# Genetic variants associated with traits and diseases

Genetic variants associated with traits and diseases
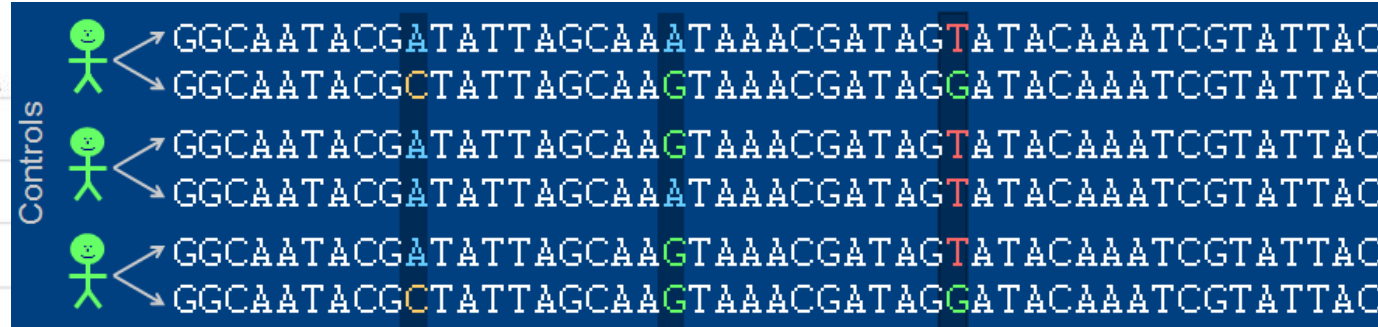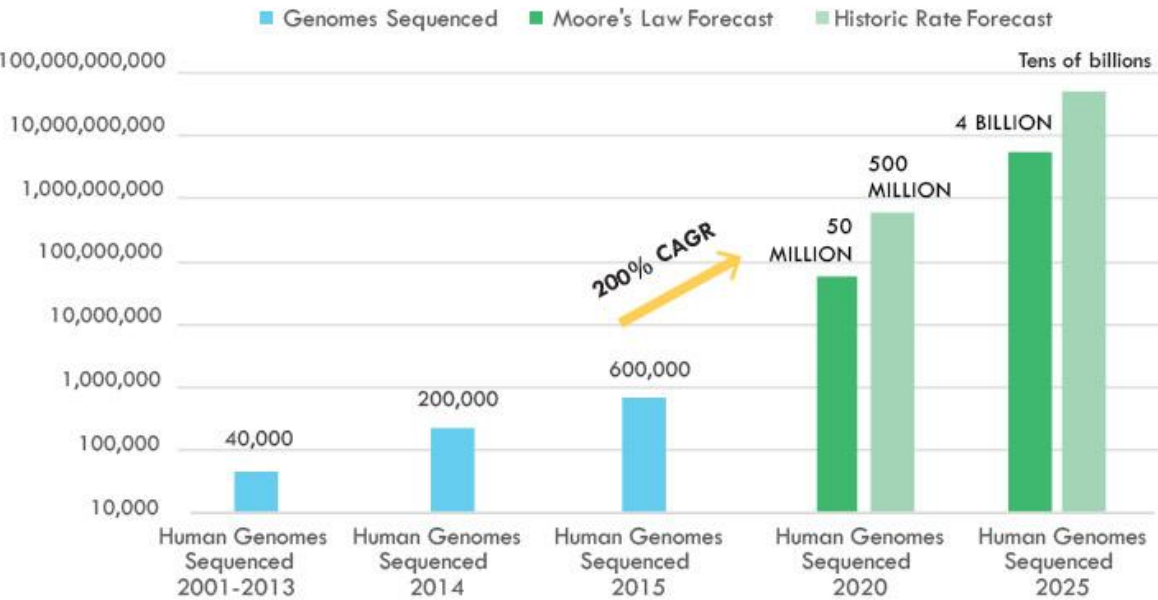
# Population sequencing to identify disease-associated genetic variants

### The Number of Human Genomes Sequenced (log scale)

■ Genomes Sequenced   ■ Moore's Law Forecast   ■ Historic Rate Forecast

Tens of billions

4 BILLION

500 MILLION

50 MILLION

200% CAGR

600,000

200,000

40,000

| Human Genomes Sequenced 2001-2013 | Human Genomes Sequenced 2014 | Human Genomes Sequenced 2015 | Human Genomes Sequenced 2020 | Human Genomes Sequenced 2025 |

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

Controls

GGCAATACGATATTAGCAAATAAACGATAGTATACAAATCGTATTAC
GGCAATACGCTATTAGCAAGTAAACGATAGGATACAAATCGTATTAC

GGCAATACGATATTAGCAAGTAAACGATAGTATACAAATCGTATTAC
GGCAATACGATATTAGCAAATAAACGATAGTATACAAATCGTATTAC

GGCAATACGATATTAGCAAGTAAACGATAGTATACAAATCGTATTAC
GGCAATACGCTATTAGCAAGTAAACGATAGGATACAAATCGTATTAC

GA II
**1.6 billion** bp per day
(2008)

GA IIx
**5 billion** bp per day
(2009)

HiSeq 2500
**60 billion** bp per day
(2012)

Oxford Nanopore technology

Images: www.illumina.com/systems
Numbers: www.politigenomics.com/next-generation-sequencing-informatics
Dates: Illumina press releases

## Millions of common and rare genetic variants found in human population

# Population sequencing to identify disease-associated genetic variants



The Number of Human Genomes Sequenced (log scale)

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

GA II
**1.6 billion** bp per day
(2008)

GA IIx
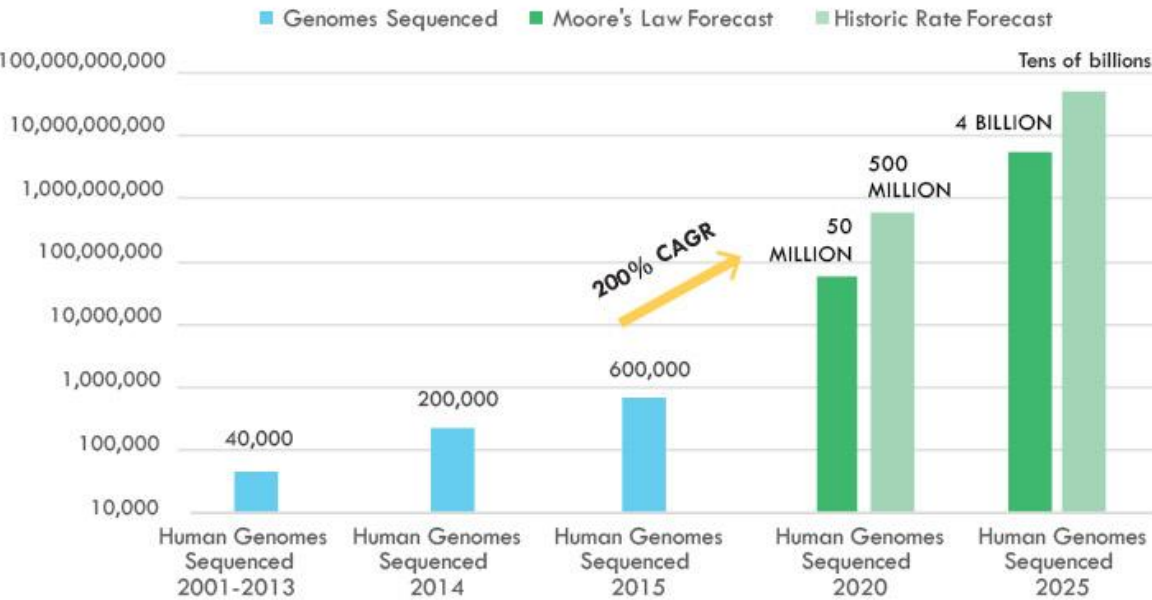**5 billion** bp per day
(2009)

HiSeq 2500
**60 billion** bp per day
(2012)

Oxford Nanopore technology

Images: www.illumina.com/systems
Numbers: www.politigenomics.com/next-generation-sequencing-informatics
Dates: Illumina press releases

Statistically significant association?

Millions of common and rare genetic variants found in human population

# Functional components of the human genome



DNA binding proteins

sequence motifs
(complex syntax)

Repressed gene

Active gene

Control elements

Protein

chromatin fiber

DNA

nucleus

nucleosome

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

# Molecular mapping of functional components of the genome

chromatin fiber

DNA

nucleus

nucleosome

Repressed gene

Active gene

Control elements

Protein

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

# Molecular mapping of functional components of the genome



chromatin fiber

DNA

nucleus

nucleosome

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

Repressed gene

Active gene

Protein

Control elements

# Molecular mapping of functional components of the genome



chromatin fiber

DNA

nucleus

nucleosome

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

Repressed gene

Active gene

Protein

Control elements

# Molecular mapping of functional components of the genome



chromatin fiber

DNA

nucleus

nucleosome

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

Repressed gene

Active gene

Protein

Control elements
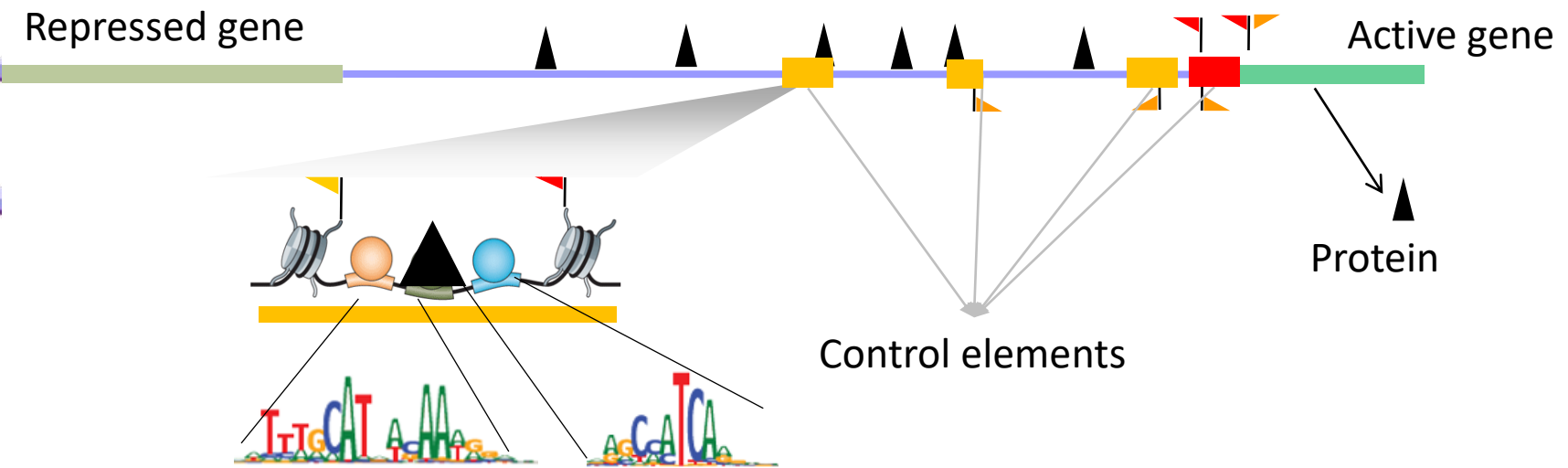
# Molecular mapping of functional components of the genome



chromatin fiber

DNA

nucleus

nucleosome

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

Repressed gene

Active gene

Protein

Control elements
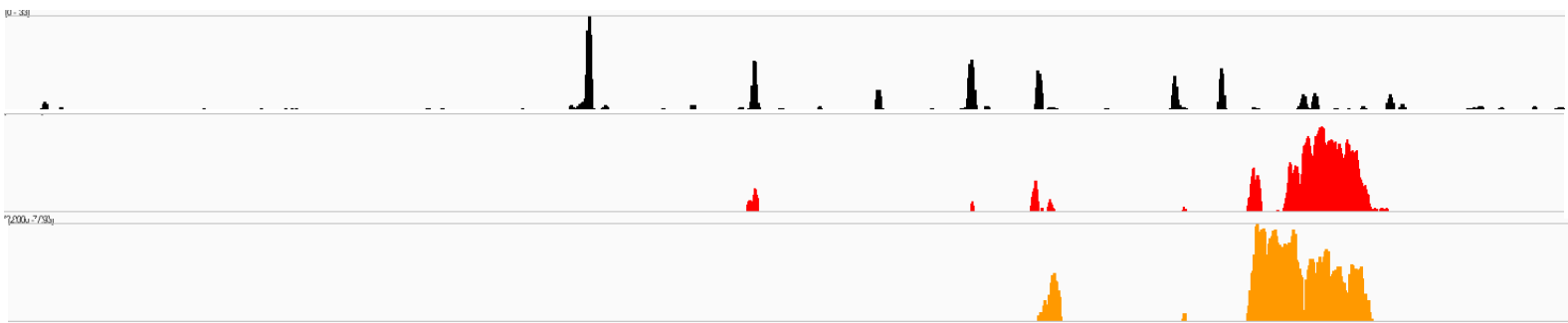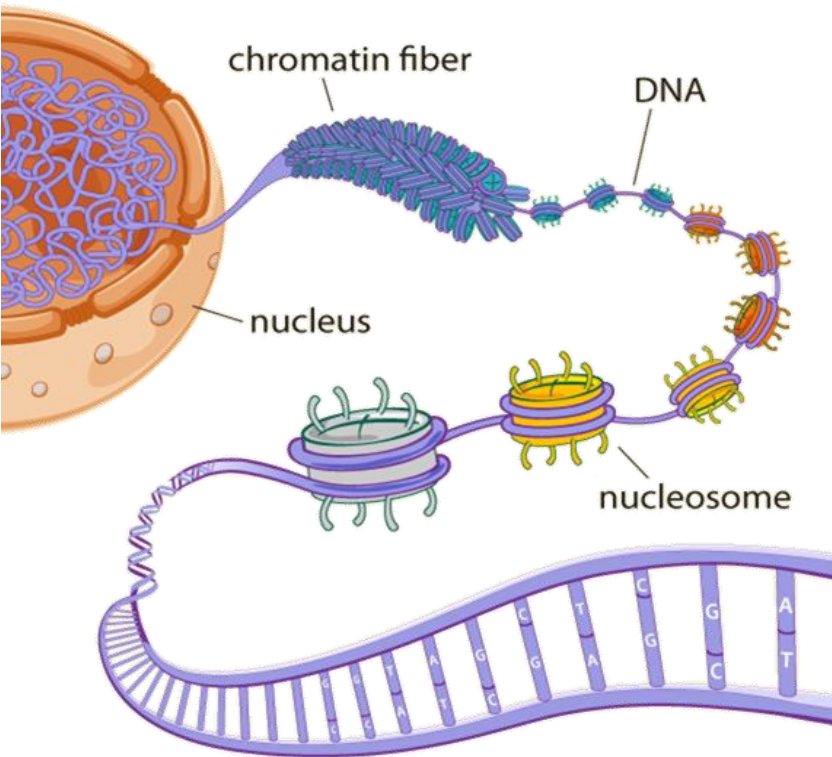
# Molecular mapping of functional components of the genome



