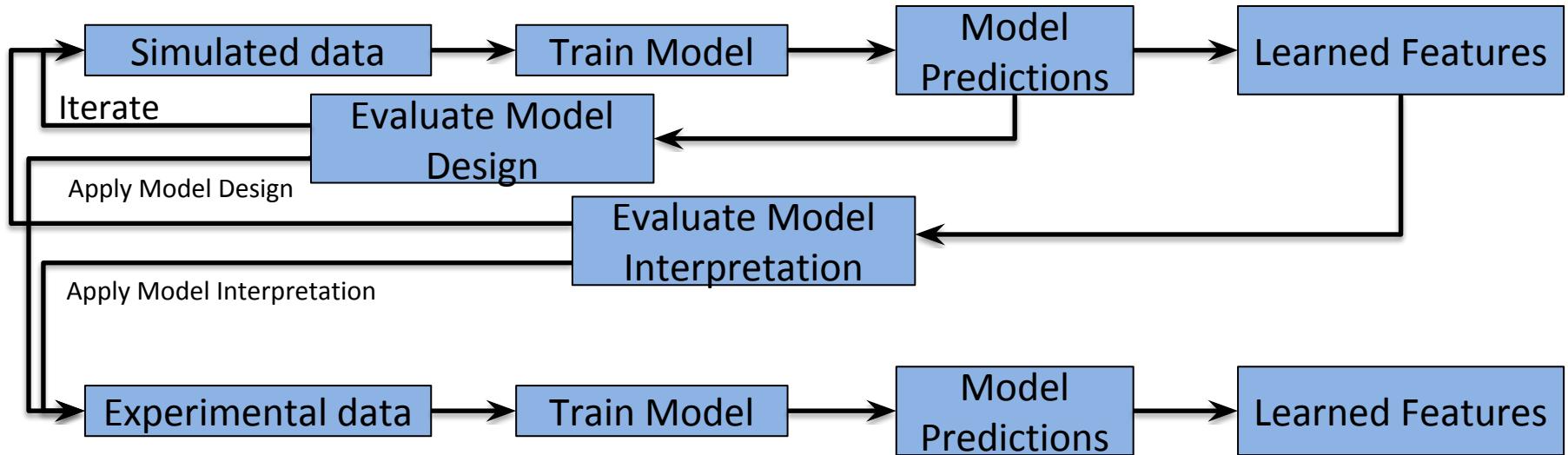
Single layer -> multi-layer -> regularized multi layer led to

- improved auROC
- Less interpretable 1st layer filters
- Worse ISM feature importance
- Better DeepLIFT feature importance

