# Many transcription factors recognize DNA shape

## Katie Pollard
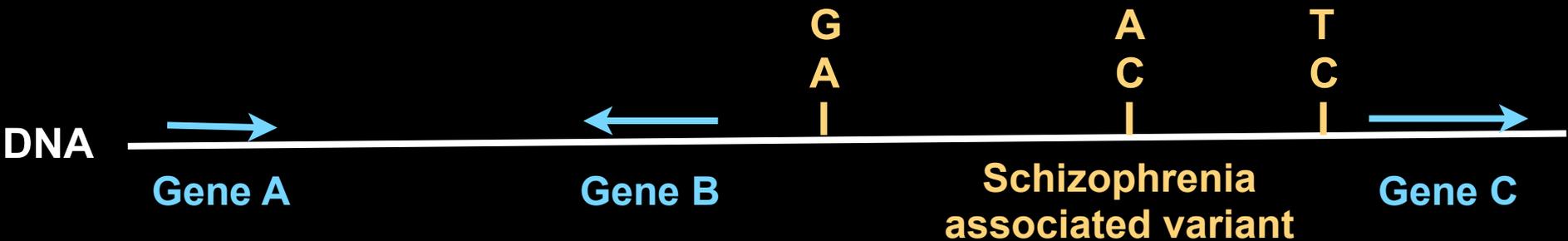
**Gladstone Institutes**

**UCSF Division of Biostatistics, Institute for Human Genetics, and Institute for Computational Health Sciences**

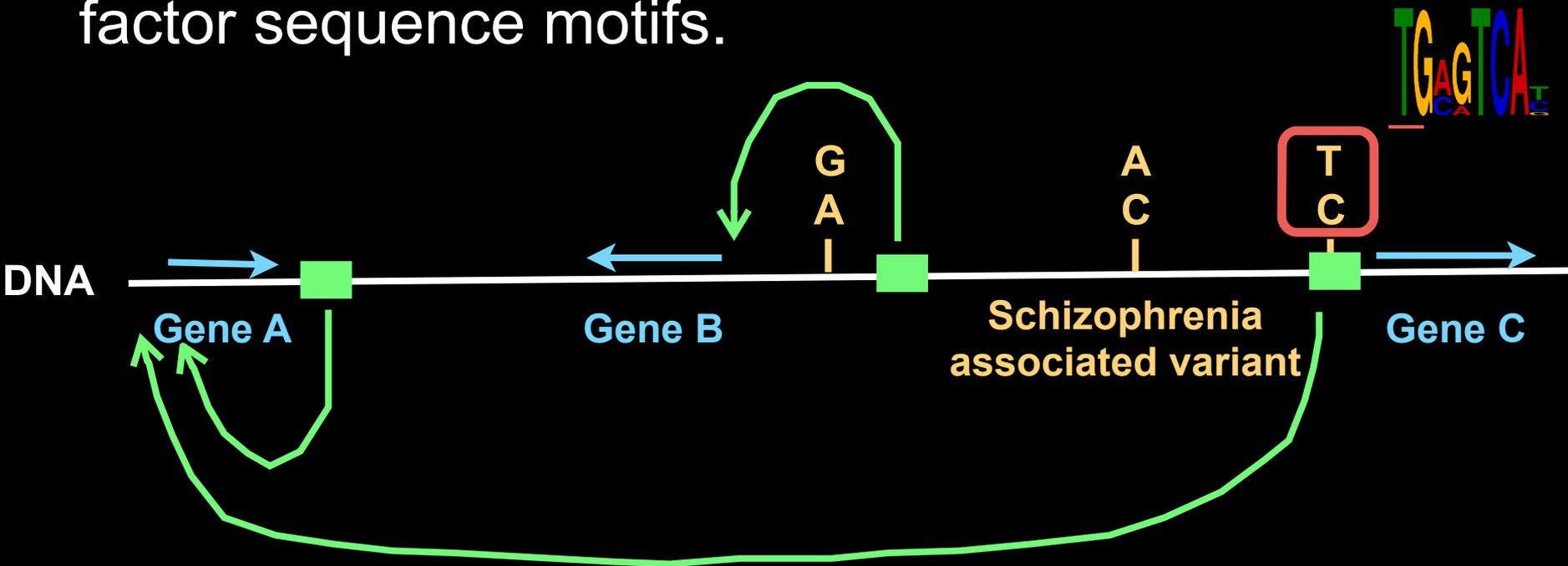ENCODE Users Meeting - Stanford, CA

June 10, 2016

# Most disease associated mutations are outside coding regions

Hypothesis 1: Non-coding variants alter transcription factor sequence motifs.
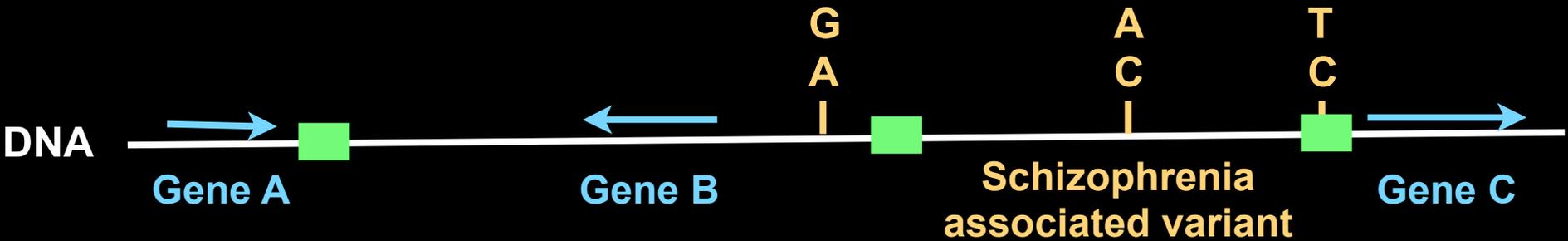