chromatin fiber

DNA

nucleus

nucleosome

Repressed gene

Active gene

Protein

Control elements

https://www.broadinstitute.org/news/1504
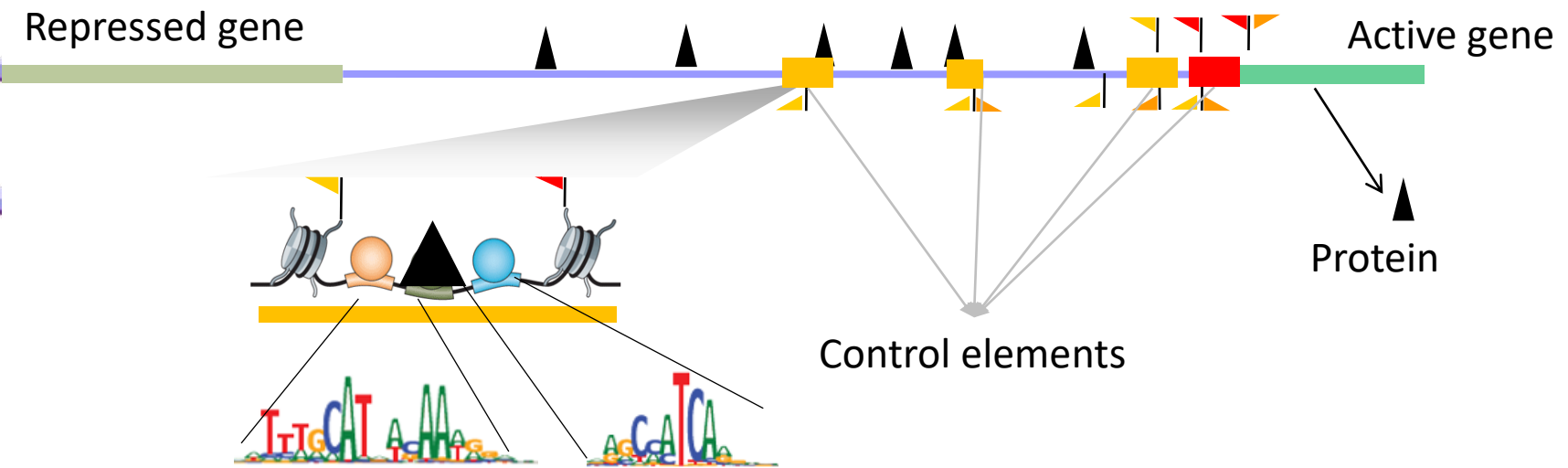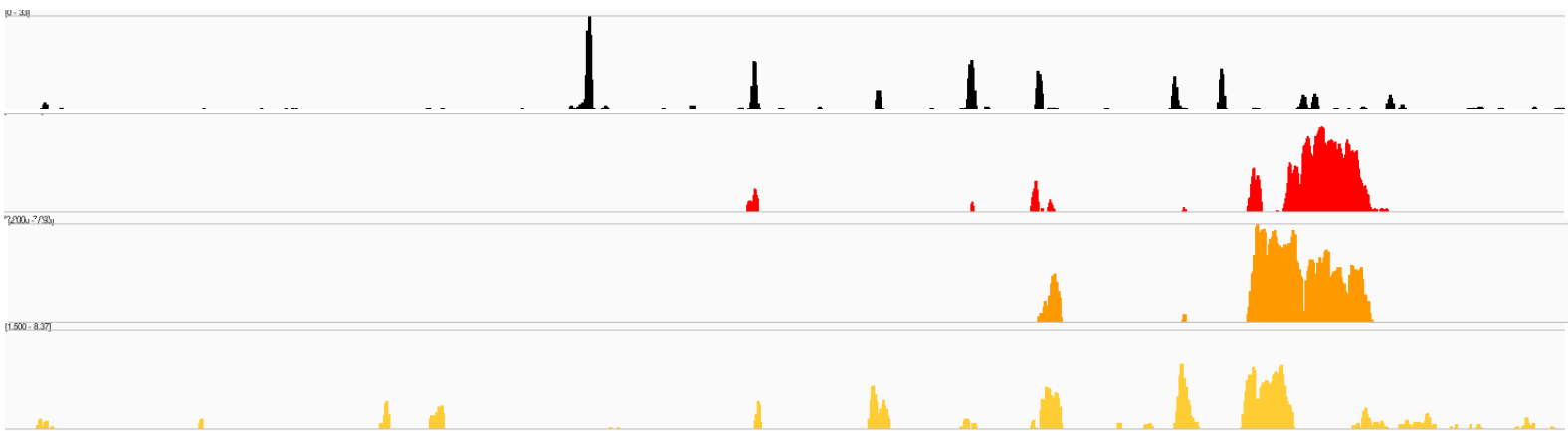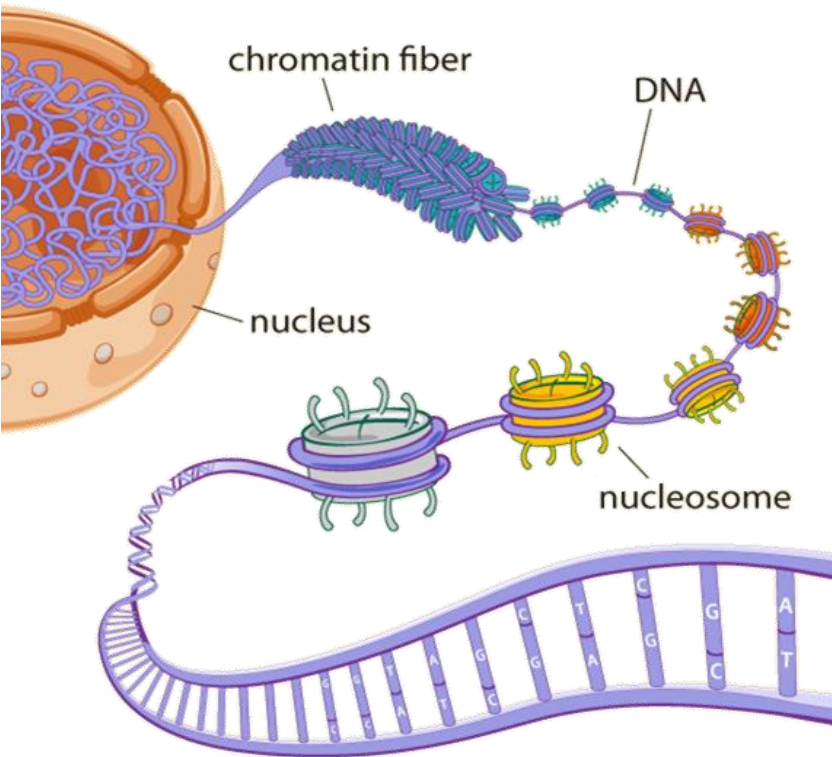
# Molecular mapping of functional components of the genome



Repressed gene

Active gene

Protein

Control elements

chromatin fiber
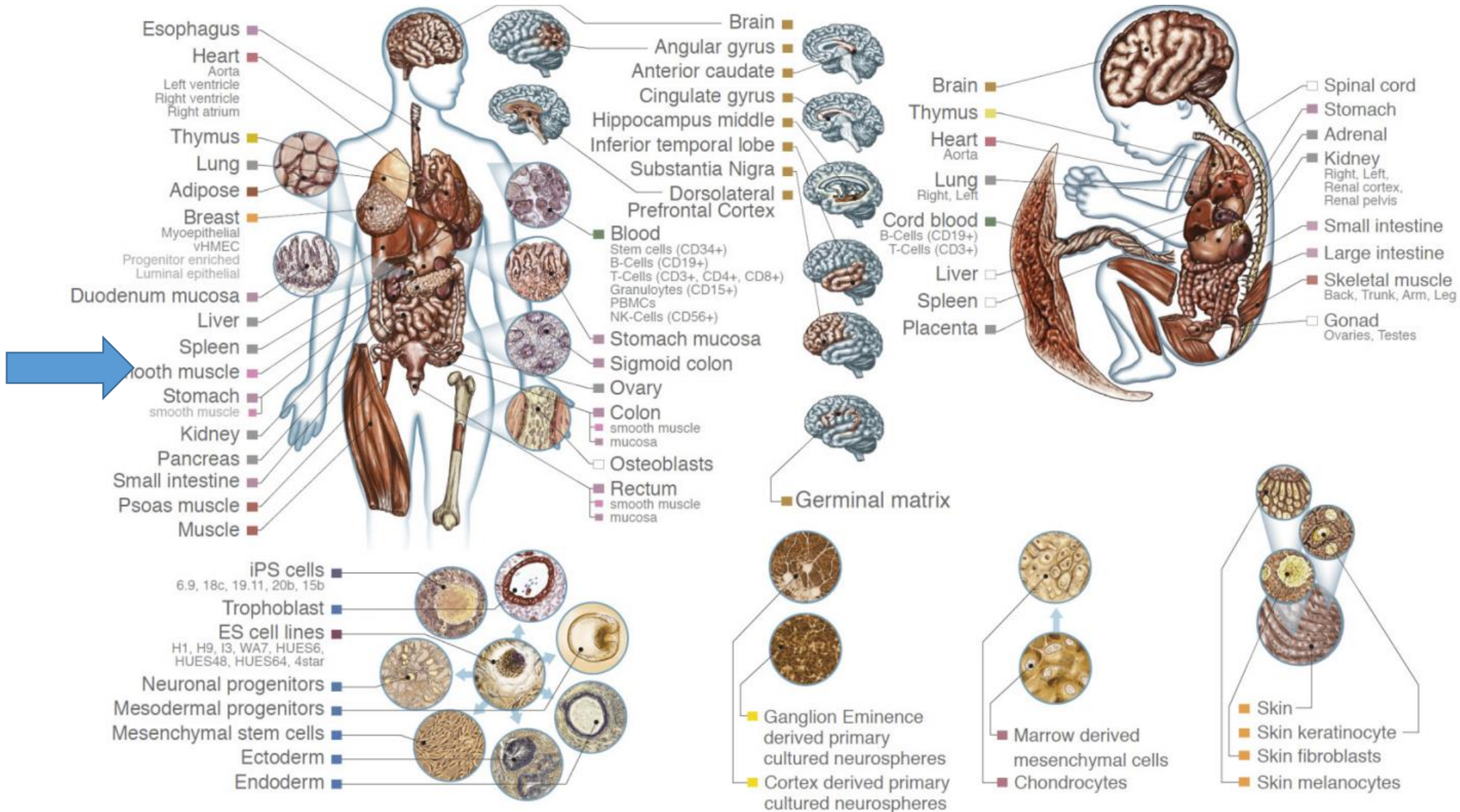
DNA

nucleus

nucleosome

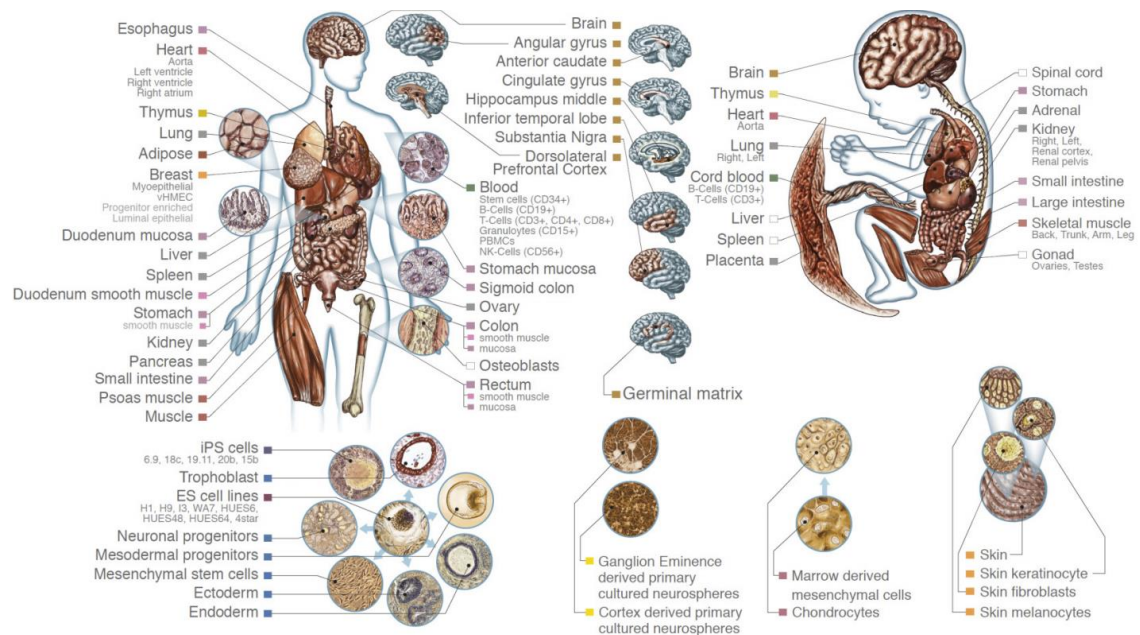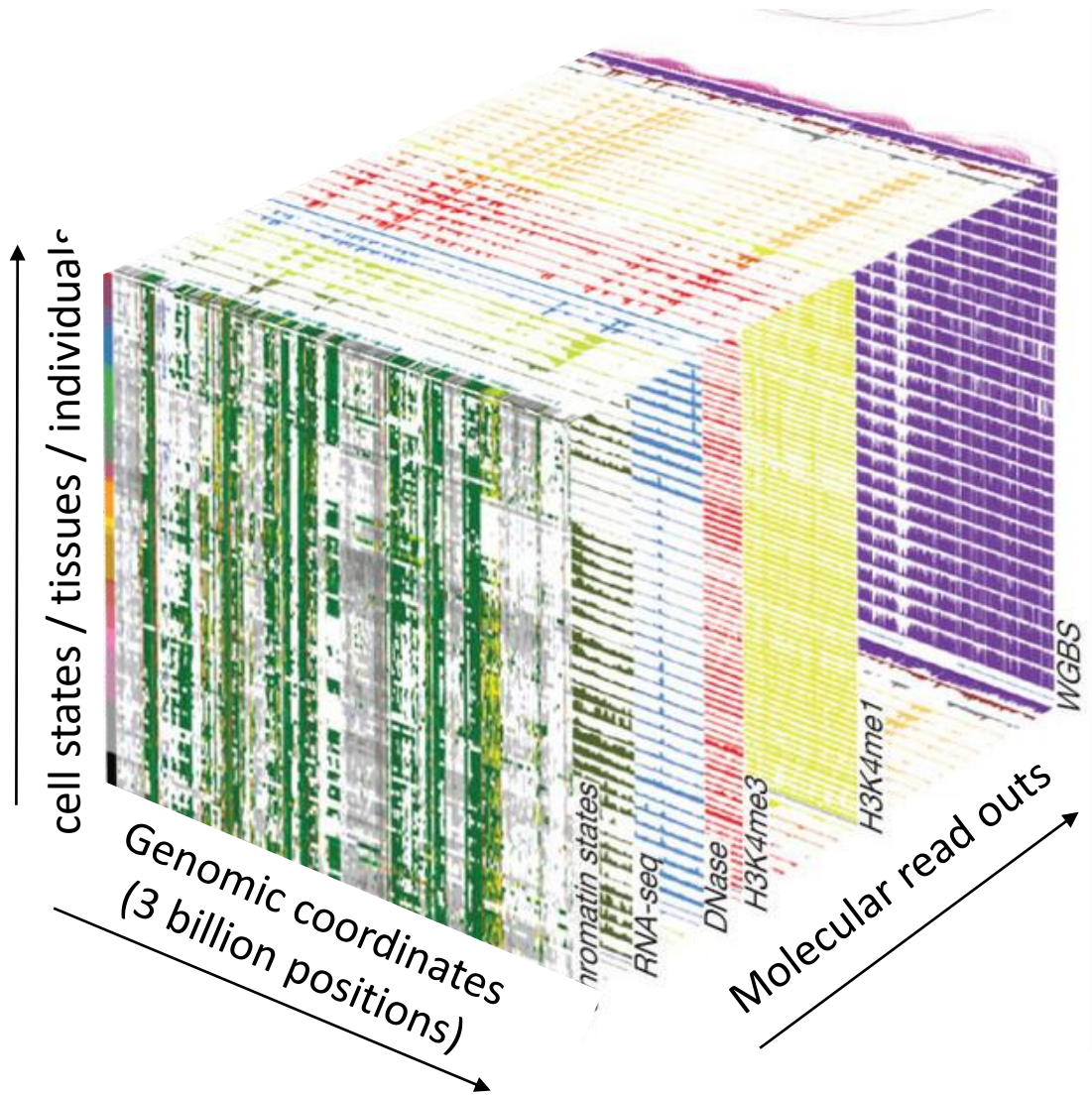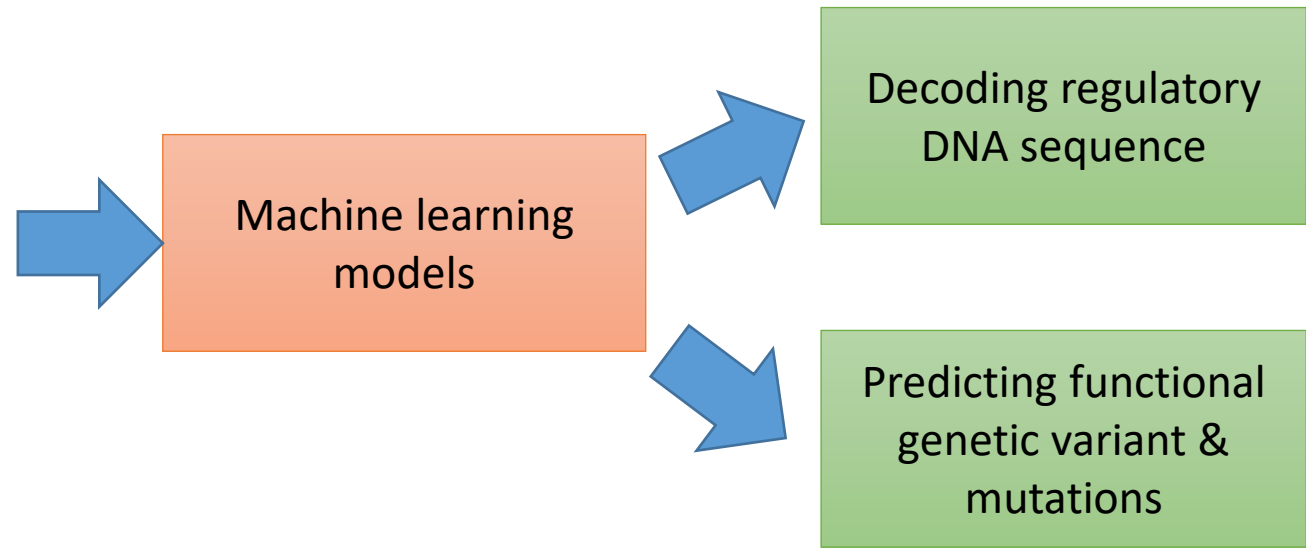https://www.broadinstitute.org/news/1504

# One genome ⇔ many cell types

ACCAGTTACGACGG
TCAGGGTACTGATA
CCCCAAACCGTTGA
CCGCATTTACAGAC
GGGGTTTGGGTTTT
GCCCCACACAGGTA
CGTTAGCTACTGGT
TTAGCAATTTACCG
TTACAACGTTTACA
GGGTTACGGTTGGG
ATTTGAAAAAAAGT
TTGAGTTGGTTTTT
TCACGGTAGAACGT
ACCTTACAAA...........