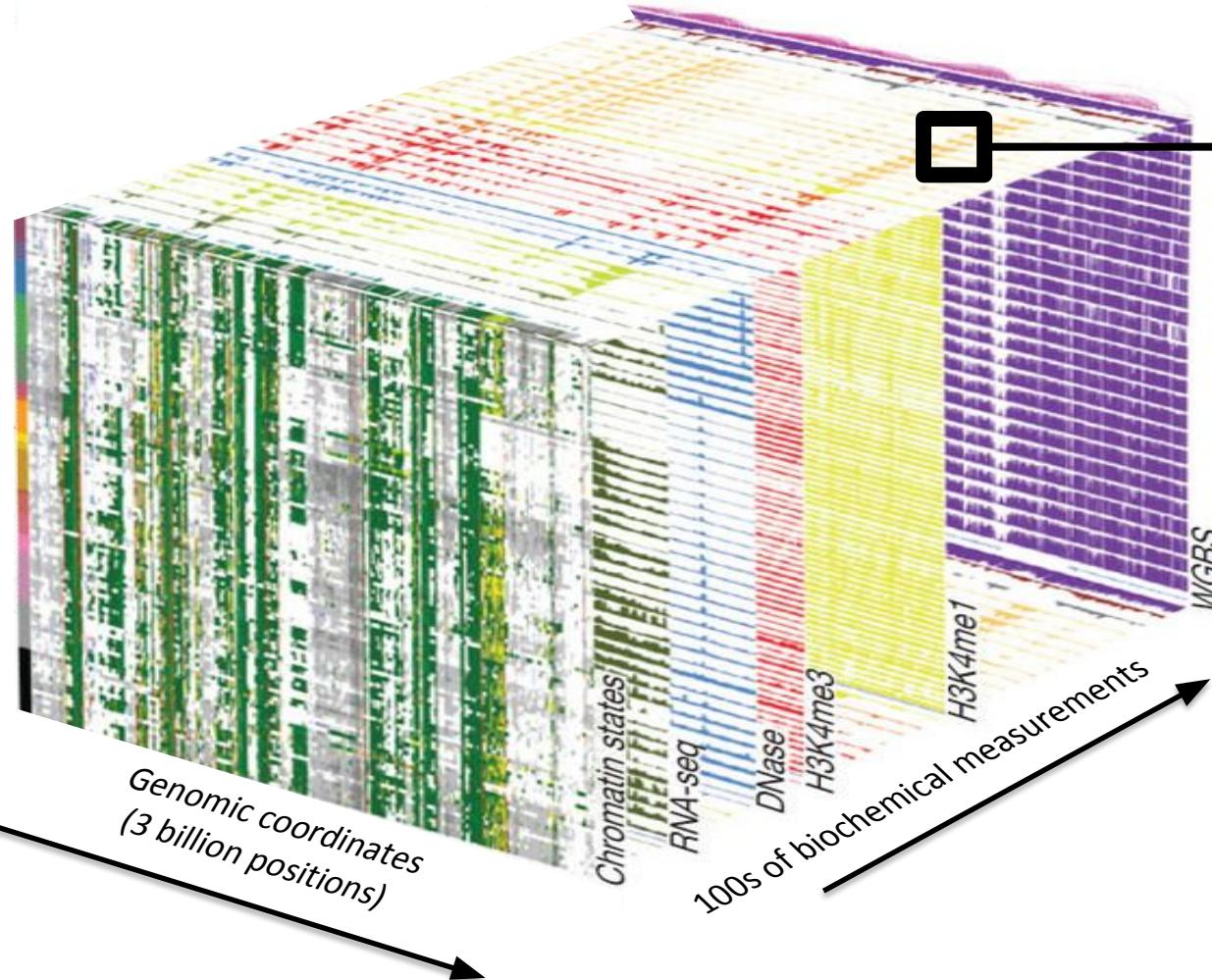
Why? What is the expected ISM behavior in the limit of the “perfect” model of this simulation?