# Most disease associated mutations are outside coding regions

Hypothesis 1: Non-coding variants alter transcription factor sequence motifs.



Approach: Map variants to correct pathways by predicting enhancers and their target genes. Score variants for changes in binding affinity.

# EnhancerFinder distinguishes biologically active enhancers



DNA

Gene A      Gene B      Schizophrenia associated variant      Gene C

G A I      A C I      T C
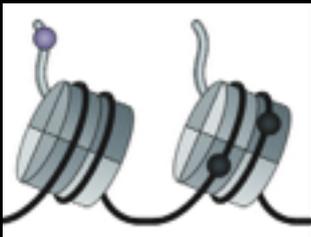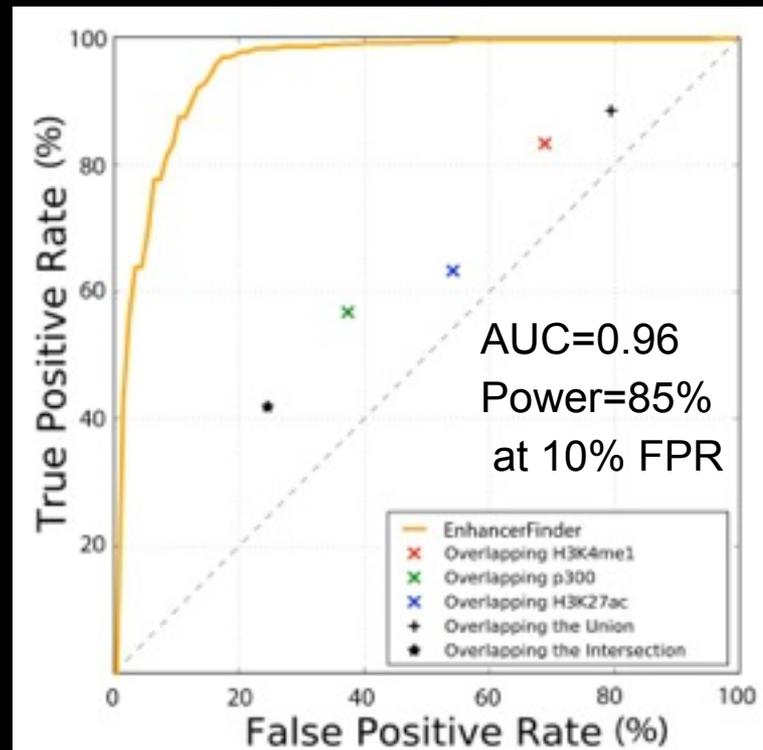
## Training Data

**VISTA Enhancers**



**Yes**      **No**

Functional Genomics



DNA Sequences

AAAA,AAAC,AAAG,AAAT,..