**100s of Cell-Types/Tissues**

cell states / tissues / individuals

Genomic coordinates (3 billion positions)

Molecular read outs

chromatin states
RNA-seq
DNase
H3K4me3
H3K4me1
WGBS

Machine learning models

Decoding regulatory DNA sequence

Predicting functional genetic variant & mutations

*Dunham, Kundaje et al. 2012 Nature*
*Kundaje et al. 2015 Nature*

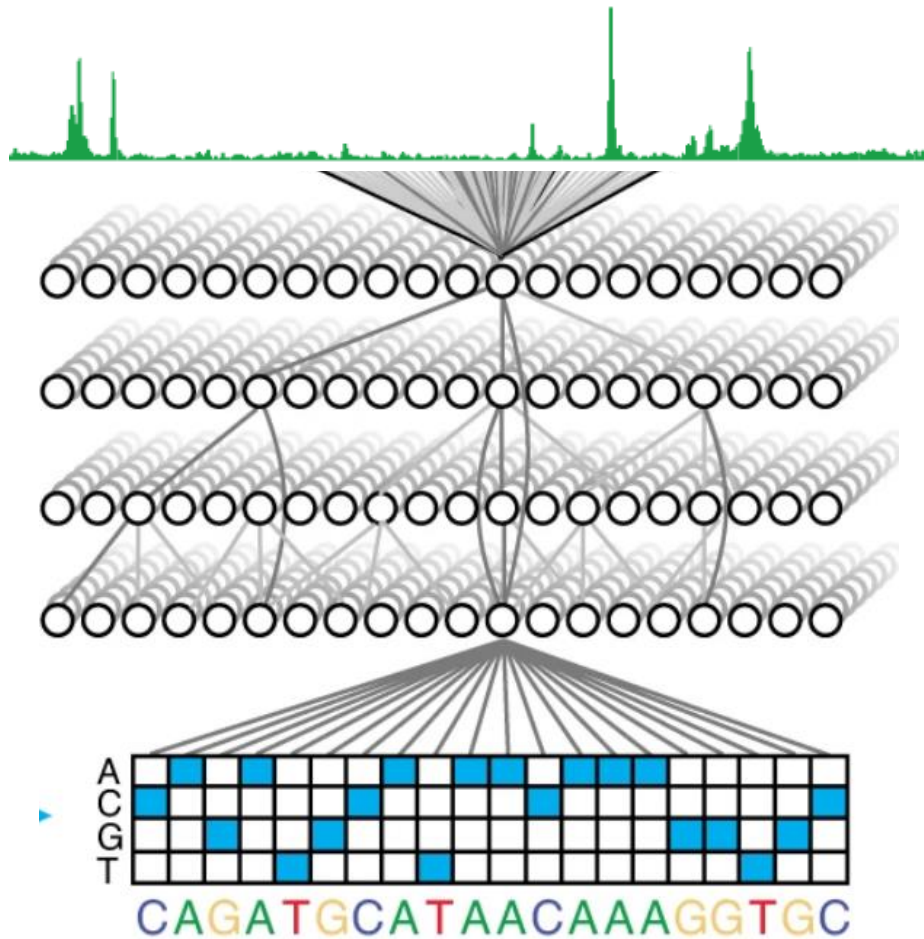# Deep learning framework for decoding regulatory DNA

Ziga Avsec

Anusri Pampari

Anna Shcherbina

Avanti Shrikumar

Alex Tseng

Surag Nair

Jacob Schreiber

**BPNet**
(maps sequence to base-resolution profiles)
One model for every expt.

*Avsec et al. 2021, Nature Genetics*
*Shrikumar et al. 2017, ICML*
*Tseng et al. 2020, NeurIPS*
*Nair et al, 2022, Bioinformatics*
*Schreiber et al. 2022, Biorxiv*

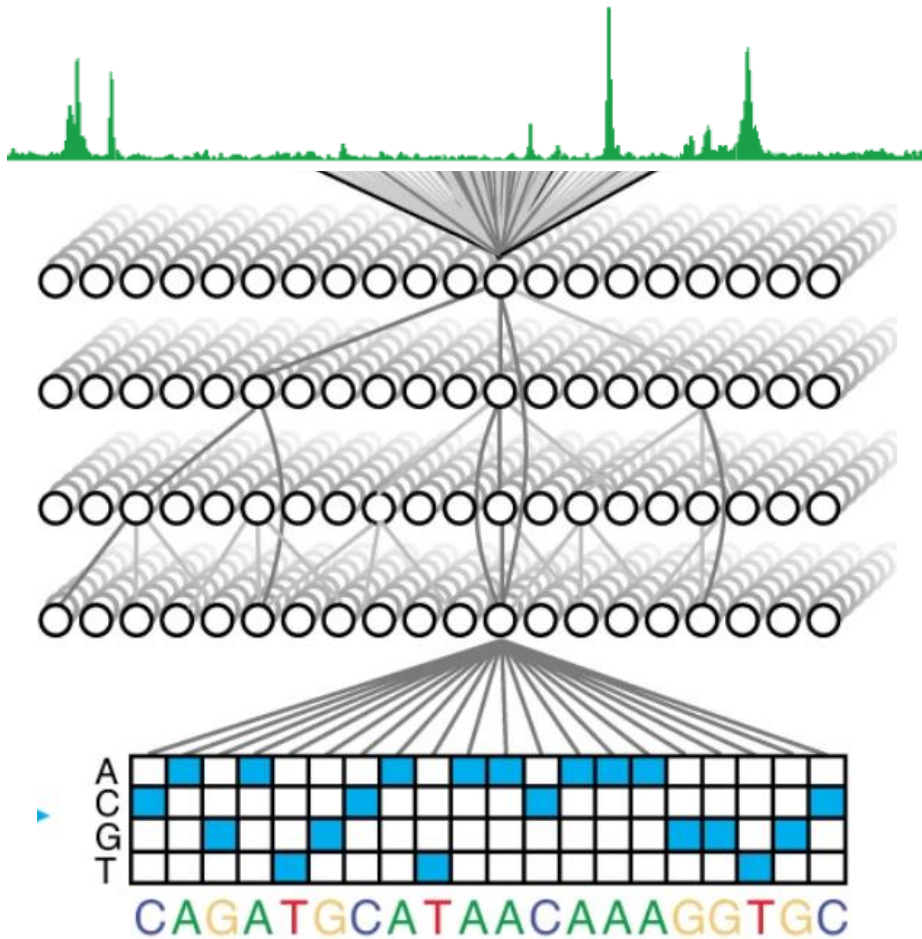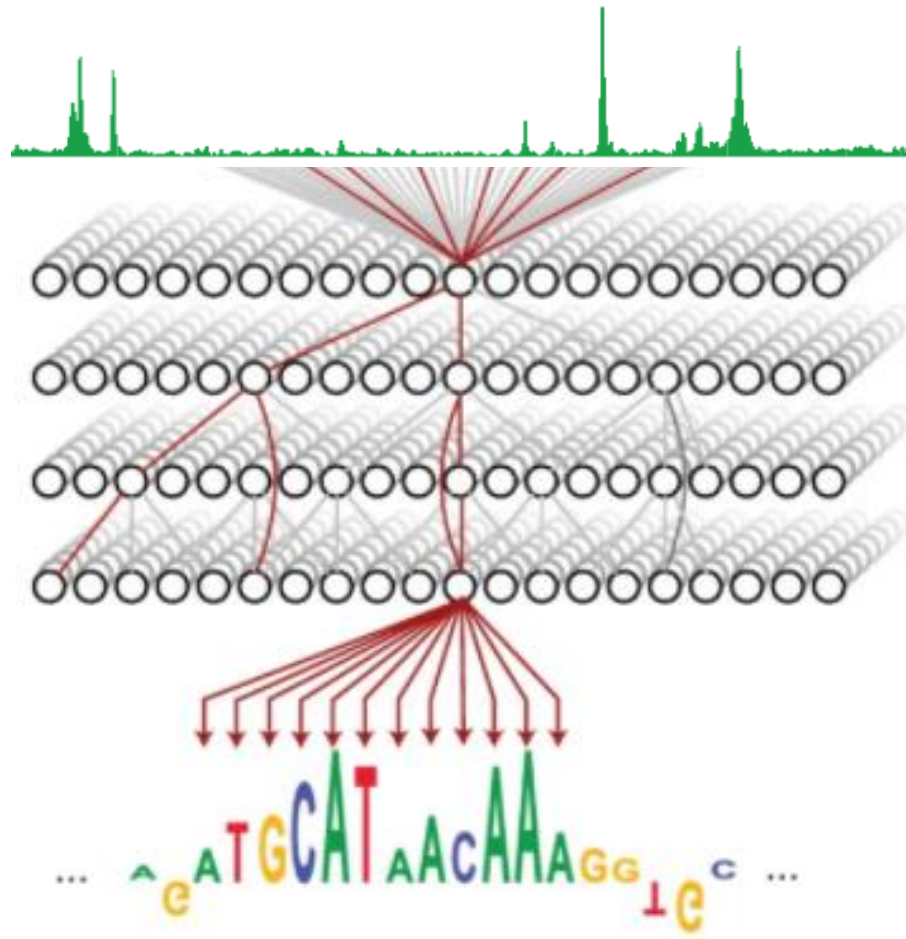# Deep learning framework for decoding regulatory DNA

Ziga Avsec

Anusri Pampari

Anna Shcherbina

Avanti Shrikumar

Alex Tseng

Surag Nair

Jacob Schreiber

*Avsec et al. 2021, Nature Genetics*
*Shrikumar et al. 2017, ICML*
*Tseng et al. 2020, NeurIPS*
*Nair et al, 2022, Bioinformatics*
*Schreiber et al. 2022, Biorxiv*

**BPNet**
(maps sequence to base-resolution profiles)
One model for every expt.

**DeepLIFT, FastISM, Yuzu, MoDISCo**
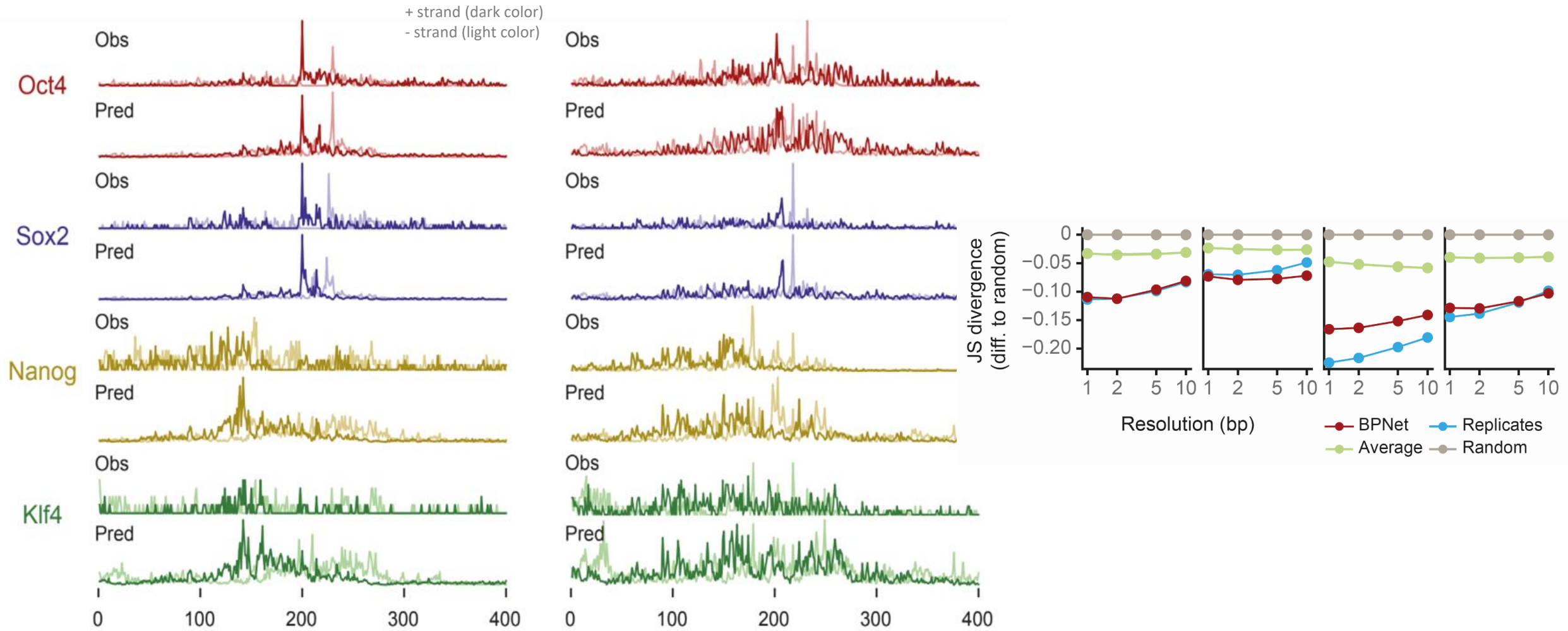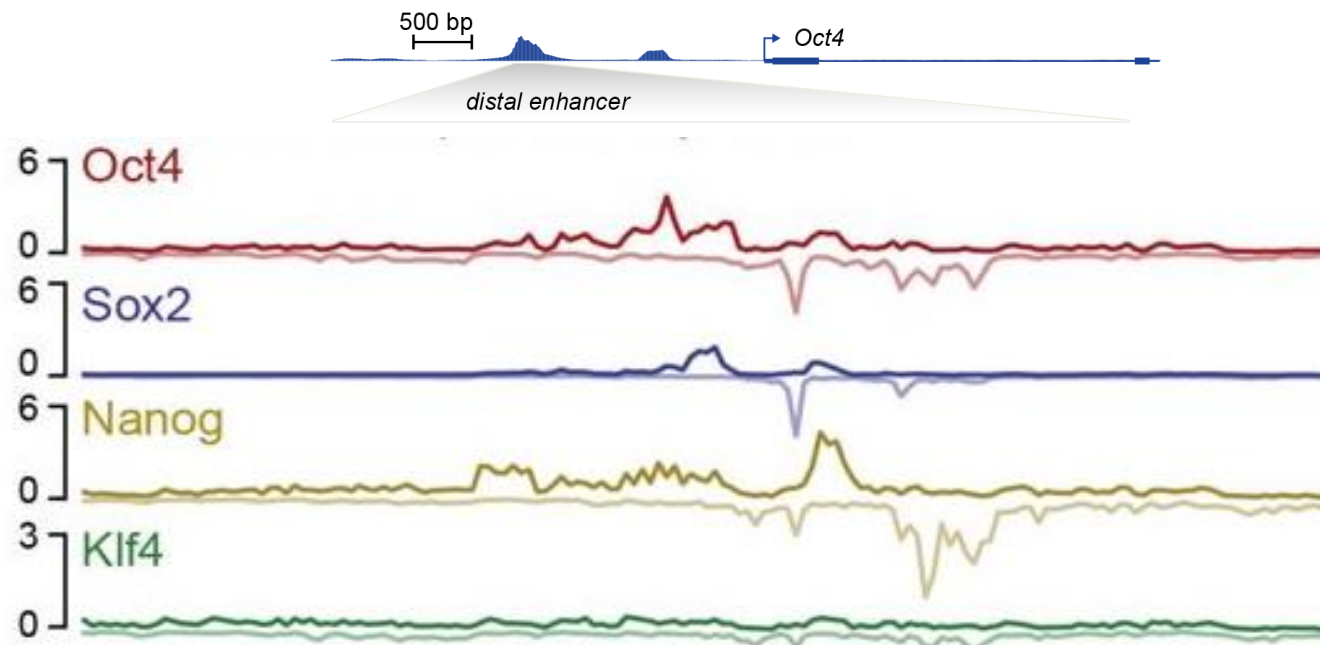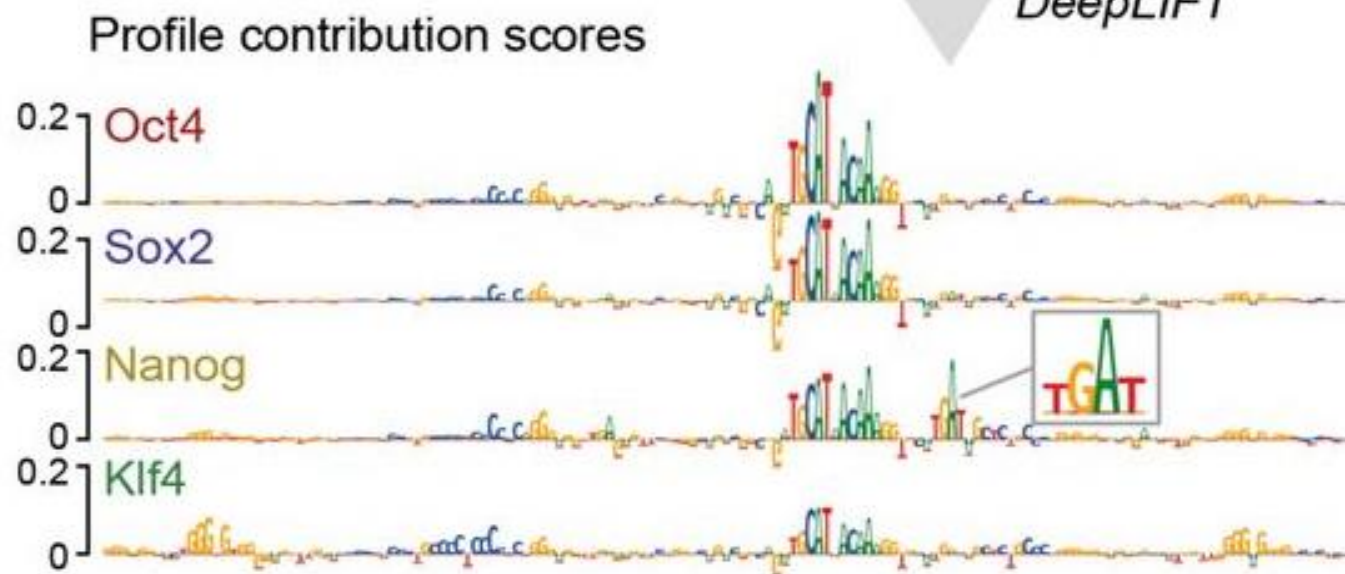(infers contribution of every base in each control sequence thru lens of model)