# The DragoNN workflow

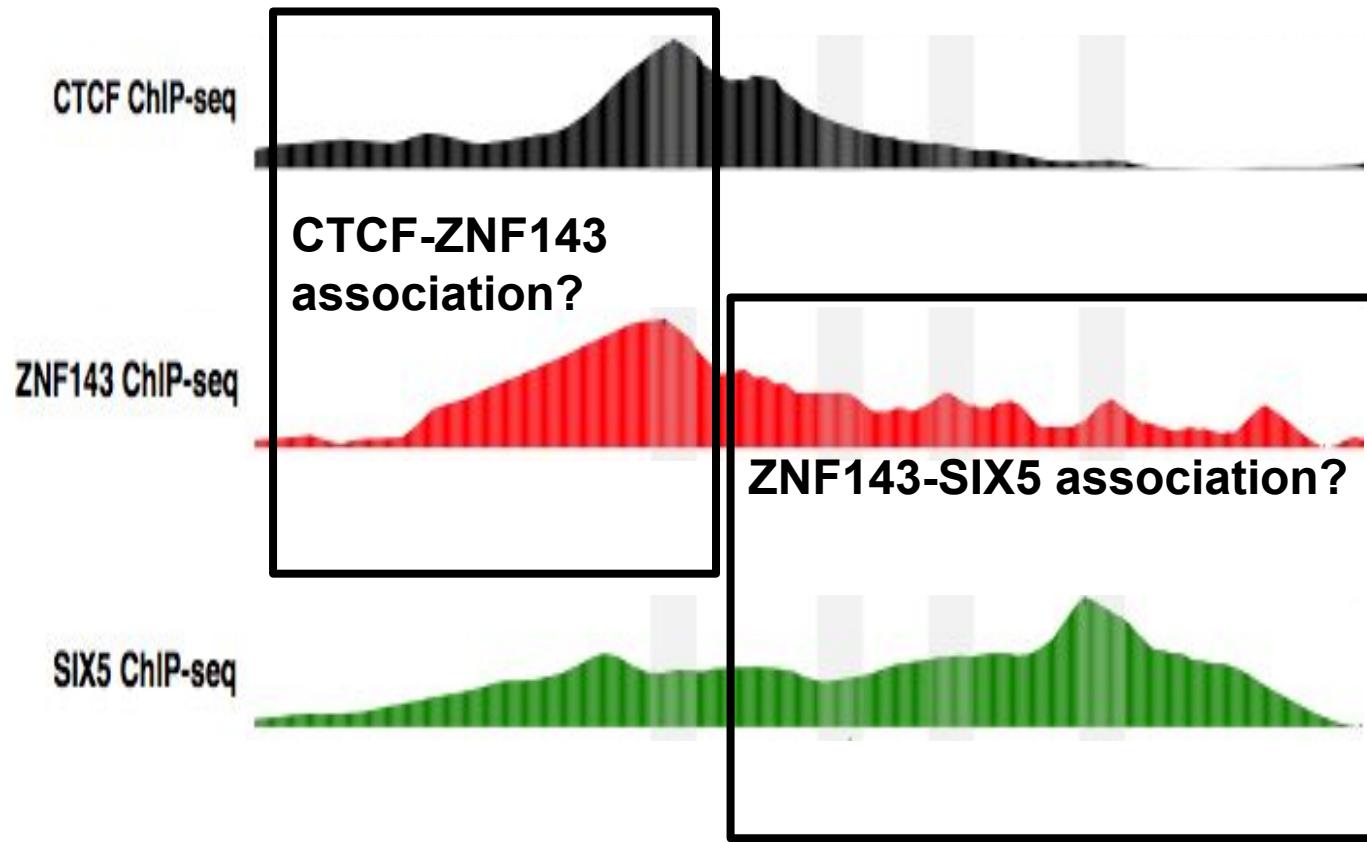
Systematic Development and Application of Model Design and Interpretation