## Performance on held out data



AUC=0.96
Power=85%
at 10% FPR

Legend:
- EnhancerFinder
- Overlapping H3K4me1
- Overlapping p300
- Overlapping H3K27ac
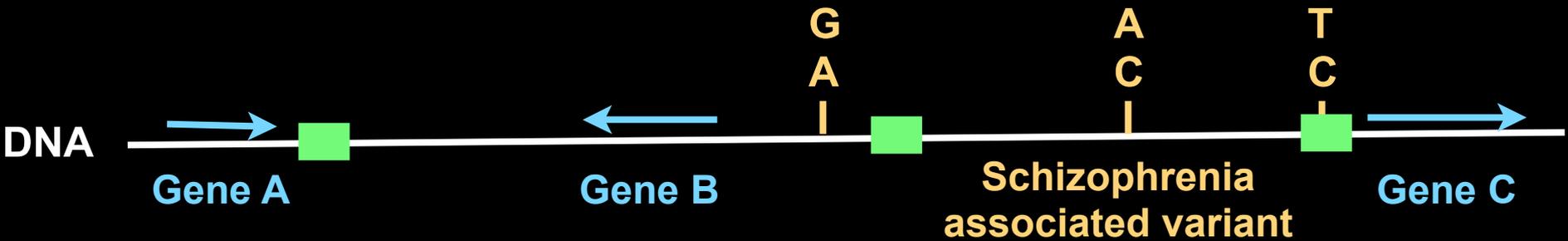- Overlapping the Union
- Overlapping the Intersection

## 80% validate in vivo



Erwin et al. (2014)
Capra et al. (2014)

# EnhancerFinder distinguishes biologically active enhancers



DNA

G
A
I
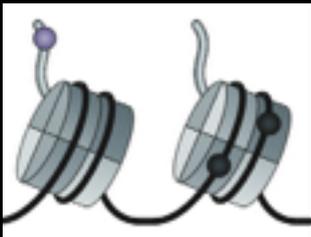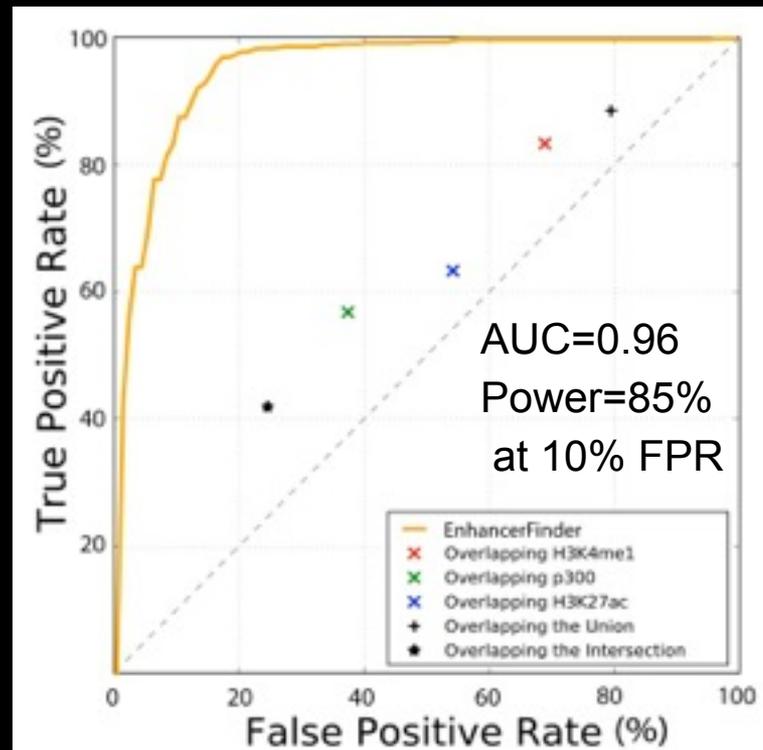
A
C
I

T
C
I

**Gene A**

**Gene B**

**Schizophrenia associated variant**

**Gene C**

**Training Data**
VISTA Enhancers

**Yes** **No**

Functional Genomics

DNA Sequences
AAAA,AAAC,AAAG,AAAT,..