# BPNet maps DNA sequence to base-resolution molecular profiles with unprecedented accuracy
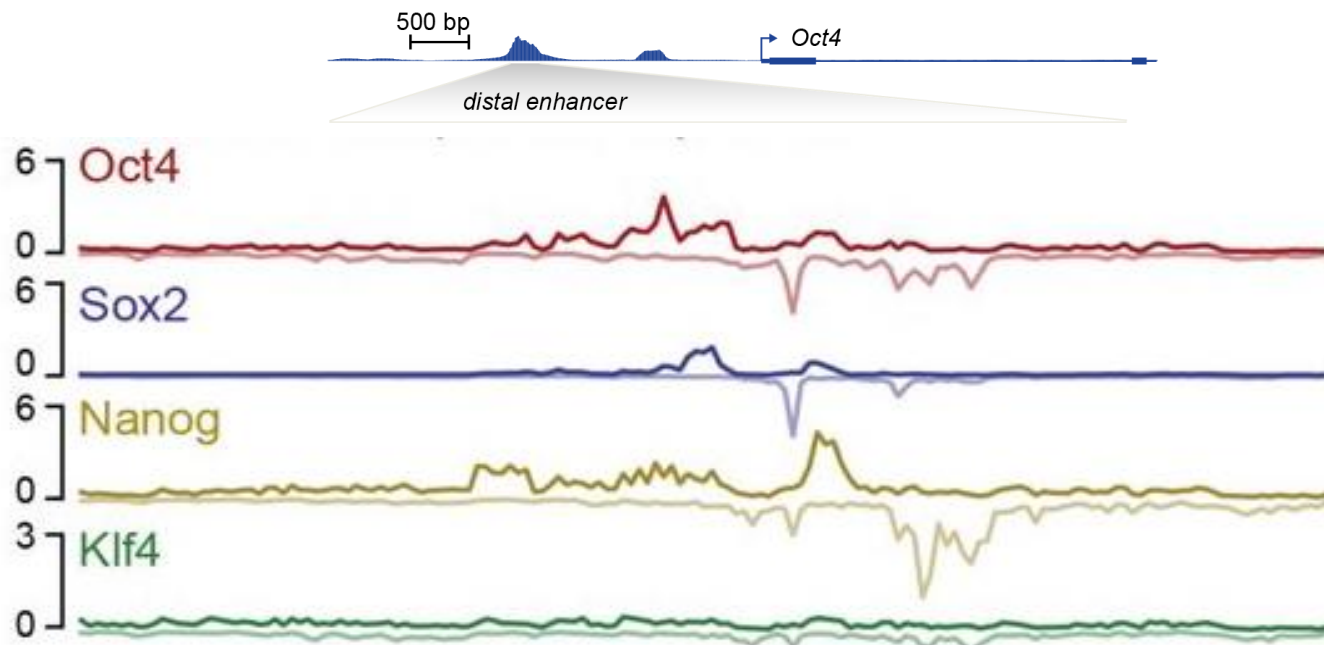# (on par with concordance between replicate experiments)



+ strand (dark color)
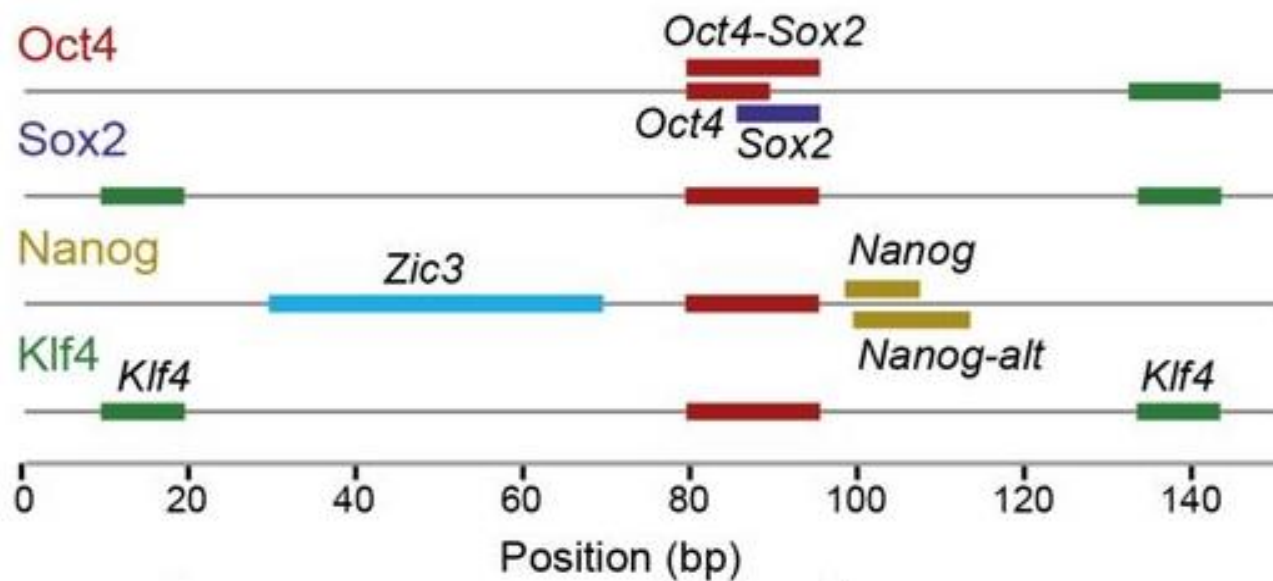- strand (light color)

500 bp | Oct4

distal enhancer

Oct4
Sox2
Nanog
Klf4

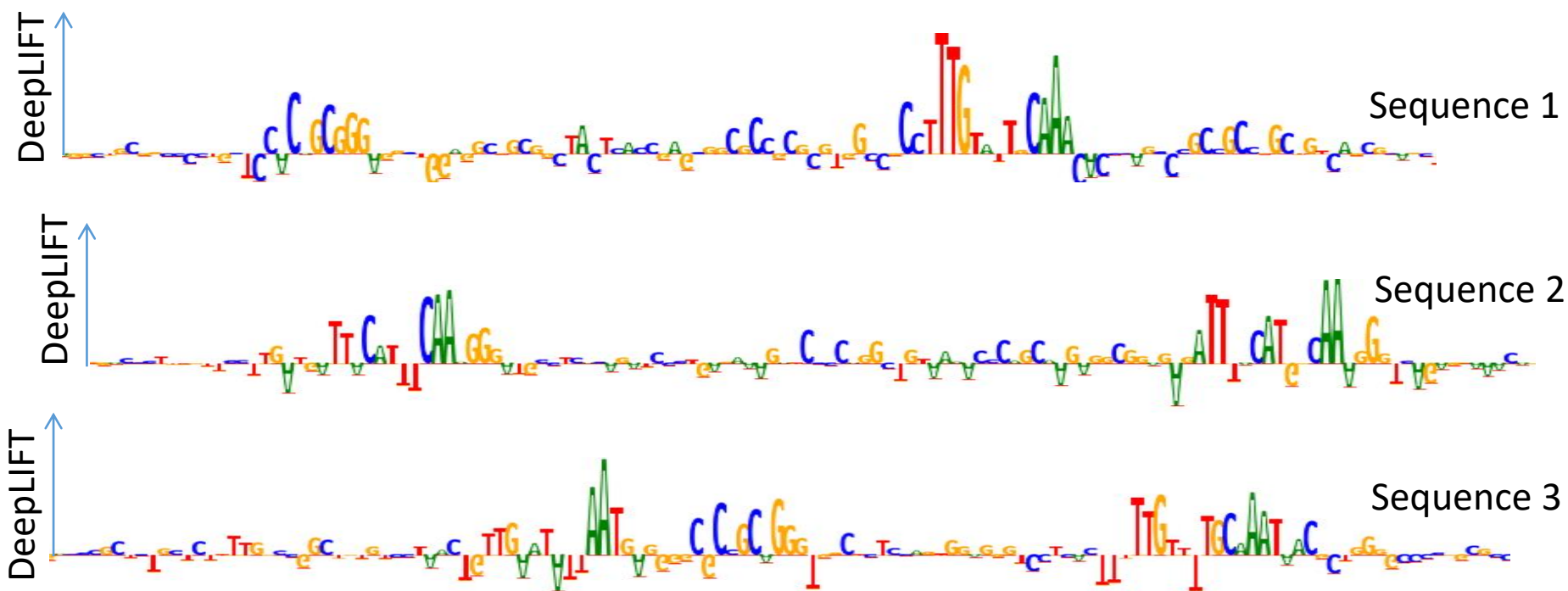DeepLIFT

Profile contribution scores

Oct4
Sox2
Nanog
Klf4

TGAT

Avanti Shrikumar

Alex Tseng

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

500 bp

Oct4

distal enhancer

Oct4

Sox2

Nanog

Klf4

DeepLIFT

Profile contribution scores

Oct4

Oct4-Sox2

Sox2

Oct4
Sox2

Nanog

Zic3

Nanog

Nanog-alt

Klf4

Klf4

Klf4

Position (bp)

Avanti Shrikumar

Alex Tseng

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

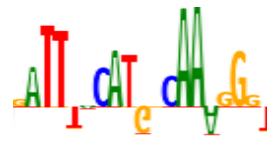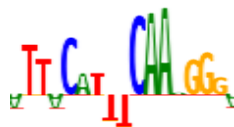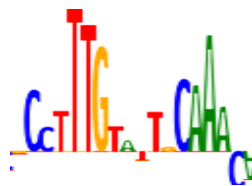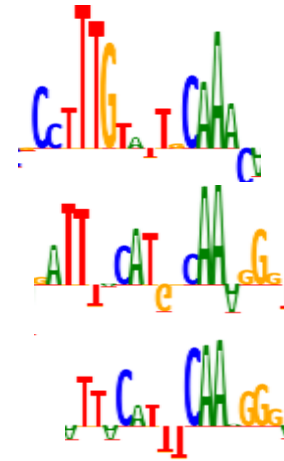# TF-MoDISCO: Consolidate predictive subsequences into non-redundant motif representations



Sequence 1

Sequence 2

Sequence 3

Avanti Shrikumar

Alex Tseng

*Shrikumar et al. 2018*
*Avsec et al. Nature Genetics 2021*

# TF-MoDISCO: Consolidate predictive subsequences into non-redundant motif representations



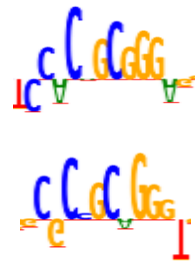Avanti Shrikumar

Alex Tseng

*Shrikumar et al. 2018*
*Avsec et al. Nature Genetics 2021*
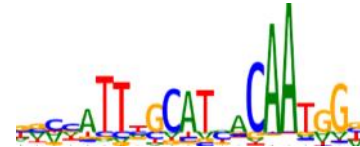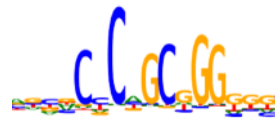
# TF-MoDISCO: Consolidate predictive subsequences into non-redundant motif representations



Avanti Shrikumar

Alex Tseng

Shrikumar et al. 2018
Avsec et al. Nature Genetics 2021

# TF-MoDISCO: Consolidate predictive subsequences into non-redundant motif representations



Avanti Shrikumar

Alex Tseng

*Shrikumar et al. 2018*
*Avsec et al. Nature Genetics 2021*
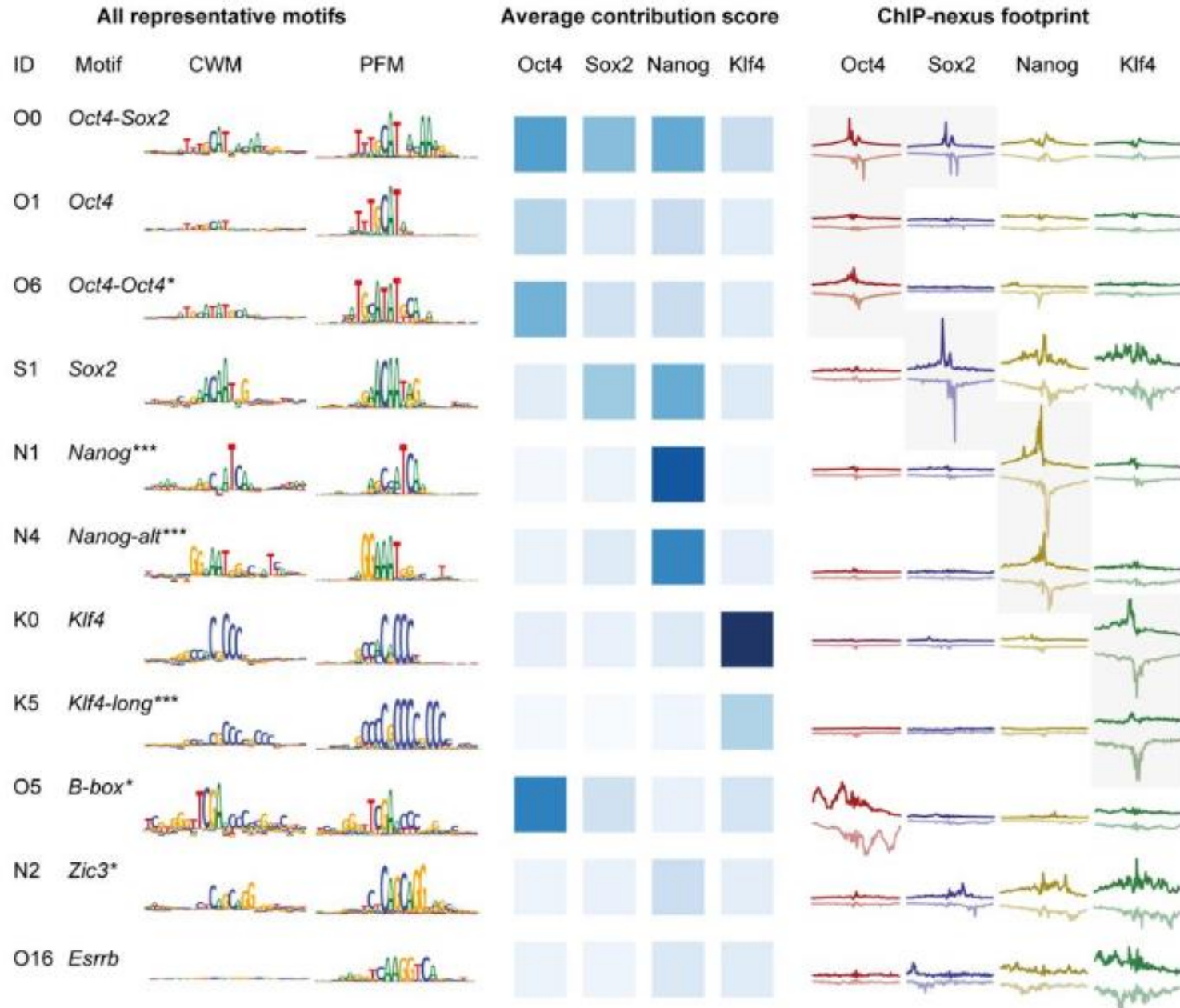
# Complex repertoire of motifs due to cooperative binding



50 motifs for 4 proteins!