# Example Application: in-vivo TF binding



# Example Application



i **CGCCCTCTGGCG**  
Identified predictive sequence

Identified  
predictive  
sequence

ii **GGGAACGTGAG**  
iii **GGGACTCGGAG**  
iv **GGGAATAGTAG**

No CTCF-SIX5  
direct association?

**CTCF-SIX5 association conditional on ZNF143 first discovered in:**

**Architecture of the human regulatory network derived from ENCODE data**

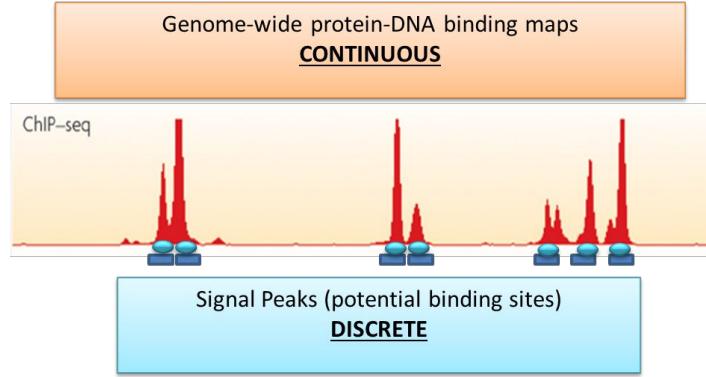
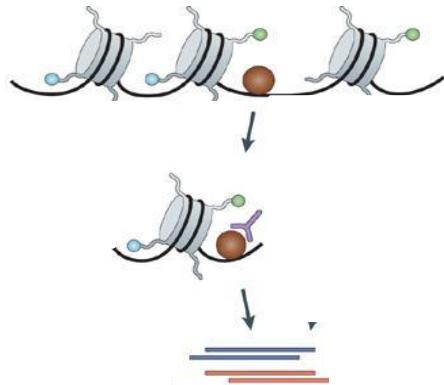
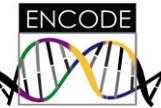
**Mark B. Gerstein, Anshul Kundaje et al.**

*Nature* **489**, 91–100 (06 September 2012) | doi:10.1038/nature11245

Received 09 December 2011 | Accepted 22 May 2012 | Published online 05 September 2012

# State of the Field: TF Binding Models

# Challenge Data: “Ground truth”



## Genome-wide in-vivo TF binding data

- **Ground truth:** High quality ChIP-seq data targeting specific TF in specific cell type
- **Peaks:** Regions of enrichment => likely binding events. Peaks ranked by enrichment and reproducibility across replicates