**Performance on held out data**

True Positive Rate (%)

AUC=0.96
Power=85%
at 10% FPR

— EnhancerFinder
× Overlapping H3K4me1
× Overlapping p300
× Overlapping H3K27ac
+ Overlapping the Union
⬠ Overlapping the Intersection

False Positive Rate (%)

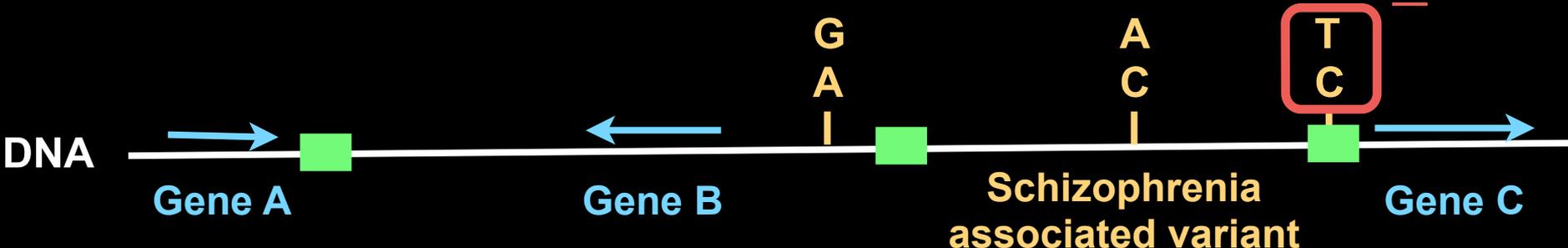**80% validate in vivo**

Erwin et al. (2014)
Capra et al. (2014)

# MotifDiverge quantifies loss/gain of TF binding sites

**DNA**

Gene A    Gene B    Schizophrenia associated variant    Gene C

G A I    A C I    T C



## Statistical model for TFBS evolution with turnover

Slide TFBS motif

Bernoulli trial $(n_A - n_B)$

Seq A    # Hits $(m_A)$

GGGAGTGTTGAGGGGGCCTGAAGGGTTCCGCTCCTCCCACCCAGGGAACCGCCATGCCACTAGTGGGCTGTCCTGGAGACTCGGGGAGAAAGCACACAGGCTGTCGGGAAAGGTGGGTCGCAGGC

Seq B    # Hits $(m_B)$

TTGGAGGATGATTAGGTTTCATAAGATGGAGGTGATGGTAAAGTGATCAATGTAAGAGAGTGTCACTTAATAACATATGGATGATTTGTATGC

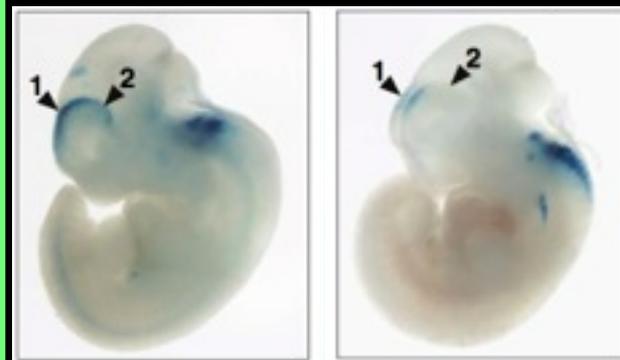Two correlated Bernoulli trials $(n_B)$

$$P(N_{xy} = n_{xy}) =$$

$$(6) \quad \begin{cases} \sum_{j=0}^{k_x - k_y} P_s(N_1 = n_{xy} - j)Bin(N_2 = j) & \text{for} \quad k_x \geq k_y \\ \sum_{j=0}^{k_y - k_x} P_s(N_1 = n_{xy} + j)Bin(N_2 = j) & \text{for} \quad k_x < k_y, \end{cases}$$