# Syntax discovery using *in-silico* perturbations

Use BPNet model as in-silico oracle to perform perturbation experiments



1)    On synthetic sequences

# Syntax discovery using *in-silico* perturbations

Use BPNet model as in-silico oracle to perform perturbation experiments



1) On synthetic sequences

# Syntax discovery using *in-silico* perturbations

Use BPNet model as in-silico oracle to perform perturbation experiments

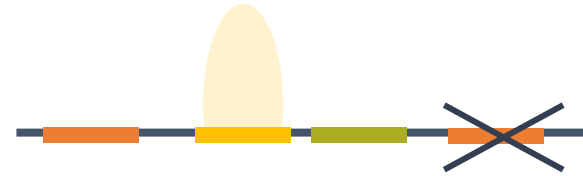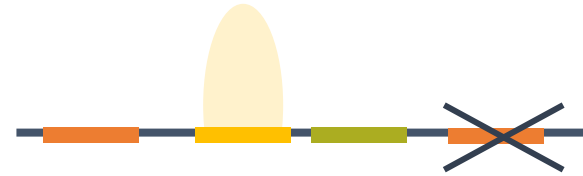1) On synthetic sequences

2) By mutating genomic sequences

# Syntax discovery using *in-silico* perturbations

Use BPNet model as in-silico oracle to perform perturbation experiments
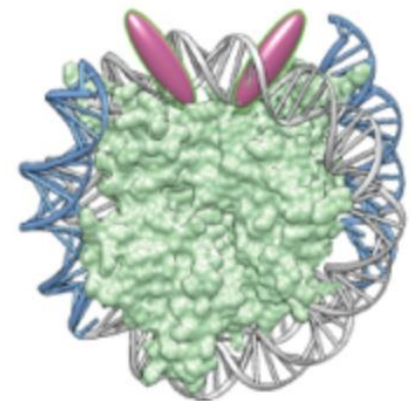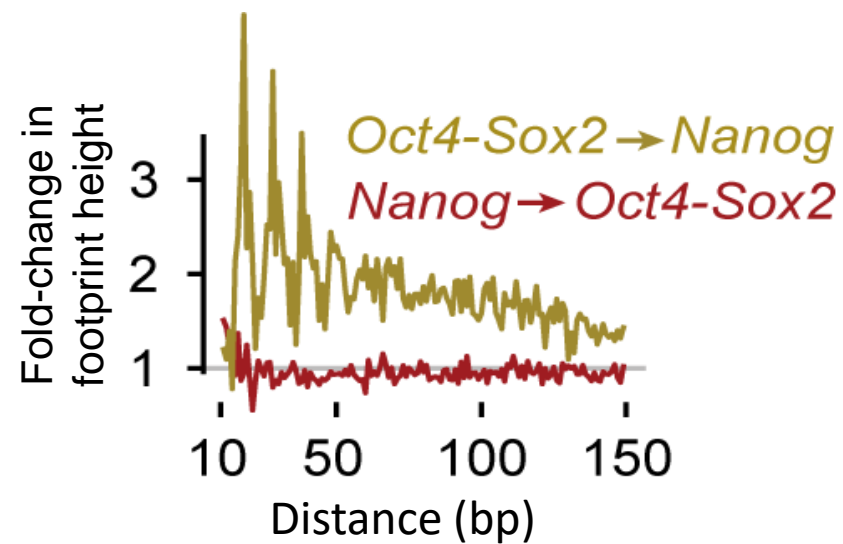
1) On synthetic sequences

2) By mutating genomic sequences

# Syntax discovery using *in-silico* perturbations

Use BPNet model as in-silico oracle to perform perturbation experiments

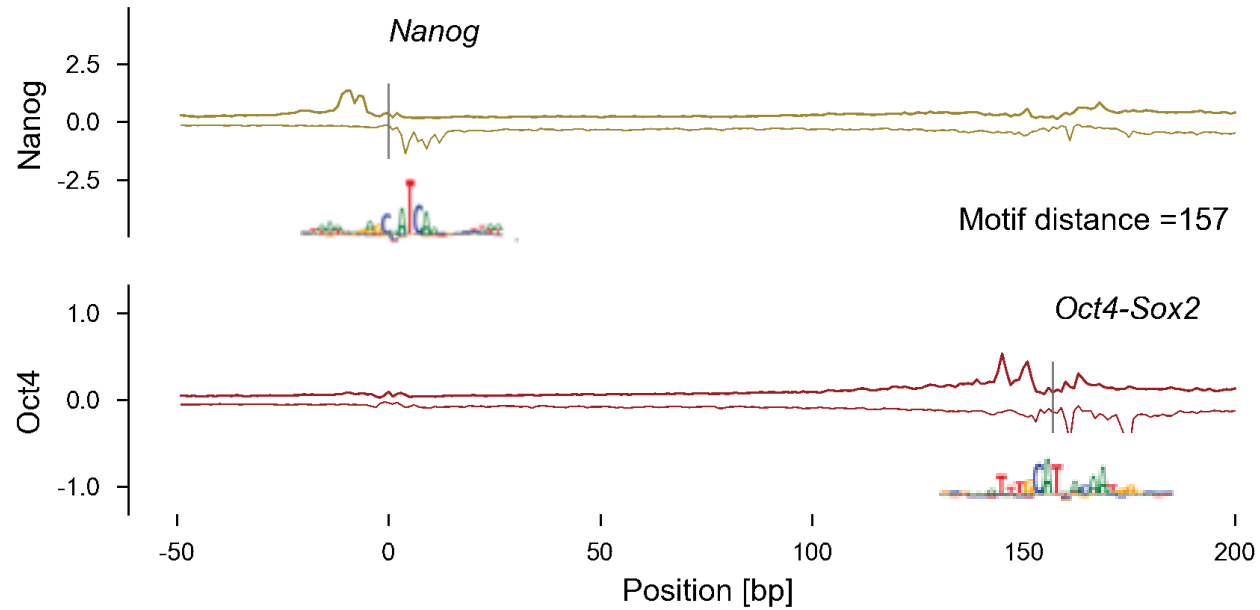1) On synthetic sequences

2) By mutating genomic sequences

# Syntax discovery using *in-silico* perturbations

Use BPNet model as in-silico oracle to perform perturbation experiments
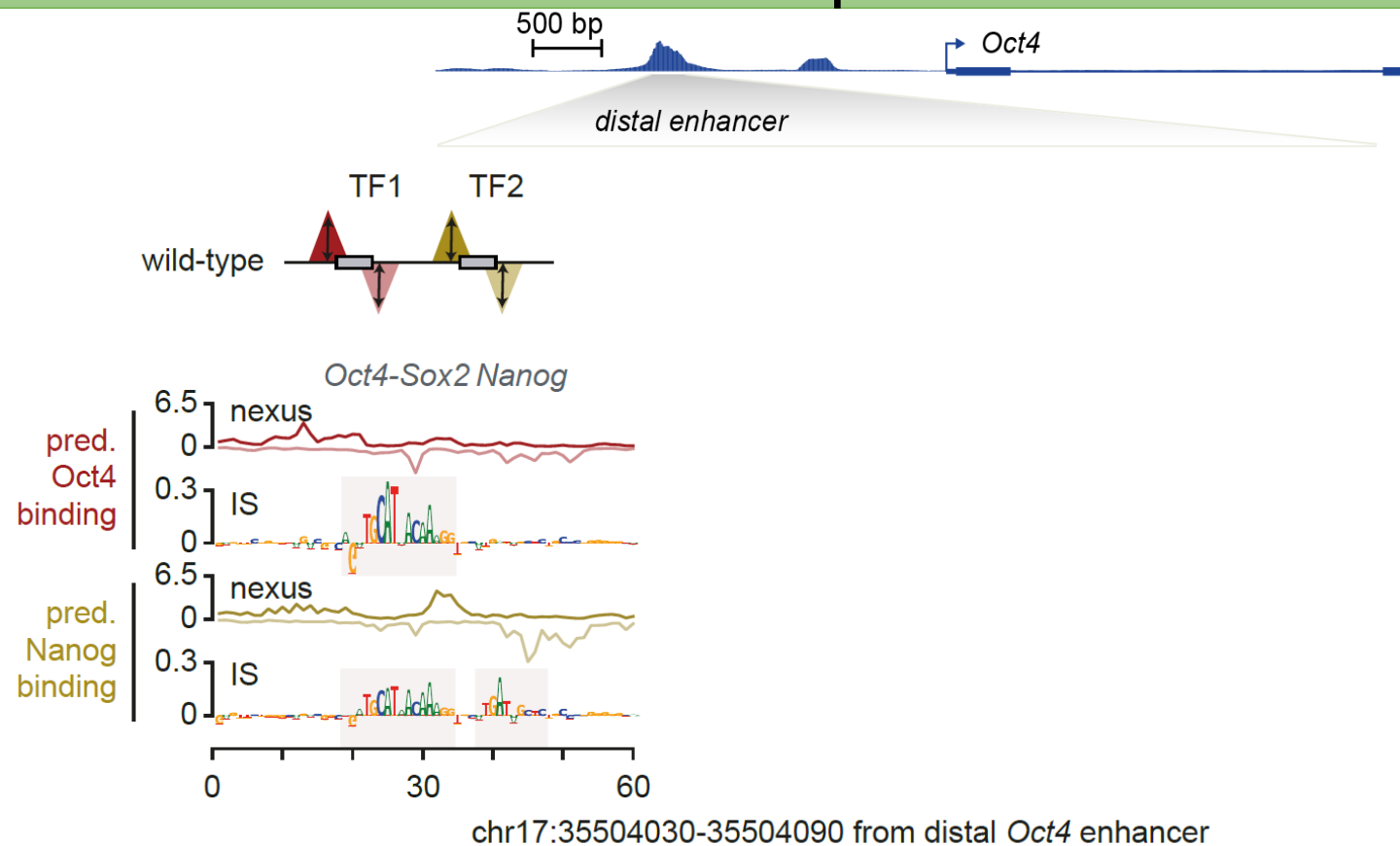
1) On synthetic sequences

*In silico* biochemistry
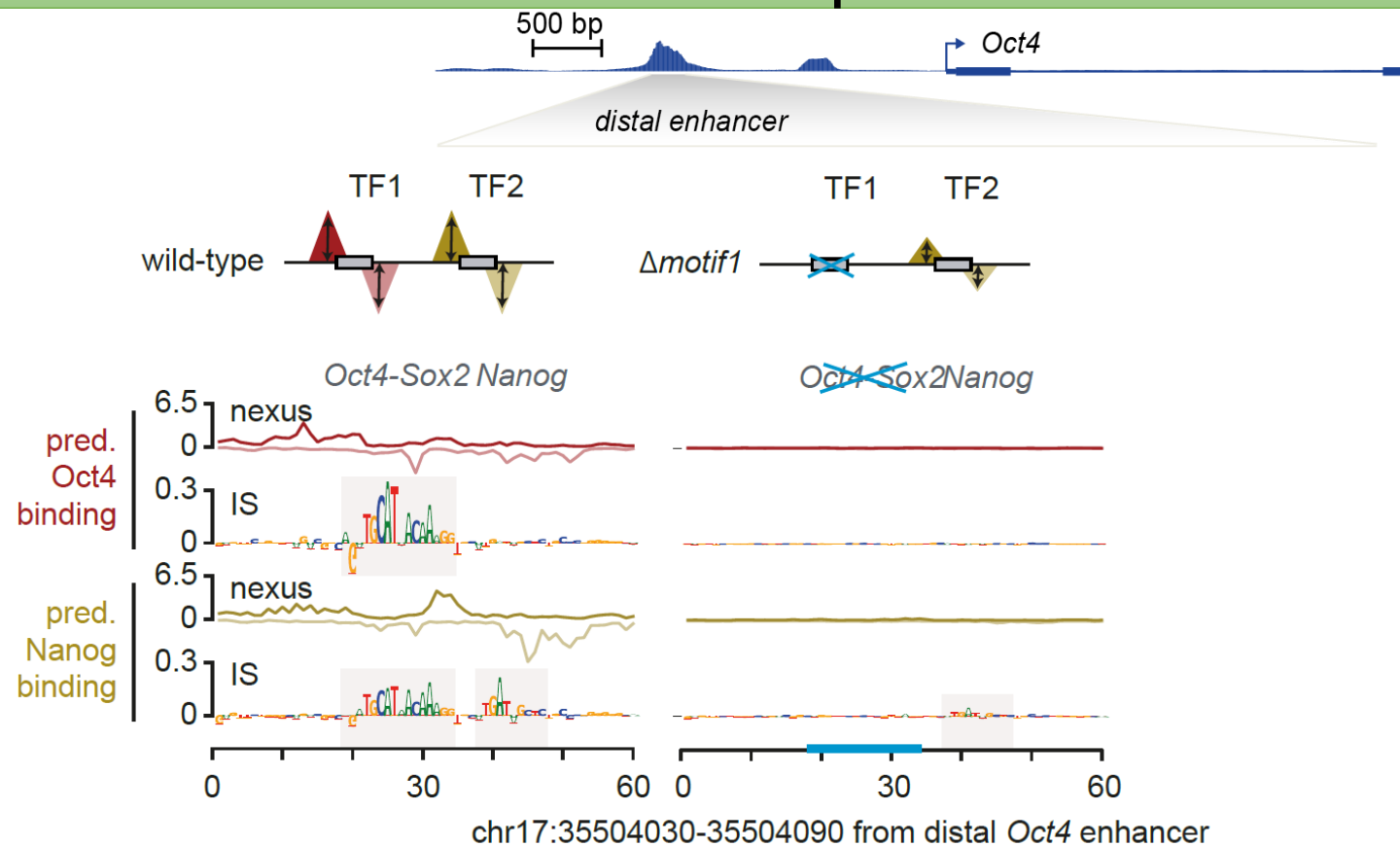
2) By mutating genomic sequences

*In silico* genetics

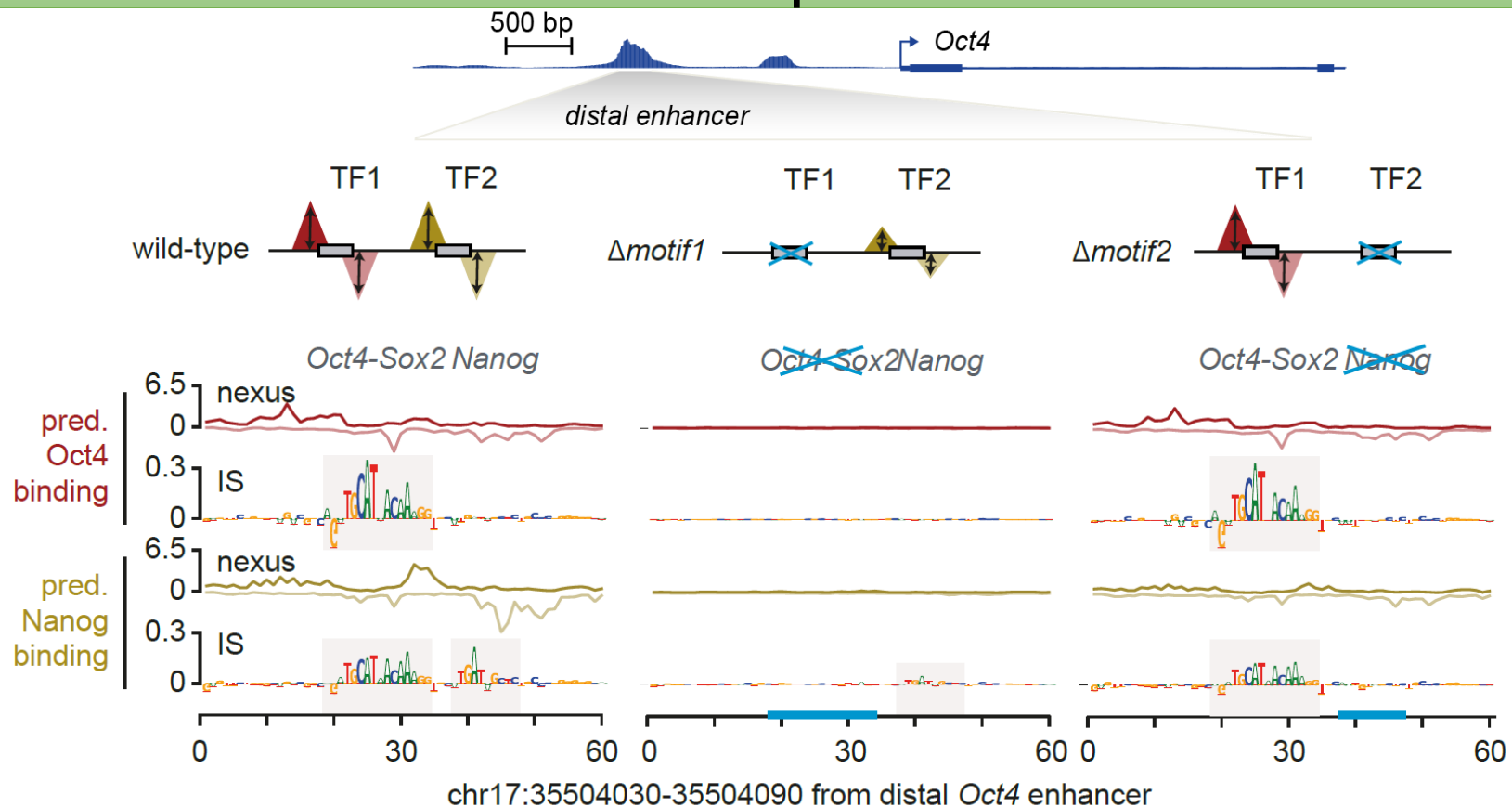# *In-silico* reporters: Designing synthetic sequences to query models to reveal syntax

chr17:35504030-35504090 from distal *Oct4* enhancer

chr17:35504030-35504090 from distal *Oct4* enhancer

# *in-silico* genome editing: Deciphering syntax by perturbing genomic sequences



chr17:35504030-35504090 from distal *Oct4* enhancer

chr17:35504030-35504090 from distal *Oct4* enhancer

# Designing CRISPR experiments to validate motif syntax



Sox2 ChIP-nexus

Predicted

Sox2 motif CCT**TT**GTTCC

Observed

Genomic position (bp)