# Challenge Data: Output variable

---

## Binary Classification problem

- Is a genomic region bound or not by TF of interest?
- Two types of TF ChIP-seq peaks called
  - **Conservative:** High-confidence, reproducible peaks labeled as **bound (B)**
  - **Relaxed:** Lower confidence, reproducible peaks labeled as **ambiguous (A)** – not used in evaluation
- Rest of the genome is labeled **unbound (U)**
- Resolution of data: 200 bp bins every 50 bp

chr	start	stop	GM12878	H1-hESC
chr10	600	800	B	B
chr10	650	850	U	U
chr10	700	900	U	U
chr10	750	950	U	U
chr10	800	1000	B	U
chr10	850	1050	U	U
chr10	900	1100	U	U
chr10	950	1150	A	B
chr10	1000	1200	U	U

# Challenge Data: Problem set up

---

## Held-out cell types

For each TF

- ChIP-seq data is provided for some **training cell types** for participants to train models
- Binding to be predicted in new **held-out cell types**

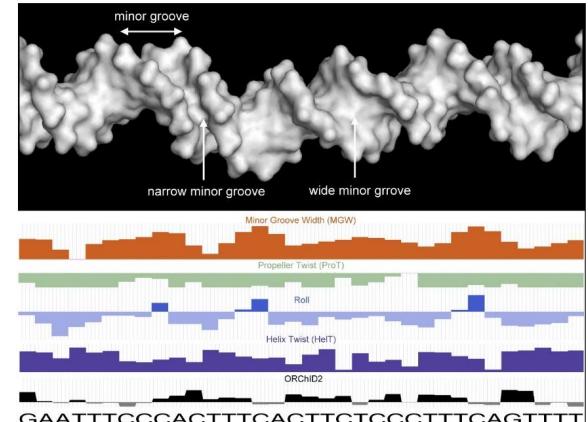
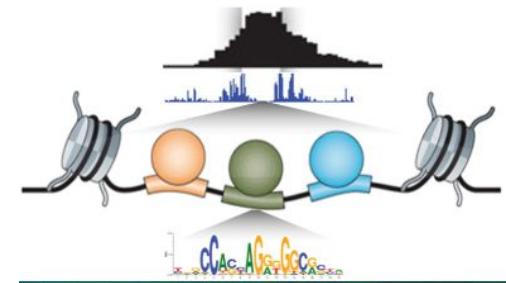
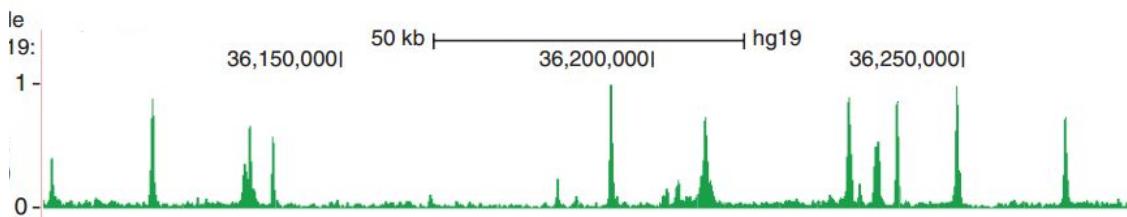
## Held-out chromosomes

- Binding data in **chr1, 21 and 8** is never provided (even in training cell types)
- Performance evaluated on these 3 held-out chromosomes

**32 diverse TFs, 14 cell types** (Training and held-out cell types are different for each TF)

# Challenge Data: Predictor variables

1. DNA Sequence features
2. Chromatin accessibility
3. Gene expression
4. Gene annotations
5. DNA shape features



(Zhou et al. 2013 Nucleic Acids Res., 41, W56-W62.)

# Challenge Questions

---

**Binary prediction task:** Predict probability (between 0 and 1) of being bound by TF for each 200 bp bin every 50 bp in specified chromosomes in specified cell type

Chromosome	Start Position	End Position	Predicted Probability
chr1	600	800	2.16e-39
chr1	650	850	5.65e-169
chr1	700	900	6.18e-37
chr1	750	950	2.70e-26
chr1	800	1000	6.40e-74