### P-value for net change in binding
- One or many TFs
- Alignment-free
- Evolutionary model
- Motif specific
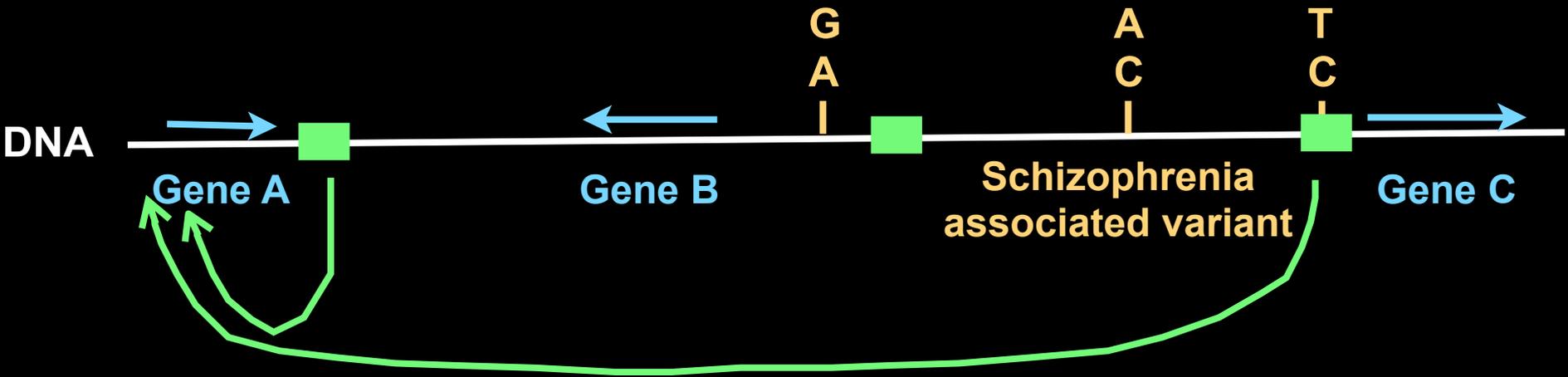
## Predicts change of function



### Detects loss/gain of function mutations with high accuracy
- Better than conservation scores
- In vivo and MPRAs in cell lines

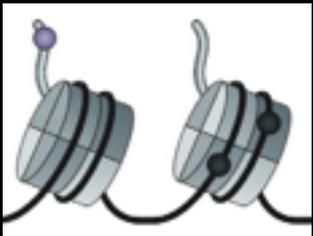Ritter et al. (2010)
Kostka et al. (2015)

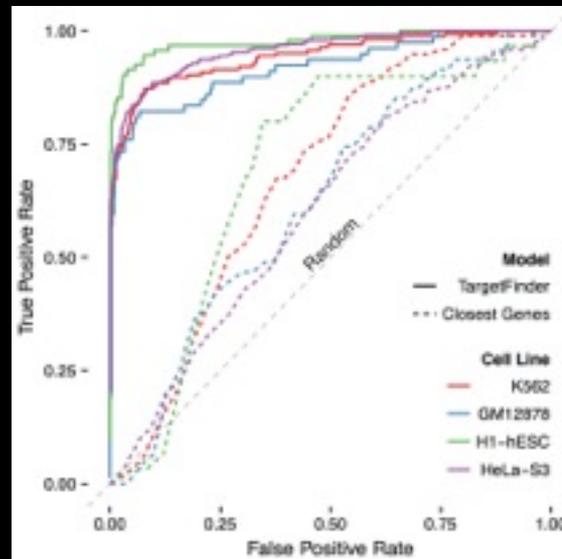# TargetFinder maps distal regulatory elements to genes



DNA

G A I     A C I     T C

Gene A     Gene B     Schizophrenia associated variant     Gene C

**Training Data**

Active enhancers
Expressed genes
Hi-C interactions

Functional Genomics

**Closest gene usually wrong**



**Reveals distinct genomic signature of looping DNA**

- Heterochromatin on loop
- Cohesin within 6Kb of enhancer and promoter but not mid-loop
- TFs bound with CTCF improve predictions

Whalen et al. (2016)

# Summary and Challenges

• Machine-learning on biologically validated enhancers identifies non-coding variants most likely to affect gene regulation <u>and</u> the targeted genes.

