# Deep learning oracles for genomic discovery

Kundaje lab
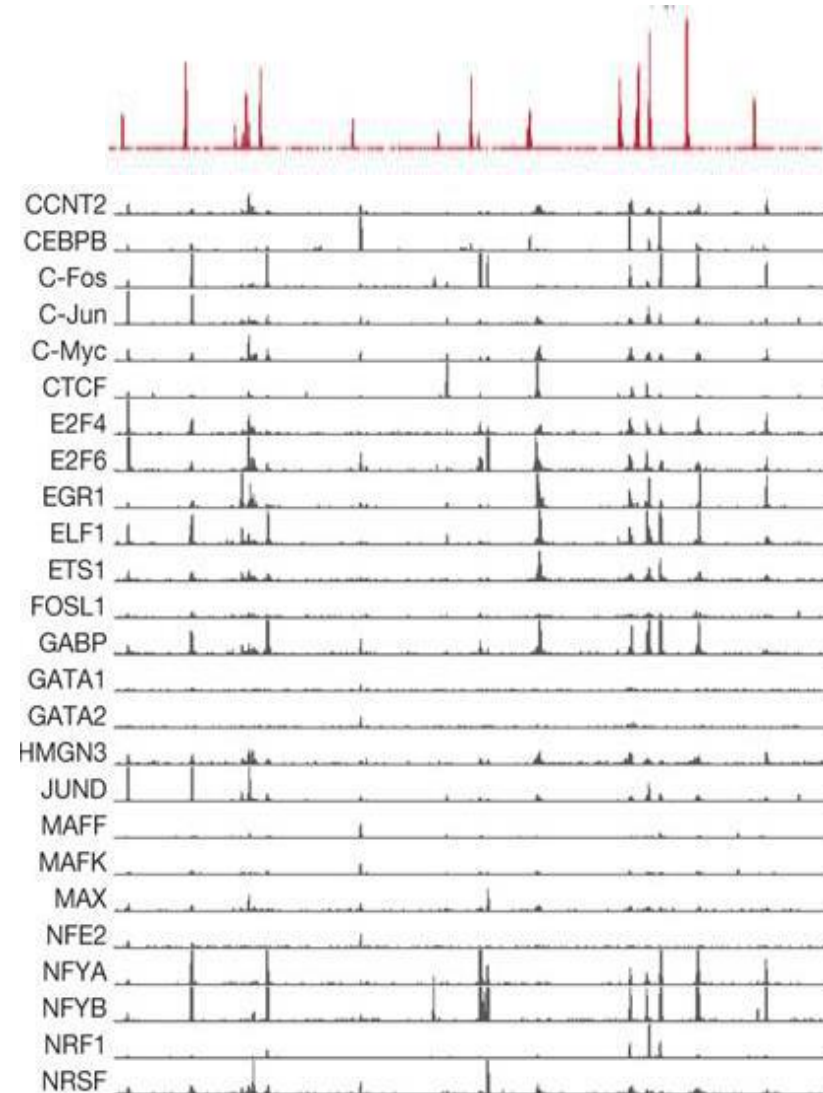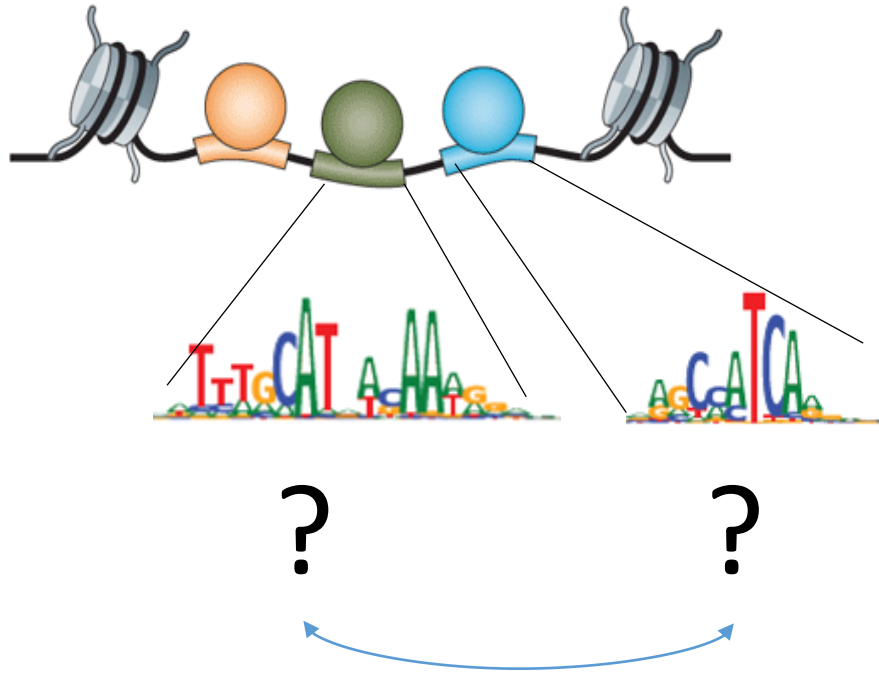
Genetics, Computer Science

Stanford University

http://anshul.kundaje.net

# Deciphering functional DNA words and their syntax in regulatory DNA



chromatin accessibility (ATAC-seq / DNase-seq)

Syntax: Rules of arrangement, preferred spacing, orientation, interactions between works

Protein-DNA binding maps (ChIP-seq, ChIP-exo)

Adapted from Thurman et al 2012

# Predictive model of regulatory DNA

Genome-wide protein-DNA binding map



…GACTTGAAACGGCATTG…
Inactive (0) (0.3)
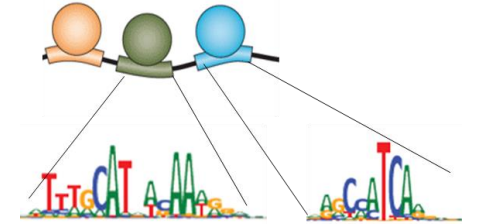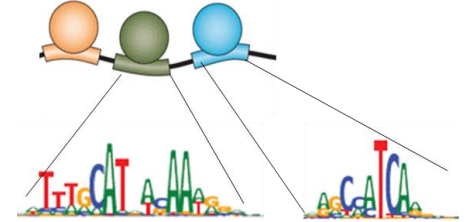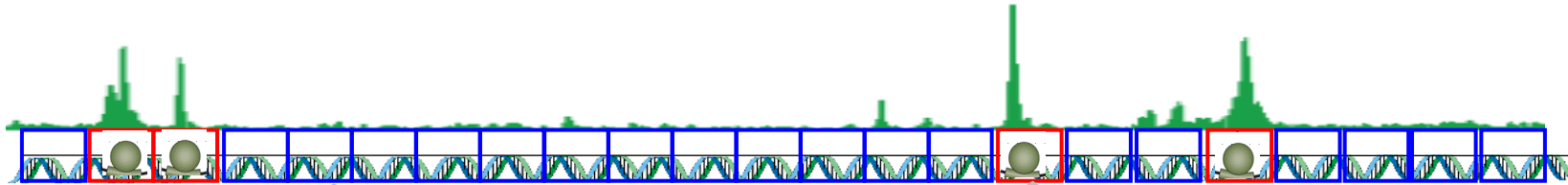
…GACAGATAATGCATTGA…
Active (+1) (20.2)

# Predictive model of regulatory DNA

Genome-wide protein-DNA binding map

...GACTTGAAACGGCATTG...
Inactive (0) (0.3)

...GACAGATAATGCATTGA...
Active (+1) (20.2)

...GACAGATAATGCATTGA...

...ACTGTCATGGATATTCT...

...GATATTCTACTGTAAG...

DNA sequences $(S_i)$

...CAACCTTGAACGGCATTG...

...GACTTGAAACGGCATTG...

...CAGTATGCATACGTGAA...

Classification
or Regression
model
$F(S_i)$

Arvey et al. 2012
Ghandi et al. 2014
Setty et al. 2015
Alipanahi et al. 2015
Zhou et al. 2015
Kelly et al. 2016, 2018
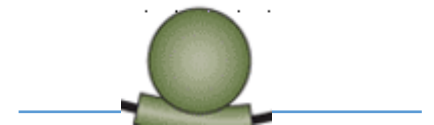Avsec et al. 2021

Class = +1 (20.2)

Class = +1 (10.6)

Class = +1 (15.8)

Measured
Labels $(Y_i)$
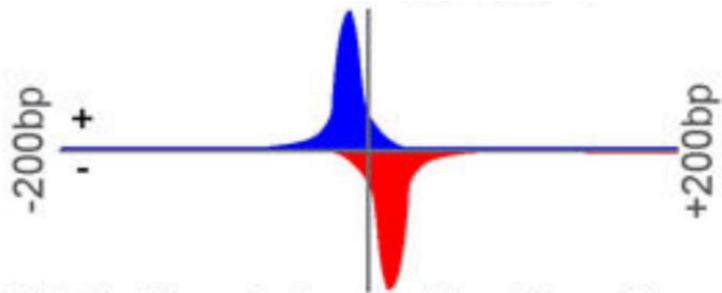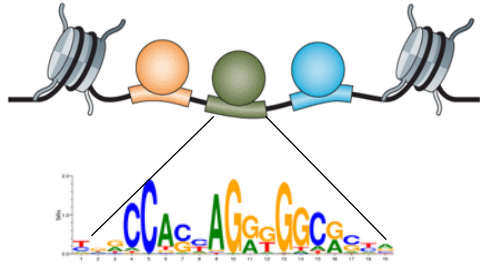
Bound

Class = 0 (0.3)

Class = 0 (1.2)

Class = 0 (3.5)

Unbound

# High-resolution 'shapes' and 'spans' of TF and chromatin profiles capture exquisite information about protein-DNA contacts



Distribution of stranded tag 5' positions around binding event (bp)
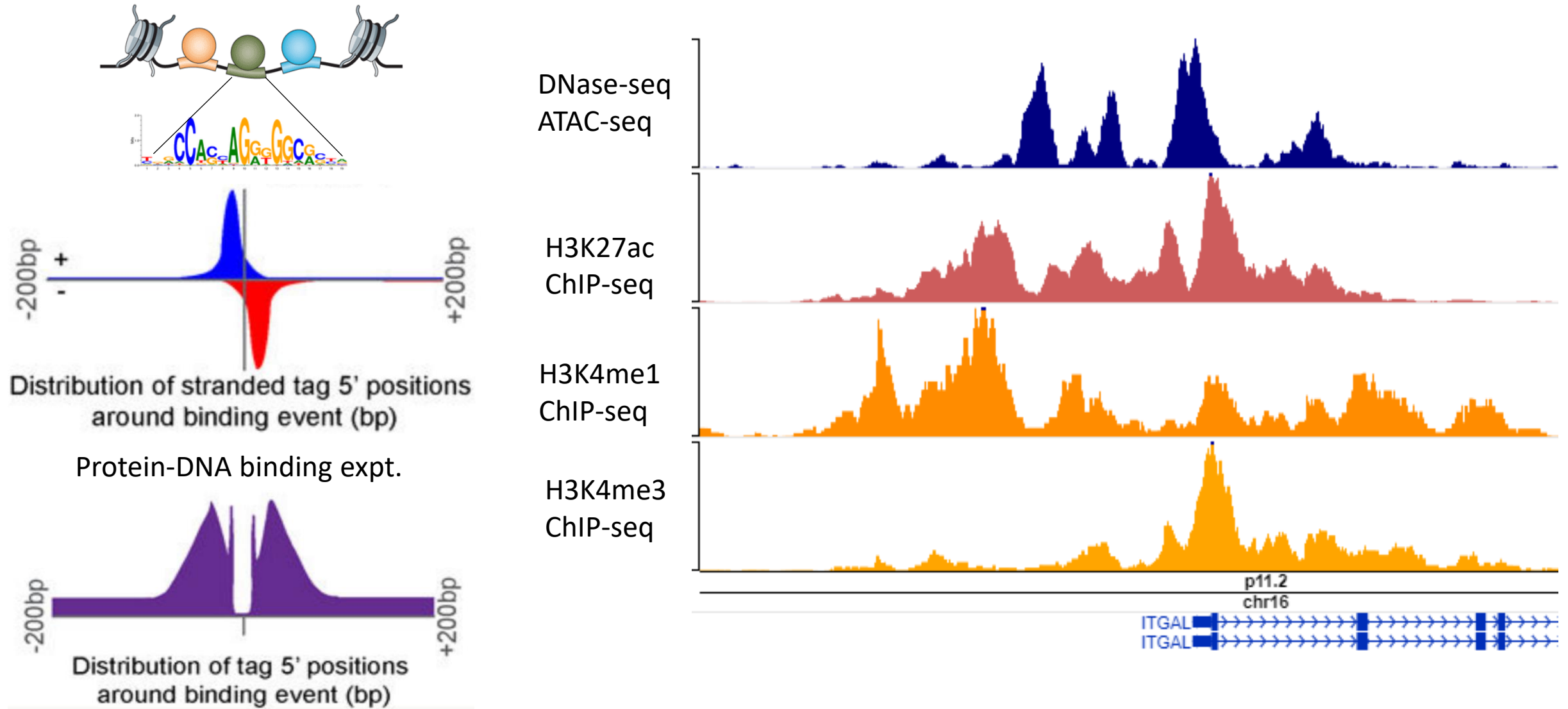
Protein-DNA binding expt.

Distribution of tag 5' positions around binding event (bp)
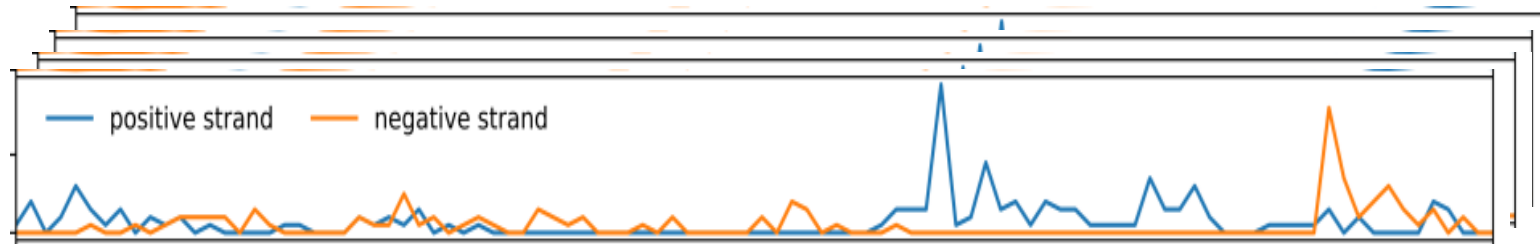
DNA accessibility experiments

# High-resolution 'shapes' and 'spans' of TF and chromatin profiles capture exquisite information about protein-DNA contacts
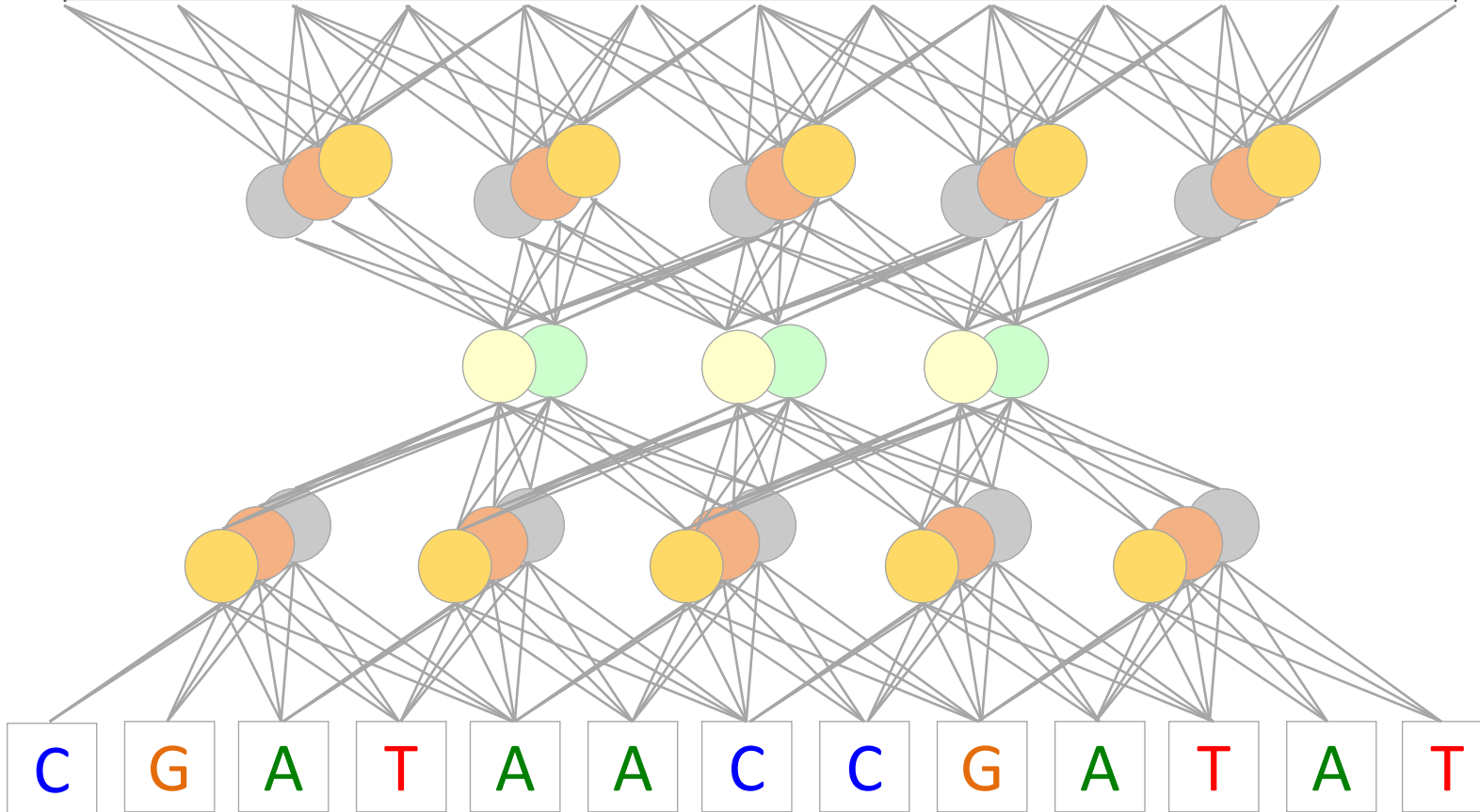


Distribution of stranded tag 5' positions around binding event (bp)

Protein-DNA binding expt.

Distribution of tag 5' positions around binding event (bp)

DNA accessibility experiments

DNase-seq
ATAC-seq

H3K27ac
ChIP-seq

H3K4me1
ChIP-seq

H3K4me3
ChIP-seq

p11.2
chr16
ITGAL
ITGAL

# BPNet : Sequence to base-res. TF binding profiles

**Total reads + base-resolution probability profile (1 kb)**



Ziga Avsec

Sequence windows (2 kb)
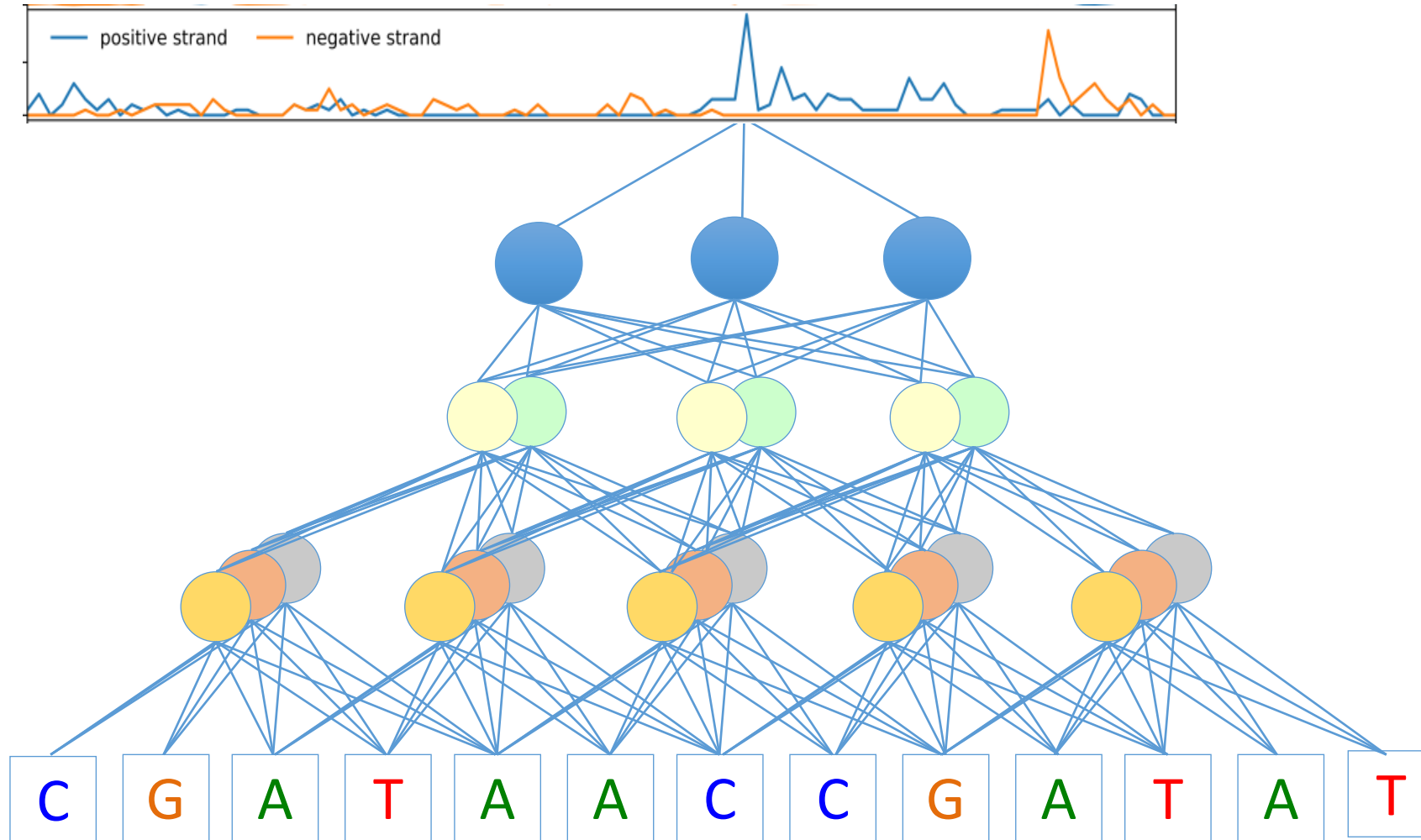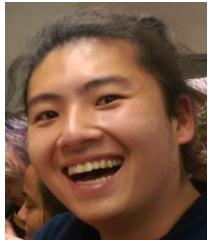
*Avsec et al. 2021 Nature Genetics*

# BPNet : Sequence to base-res. TF binding profiles

**Total reads + base-resolution probability profile (1 kb)**



Ziga Avsec

— positive strand   — negative strand

**Assay bias/control track**

C G A T A A C C G A T A T

Sequence windows (2 kb)

*Avsec et al. 2021 Nature Genetics*

# BPNet predicts base resolution protein-DNA binding profiles with unprecedented accuracy (on par with replicate concordance)

Oct4, Sox2, Nanog and Klf4 in mESCs

*Julia Zeitlinger lab*

# DeepLIFT: Inferring predictive nucleotides in any sequence



Avanti Shrikumar

Alex Tseng

Shrikumar et al. 2017 ICML
Shrikumar et al. 2019 ISMB
Tseng et al. 2020 NeurIPS
Greenside et al. 2018, ECCB