# Main Challenge Question

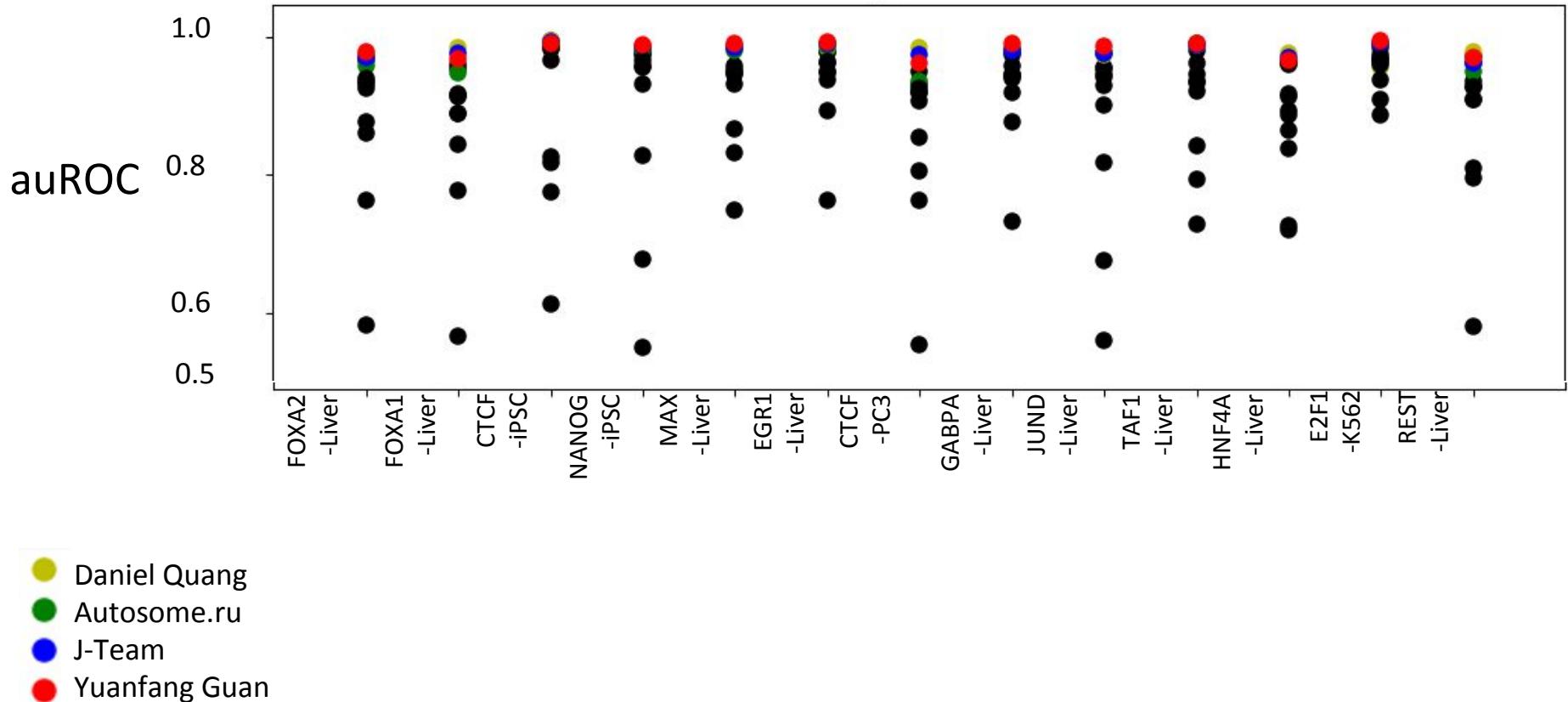
---

Across cell-type prediction

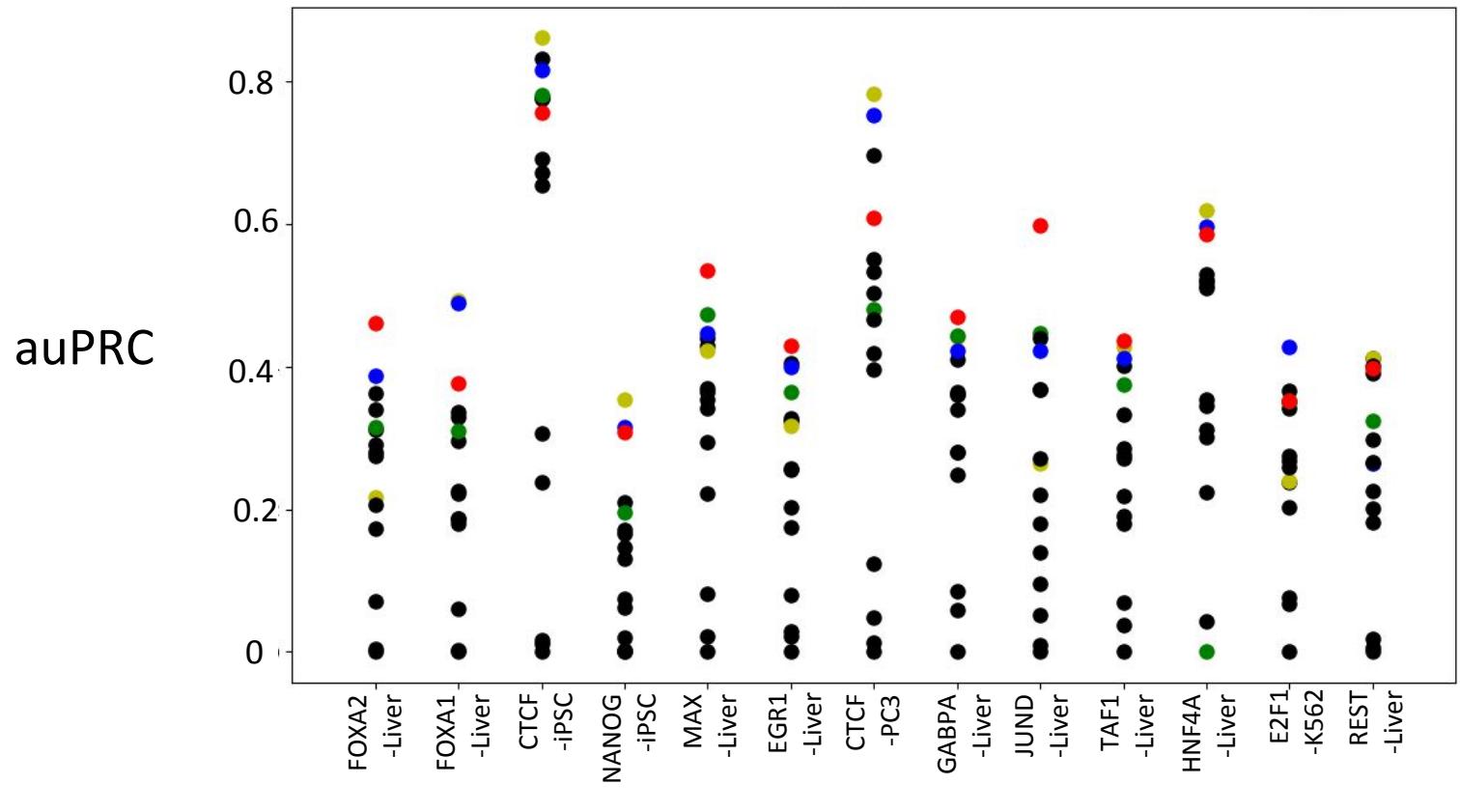
Predict binding in hidden cell types on hidden chromosomes

(Practical definition of “prediction”)