   – Massive integration of functional genomics data enables cell type specific predictions

   – Many enhancer-like regions are minimally active and not consistently looping to a target gene
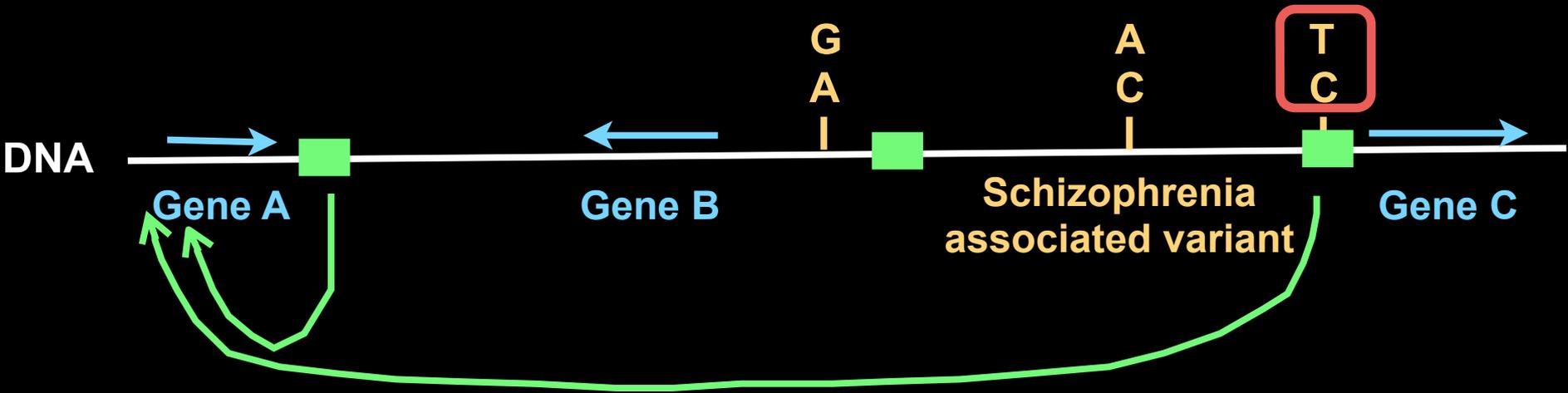
# Summary and Challenges

• Machine-learning on biologically validated enhancers identifies non-coding variants most likely to affect gene regulation <u>and</u> the targeted genes.

- Massive integration of functional genomics data enables cell type specific predictions

- Many enhancer-like regions are minimally active and not consistently looping to a target gene

• But much remains to be explained…

- Functional variants outside enhancers

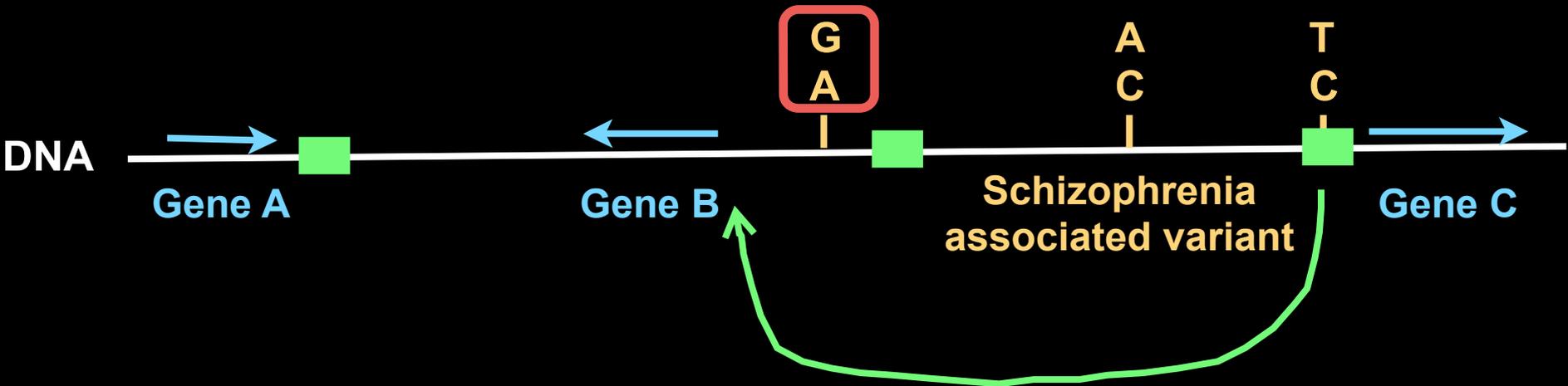# TargetFinder maps distal regulatory elements to genes