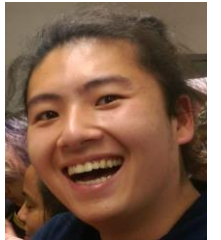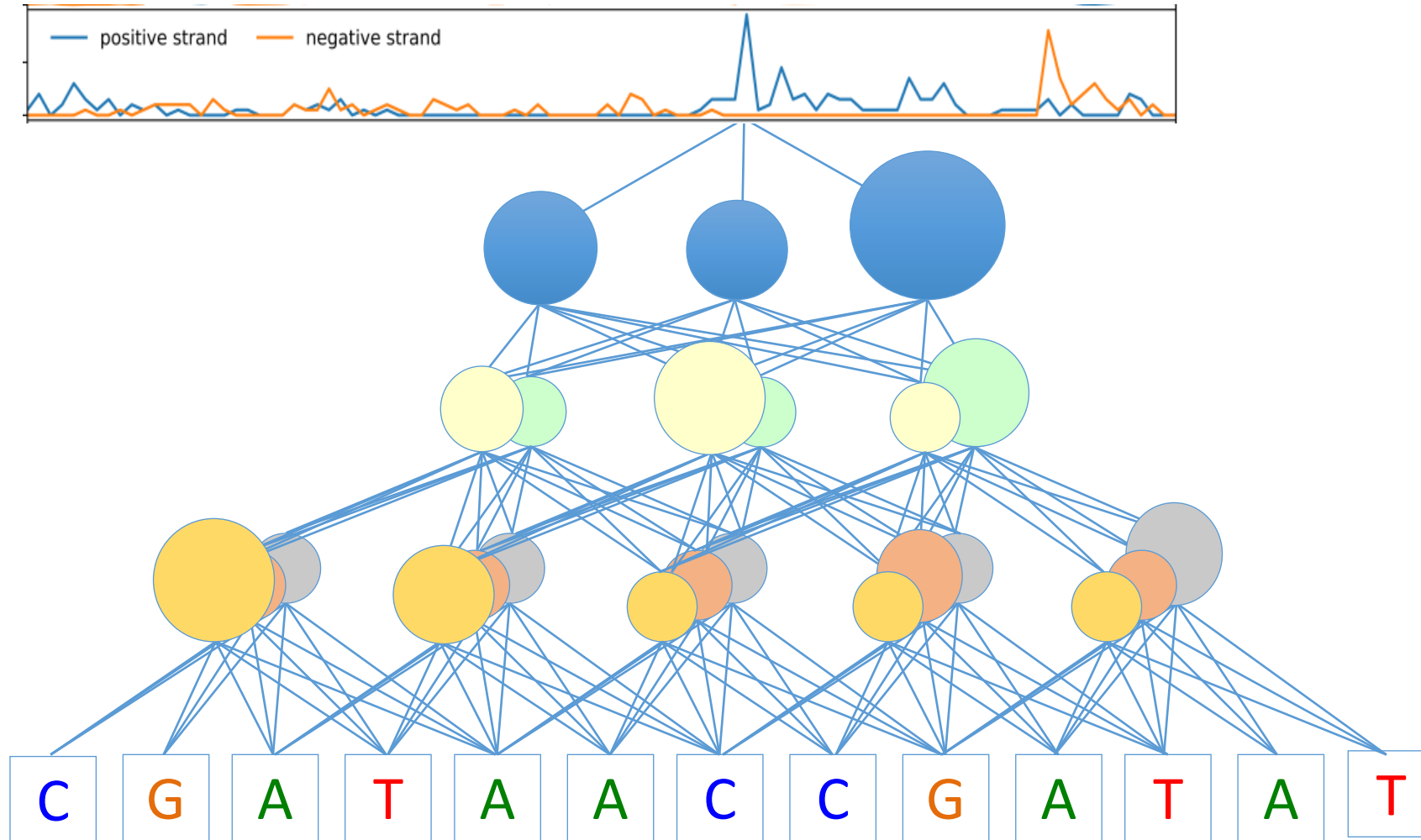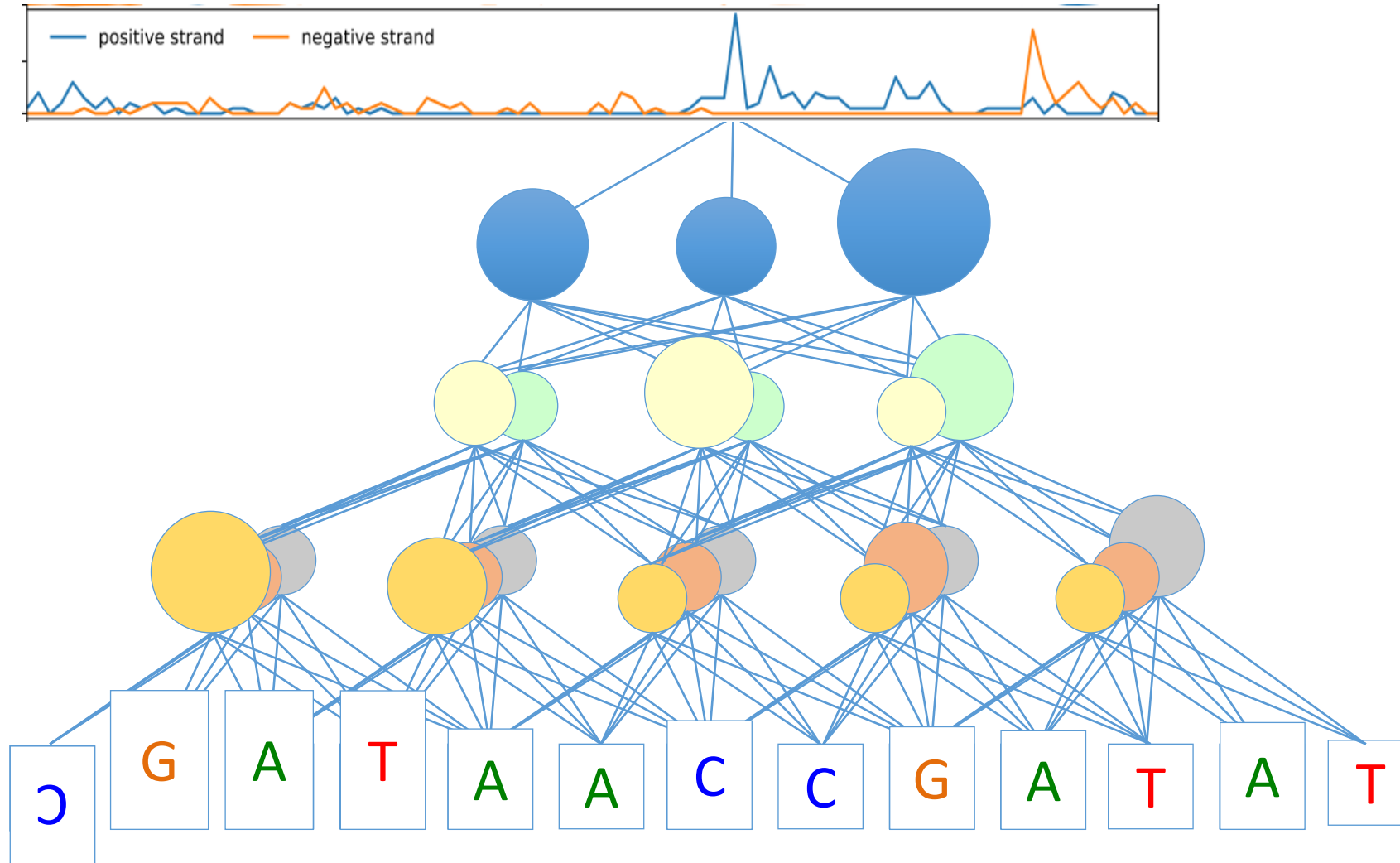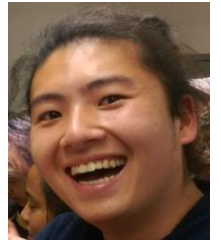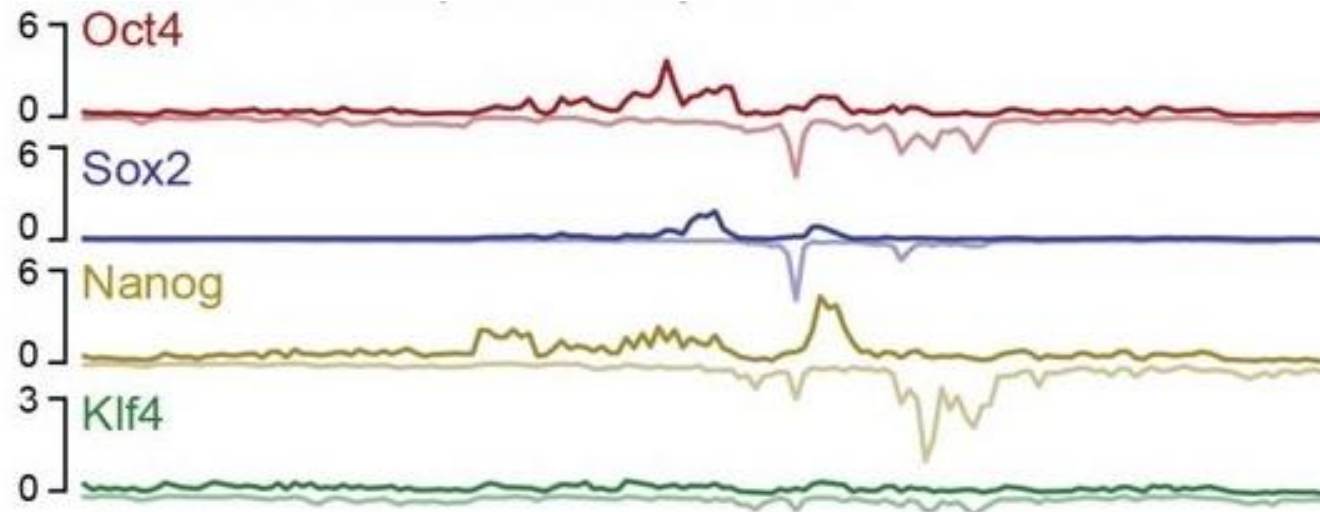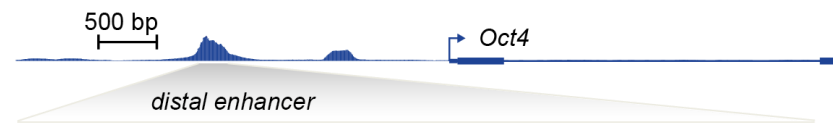# DeepLIFT: Inferring predictive nucleotides in any sequence
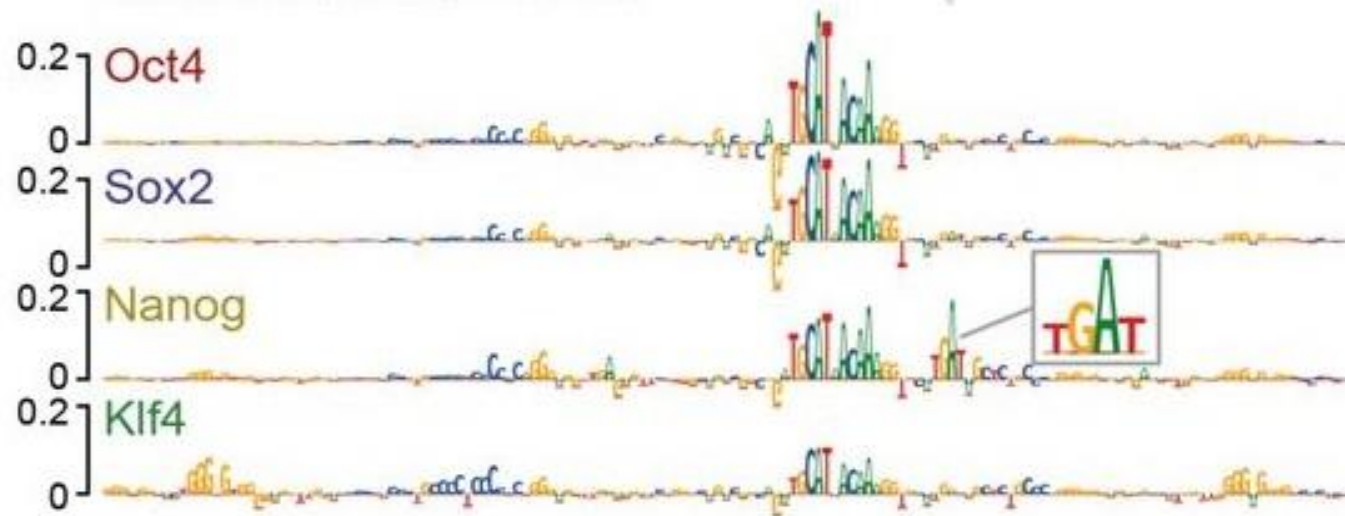


Avanti Shrikumar

Alex Tseng

*Shrikumar et al. 2017 ICML*
*Shrikumar et al. 2019 ISMB*
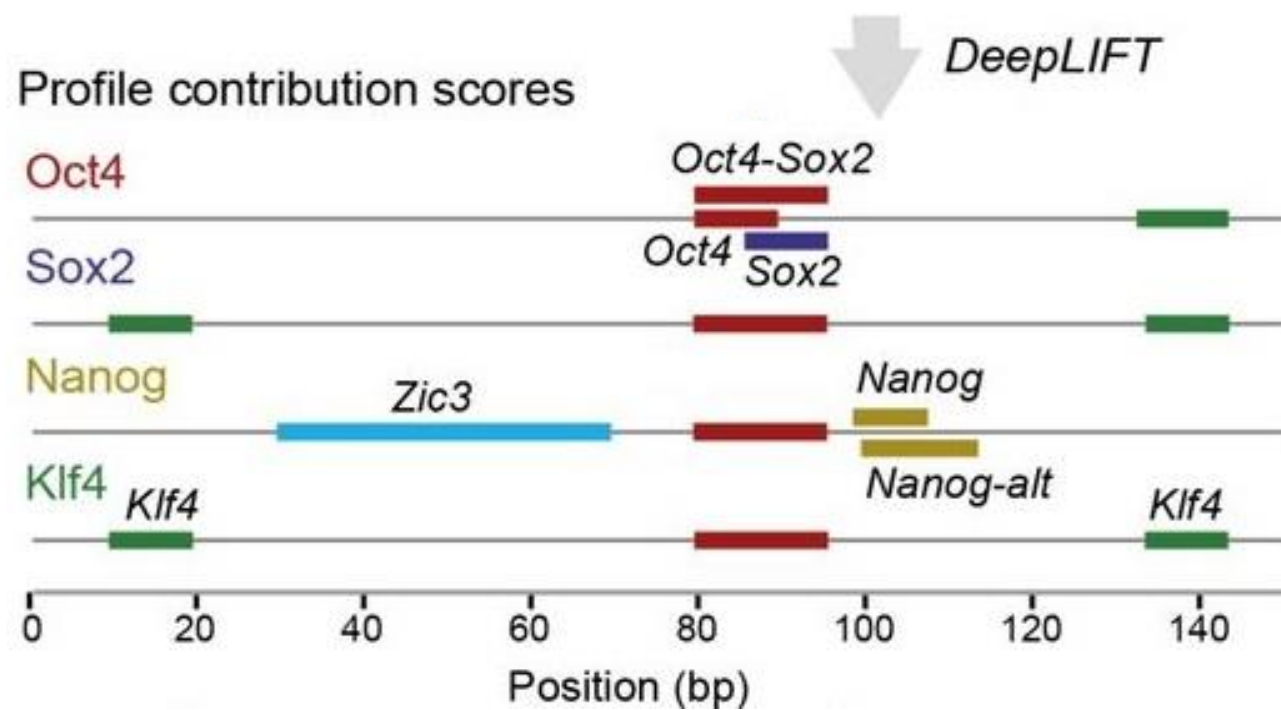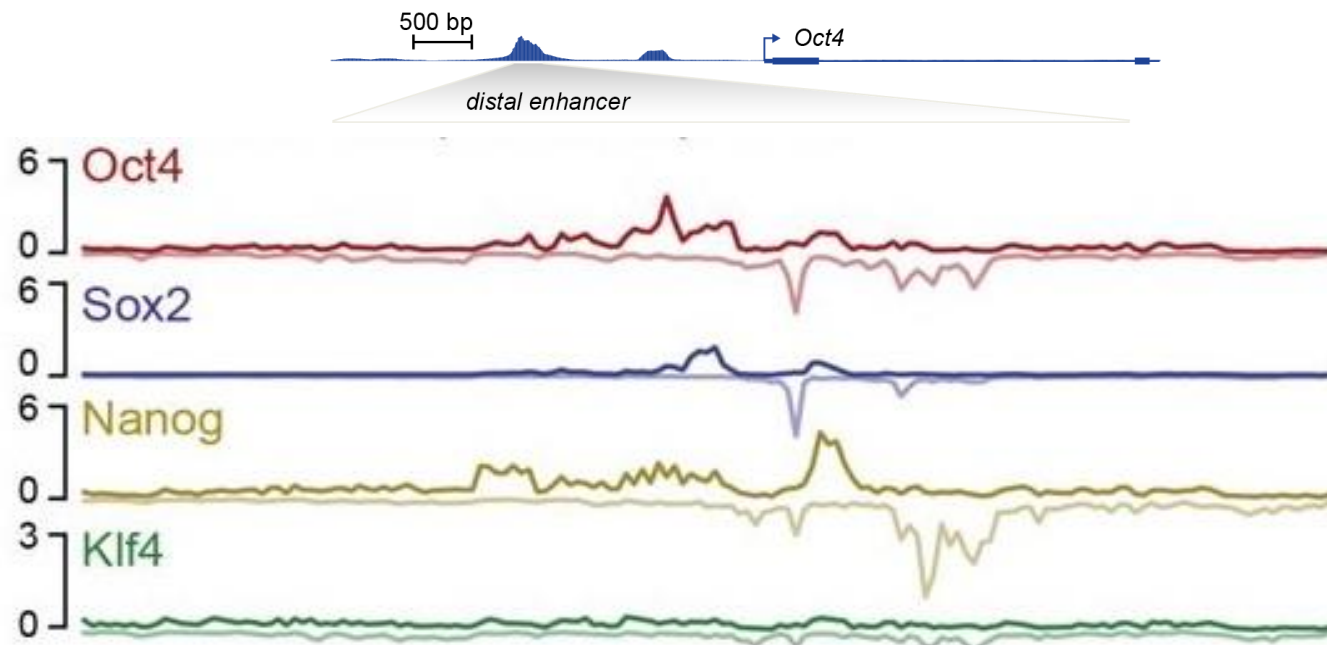*Tseng et al. 2020 NeurIPS*
*Greenside et al. 2018, ECCB*

# DeepLIFT: Inferring predictive nucleotides in any sequence

Avanti Shrikumar

Alex Tseng

*Shrikumar et al. 2017 ICML*
*Shrikumar et al. 2019 ISMB*
*Tseng et al. 2020 NeurIPS*
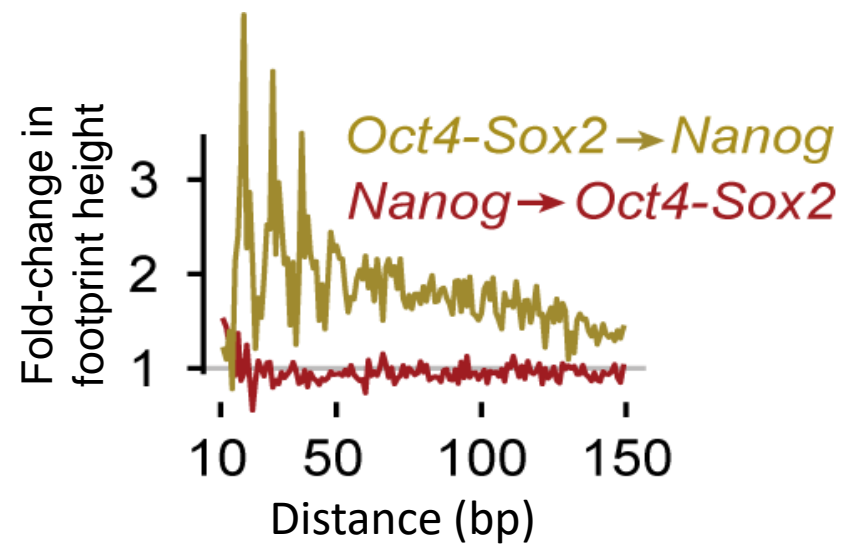*Greenside et al. 2018, ECCB*

# DeepLIFT: Inferring predictive nucleotides in any sequence



Avanti Shrikumar

Alex Tseng

*Shrikumar et al. 2017 ICML*
*Shrikumar et al. 2019 ISMB*
*Tseng et al. 2020 NeurIPS*
*Greenside et al. 2018, ECCB*

# DeepLIFT: Inferring predictive nucleotides in any sequence



Avanti Shrikumar

Alex Tseng

*Shrikumar et al. 2017 ICML*
*Shrikumar et al. 2019 ISMB*
*Tseng et al. 2020 NeurIPS*
*Greenside et al. 2018, ECCB*

mESCs

distal enhancer

500 bp

Oct4

DeepLIFT

Profile contribution scores

mESCs

Oct4
Sox2
Nanog
Klf4

distal enhancer

DeepLIFT

Profile contribution scores

Oct4

Oct4-Sox2

Sox2

Oct4 Sox2

Nanog

Zic3

Nanog

Nanog-alt

Klf4

Klf4

Klf4

Position (bp)

# Deciphering syntax dependent TF cooperativity with synthetic designed sequences

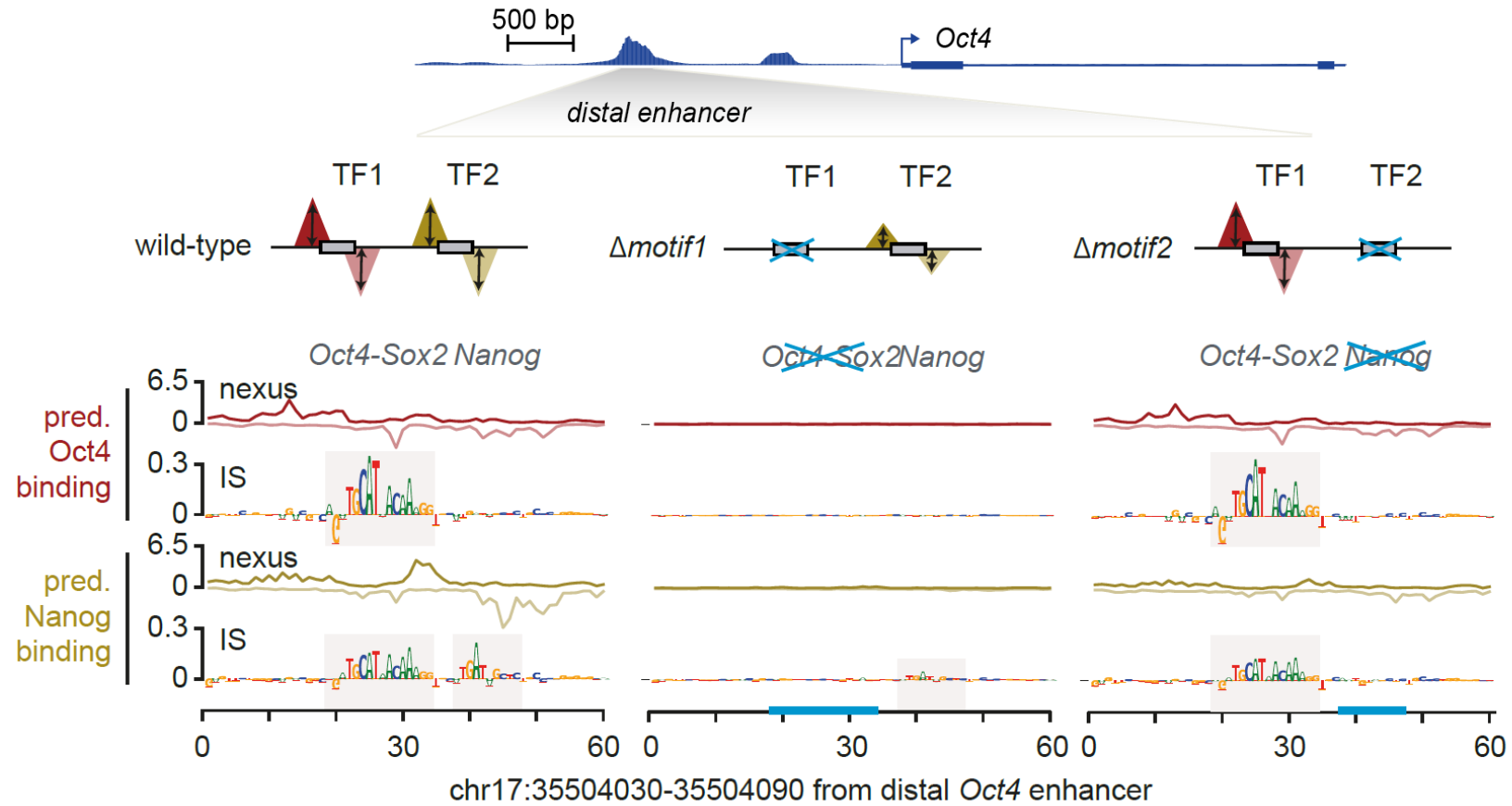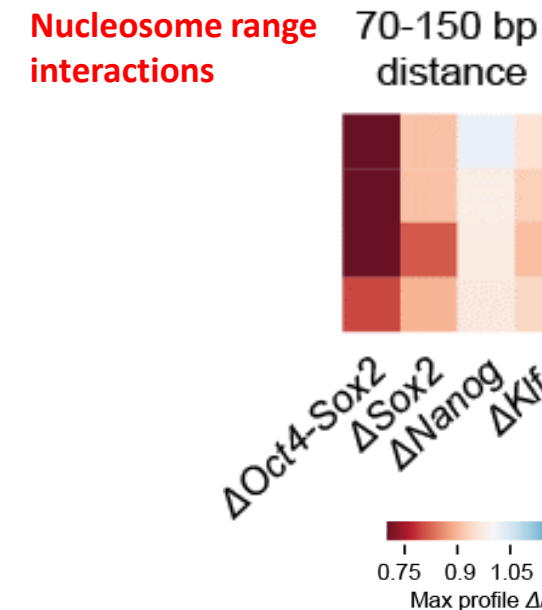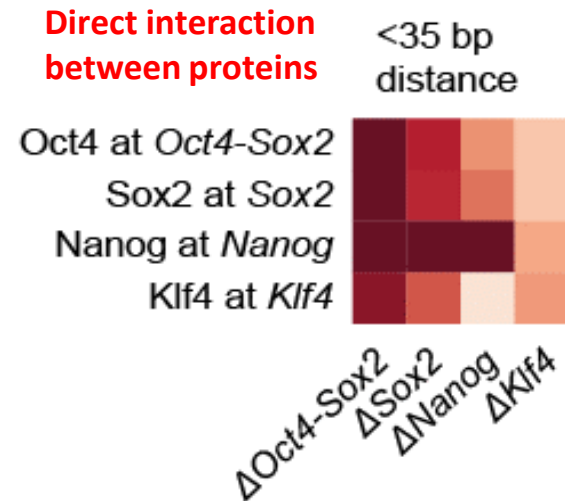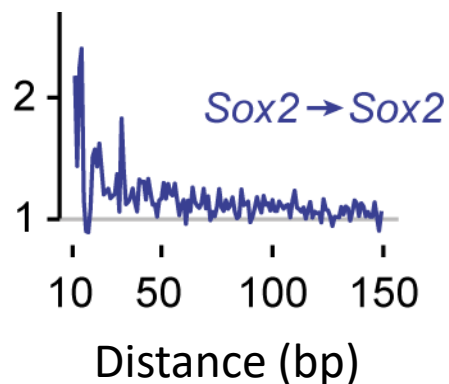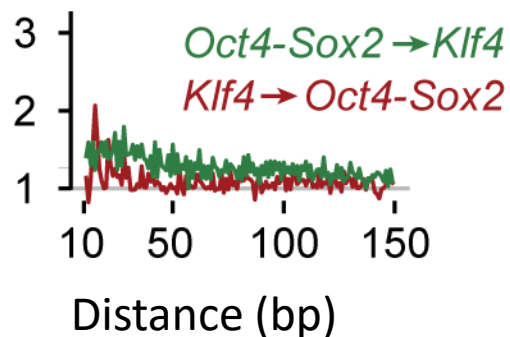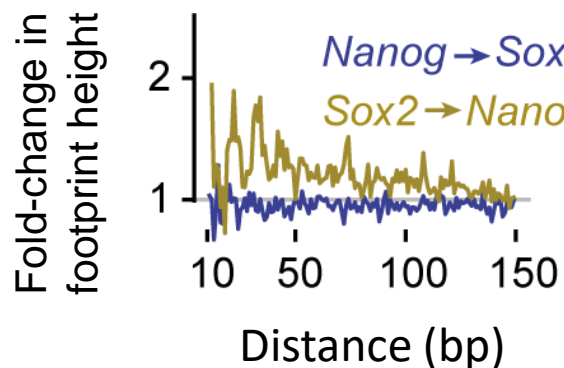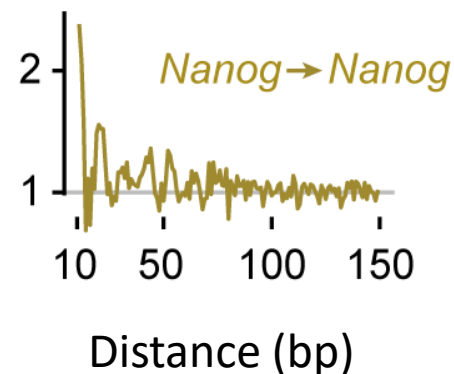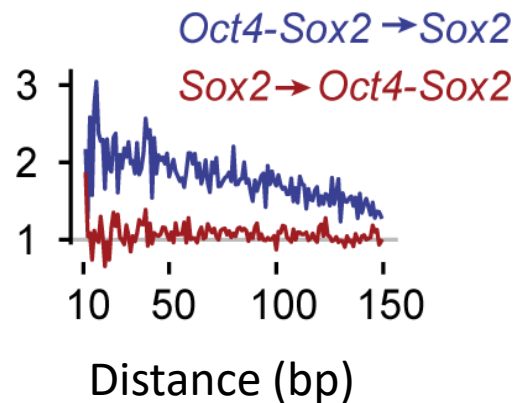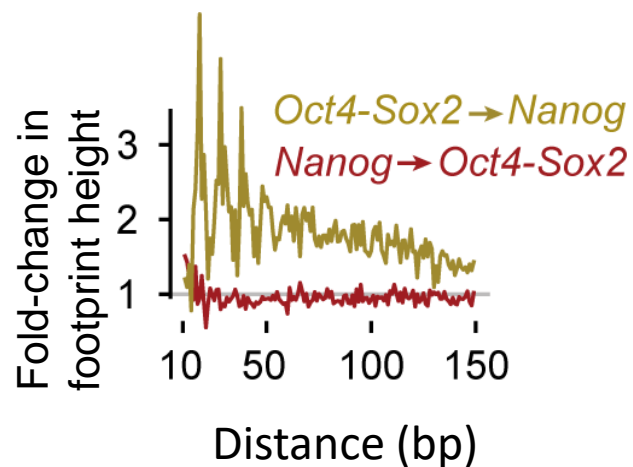# Deciphering syntax dependent TF cooperativity with *in-silico* genome editing