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

# Designing CRISPR experiments to validate motif syntax



Sox2 ChIP-nexus

Predicted

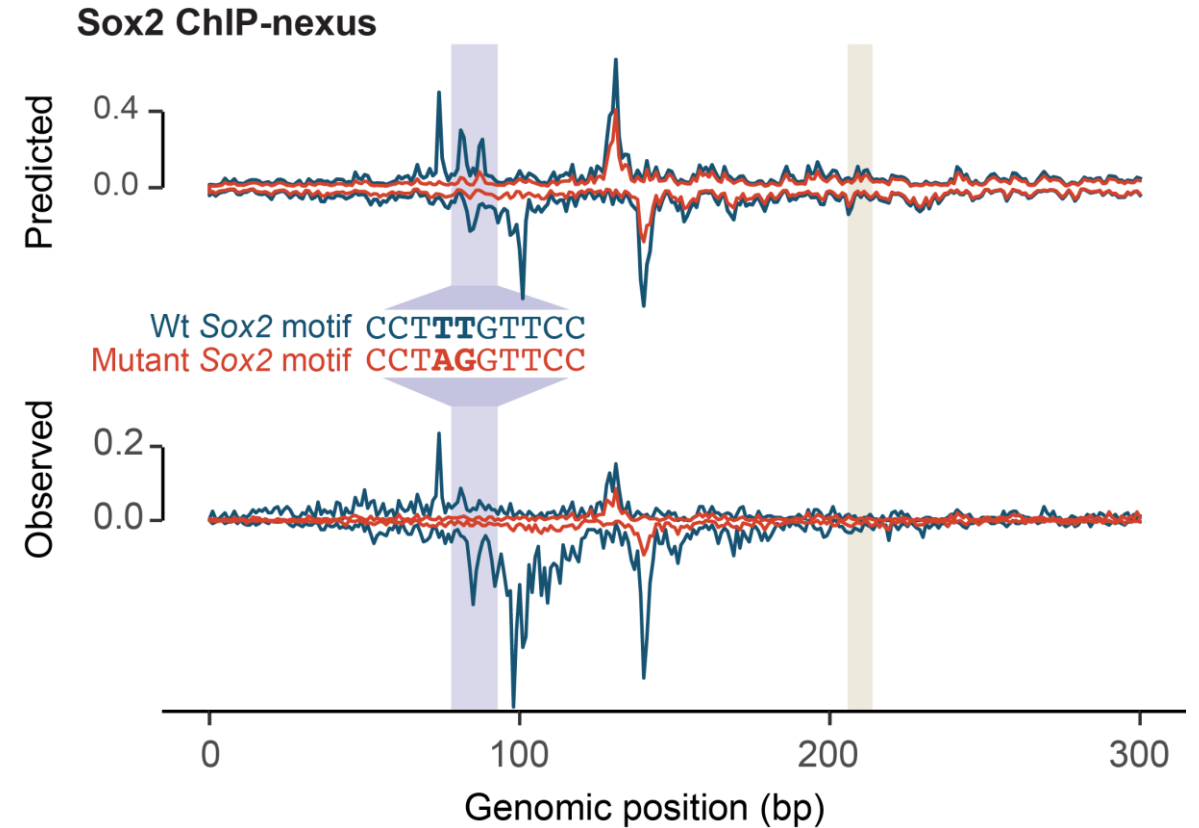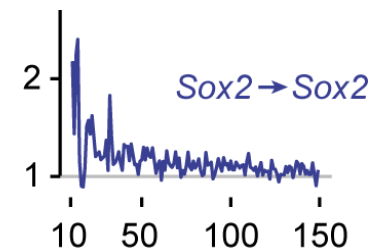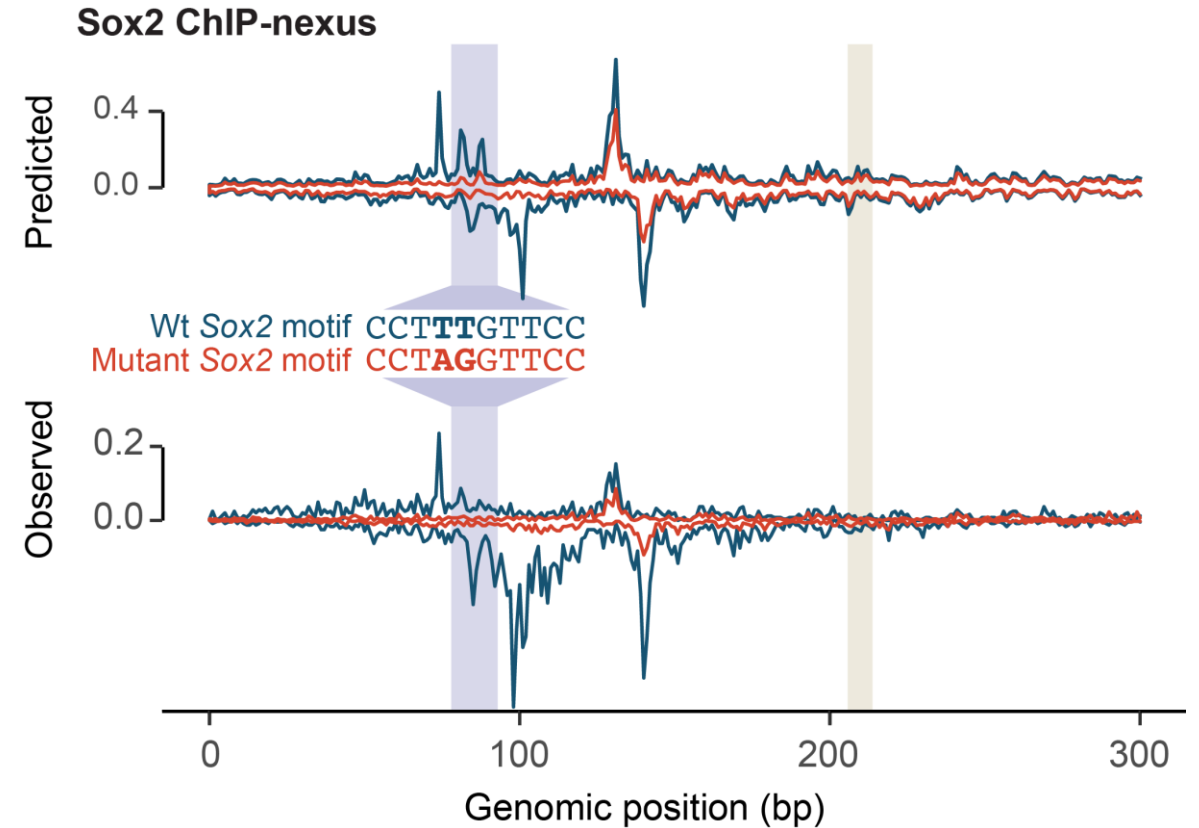Wt *Sox2* motif CCT**TT**GTTCC
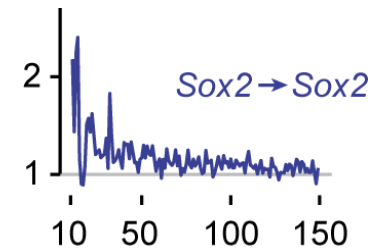Mutant *Sox2* motif CCT**AG**GTTCC

Observed

Genomic position (bp)

Julia Zeitlinger, Sabrina Krueger, Melanie Weilert

Sox2 ChIP-nexus

Julia Zeitlinger, Sabrina Krueger, Melanie Weilert

Sox2 ChIP-nexus

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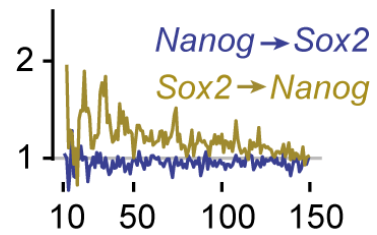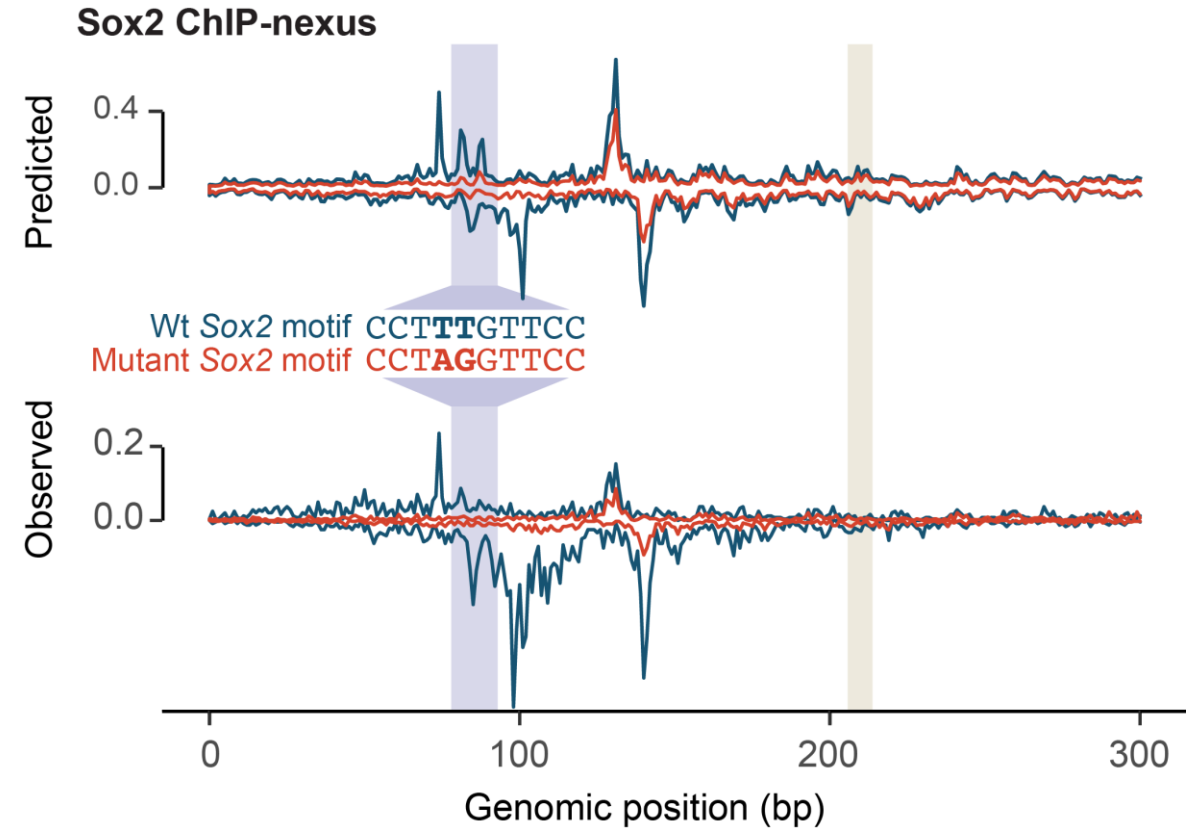Mutant *Sox2* motif CCT**AG**GTTCC

Nanog→Sox2
Sox2→Nanog

Sox2→Sox2

Julia Zeitlinger, Sabrina Krueger, Melanie Weilert

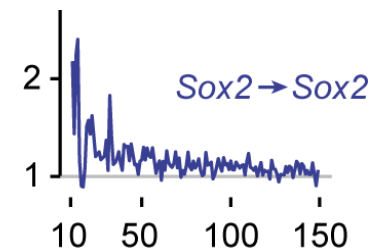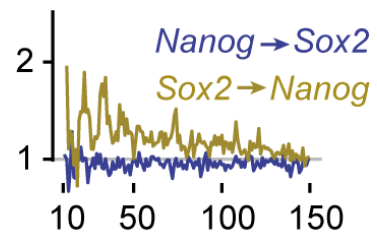# Designing CRISPR experiments to validate motif syntax



**Nanog ChIP-nexus**

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

Nanog → Sox2
Sox2 → Nanog

**Sox2 ChIP-nexus**

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

Sox2 → Sox2

Julia Zeitlinger, Sabrina Krueger, Melanie Weilert

**Sox2 ChIP-nexus**

Predicted

0.4

0.0

*Sox2* motif  CCT**TT**GTTCC

Observed

0.2

0.0

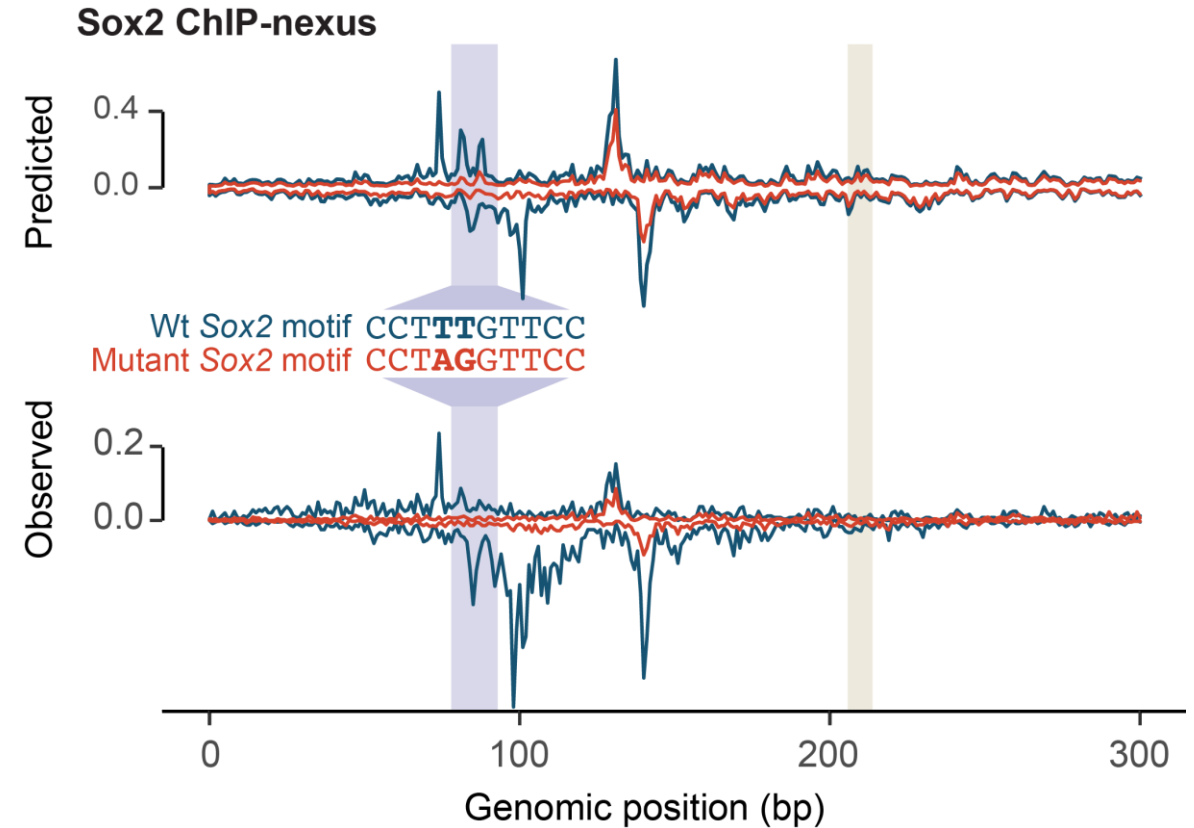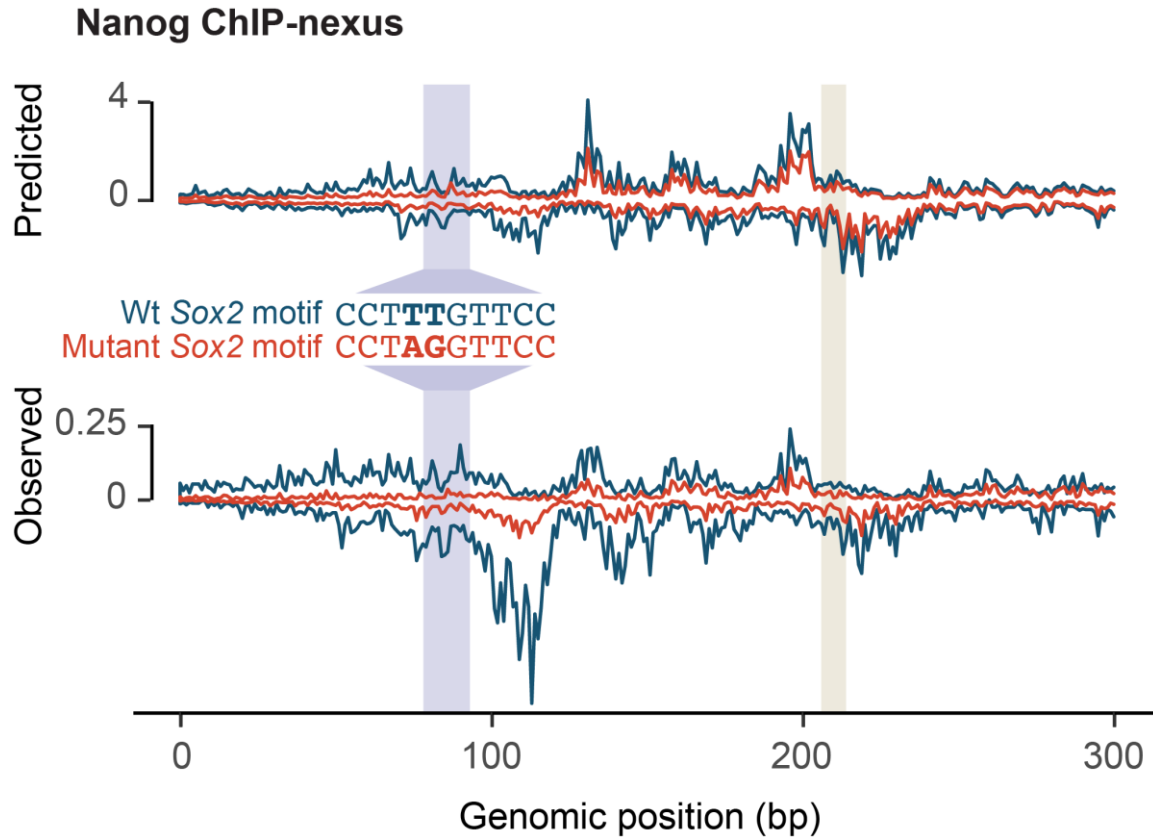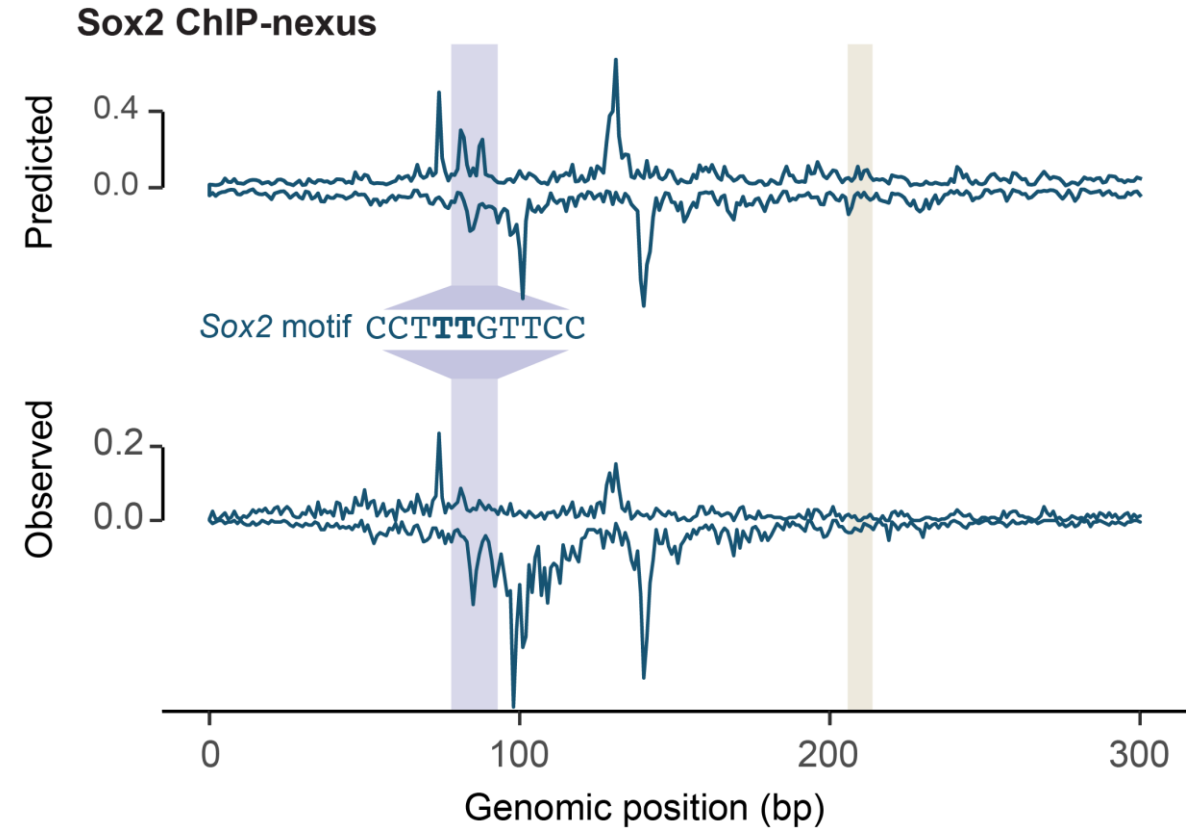0 100 200 300
Genomic position (bp)