# auROC across all TF/cell type



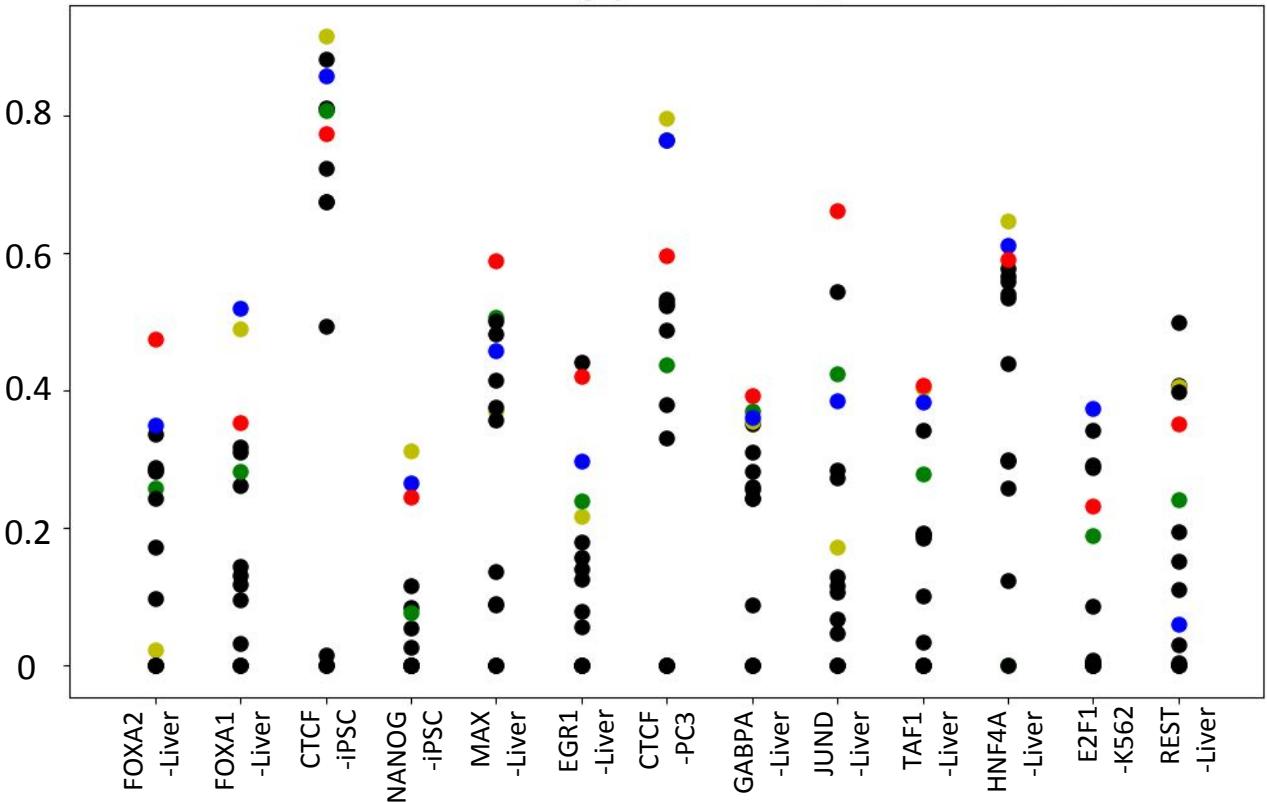
# auPRC across all TF/cell type



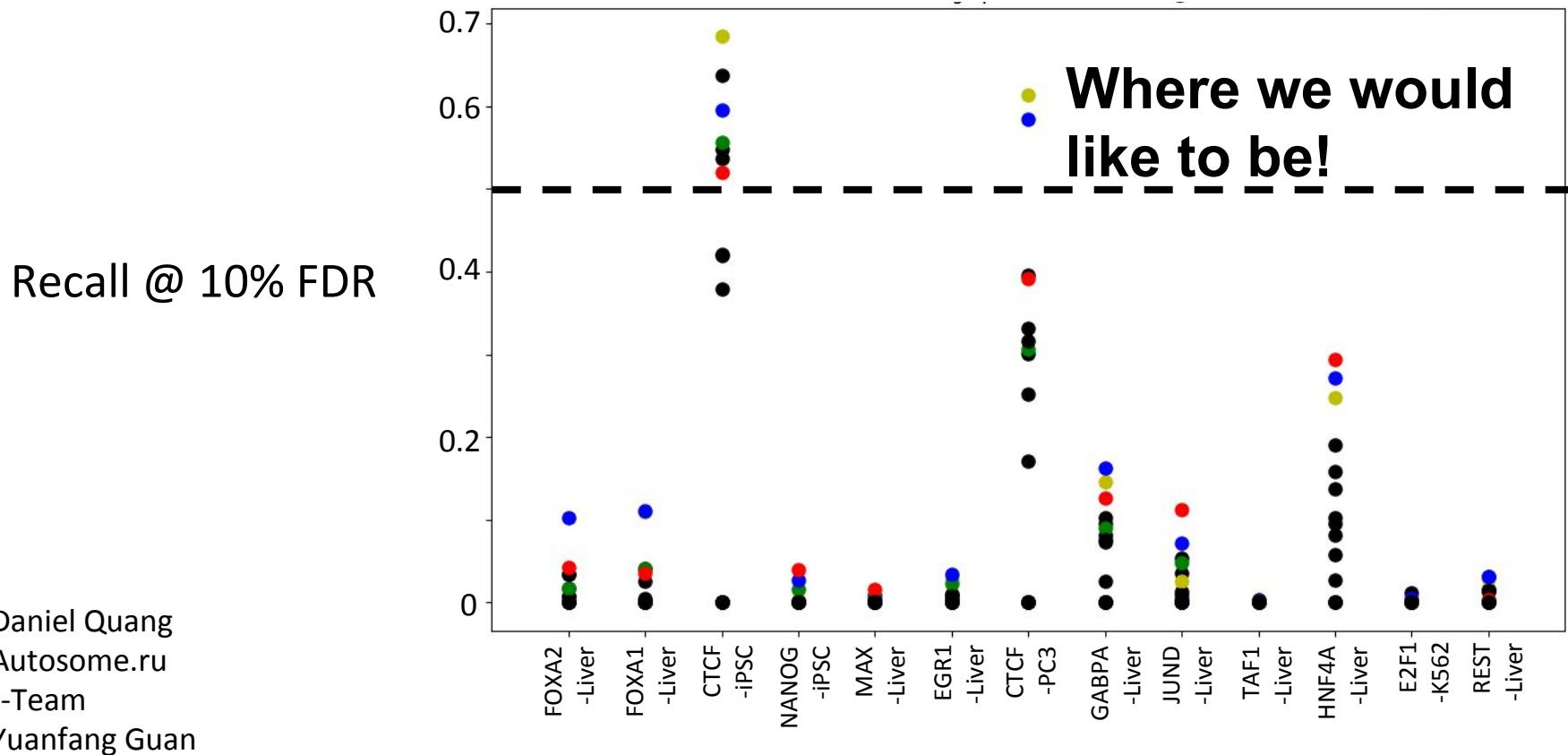
# Recall @ 50%FDR across all TF/cell type