500 bp

Oct4

distal enhancer

TF1    TF2

wild-type

*Oct4-Sox2 Nanog*

pred. Oct4 binding

6.5
0
nexus

0.3
0
IS

pred. Nanog binding

6.5
0
nexus

0.3
0
IS

0    30    60

chr17:35504030-35504090 from distal *Oct4* enhancer

chr17:35504030-35504090 from distal *Oct4* enhancer

# Deciphering syntax dependent TF cooperativity with *in-silico* genome editing



chr17:35504030-35504090 from distal *Oct4* enhancer

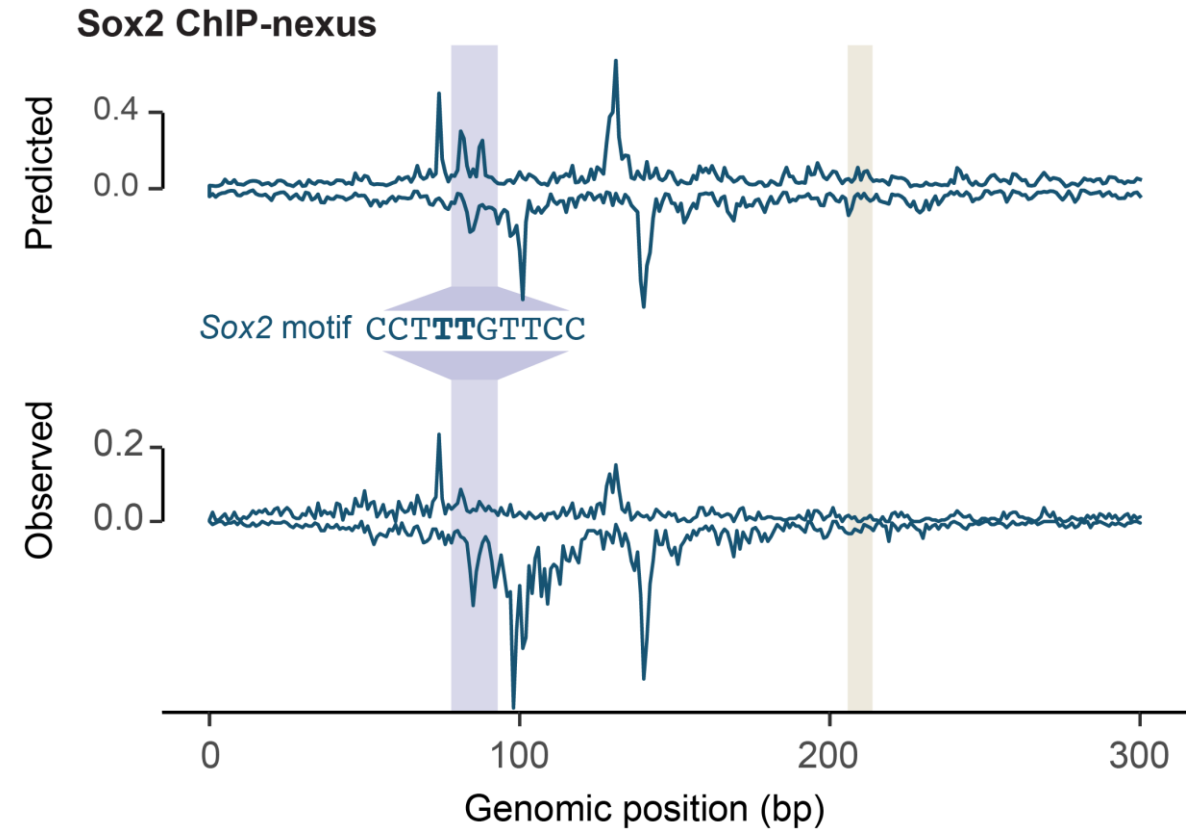Deciphering syntax dependent TF cooperativity with *in-silico* genome editing

# Distance dependent motif syntax rules of asymmetric directional cooperativity
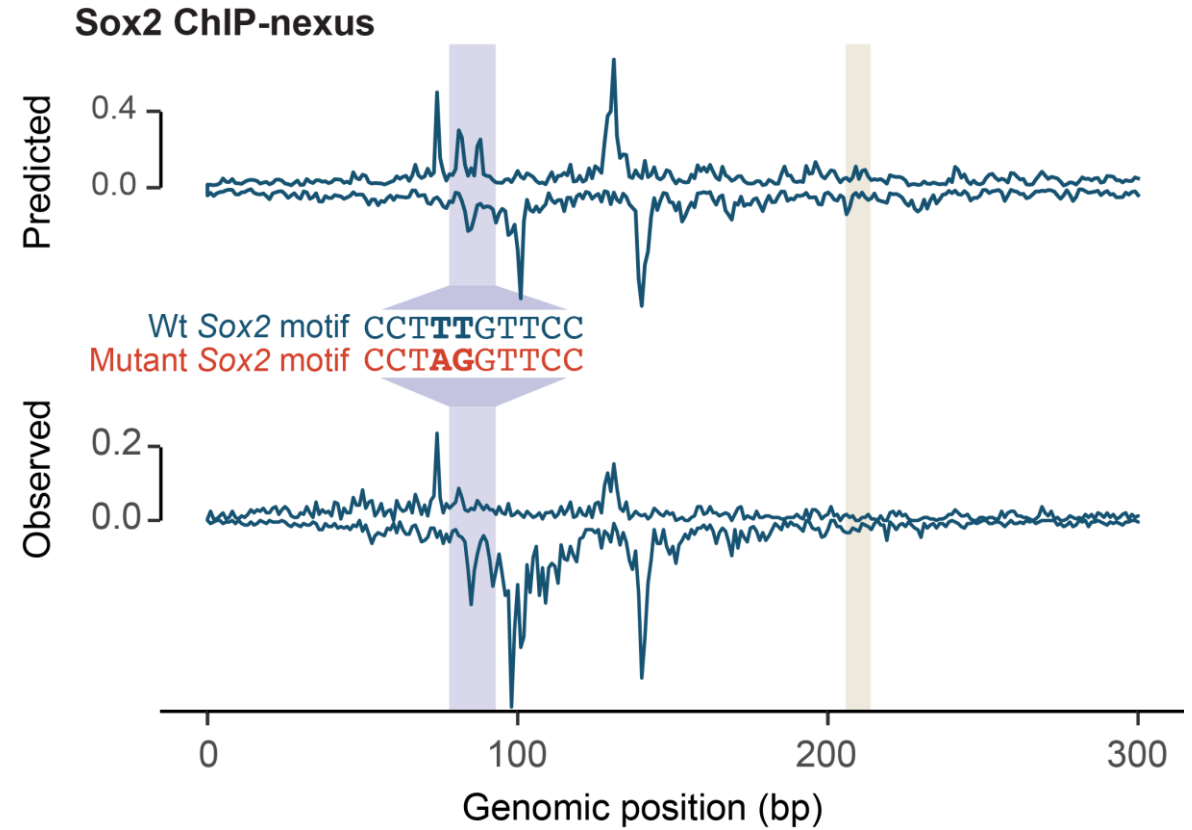


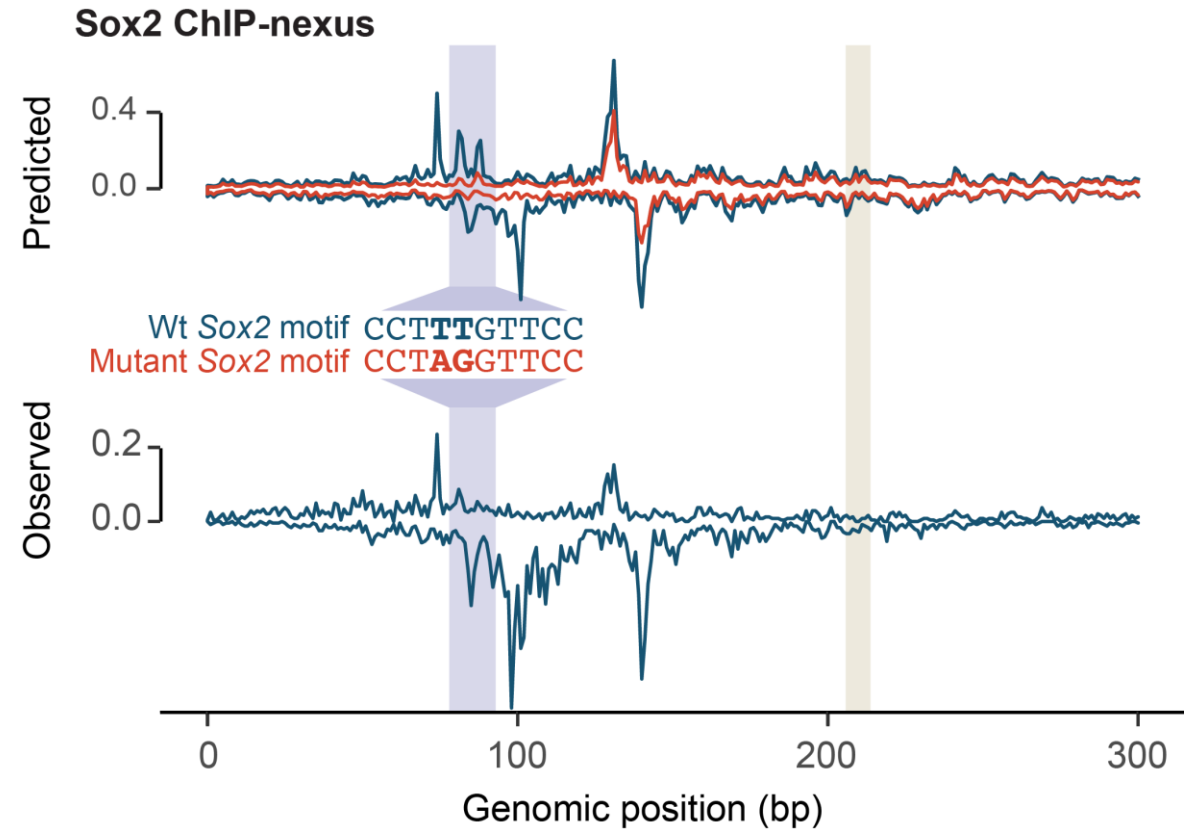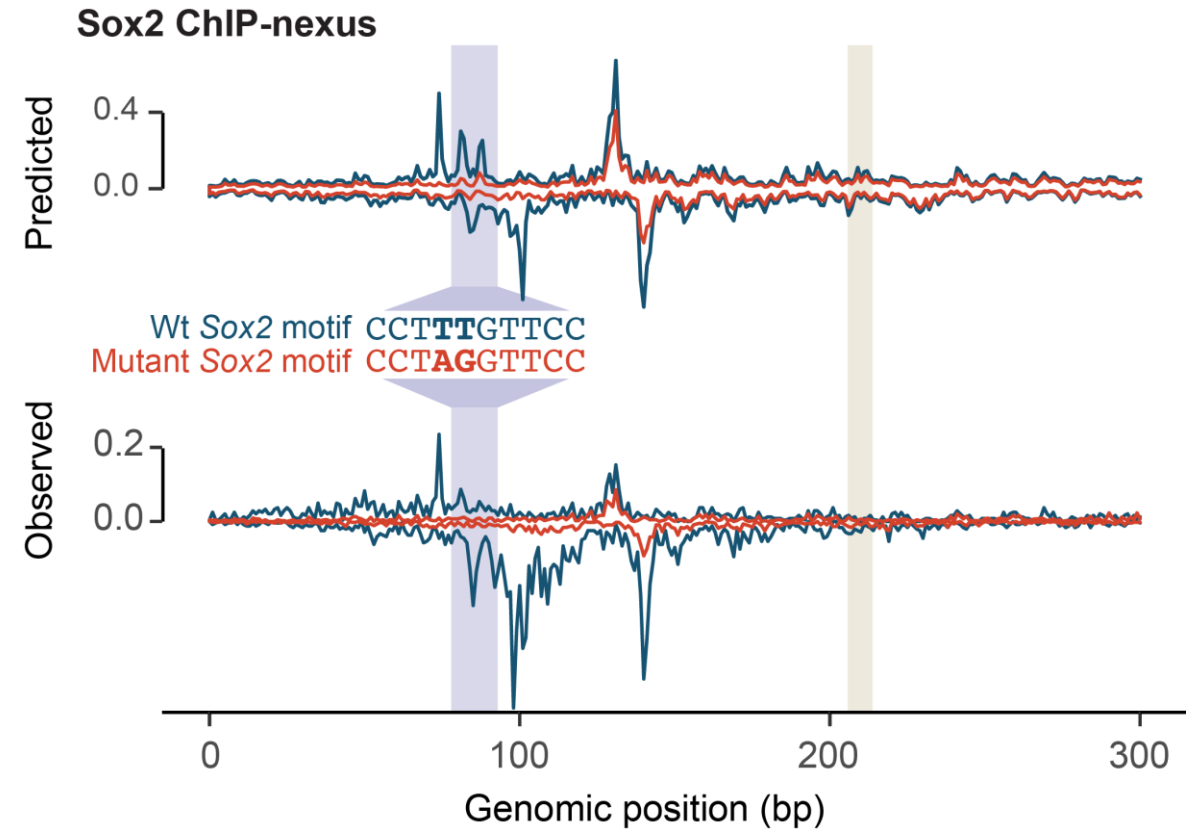*Validated with CRISPR/Cas9 syntax editing experiments*

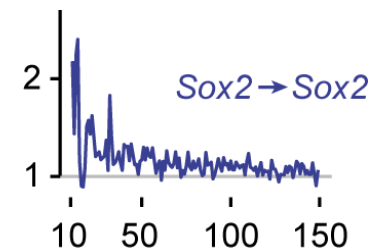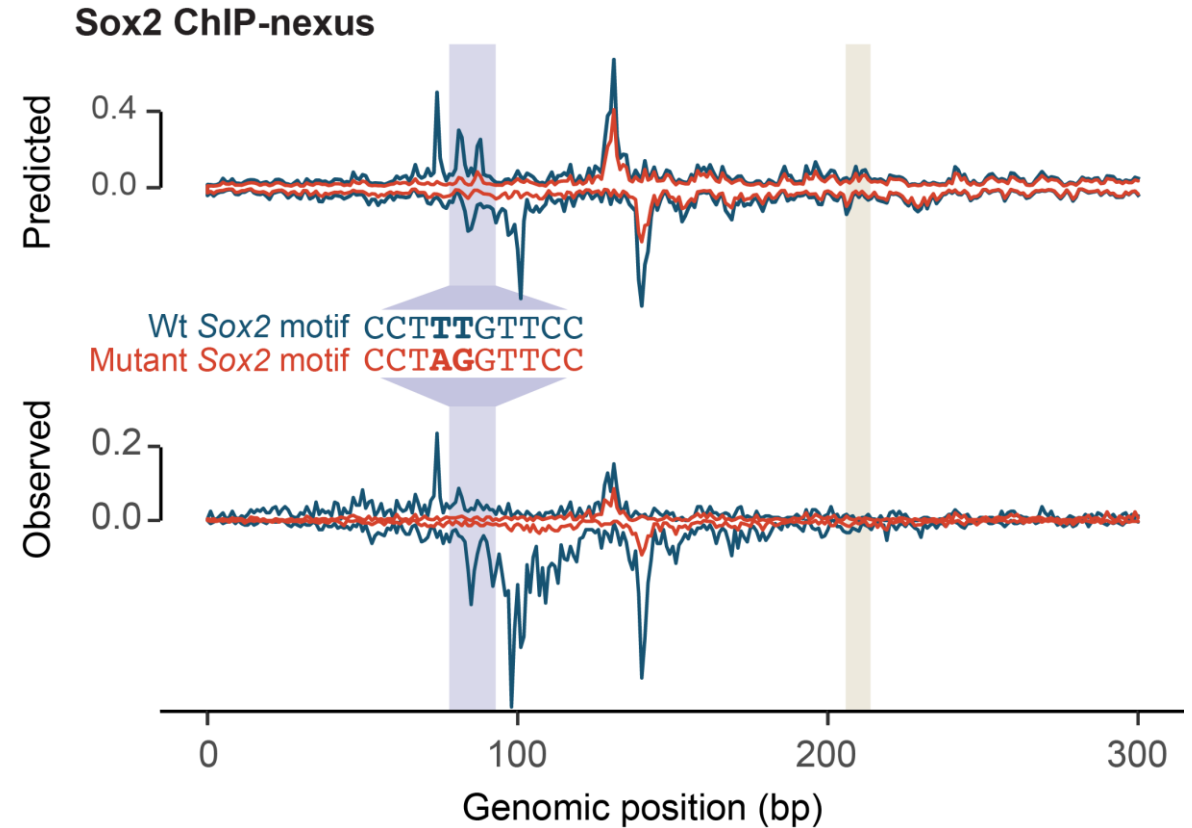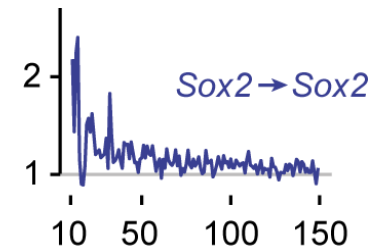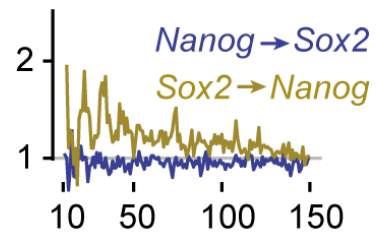# CRISPR mutations validate motif syntax Nanog <> Sox2



Julia Zeitlinger, Sabrina Krueger, Melanie Weilert

# CRISPR mutations validate motif syntax Nanog <> Sox2



Julia Zeitlinger, Sabrina Krueger, Melanie Weilert

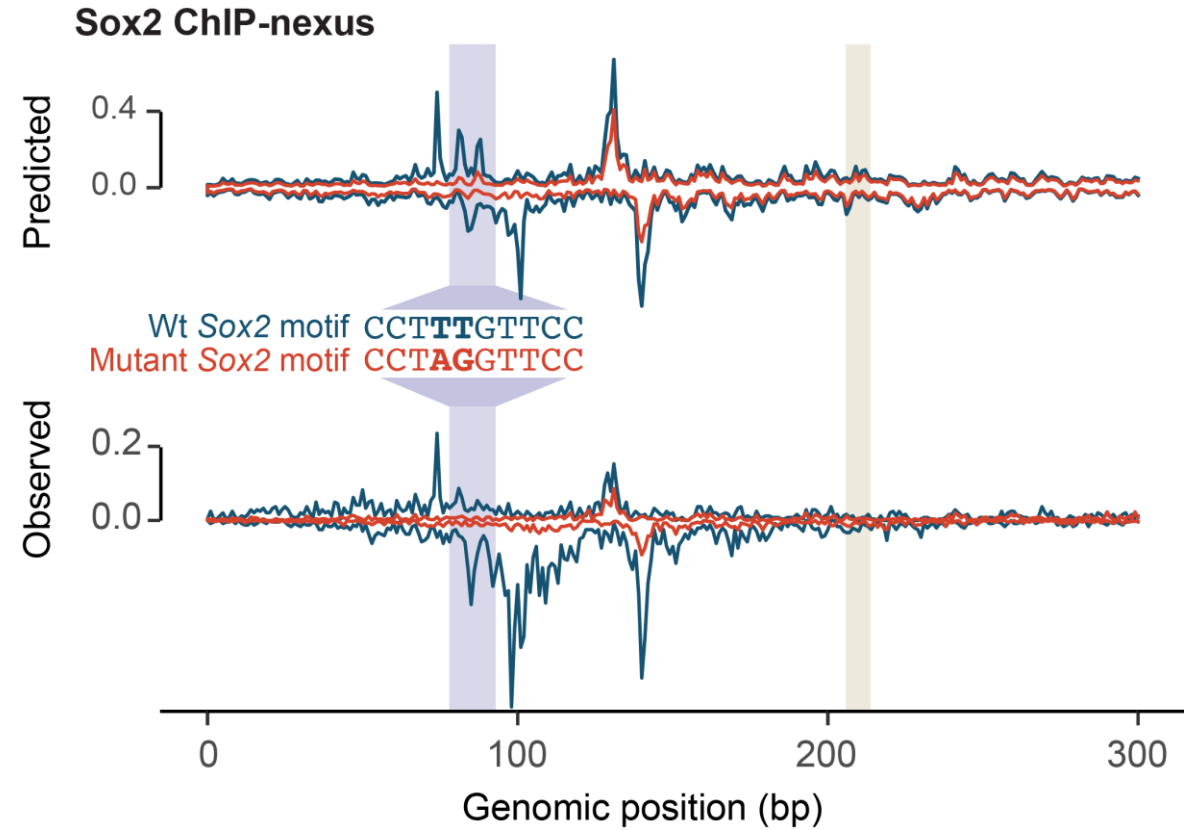**Sox2 ChIP-nexus**

Wt *Sox2* motif CCT**TT**GTTCC
Mutant *Sox2* motif CCT**AG**GTTCC

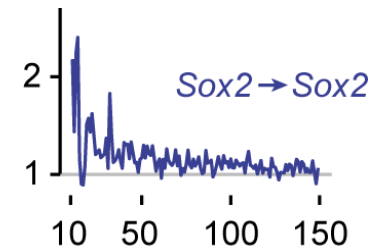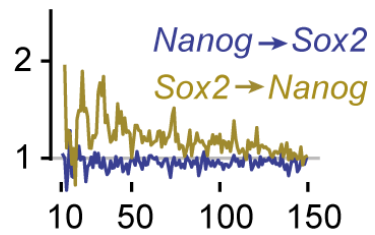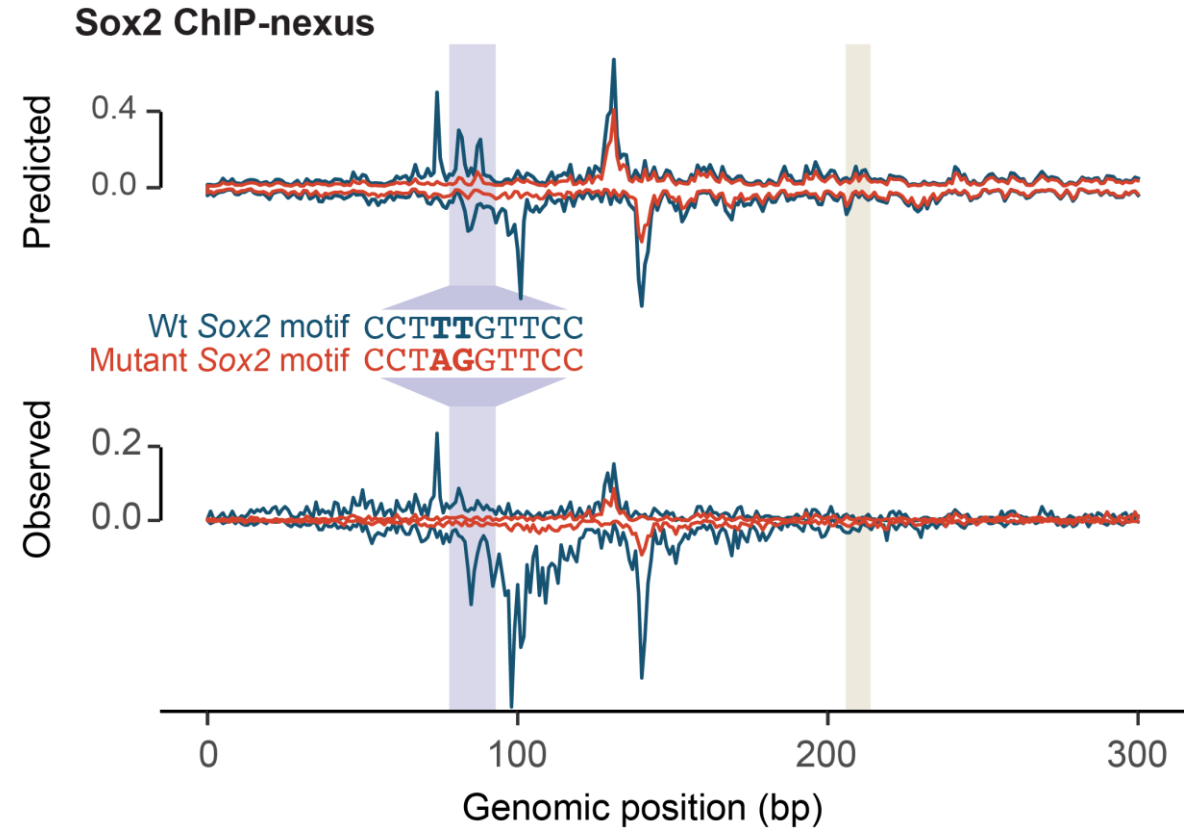Julia Zeitlinger, Sabrina Krueger, Melanie Weilert