Julia Zeitlinger, Sabrina Krueger, Melanie Weilert

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

Sox2 ChIP-nexus

Julia Zeitlinger, Sabrina Krueger, Melanie Weilert

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

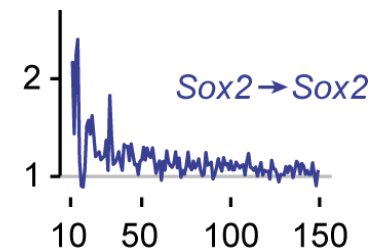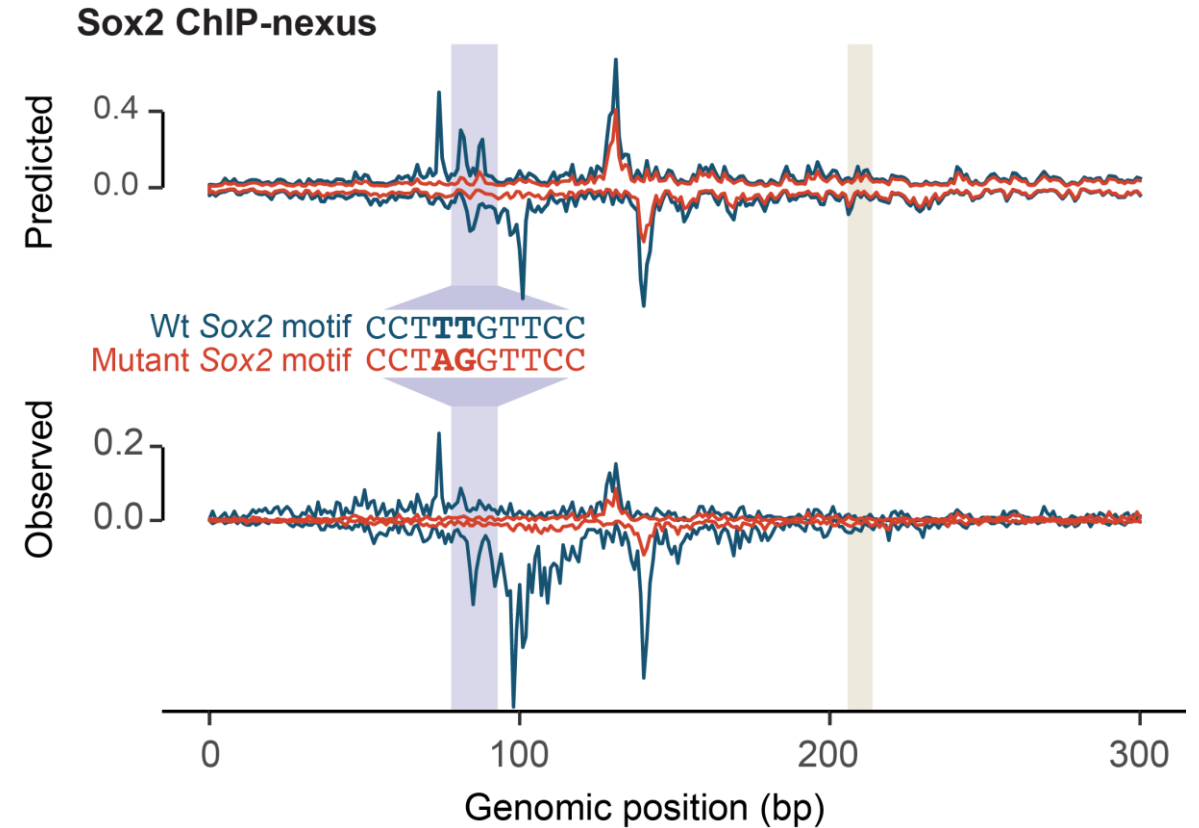Designing CRISPR experiments to validate motif syntax

**Nanog ChIP-nexus**

Predicted

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

Observed

Nanog → Sox2
Sox2 → Nanog

**Sox2 ChIP-nexus**

Predicted

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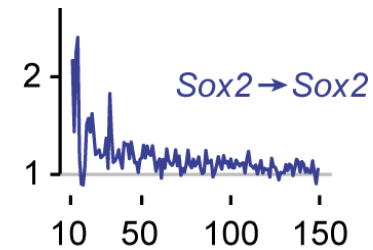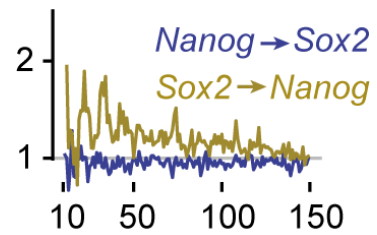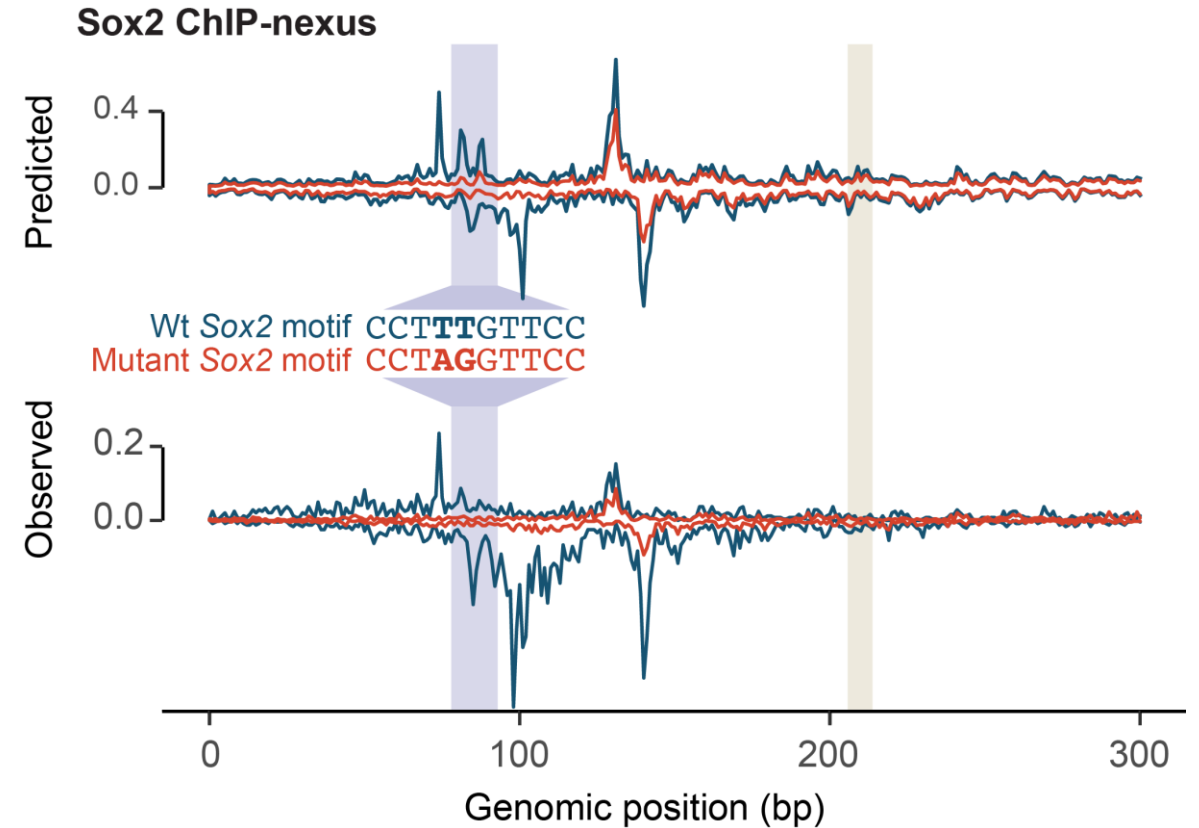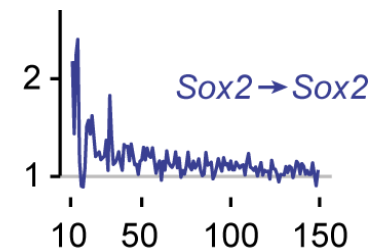Mutant *Sox2* motif CCT**AG**GTTCC

Observed

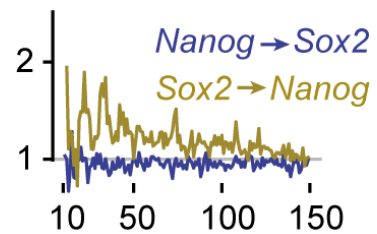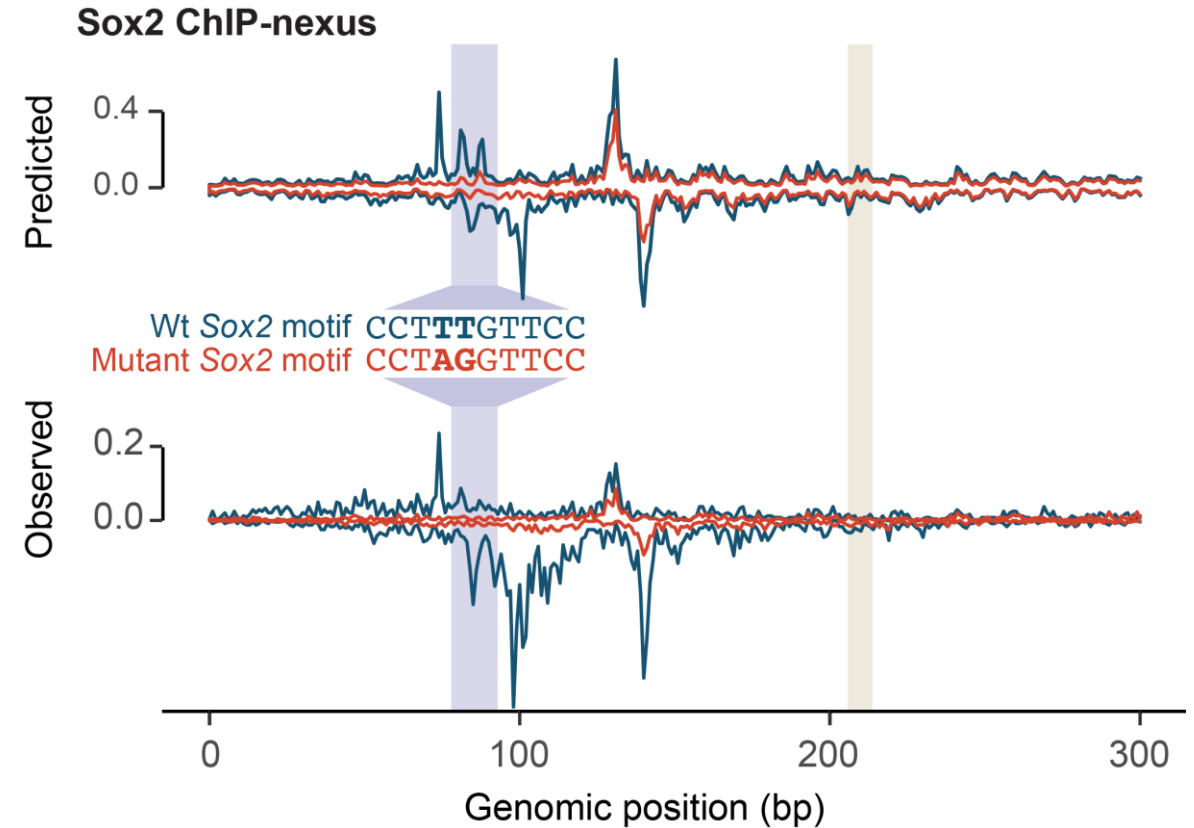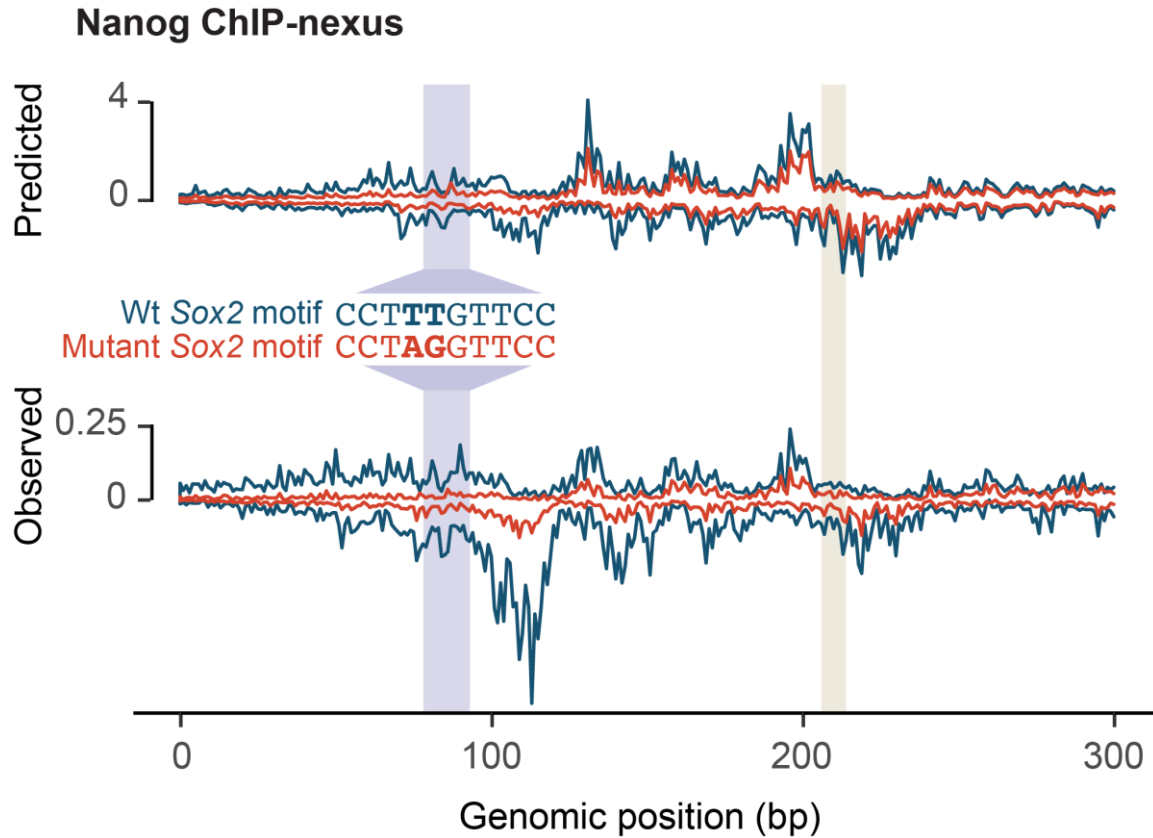Genomic position (bp)