Recall @ 50% FDR

- Daniel Quang
- Autosome.ru
- J-Team
- Yuanfang Guan



# Recall @ 10%FDR across all TF/cell type





# WHAT'S NEXT

# WHAT'S NEXT

- Access tutorial on github:  
[https://github.com/kundajelab/dragonnn/blob/GTC/examples/GTC\\_workshop\\_tutorial.ipynb](https://github.com/kundajelab/dragonnn/blob/GTC/examples/GTC_workshop_tutorial.ipynb)
- Follow DragoNN repo for discussions and more tutorials:  
<https://github.com/kundajelab/dragonnn/>
- Reach out to NVIDIA and the Deep Learning Institute
- Reach out to Johnny Israeli and Anshul Kundaje

# WHAT'S NEXT

## TAKE SURVEY

...for the chance to win an NVIDIA SHIELD TV.

Check your email for a link.

## ACCESS ONLINE LABS

Check your email for details to access more DLI training online.

## ATTEND WORKSHOP

Visit [www.nvidia.com/dli](http://www.nvidia.com/dli) for workshops in your area.

## JOIN DEVELOPER PROGRAM

Visit <https://developer.nvidia.com/join> for more.

# GPU TECHNOLOGY CONFERENCE

May 8 - 11, 2017 | Silicon Valley | #GTC17  
[www.gputechconf.com](http://www.gputechconf.com)



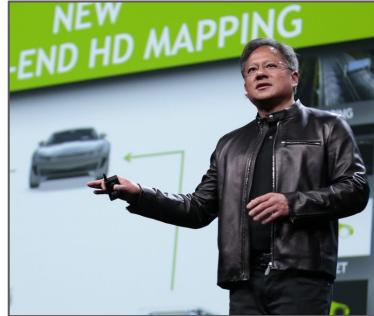
## CONNECT

Connect with technology experts from NVIDIA and other leading organizations



## LEARN

Gain insight and valuable hands-on training through hundreds of sessions and research posters



## DISCOVER

See how GPUs are creating amazing breakthroughs in important fields such as deep learning and AI



## INNOVATE

Hear about disruptive innovations from startups

# backup slides

# Performance Measures

		Predicted			
		+	-	Sensitivity (recall)	False negative rate
Actual	+	TP Type I error	FN Type II error	TP/●	FN/●
	-	FP Type I error	TN	False positive rate FP/●	Specificity TN/●
		Precision TP/	False omission rate FN/	<i>Lever et al. 2016 Nature Methods</i>	
		FDR FP/	Negative predictive value TN/		

- **Area under Receiver Operating Curve (auROC):** The area under the curve of Recall vs. False positive rate.
- **Area under Precision-Recall Curve (auPRC):** The area under the curve of Precision (1- False discovery rate) vs. Recall.
- **Recall @ X% FDR (10%, 50%):** The fraction of actual bound sites predicted as positives when controlling the False Discovery Rate (1-Precision) to be under X%.