Julia Zeitlinger, Sabrina Krueger, Melanie Weilert

# CRISPR mutations validate motif syntax Nanog <> Sox2



Sox2 ChIP-nexus

Predicted

Wt *Sox2* motif CCT**TT**GTTCC
Mutant *Sox2* motif CCT**AG**GTTCC

Observed

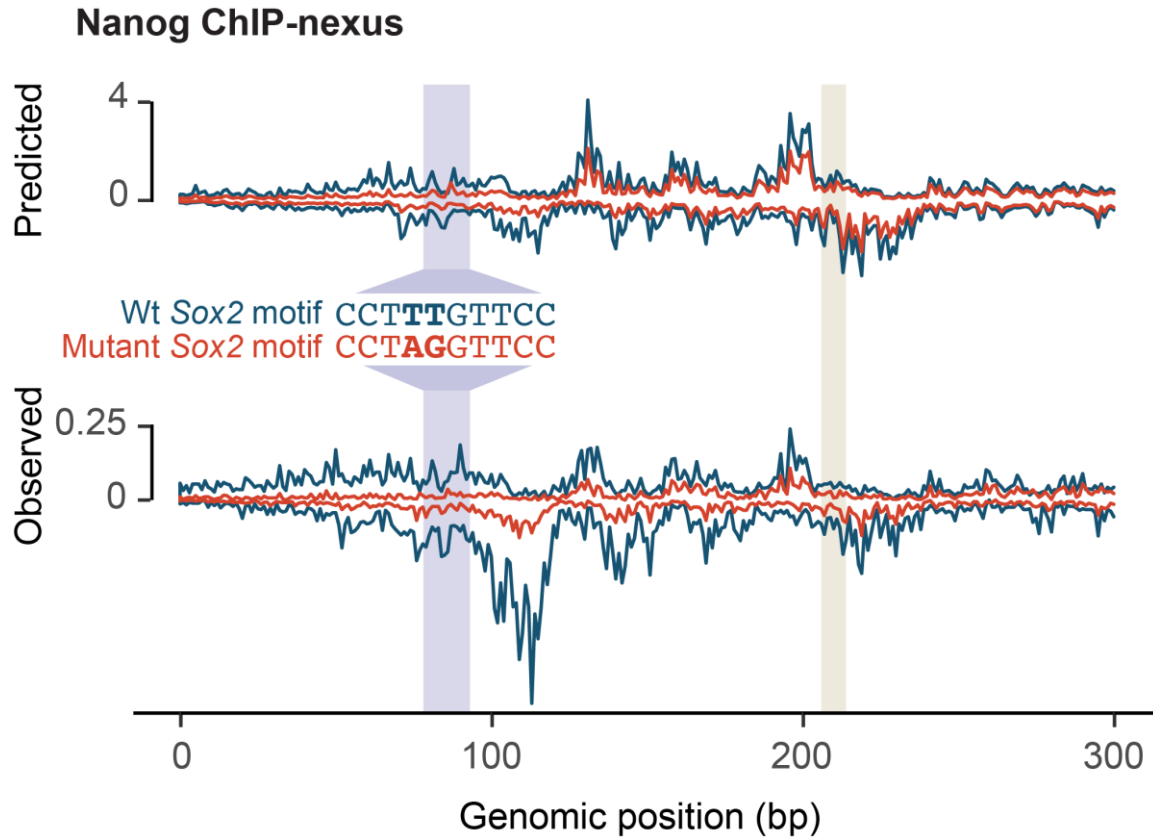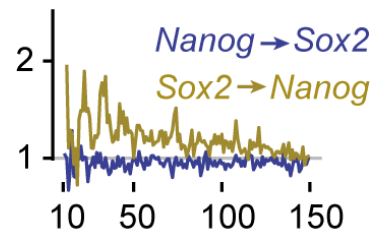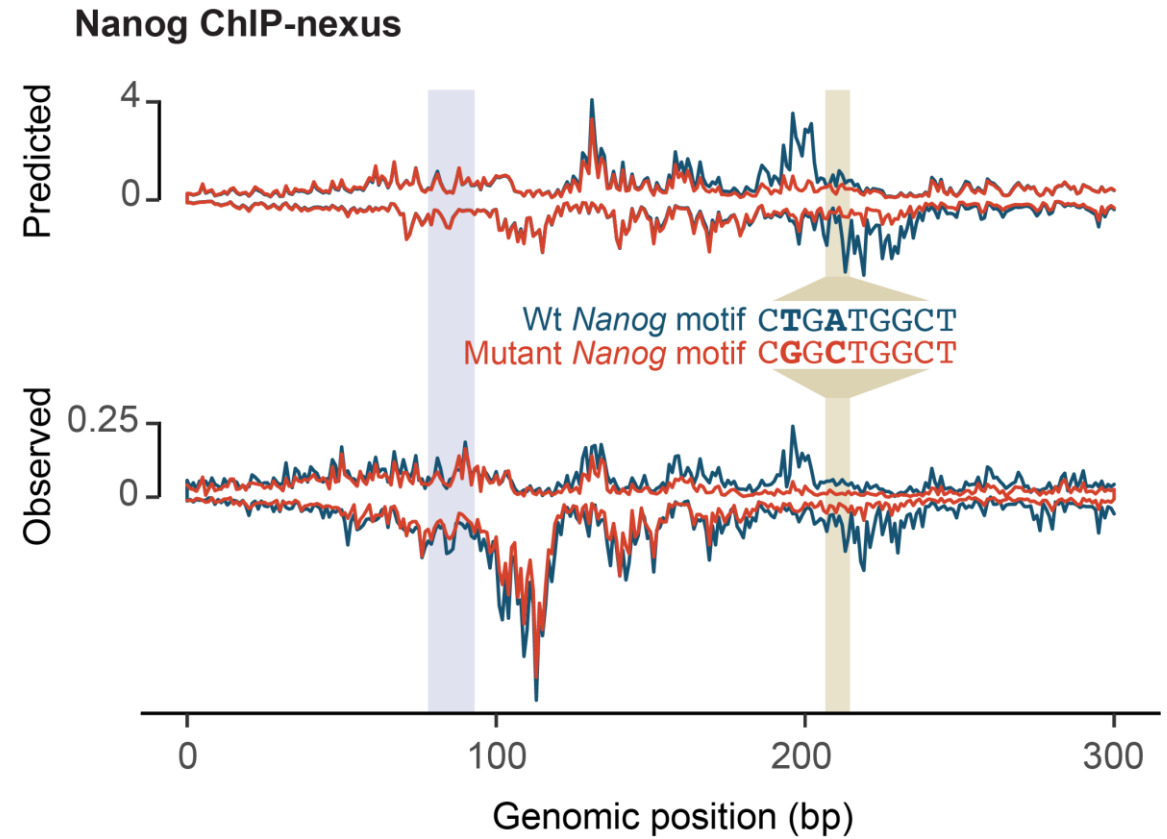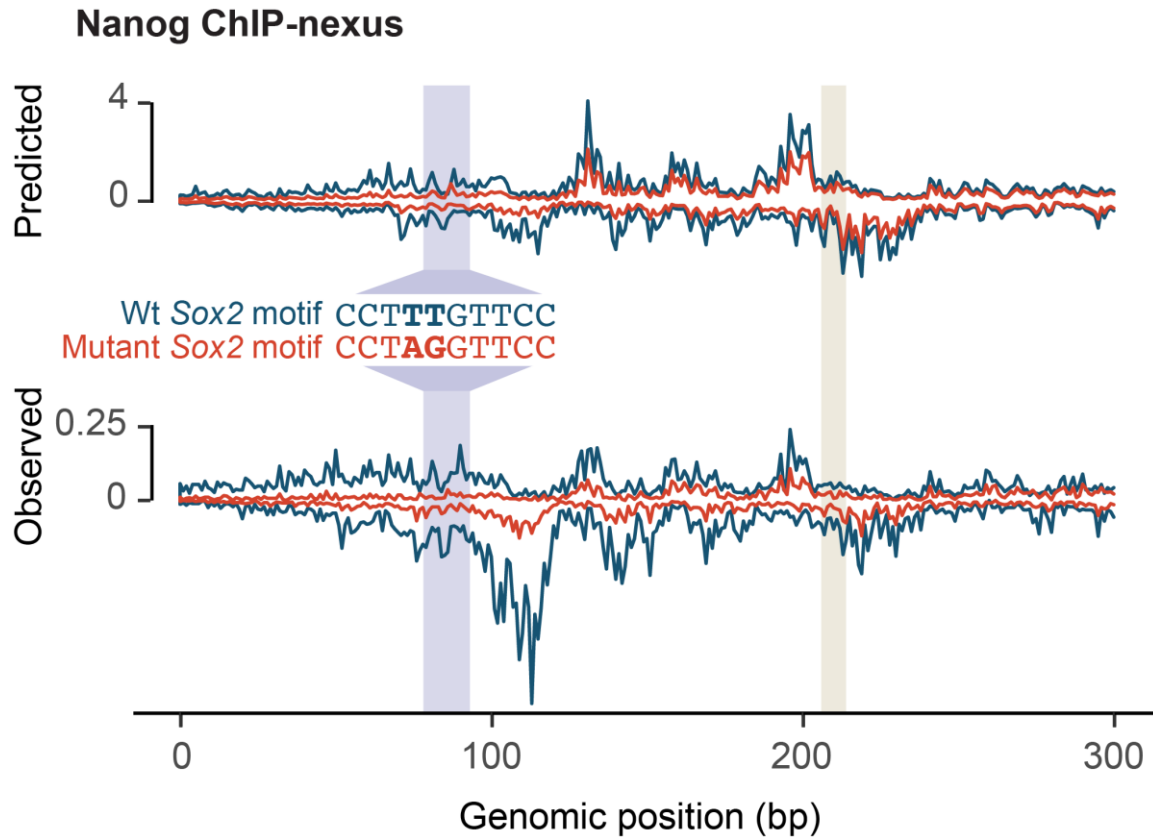Genomic position (bp)

*Sox2 → Sox2*

Julia Zeitlinger, Sabrina Krueger, Melanie Weilert

# CRISPR mutations validate motif syntax Nanog <> Sox2



Sox2 ChIP-nexus

Wt *Sox2* motif CCT**TT**GTTCC
Mutant *Sox2* motif CCT**AG**GTTCC

Nanog→Sox2
Sox2→Nanog

Sox2→Sox2

Genomic position (bp)

Julia Zeitlinger, Sabrina Krueger, Melanie Weilert
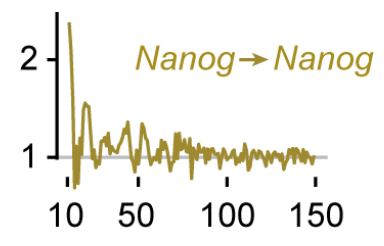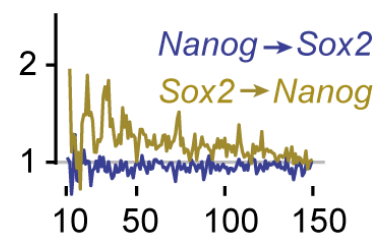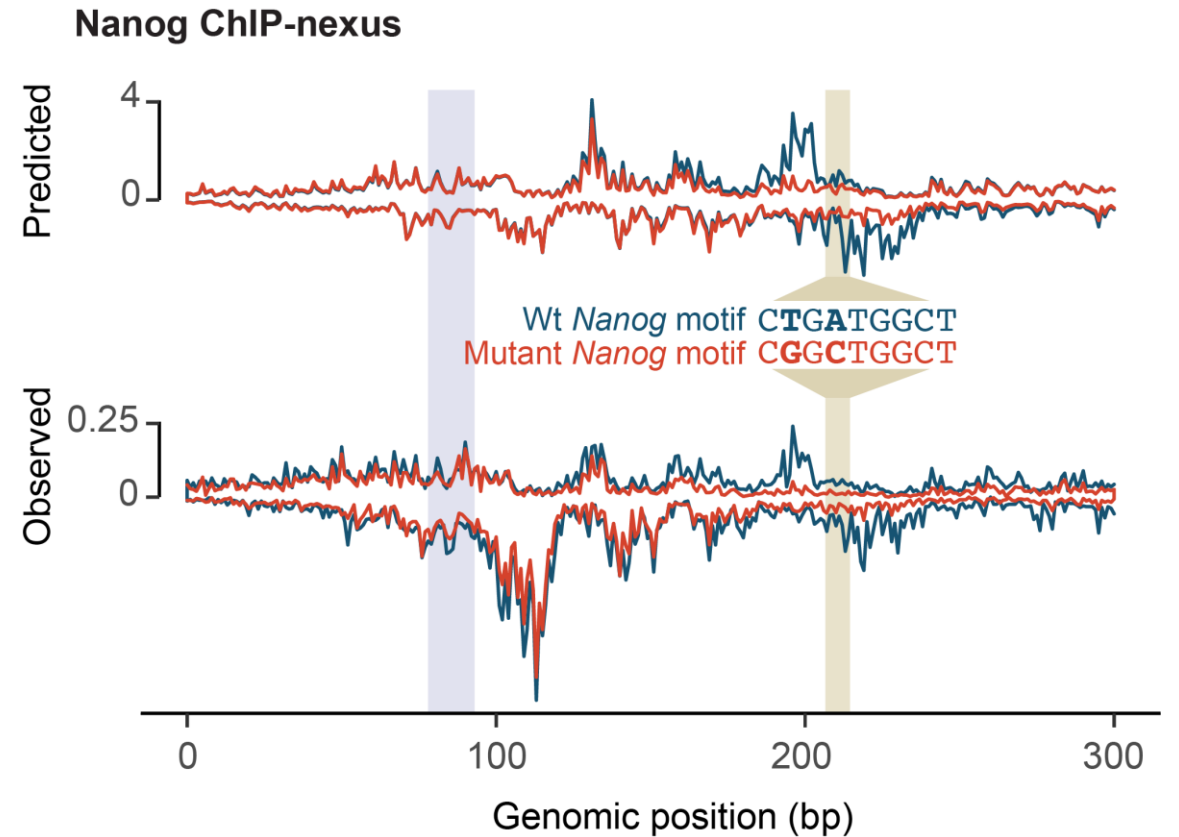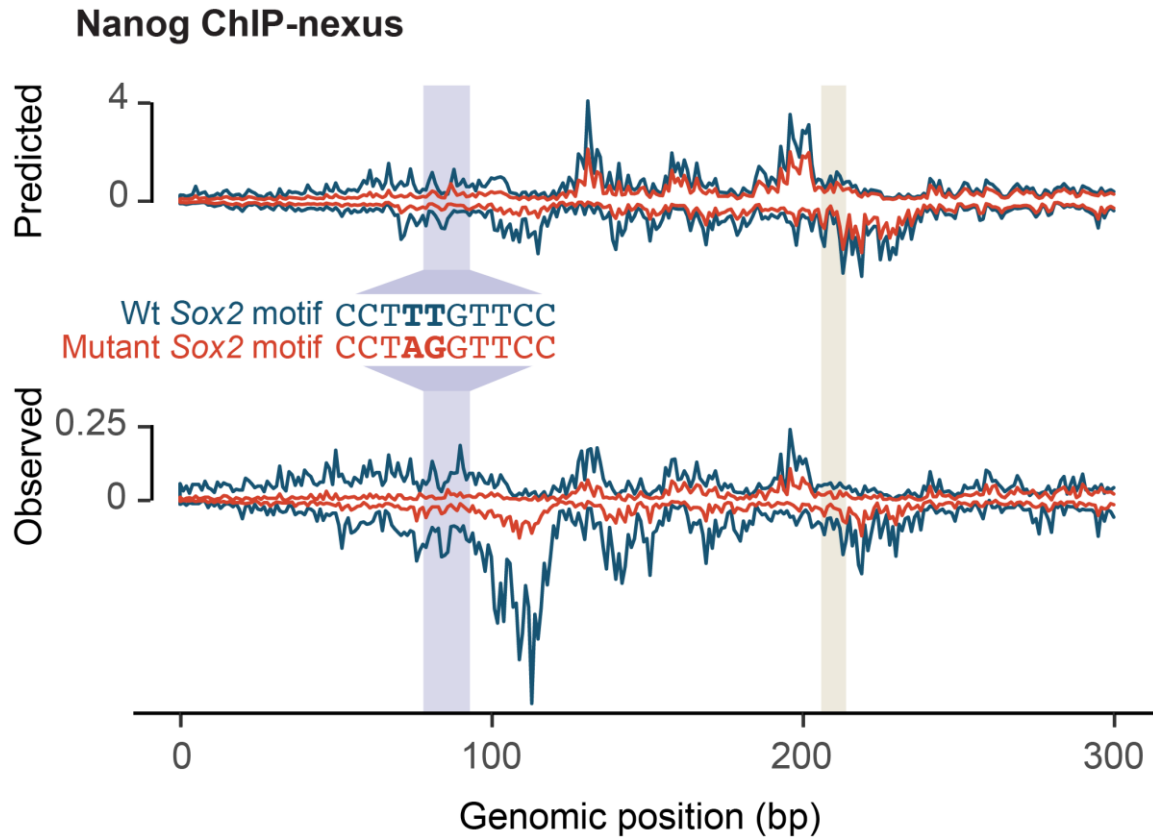
# CRISPR mutations validate motif syntax Nanog <> Sox2



Julia Zeitlinger, Sabrina Krueger, Melanie Weilert

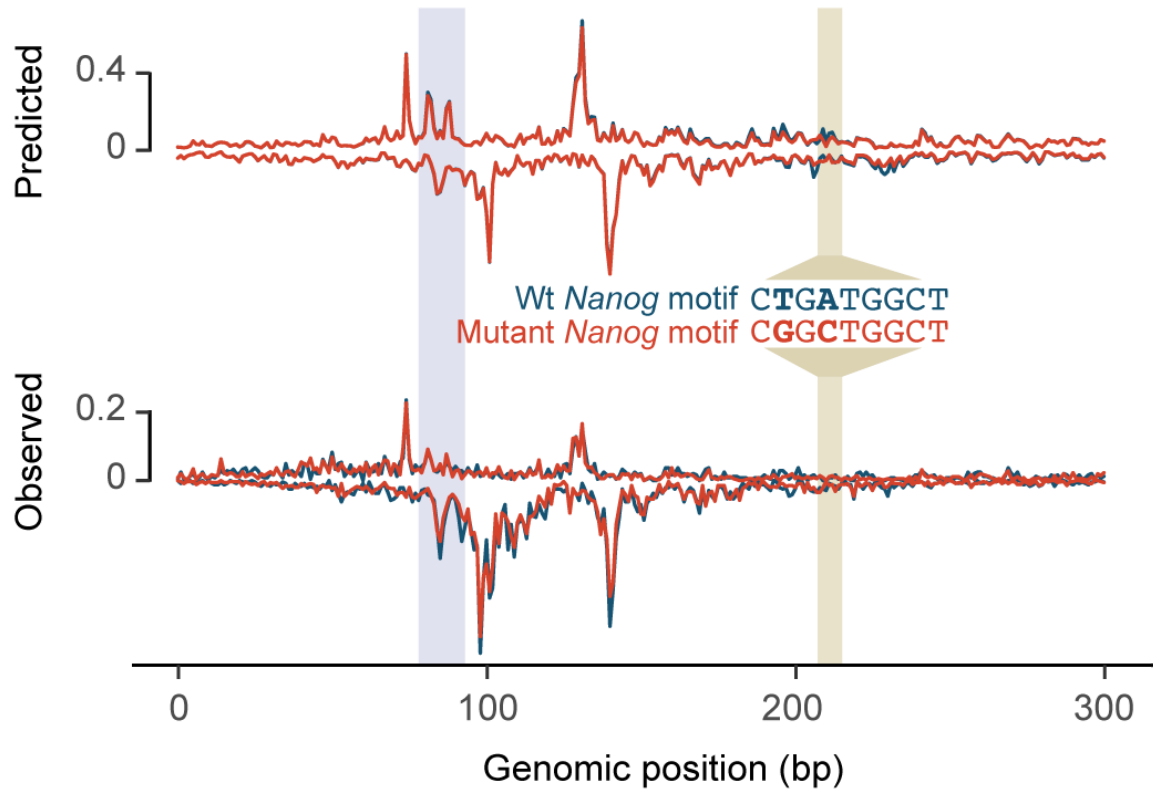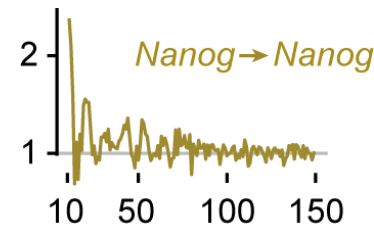# CRISPR mutations validate motif syntax Nanog <> Sox2



Julia Zeitlinger, Sabrina Krueger, Melanie Weilert

# CRISPR mutations validate motif syntax Nanog <> Sox2



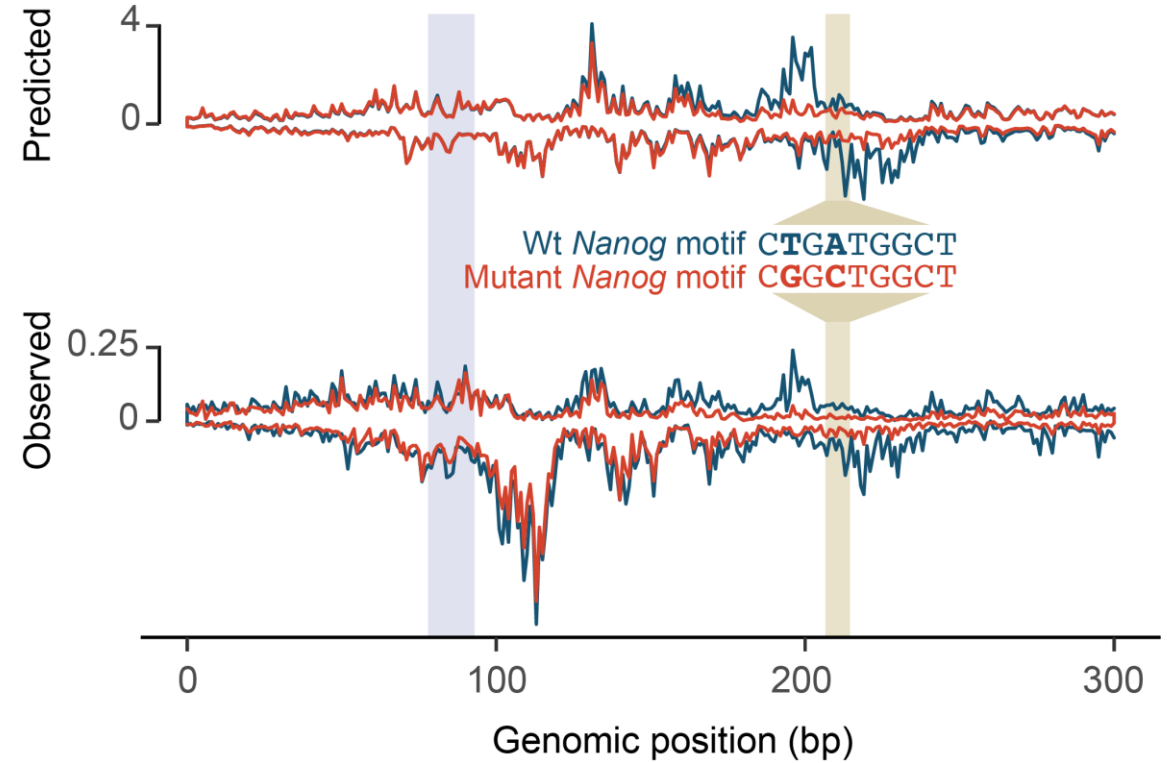Julia Zeitlinger, Sabrina Krueger, Melanie Weilert

Sox2 ChIP-nexus

Nanog ChIP-nexus

Wt *Nanog* motif **CTGA**TGGCT
Mutant *Nanog* motif C**GGC**TGGCT

Nanog → Sox2
Sox2 → Nanog

Nanog → Nanog

Julia Zeitlinger, Sabrina Krueger, Melanie Weilert
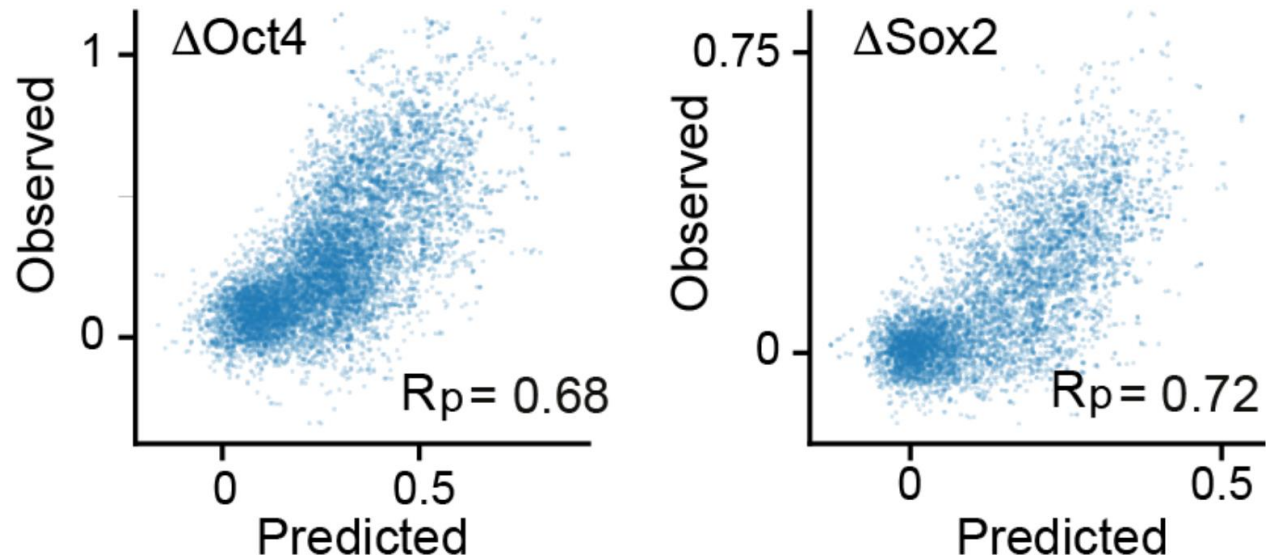
# Binding syntax is predictive of differential accessibility after TF depletion & reporter expression
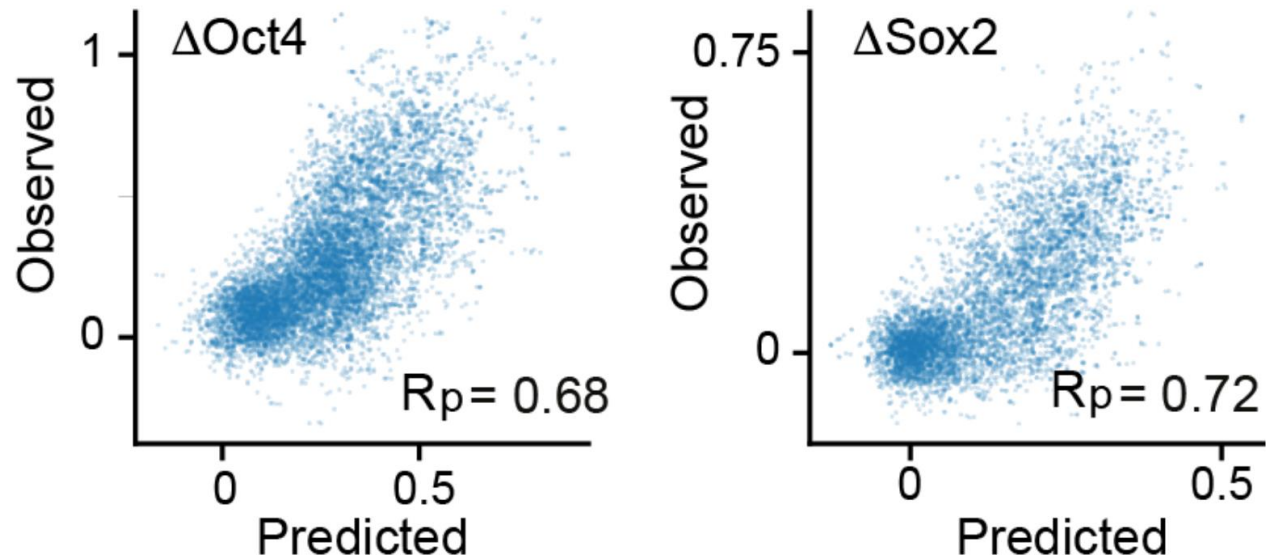