Sox2 → Sox2

Julia Zeitlinger, Sabrina Krueger, Melanie Weilert

# In-silico mutagenesis: Predict effect of genetic variant on molecular activity

Predicted molecular profile of protein-DNA binding



ref=C

1 Kb

PredictedSignal

……ACTGAT**C**GCAATCG…….
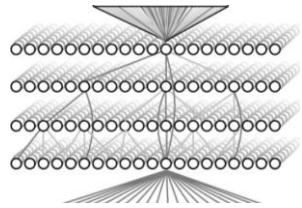
# In-silico mutagenesis: Predict effect of genetic variant on molecular activity

Predicted molecular profile of protein-DNA binding



ref=C
alt=G
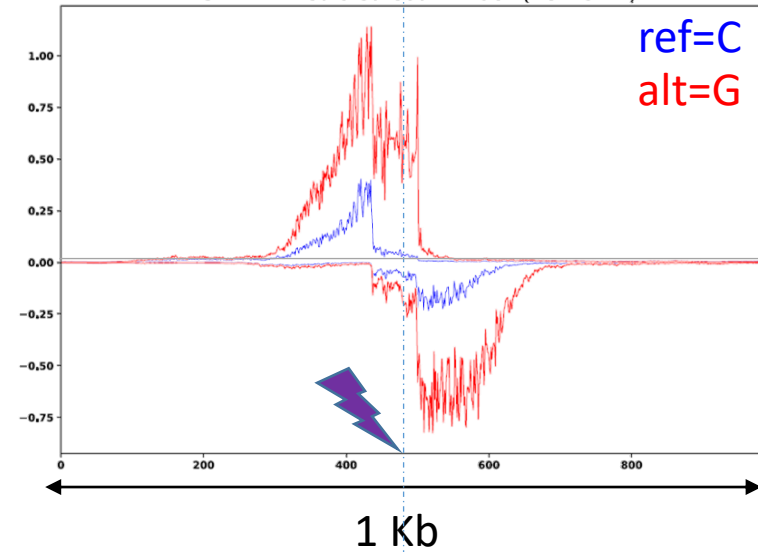
1 Kb

ΔPredictedSignal

……ACTGAT**C**GCAATCG…….

……ACTGAT**G**GCAATCG…….

Predicted molecular profile of protein-DNA binding



ref=C
alt=G

1 Kb

ref=C

alt=G

200 bp

Sequence binding motifs of SPI1 DNA binding protein

# Genetic loci associated with Alzheimer's disease



*(Lambert et al., Nat. Genet., 2013)*

Genomic position ⟶

Statistical significance of association

*Illustration by Bob Morreale, American Health Assistance Foundation*

# Molecular profiling of cell types in the brain



Corces et al. 2020, Nature Genetics

# Molecular profiling of cell types in the brain



*Corces et al. 2020, Nature Genetics*

$\Delta Predicted\ signal$

ML model

......ACTGATCG**A**AATCG.......

......ACTGATCG**C**AATCG.......
Risk

chr11:85599000-86331000 - Alzheimer's Disease rs1237999 - PICALM Locus

scATAC-seq

Excitatory Neurons
Inhibitory Neurons
Dopa. Neurons
Microglia
Oligos
Astrocytes
OPCs

Predicting and interpreting causal AD variants

Predicting and interpreting causal AD variants

# Predicting and interpreting causal AD variants

# Predicting and interpreting causal AD variants

ΔPredicted signal

ML model

.......ACTGATCG**A**AATCG.......

.......ACTGATCG**C**AATCG.......

Risk

Anna Shcherbina

Soumya Kundu

chr11:85599000-86331000 - Alzheimer's Disease rs1237999 - PICALM Locus

Excitatory Neurons

Inhibitory Neurons

scATAC-seq

rs1237999 - G to A

Importance Scores

Effect

Non-effect

Delta

FOS Motif

rs1237999 – G to A

Variant (Effect) Counts (WASP)

Reference (Non-effect) Counts (WASP)

Region
- CAUD
- HIPP
- MDFG
- PARL
- PTMN
- SMTG
- SUNI
- Below Thresh.

scATAC Coaccessibility

Genes

DLG2  CCDC89  PICALM  EED
TMEM126A  CCDC83  MIR6755
SYTL2  CREBZF  C11orf73

85640000   85760000   85880000   86000000   86120000   86240000

# Predicting and interpreting causal AD variants



Genetic variant rs1237999 disrupts a sequence motif of the FOS protein in a control element of the PICALM gene in oligodendrocyte cells in the brain

# Predicting *de-novo* mutations in Autism



Human cerebral cortex

Lakshman Sundaram

*with Greenleaf, Pasca labs*  *Trevino et al. 2021, Cell*

**Prediction: Mutation disrupts NFIA motif in control element of NFIA gene in glutamatergic neurons**

# Predicting *de-novo* mutations in Autism



Laksshman Sundaram

Human cerebral cortex

PCW24
PCW21
PCW20
PCW16

scATAC-seq    scRNA-seq

RG
oRG
tRG/Ependymal
oIPC/Astrocytes
OPC/Oligos
nIPC
Migrating GluN1
Migrating GluN2
GluN1
GluN2
CGE IN
MGE IN
Microglia
Pericytes
Endothelial cells

CP
SP
oSVZ
SVZ
VZ

*with Greenleaf, Pasca labs*    *Trevino et al. 2021, Cell*

**Prediction: Mutation disrupts NFIA motif in control element of NFIA gene in glutamatergic neurons**

# Predicting *de-novo* mutations in congenital heart disease



Human Fetal Heart
6 PCW    8 PCW    19 PCW

chr1:201352761–201402762

C1 eCM
C2 aCM
C3 vCM
C4 OFT
C5 FB1
C6 CFP
C7 FB2
C8 preCF
C9 CF
C10 preSMC
C11 SMC
C12 PC
C13 EPC
C14 NC
C15 Endo1
C16 Endo2
C17 IEC
C18 aEC
C19 Cap
C20 vEC

Laksshman Sundaram

*Mo Ameen*

*with Wang, Karakikes, Quertermous, Greenleaf labs*

# Predicting *de-novo* mutations in congenital heart disease



**Prediction: Mutation disrupts ETV motif in control element active in arterial endothelial cells**

*with Wang, Karakikes, Quertermous, Greenleaf labs*

# Predicting *de-novo* mutations in congenital heart disease



Human Fetal Heart

6 PCW    8 PCW    19 PCW

chr1:201352761–201402762

C1 eCM
C2 aCM
C3 vCM
C4 OFT
C5 FB1
C6 CFP
C7 FB2
C8 preCF
C9 CF
C10 preSMC
C11 SMC
C12 PC
C13 EPC
C14 NC
C15 Endo1
C16 Endo2
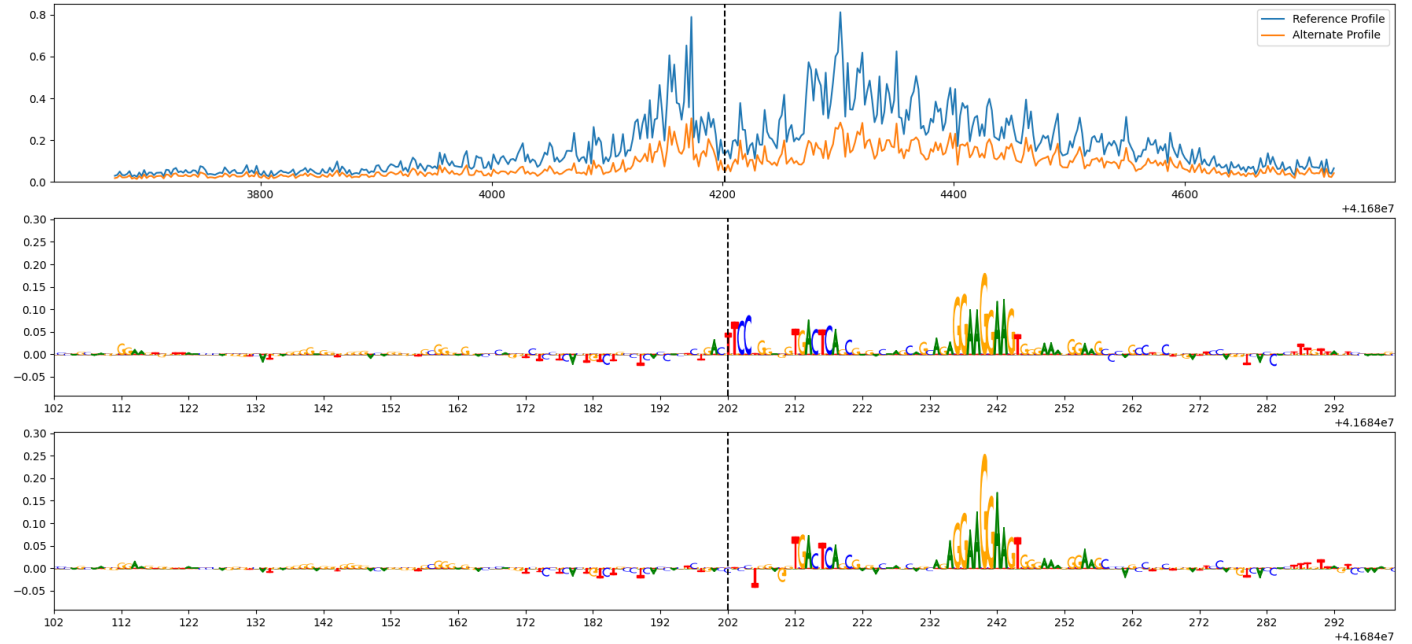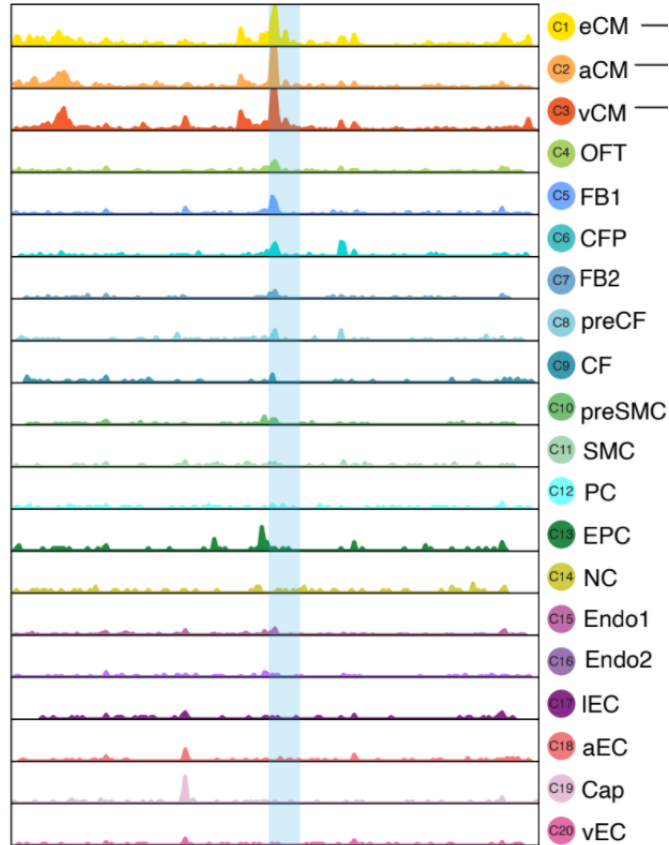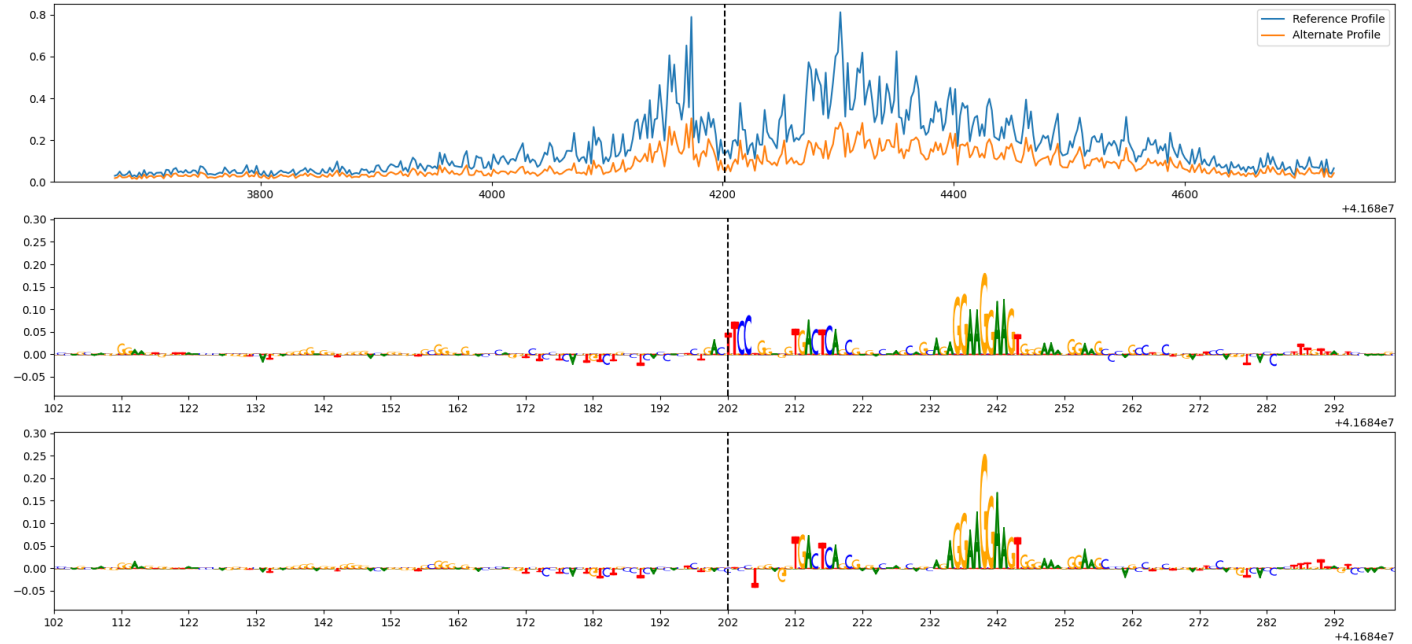C17 IEC
C18 aEC
C19 Cap
C20 vEC

Laksshman Sundaram

Mo Ameen

Reference Profile
Alternate Profile

**Prediction: Mutation disrupts ETV motif in control element active in arterial endothelial cells**

Relative JARID2 expression to WT (Normalized to ACTB)    WT-EC    JKO-EC    **

Relative NFATC1 expression to WT (Normalized to ACTB)    WT-EC    NKO-EC    ***

Relative TFAP2A expression to WT (Normalized to ACTB)    WT-NC    TKO-NC    ****

**CRISPR/Cas9 experiments validate downstream target genes**

*with Wang, Karakikes, Quertermous, Greenleaf labs*

# Democratizing ML for genomics: http://kipoi.org/



**Kipoi: Model zoo for genomics**

## Numbers

**# of models:** 1709

**# of model groups:** 16

**# of contributors:** 6

**# of model groups supporting postprocessing:**

- **Variant effect prediction:** 11/16

## Model groups by tag



Kipoi (pronounce: kípi; from the Greek κήποι: gardens) is an API and a repository of ready-to-use trained models for regulatory genomics. It currently contains 1709 different models, covering canonical predictive tasks in transcriptional and post-transcriptional gene regulation. Kipoi's API is implemented as a python package (github.com/kipoi/kipoi) and it is also accessible from the command line or R.

*Avsec et al. 2019 Nature Biotech*

## Summary

- Large-scale molecular profiling datasets => decipher genome function

- Neural networks can map DNA sequence to molecular profiles with unprecedented accuracy

- Models can be interpreted to decipher functional DNA letters, words and syntax

- Models can be used to decipher disease-associated mutations

- Predictions are validated by genome editing experiments

- Predictions can provide clues for therapeutic interventions

# Kundaje lab

STANFORD COMPUTER SCIENCE

STANFORD SCHOOL OF MEDICINE | Genetics

Daniel Kim (BMI)

Kelly Cochran (CS)

Soumya Kundu (CS)

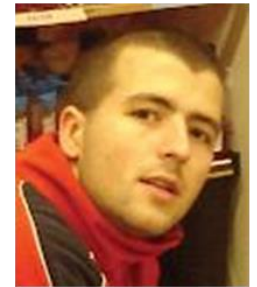Surag Nair (CS)

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)

## Collaborators