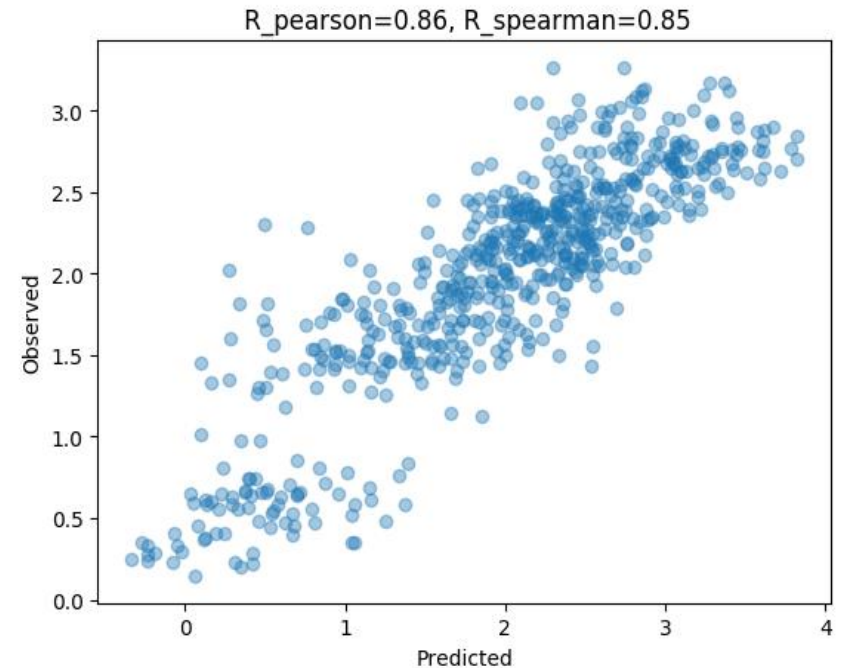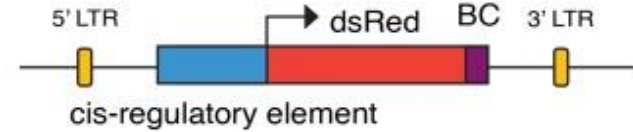ATAC-seq log fold-change loss after TF depletion

*(Independent previously published data from Friman et al. 2019)*

# Binding syntax is predictive of differential accessibility after TF depletion & reporter expression



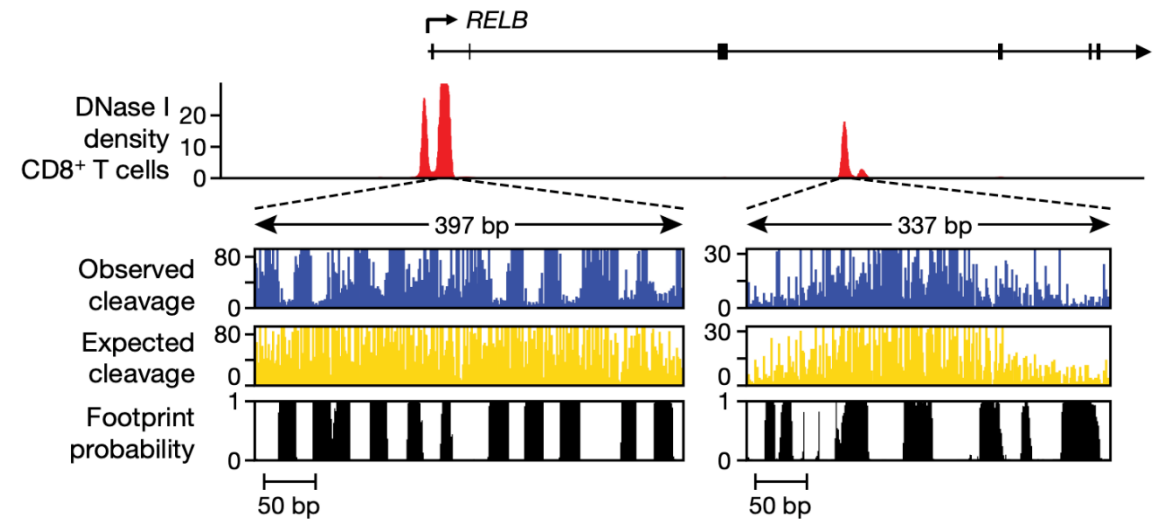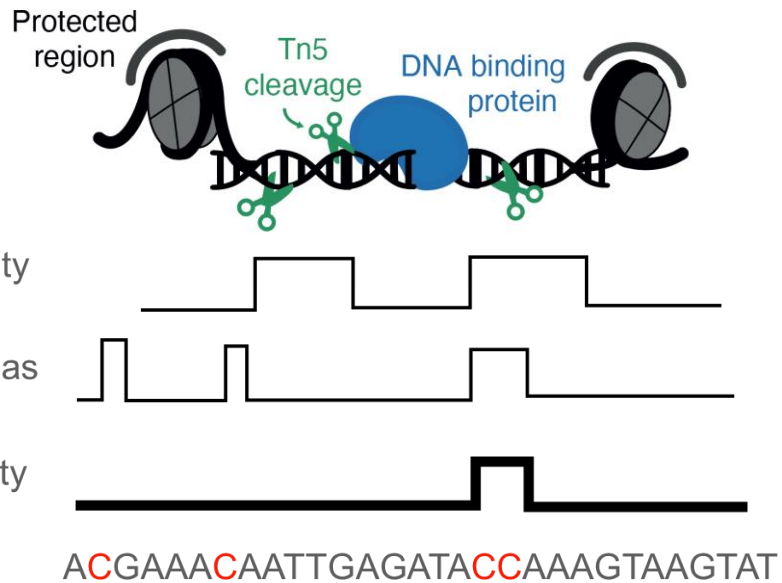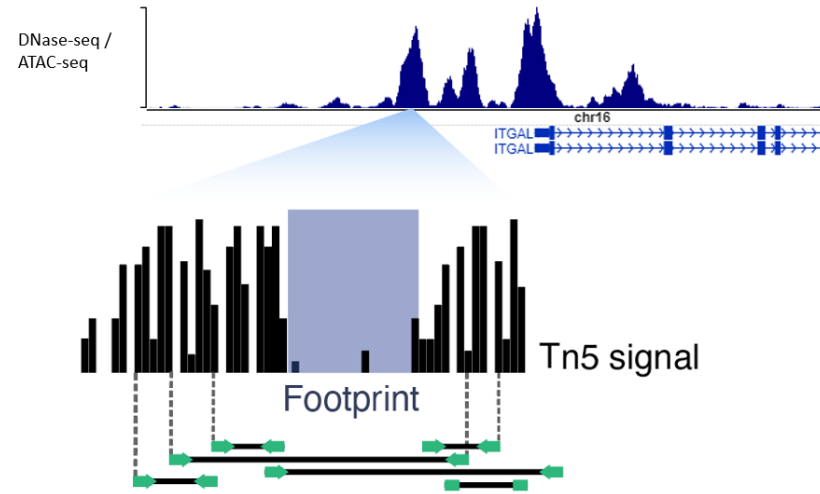ATAC-seq log fold-change loss after TF depletion

*(Independent previously published data from Friman et al. 2019)*
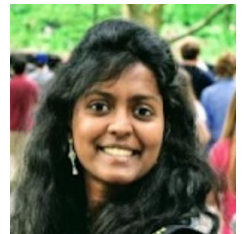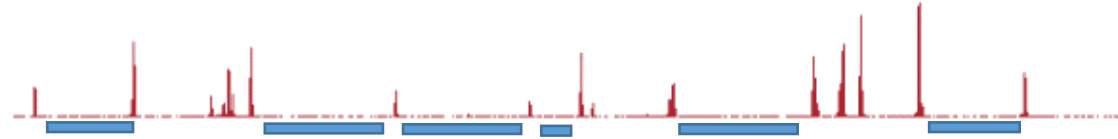
*(Independent published MPRA data from King, Maricque, Cohen 2018)*

# Modeling ATAC-seq / DNase-seq profiles (enzyme bias affects footprints)

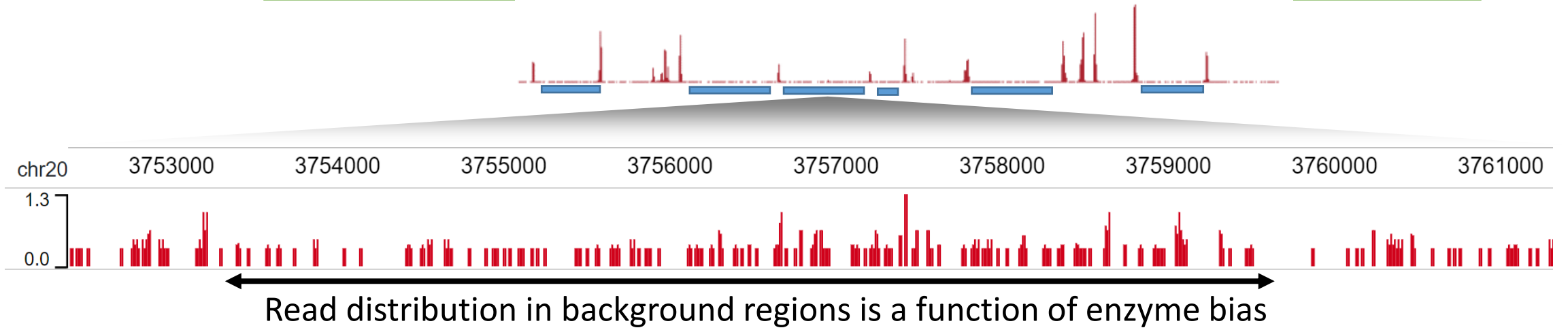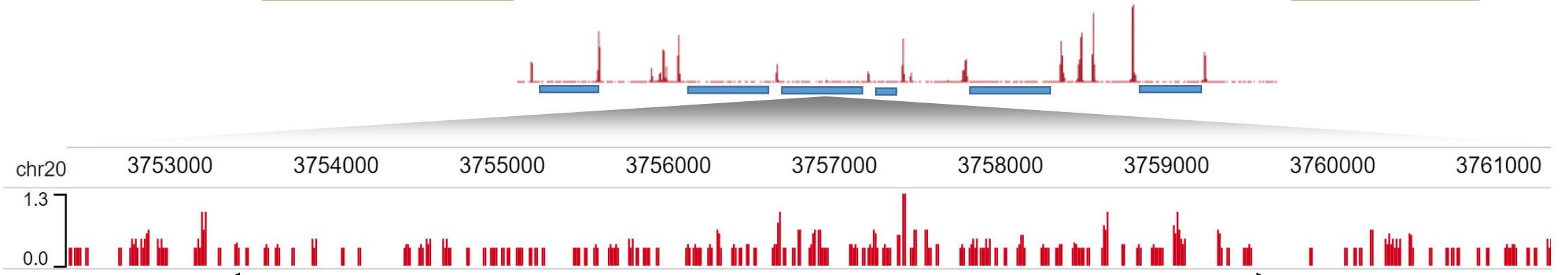# How to estimate Tn5 / DNase bias?
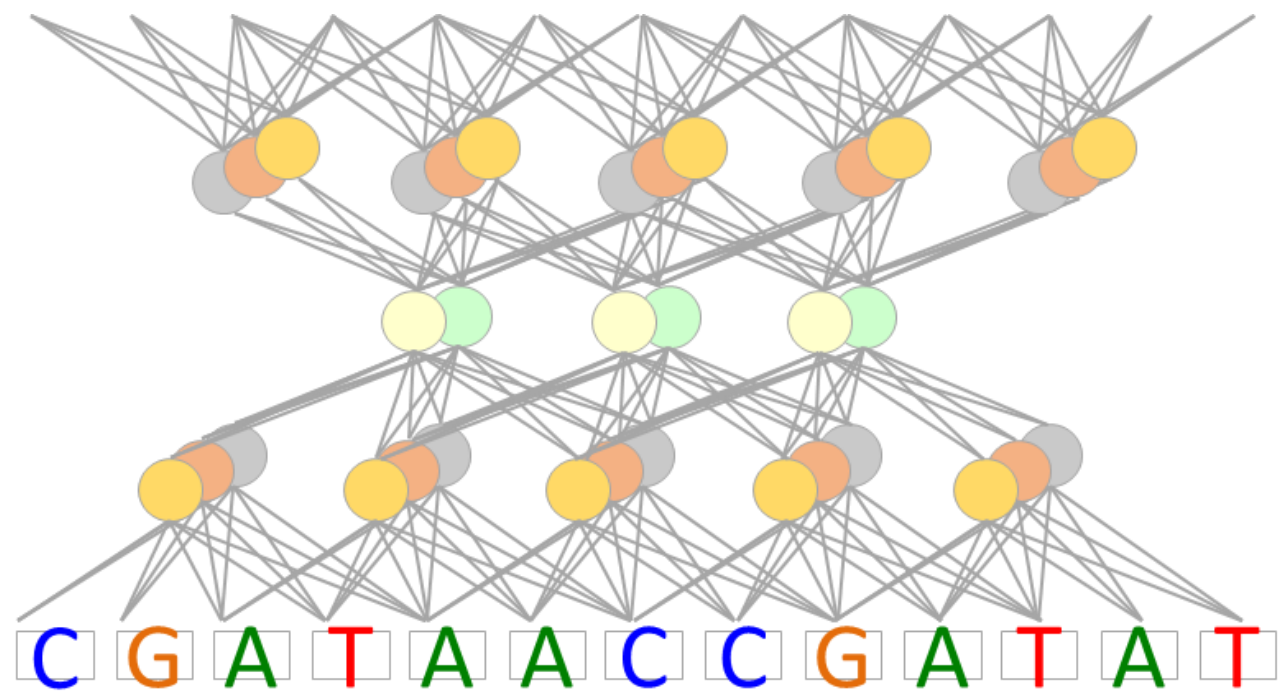


Anusri Pampari

Anna Shcherbina

# How to estimate Tn5 / DNase bias?



Read distribution in background regions is a function of enzyme bias

Anusri Pampari

Anna Shcherbina

# How to estimate Tn5 / DNase bias?



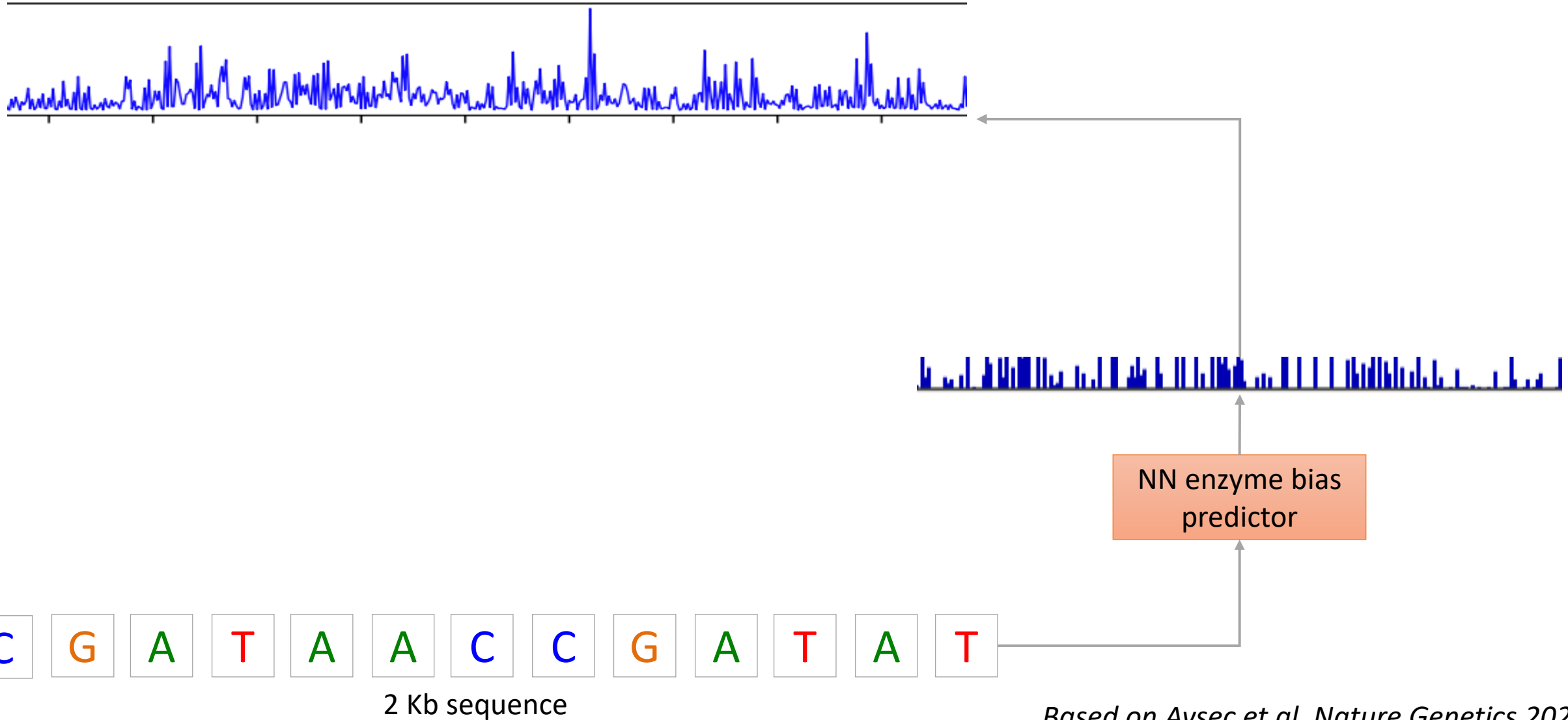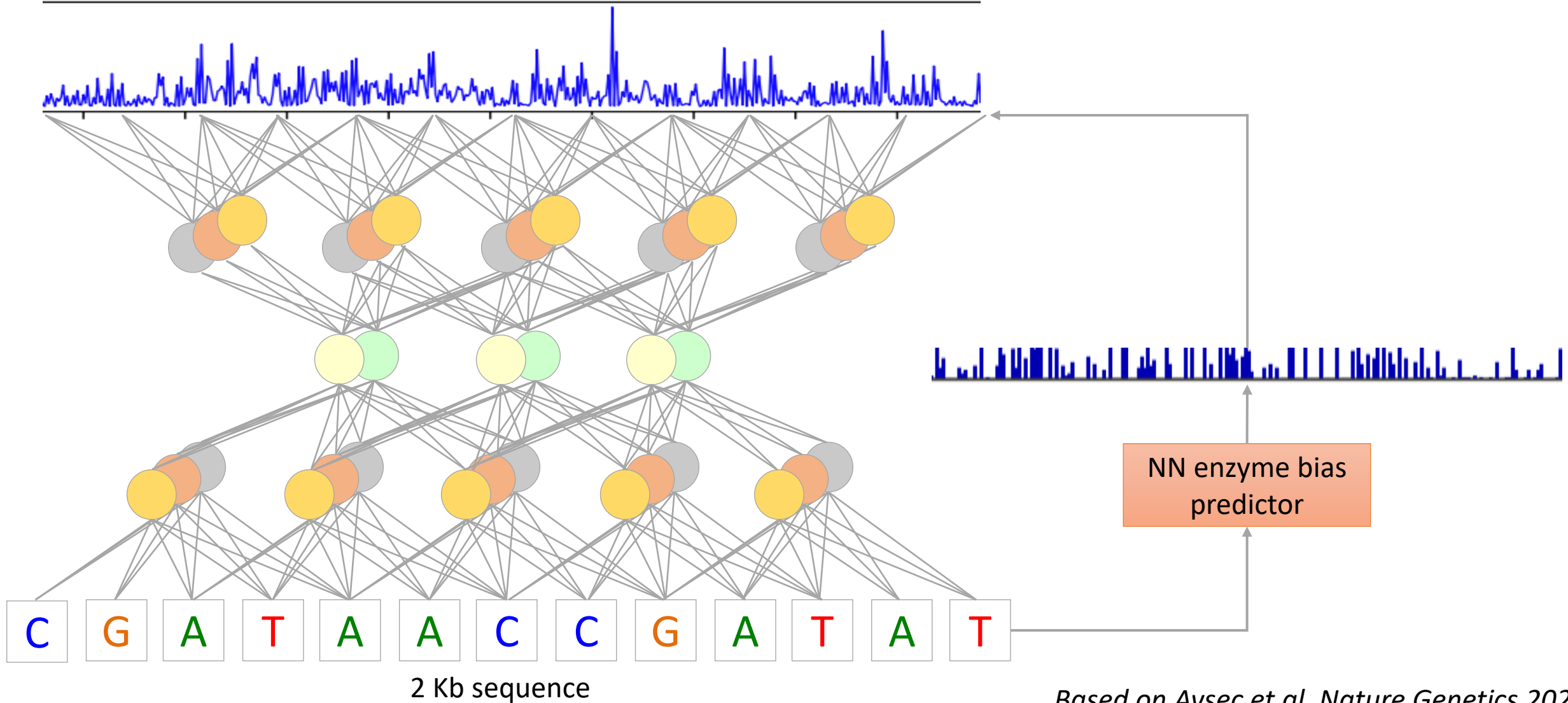Read distribution in background regions is a function of enzyme bias

Anusri Pampari

Anna Shcherbina

# ChromBPNet: Sequence to base-res chromatin accessibility profiles

total Tn5/DNase insertions (1 kb)
base-resolution probability profile (1 kb)

NN enzyme bias predictor

C G A T A A C C G A T A T

2 Kb sequence

*Based on Avsec et al. Nature Genetics 2021*

# ChromBPNet: Sequence to base-res chromatin accessibility profiles



total Tn5/DNase insertions (1 kb)
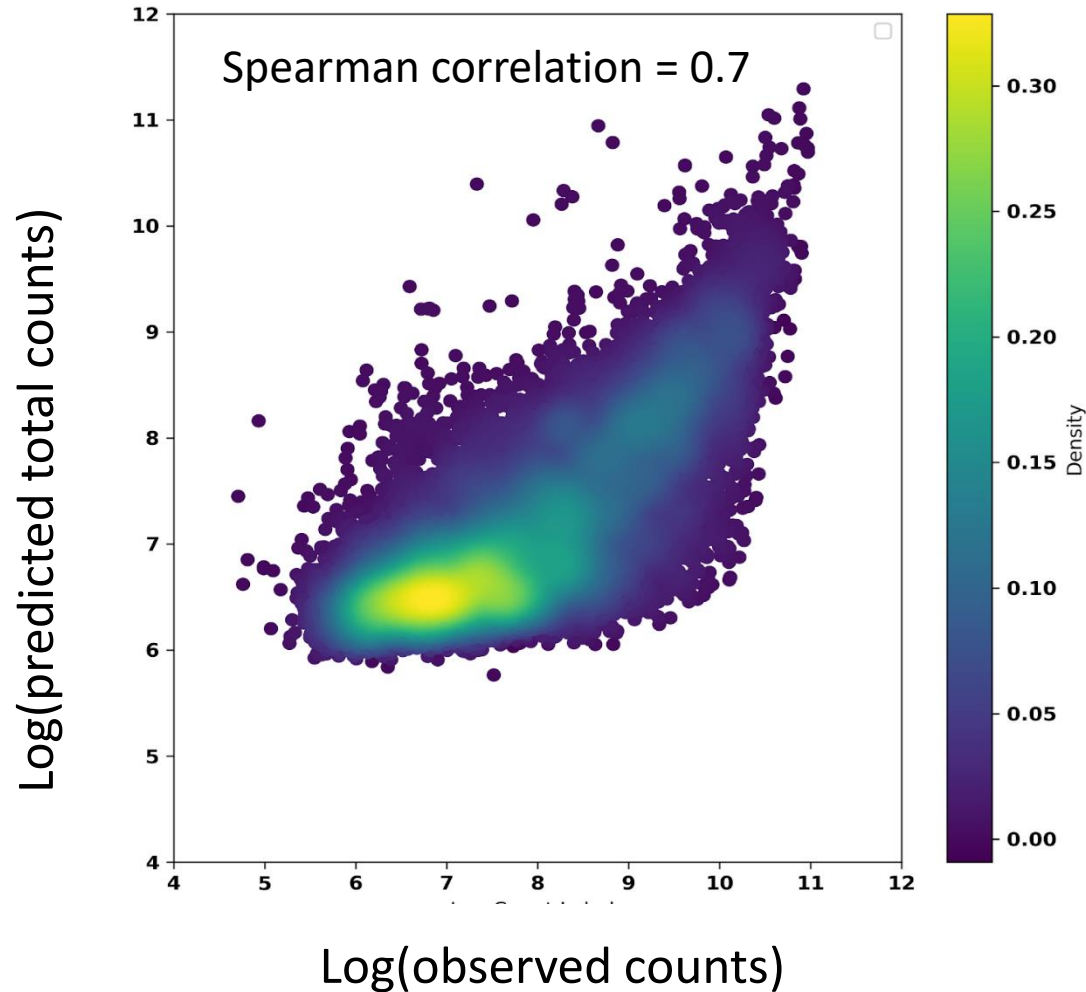base-resolution probability profile (1 kb)

NN enzyme bias predictor

C G A T A A C C G A T A T

2 Kb sequence

*Based on Avsec et al. Nature Genetics 2021*

Total counts prediction performance

# Prediction performance (held-out chromosomes)



Total counts prediction performance

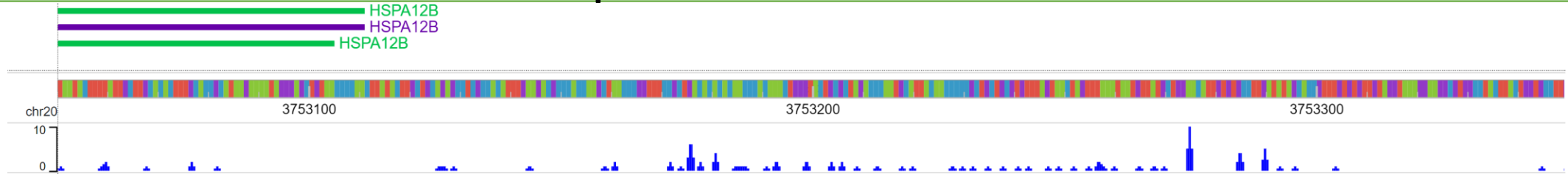Spearman correlation = 0.7

Log(predicted total counts)

Log(observed counts)

Profile prediction performance

Best limit

Observed vs. predicted profile

Worst limit

**Predicted vs Labels**
**Pseudoreps**
**Labels vs Shuffled Labels**
**Labels vs Mean Summit-Centered Label**
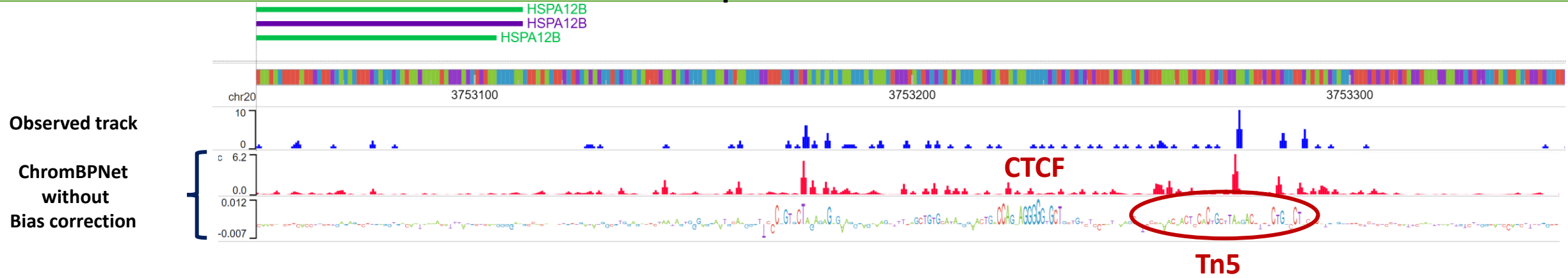
Jensen-Shannon Distance

# Denoised base-resolution bias-corrected chromatin accessibility footprints & de-biased sequence features

# Denoised base-resolution bias-corrected chromatin accessibility footprints & de-biased sequence features
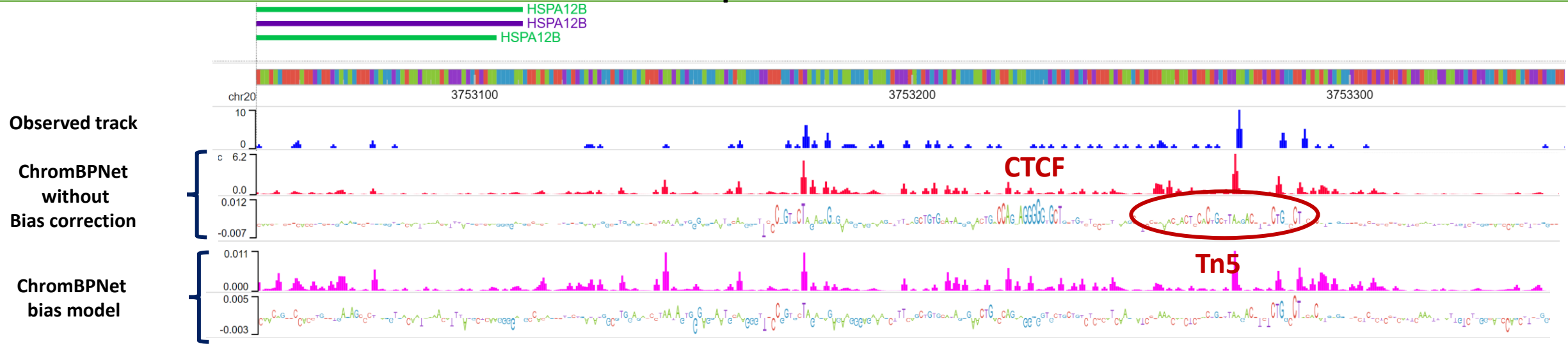
HSPA12B
HSPA12B
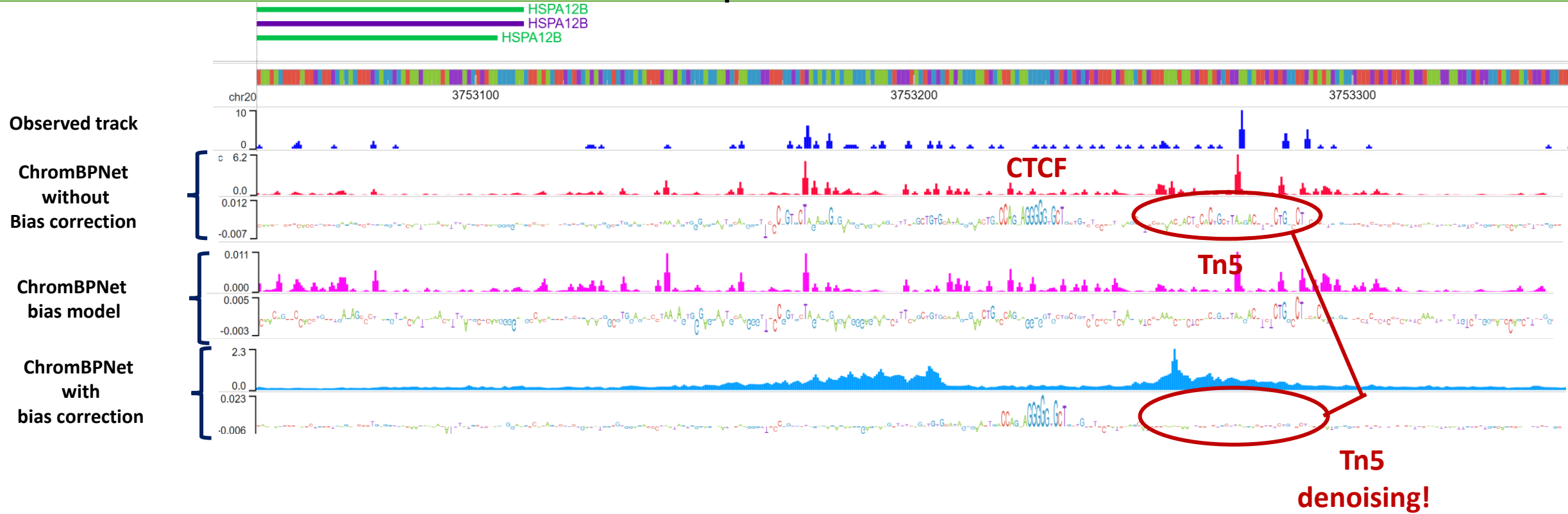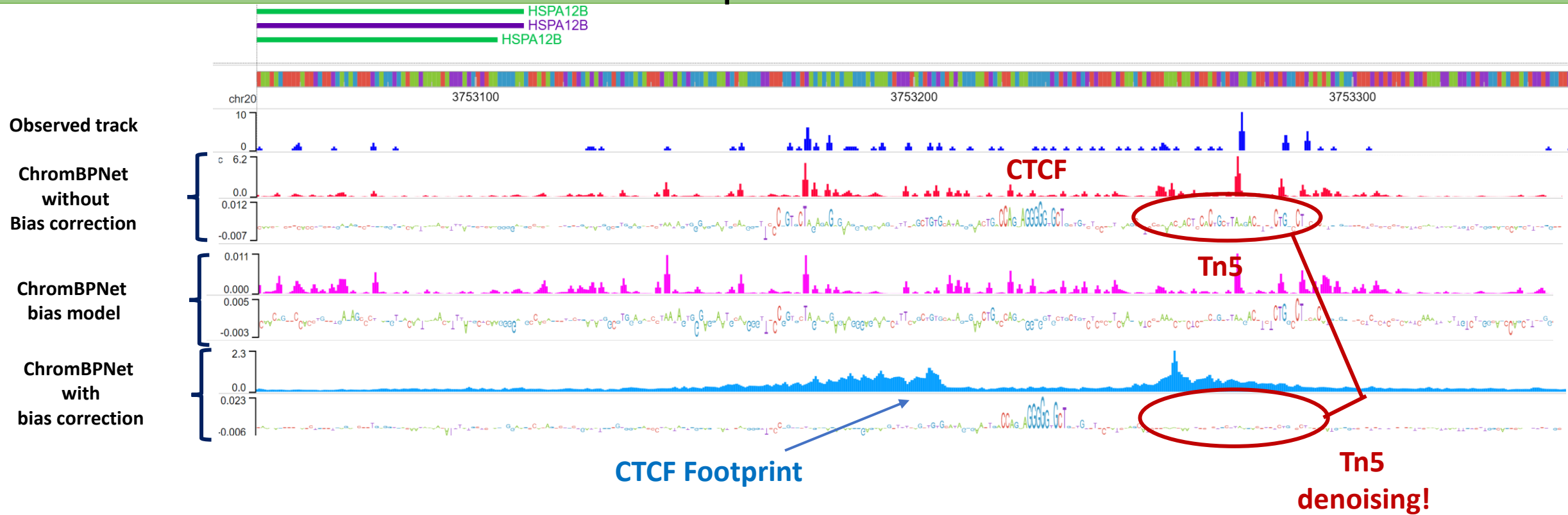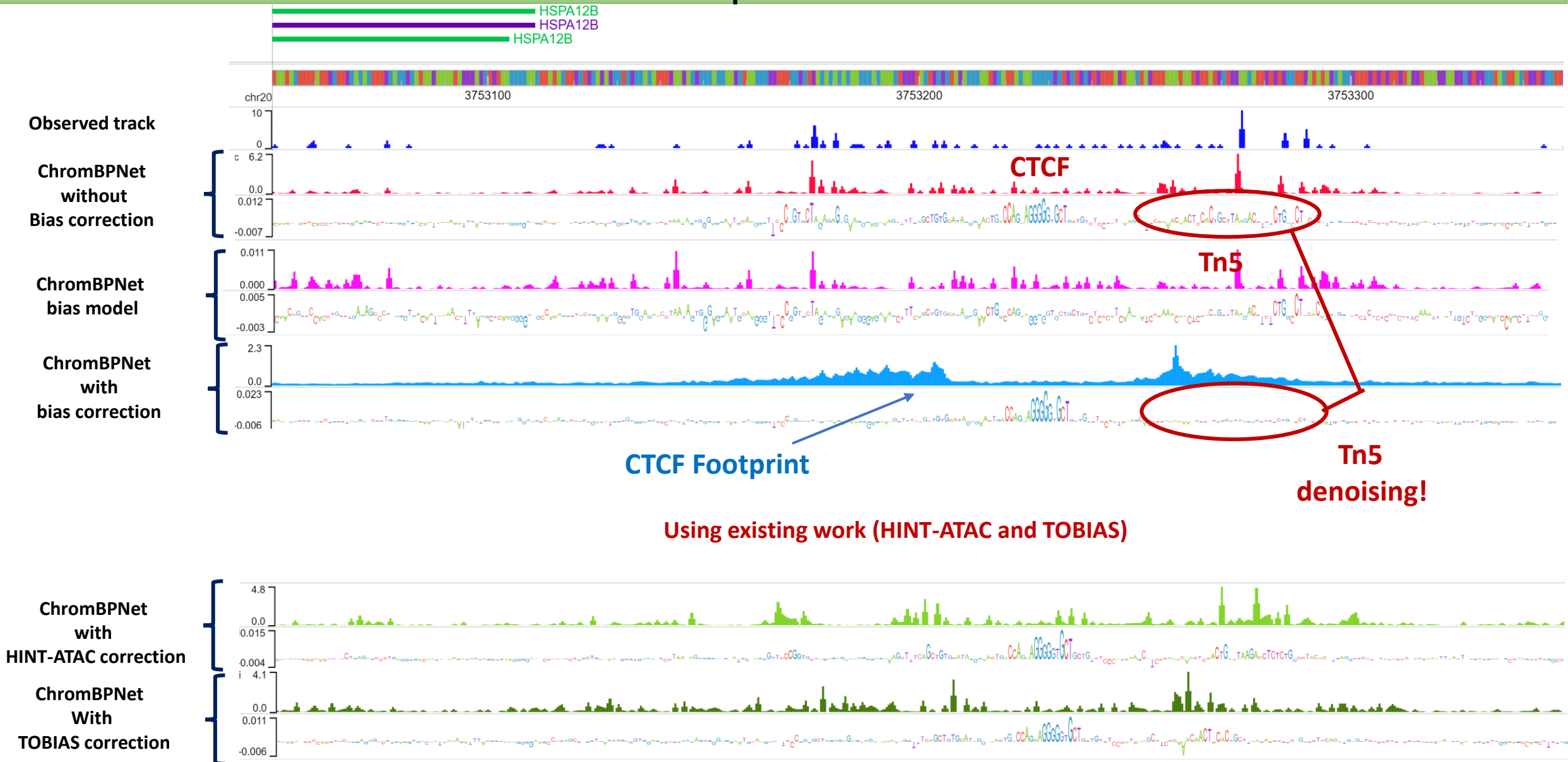HSPA12B

chr20       3753100       3753200       3753300

**Observed track**

10

0

# Denoised base-resolution bias-corrected chromatin accessibility footprints & de-biased sequence features

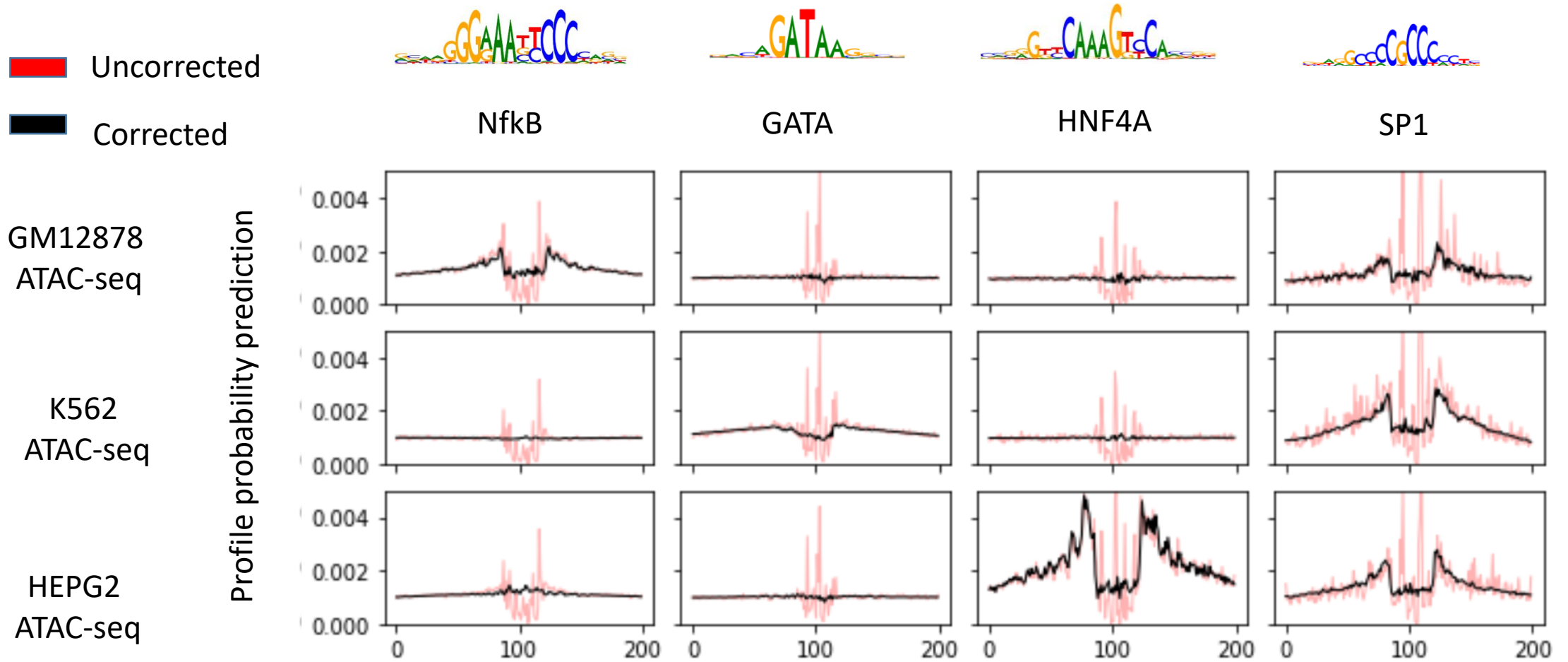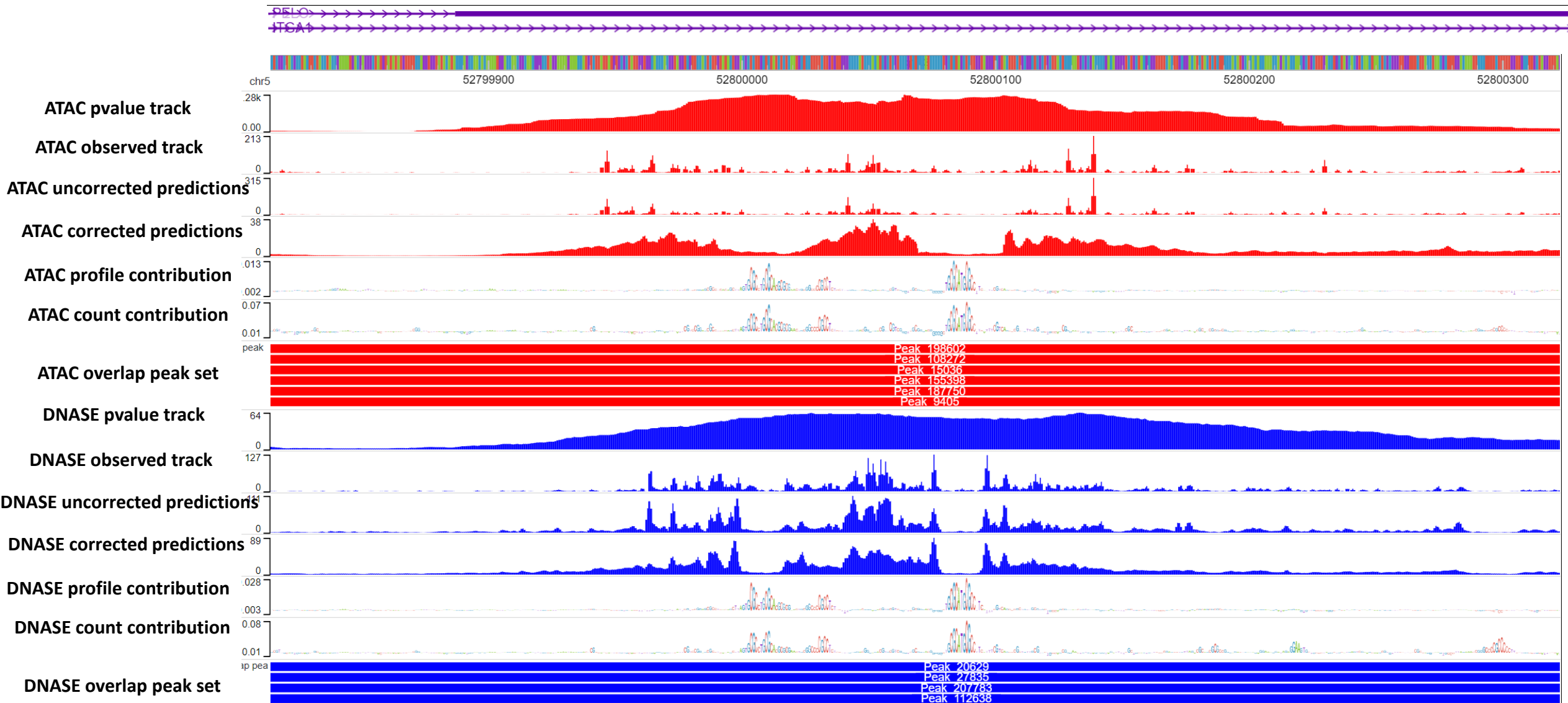# Denoised base-resolution bias-corrected chromatin accessibility footprints & de-biased sequence features

# Denoised base-resolution bias-corrected chromatin accessibility footprints & de-biased sequence features

# Denoised base-resolution bias-corrected chromatin accessibility footprints & de-biased sequence features

# Denoised base-resolution bias-corrected chromatin accessibility footprints & de-biased sequence features

# ChromBPNet can predict marginal footprints of cell-type specific TFs



200bp surrounding the motif insertion site in 10K random non-peak seqeunces

# Similar sequence syntax derived from DNase-seq and ATAC-seq data
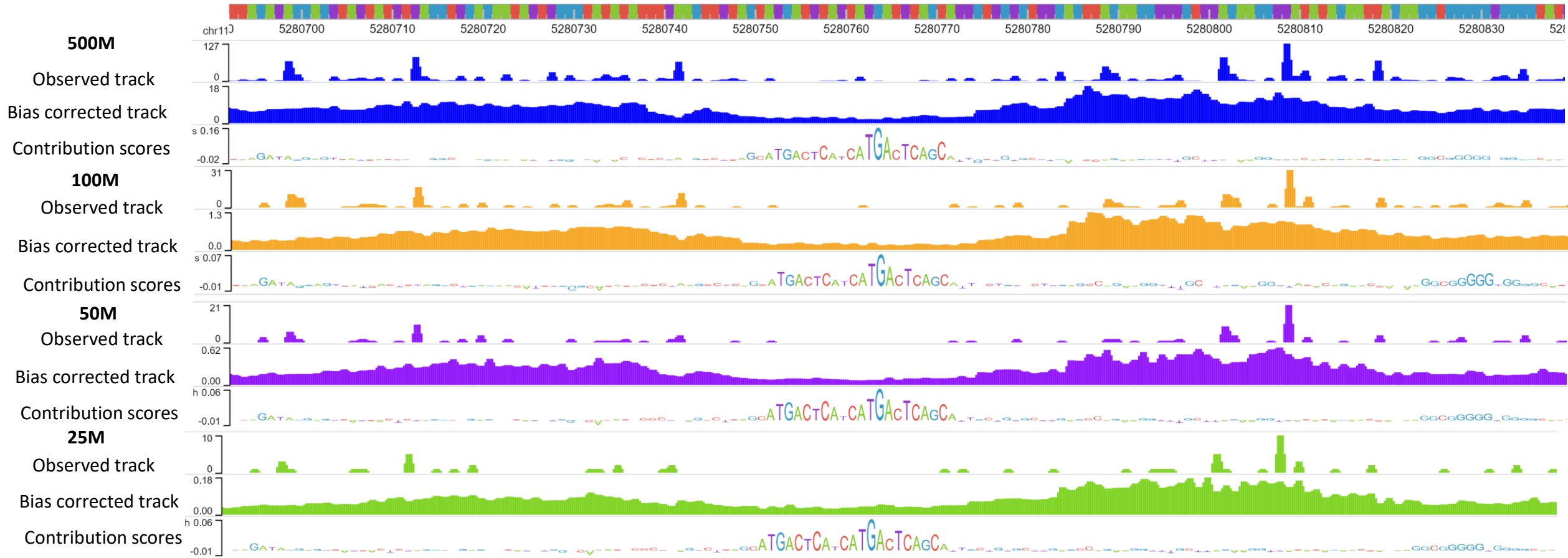
# High fidelity denoising, imputation and interpretations at different read coverages

Beta-globin locus in K562

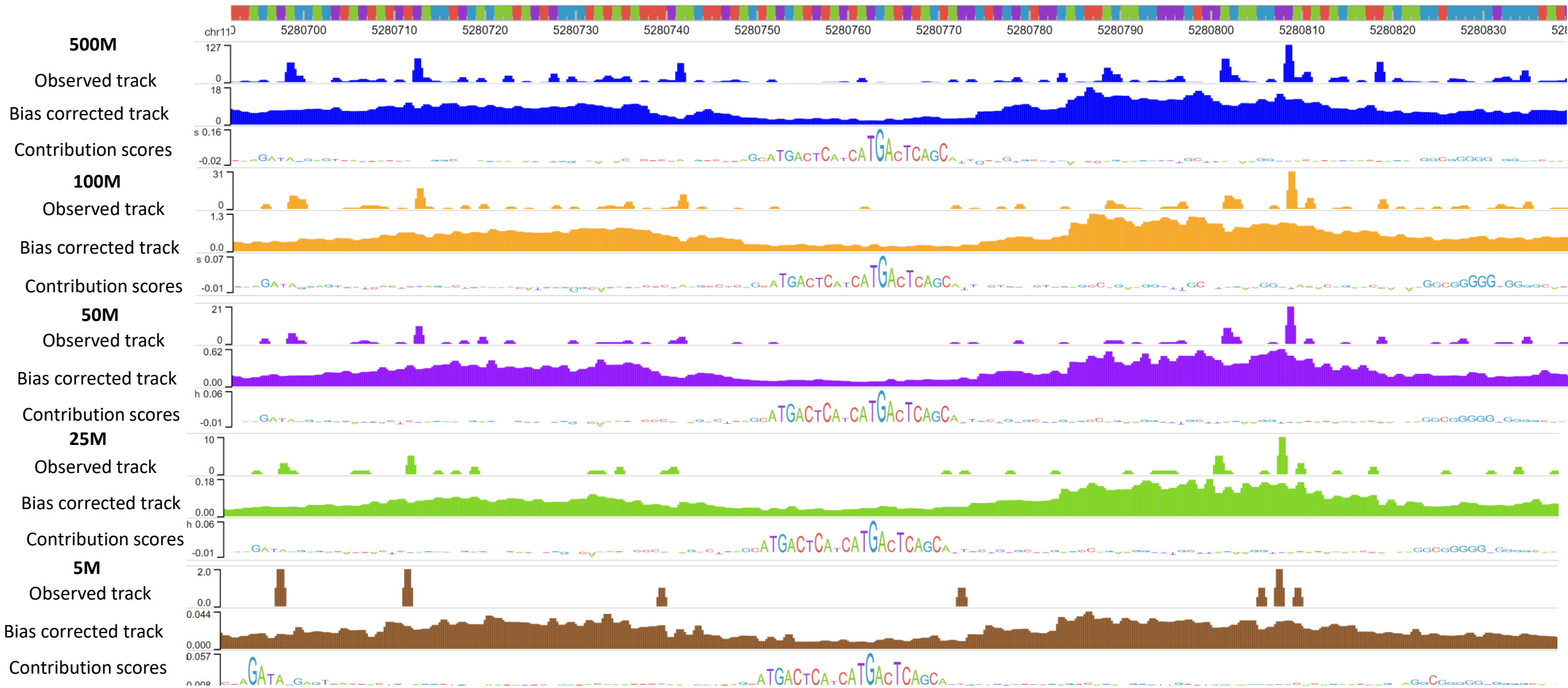# High fidelity denoising, imputation and interpretations at different read coverages

Beta-globin locus in K562

# High fidelity denoising, imputation and interpretations at different read coverages
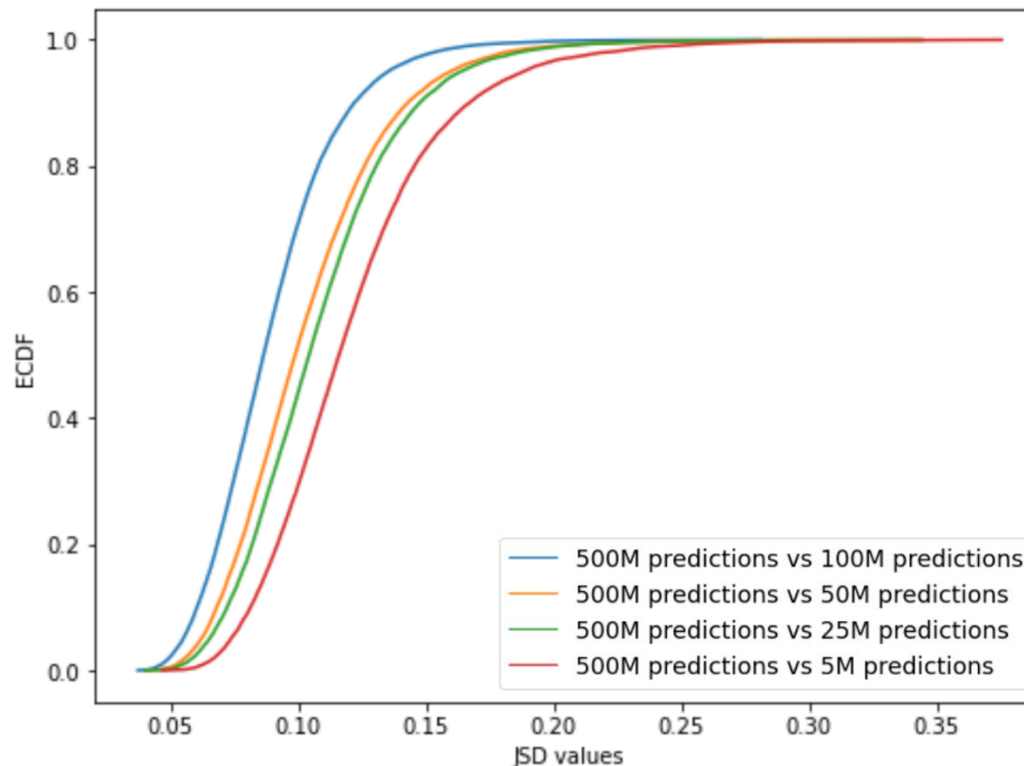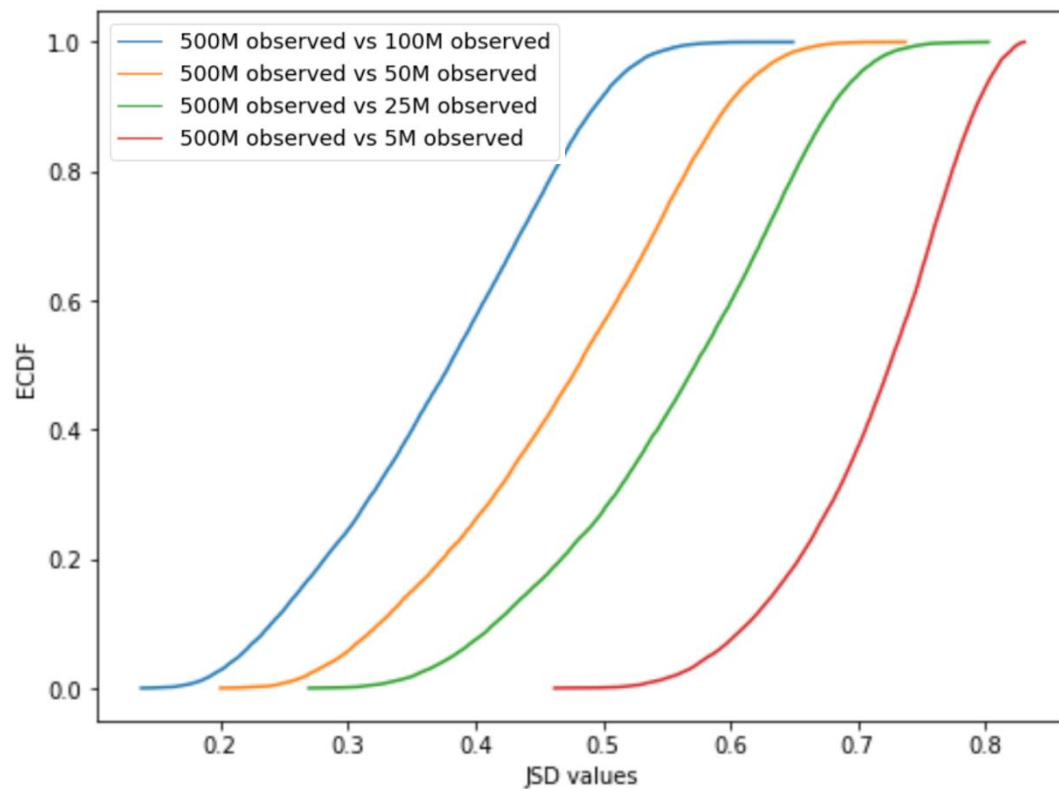
Beta-globin locus in K562

# High fidelity denoising, imputation and interpretations at different read coverages

Beta-globin locus in K562

# High fidelity denoising, imputation and interpretations at different read coverages

## Beta-globin locus in K562

# High fidelity denoising, imputation and interpretations at different read coverages

Beta-globin locus in K562

# ChromBPNet predicted tracks are substantially similar compared to observed tracks at different read depths

Using 500M as ground truth we compare degradation in signal quality at different read depths

ChromBPNet predicts substantially similar profiles compared to the observed tracks

# High fidelity marginal footprinting in K562 at different read depths

High fidelity marginal footprinting in K562 at different read depths

CTCF

SPI1

NFYB

GABPA

Profile probability prediction

500M

200bp surrounding the motif insertion site

High fidelity marginal footprinting in K562 at different read depths
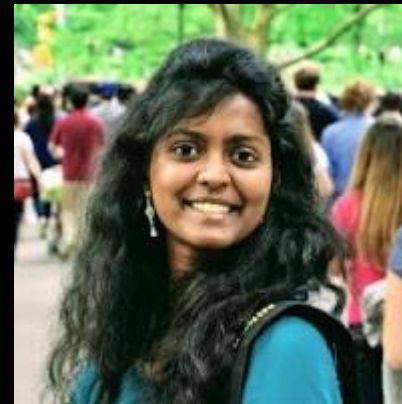
200bp surrounding the motif insertion site

High fidelity marginal footprinting in K562 at different read depths

200bp surrounding the motif insertion site

High fidelity marginal footprinting in K562 at different read depths

High fidelity marginal footprinting in K562 at different read depths

CTCF

SPI1

NFYB

GABPA

Profile probability prediction

500M
100M
50M
25M
5M

200bp surrounding the motif insertion site

# Model driven prioritization and interpretation of non-coding genetic variation

Anna Shcherbina

Soumya Kundu

Anusri Pampari

Laksshman Sundaram

# Large proportion of disease-associated genetic loci are non-coding

Benign

…….ACTGATCG**C**AATCG…….

…….ACTGATCG**G**AATCG…….

Risk

# Large proportion of disease-associated genetic loci are non-coding

Control elements
(**Non-coding variants**)

Gene
(Coding variant)

Benign

…….ACTGATCG**C**AATCG…….

…….ACTGATCG**G**AATCG…….

Risk

■ Coding   ■ Non-coding

# BPNet/ChromBPNet can predict variants influencing regulatory activity



Predicted SPI1 protein-DNA binding

Predicted chromatin accessibility

Predicted histone mark (H3K27ac)

# BPNet/ChromBPNet can predict variants influencing regulatory activity



Predicted SPI1 TF ChIP-seq

Predicted DNase-seq

Predicted H3K27ac ChIP-seq

# BPNet/ChromBPNet can interpret variants influencing regulatory activity



Predicted SPI1 TF ChIP-seq

Predicted DNase-seq

Predicted H3K27ac ChIP-seq

1 Kb

1 Kb
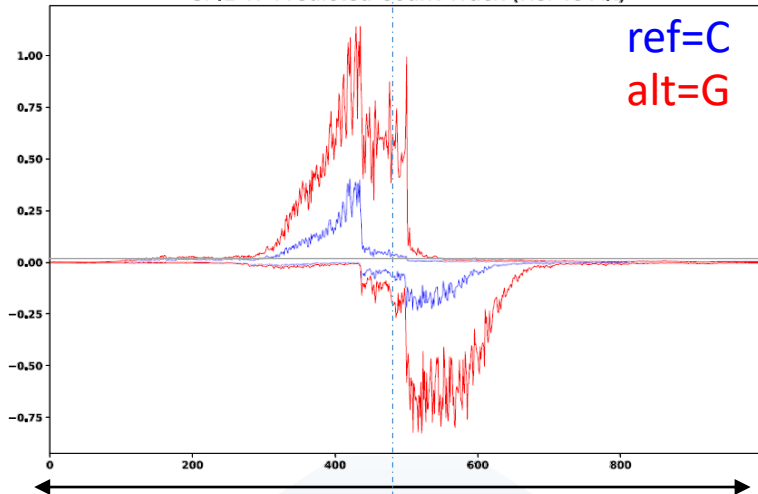
6 Kb

200 bp

200 bp

200 bp

Model interpretation predicts sequence drivers

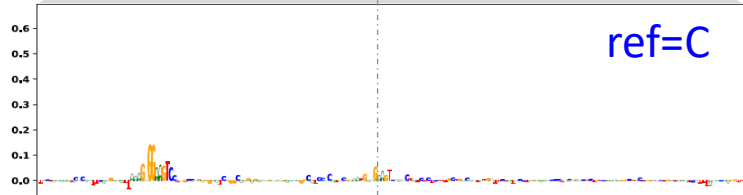# BPNet/ChromBPNet can interpret variants influencing regulatory activity
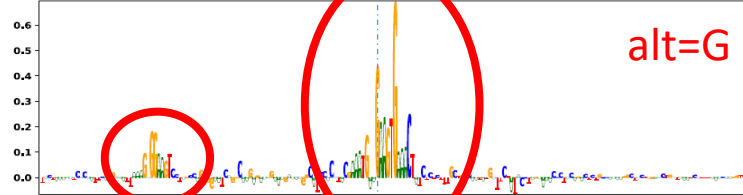


Predicted SPI1 TF ChIP-seq

Predicted DNase-seq

Predicted H3K27ac ChIP-seq

SPI1 motifs

Model interpretation predicts sequence drivers

# Variant Effect Scoring with ChromBPNet

- ChromBPNet has two heads counts and profiles

- Variant effect scoring with counts head

$$log(counts_{alt}) - log(counts_{ref})$$
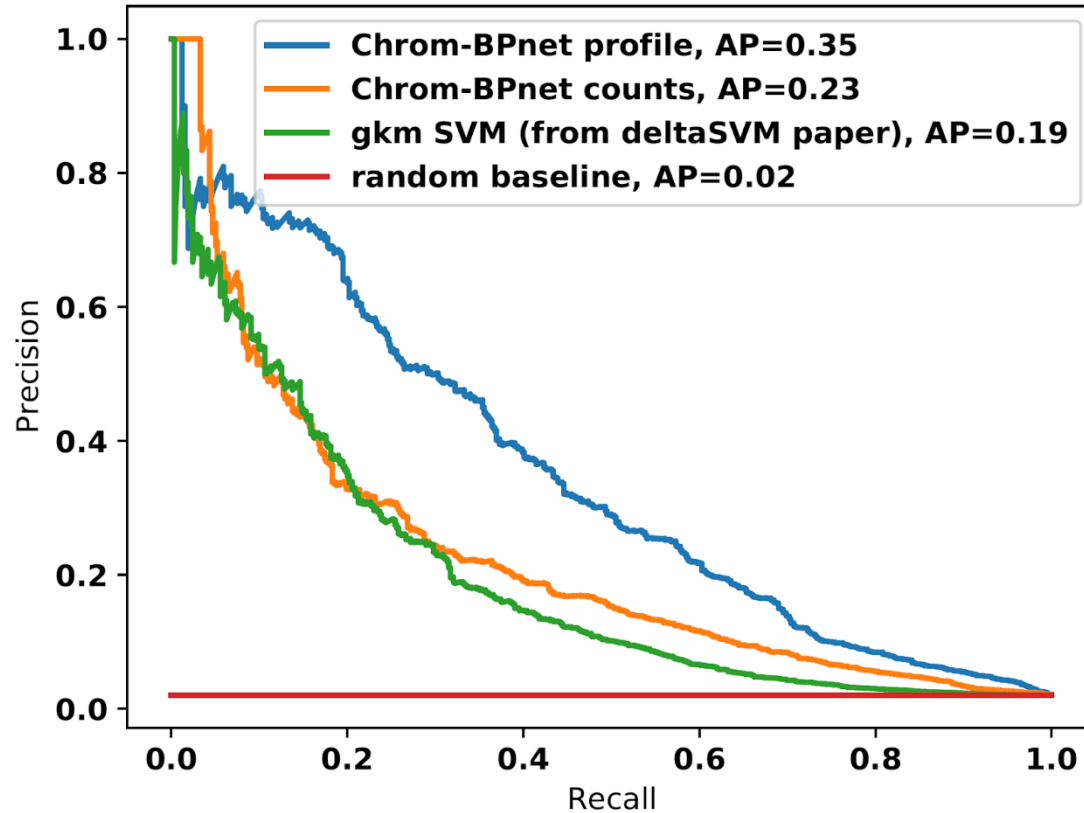
- Variant effect scoring with profile head

$$JensenShanon(Profile_{alt,} Profile_{ref}) * Sign(log(counts_{alt}) - log(counts_{ref}))$$

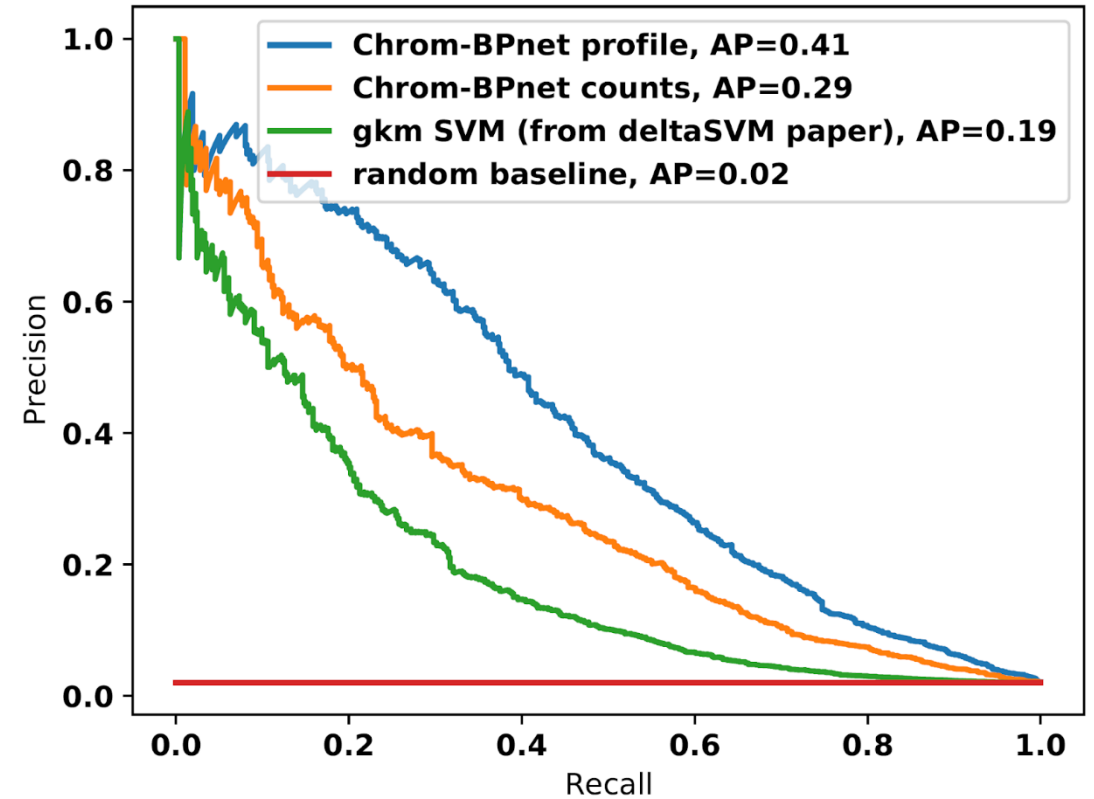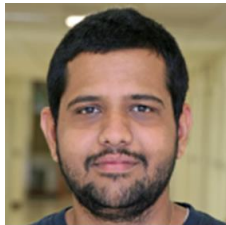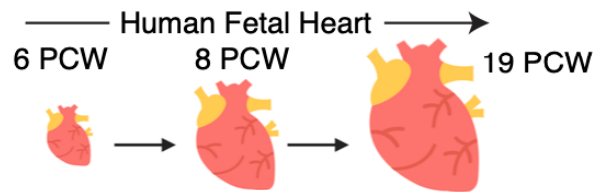# Dnase-seq ChromBPNet outperforms deltaSVM for predicting dsQTLs in LCLs

## GM12878 DNASE-seq model
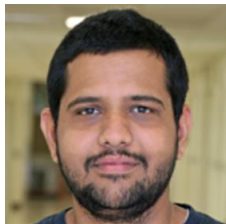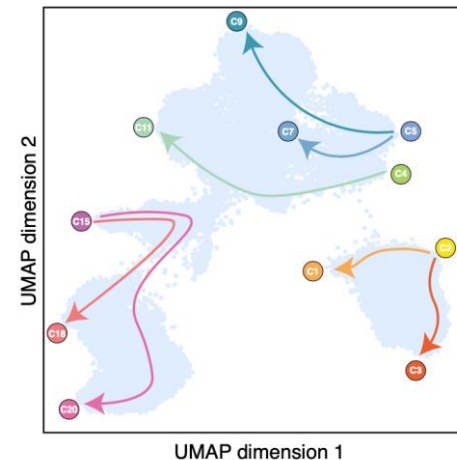
### 85M read depth



## GM12878 ATAC-seq model

### 175M read depth



*dsQTLs: Degner et al 2012*

# Single cell chromatin dynamics during human cardiogenesis



Human Fetal Heart

6 PCW    8 PCW    19 PCW

Sundaram*, Ameen*, et al. In review

Laksshman
Sundaram

# Single cell chromatin dynamics during human cardiogenesis



Myocardium
Atrial Cardiomyocytes
Ventricular Cardiomyocytes
Early Cardiac Fibroblast
Cardiac Fibroblast Progenitors
Cardiac Fibroblast
Endocardial Cushion
Late Endocardial Cushion
OFT SMC
Vasculature Development
vSMC
Pericytes
Neural Crest
Undifferentiated Epicardium
Endocardium
Transitioning Endocardium
Lymph Endothelium
Arterial Endothelium
Capillary Endothelium
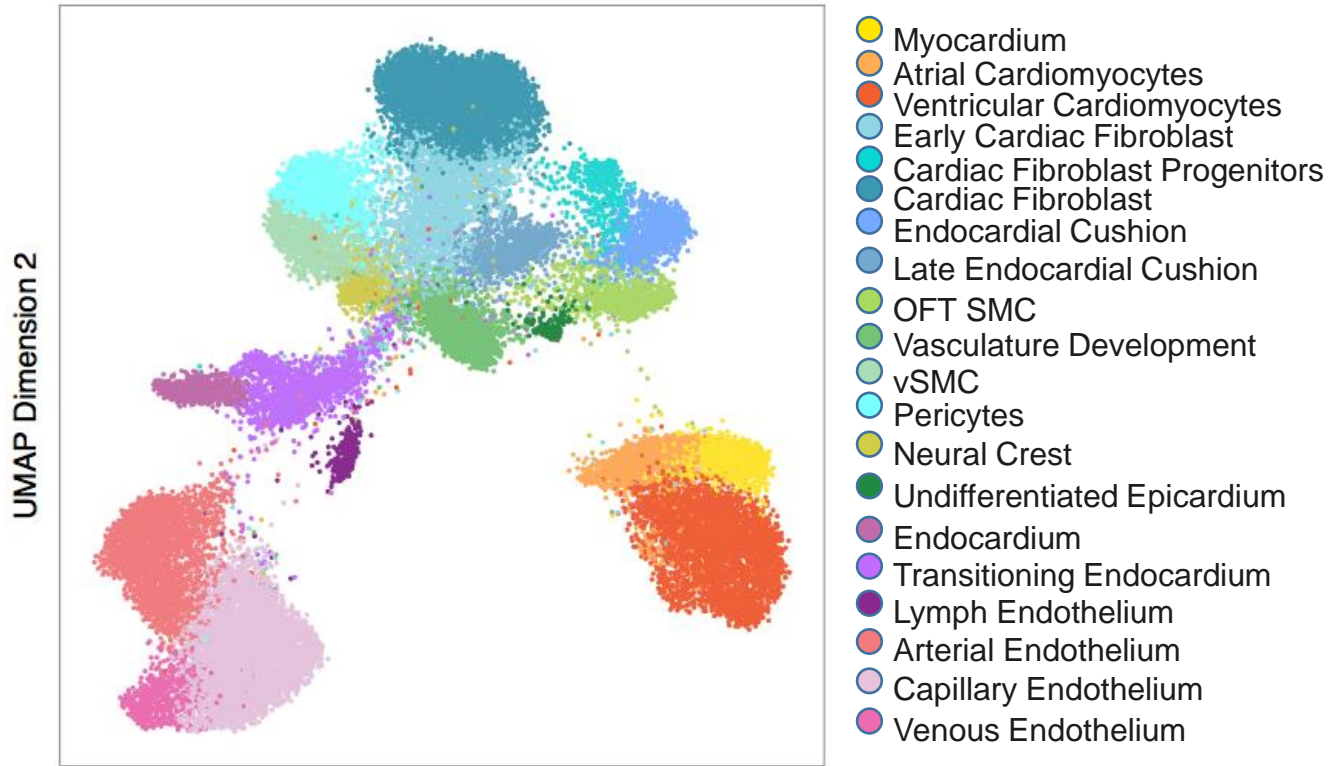Venous Endothelium
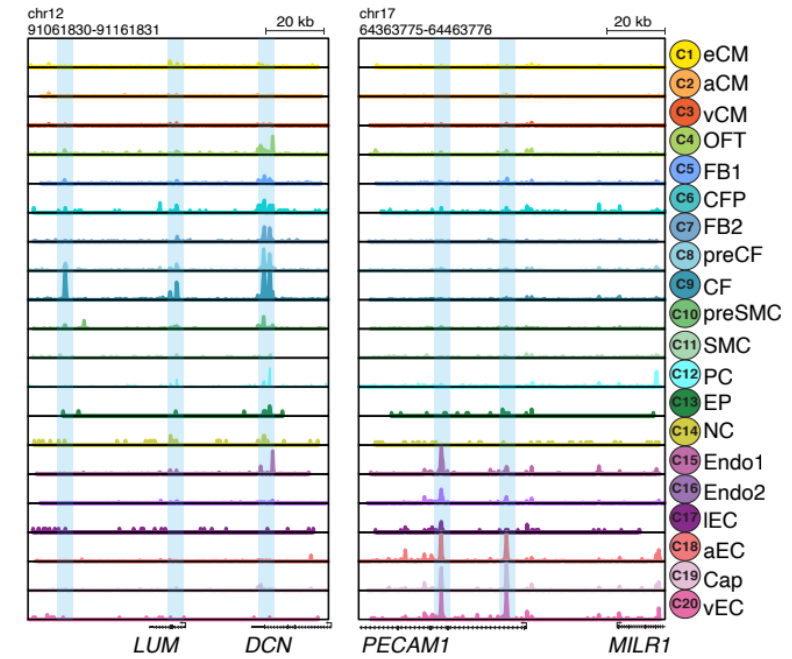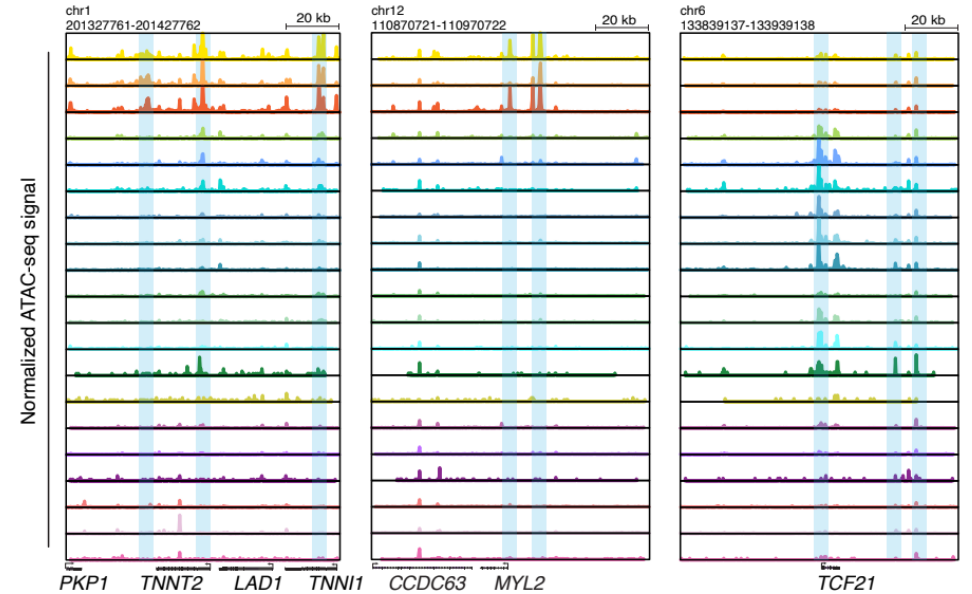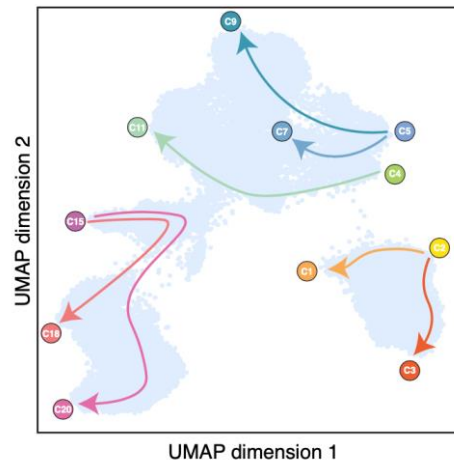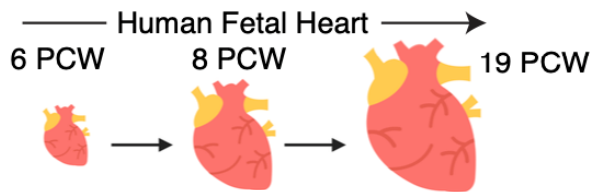
Sundaram*, Ameen*, et al. In review

Lakshman Sundaram

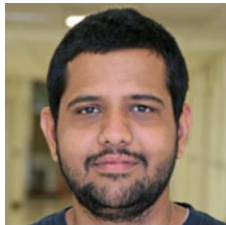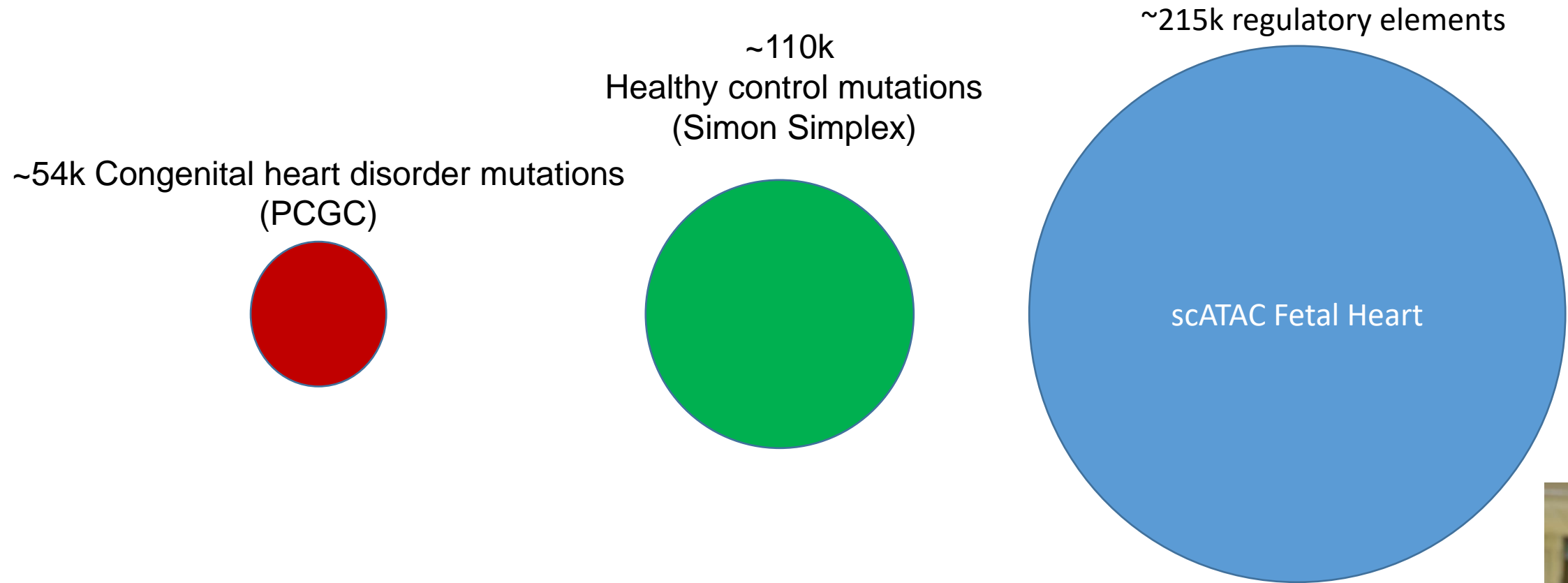# Single cell chromatin dynamics during human cardiogenesis



Sundaram*, Ameen*, et al. In review

# Prioritizing de-novo mutations in congenital heart disease with cell-type resolved regulatory map of fetal heart



~54k Congenital heart disorder mutations (PCGC)

~110k Healthy control mutations (Simon Simplex)
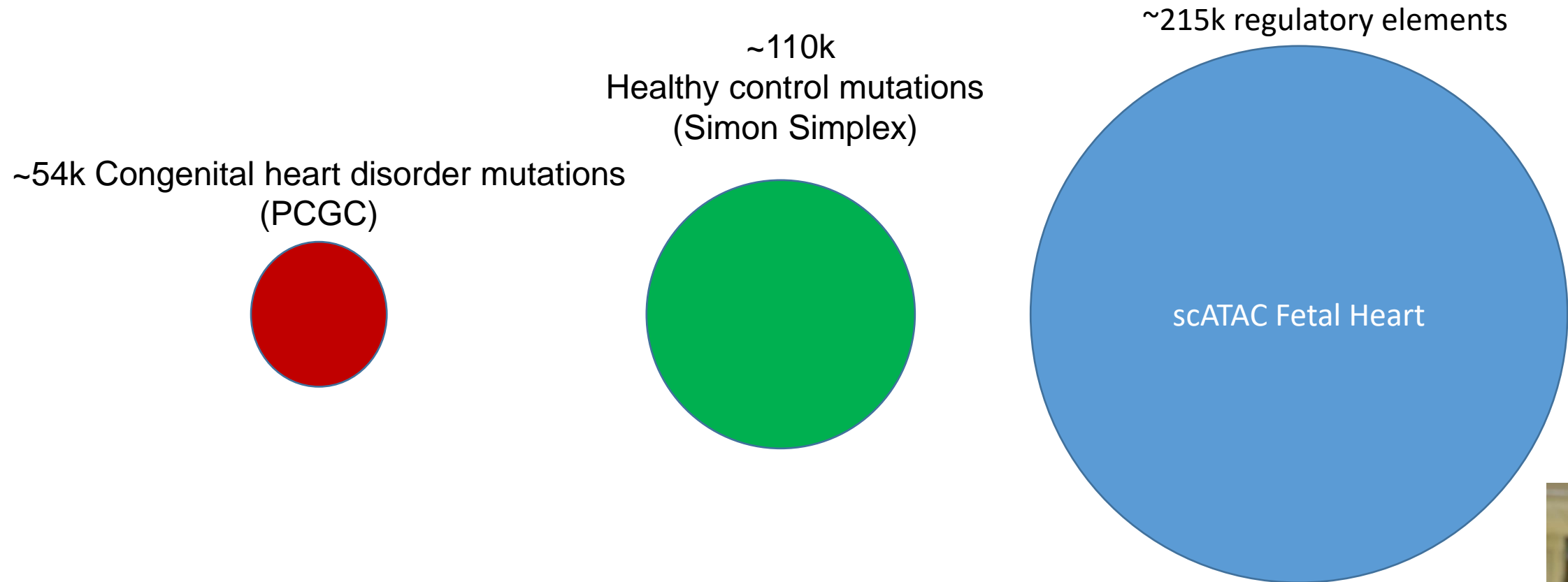
~215k regulatory elements

scATAC Fetal Heart

Laksshman Sundaram

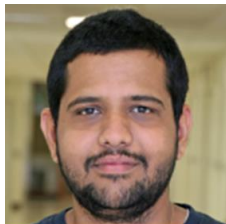# Prioritizing de-novo mutations in congenital heart disease with cell-type resolved regulatory map of fetal heart



~215k regulatory elements

~110k
Healthy control mutations
(Simon Simplex)

~54k Congenital heart disorder mutations
(PCGC)

scATAC Fetal Heart

Laksshman
Sundaram

**No enrichment of CHD mutations in all/cell type resolved scATAC-seq peaks!**

# Prioritizing mutations with cell-type resolved ChromBPNet models

Laksshman Sundaram
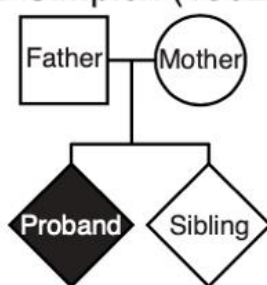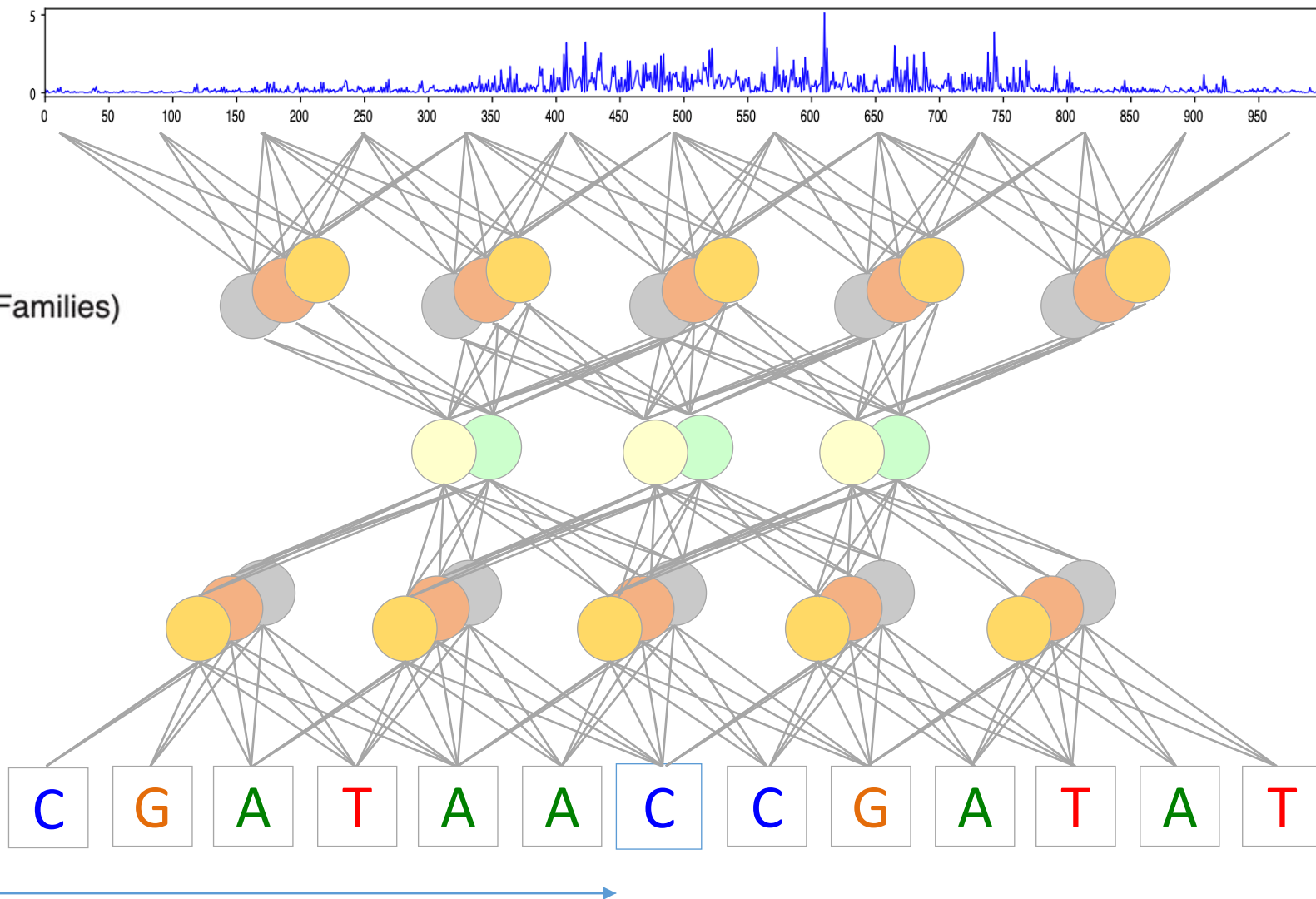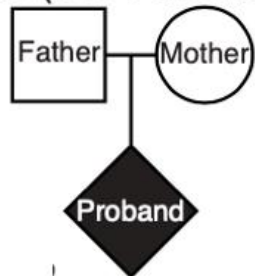
PCGC (750 Families)
Father — Mother
Proband
**Cases**

ASD Simplex (1902 Families)
Father — Mother
Proband    Sibling
**Controls**

*De novo* non-coding mutations

C G A T A A A ✗ C G A T A T
A

# Eg: CHD case mutation affecting accessibility of enhancer in Art/Cap endothelial cells
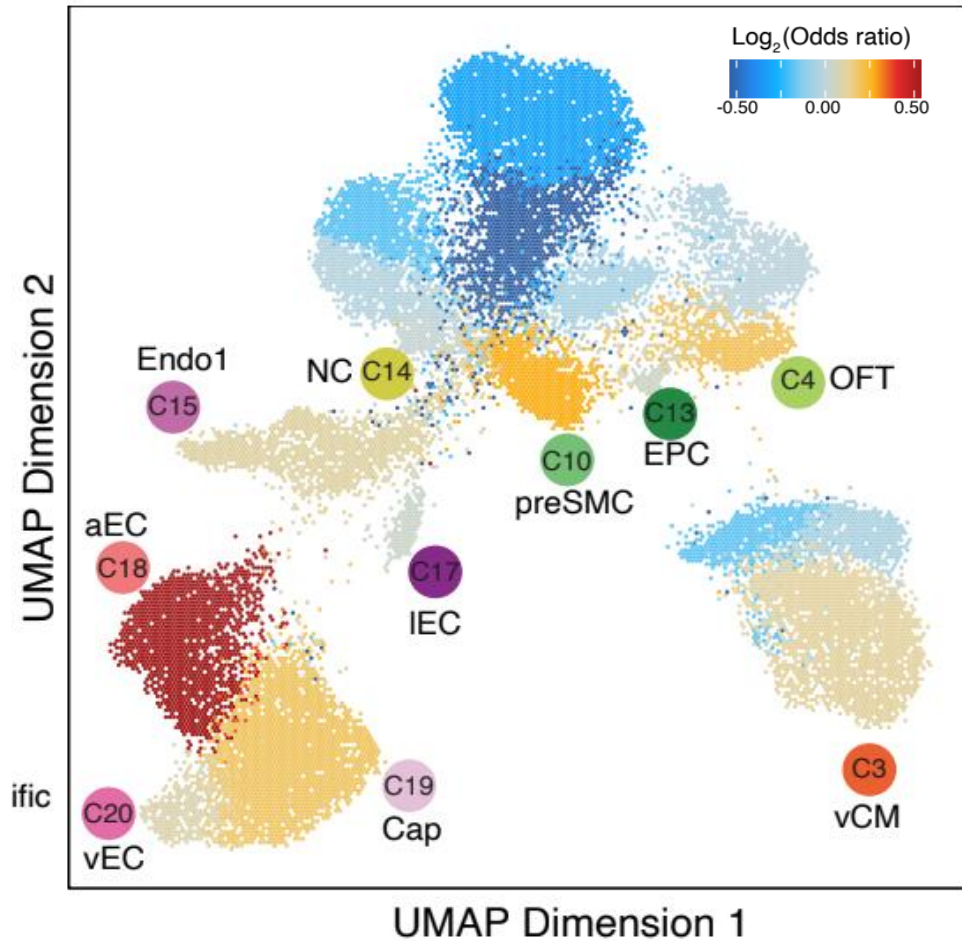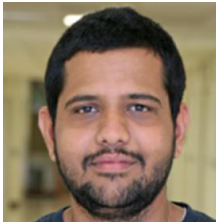


Mutation disrupts an ETS/ELK/ETV family motif

Laksshman Sundaram

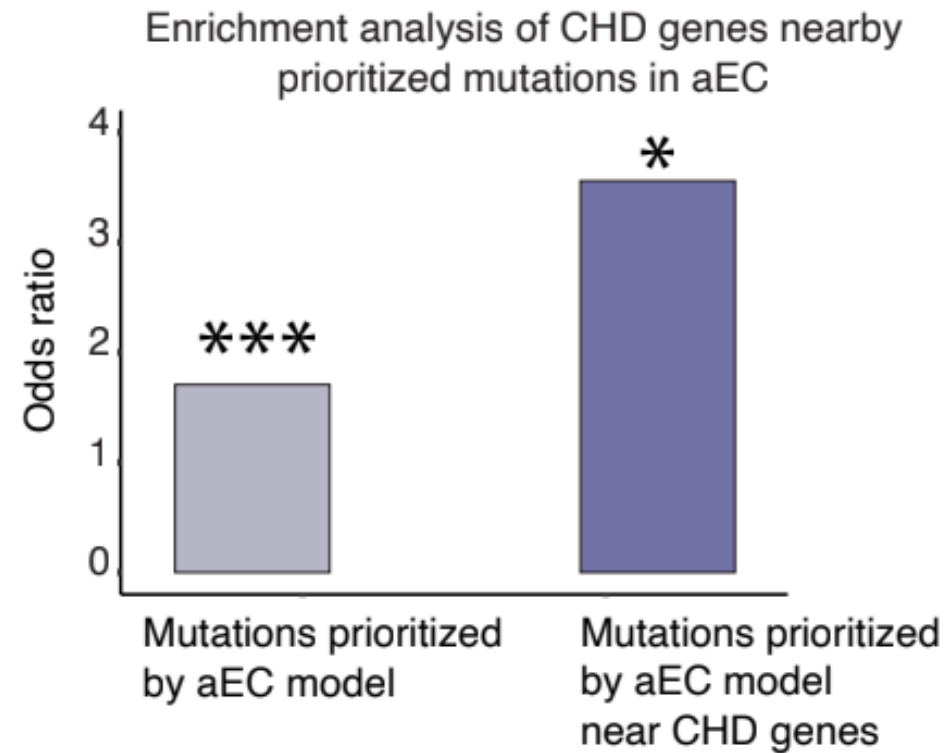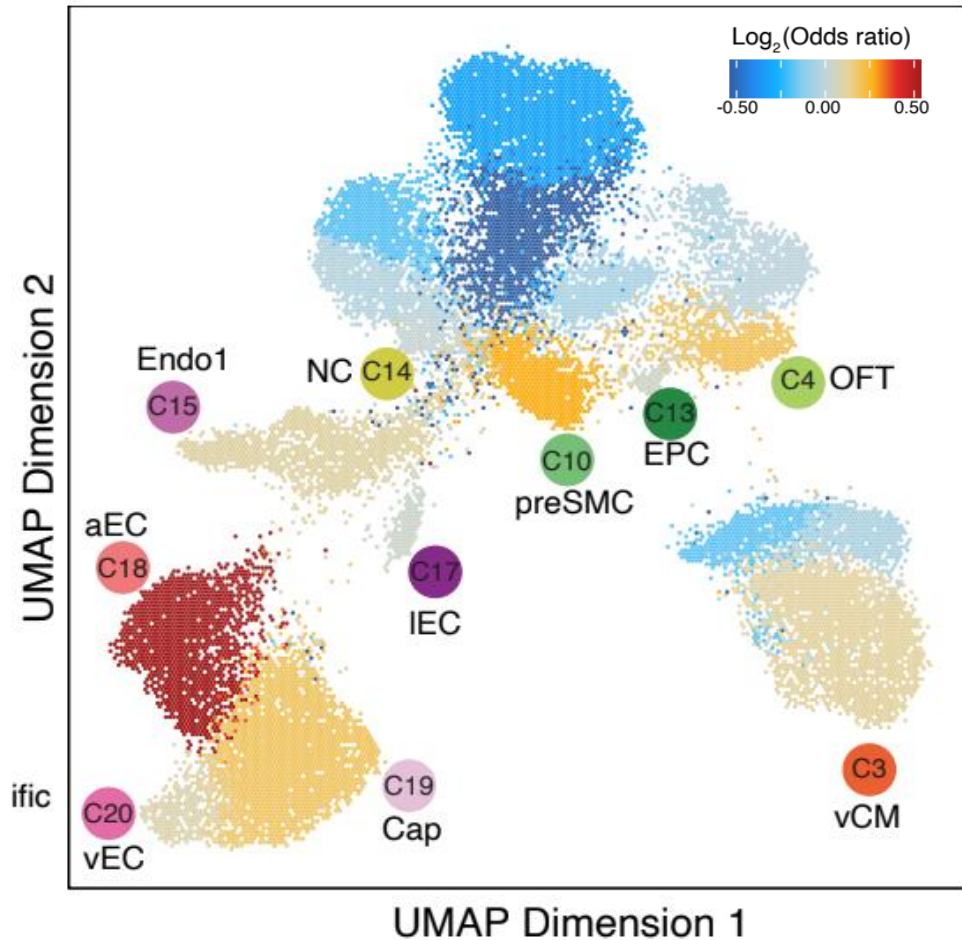# Cell states enriched for prioritized *de novo* non-coding mutations in CHD



**Arterial & Capillary endothelial cells** are most significantly enriched for CHD mutations (structural defects)

Laksshman Sundaram

# Cell states enriched for prioritized *de novo* non-coding mutations in CHD



**Arterial & Capillary endothelial cells** are most significantly enriched for CHD mutations (structural defects)
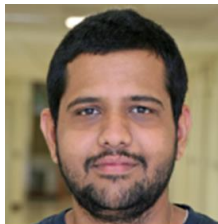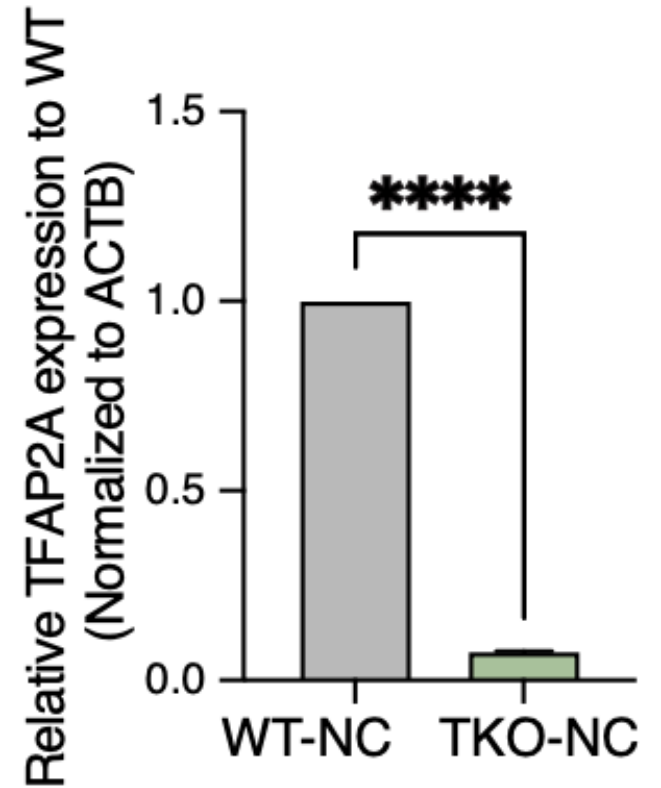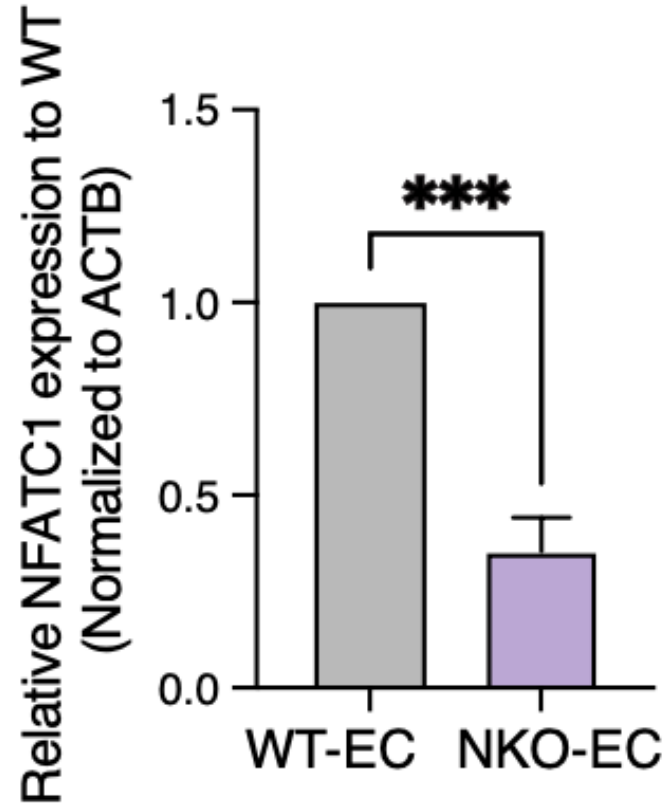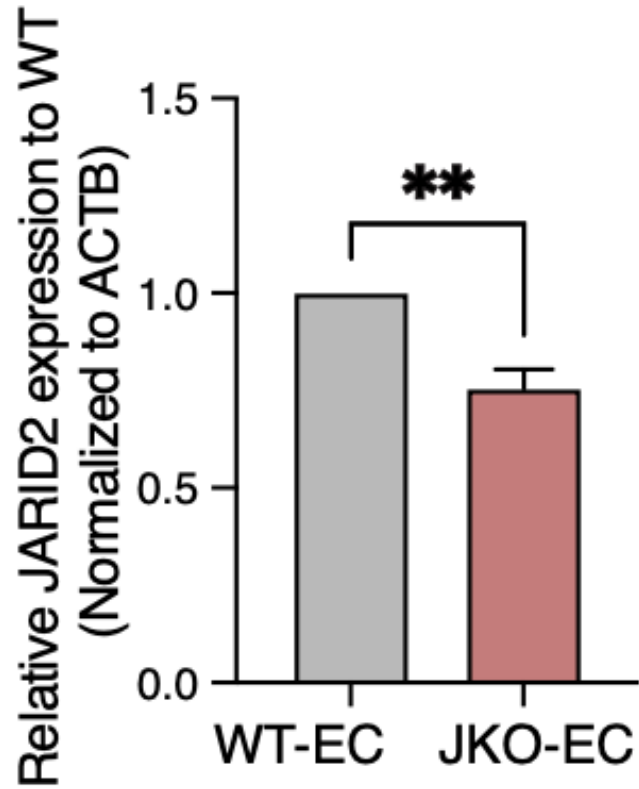
Laksshman Sundaram

CRISPR experiments confirm downstream gene targets of enhancers containing prioritized CHD mutations

## Summary

- Base-resolution neural networks can learn very accurate models of regulatory DNA sequence from bulk and single cell regulatory profiling experiments

- Can be queried to decipher novel subtle sequence syntax properties

- Can be used to decipher regulatory genetic variation

- Can be used to prioritize likely causal variants in GWAS loci and de-novo non-coding mutations

- Can be used to design precise genome editing experiments

- <u>Foundation of *in-silico* platforms for biological discovery, hypothesis generation & model-driven iterative expt. design</u>

# Kundaje lab

Daniel Kim (BMI)

Kelly Cochran (CS)

Soumya Kundu (CS)

Surag Nair (CS)

Maxim Zaslavsky (CS)

Vivek Ramalingam (Postdoc)

Caleb Lareau (Postdoc)

Akshay Balsubramani (Postdoc)

Georgi Marinov (Postdoc)

Alex Tseng (CS)

Amr Alexandari (CS)

Abhimanyu Banerjee (Physics)

Laksshman Sundaram (CS)

Anusri Pampari (CS)

Kristy Mualim (Bioinformatician)

Jacob Schreiber (Postdoc)

Mahfuza Sharmin (Postdoc)
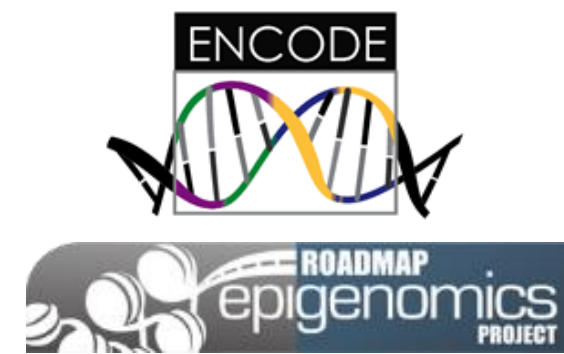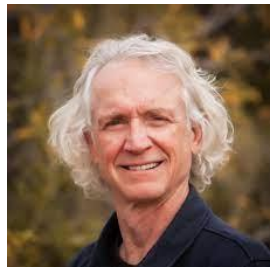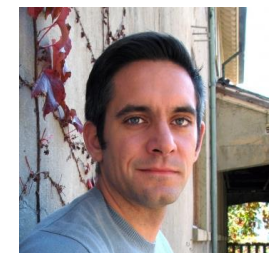
Eran Kotler (Postdoc)

Zahoor Zafrulla (ML engineer)

Jin Wook Lee (Software engineer)

# Collaborator labs