Hypothesis 2: Non-coding variants alter binding sites of structural proteins and chromatin modifiers.

**Reveals distinct genomic signature of looping DNA**

- Heterochromatin on loop
- Cohesin within 6Kb of enhancer and promoter but not mid-loop
- TFs bound with cohesin improve predictions

Whalen et al. (2016)

# TargetFinder maps distal regulatory elements to genes

DNA

Gene A    Gene B    Schizophrenia associated variant    Gene C
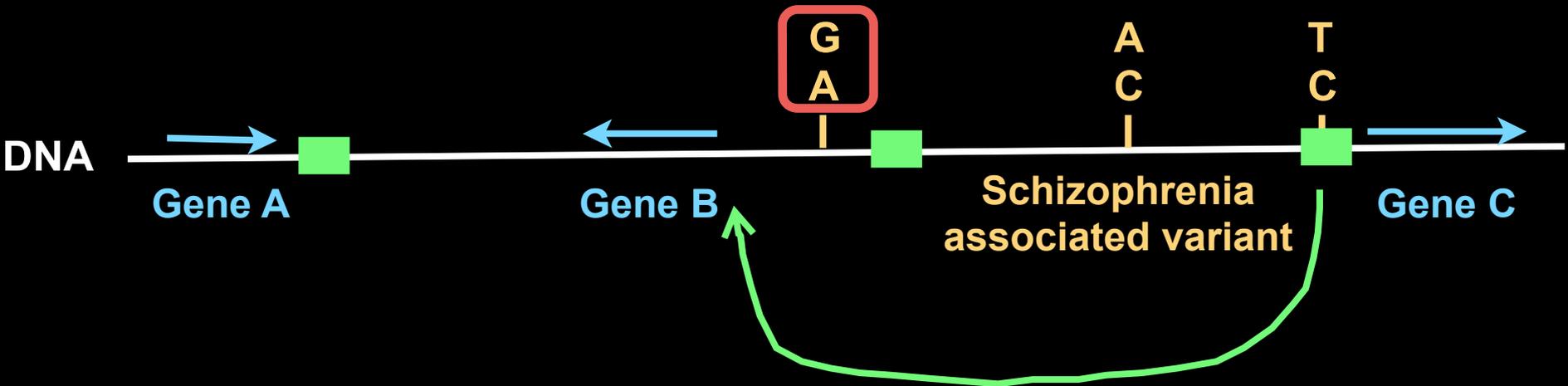
G A    A C    T C

Hypothesis 2: Non-coding variants alter binding sites of structural proteins and chromatin modifiers.

**Reveals distinct genomic signature of looping DNA**

- Heterochromatin on loop
- Cohesin within 6Kb of enhancer and promoter but not mid-loop
- TFs bound with cohesin improve predictions

Whalen et al. (2016)

# TargetFinder maps distal regulatory elements to genes

DNA

Gene A

Gene B

Schizophrenia associated variant

Gene C

G
A

A
C

T
C

Hypothesis 2: Non-coding variants alter binding sites of structural proteins and chromatin modifiers.

Approach: CRISPR edit sites identified by TargetFinder, then test chromatin and expression.

**Reveals distinct genomic signature of looping DNA**

- Heterochromatin on loop

- Cohesin within 6Kb of enhancer and promoter but not mid-loop

- TFs bound with cohesin improve predictions

Whalen et al. (2016)

# Summary and Challenges

• Machine-learning on biologically validated enhancers identifies non-coding variants most likely to affect gene regulation and the targeted genes.

- Massive integration of functional genomics data enables cell type specific predictions
- Many enhancer-like regions are minimally active and not consistently looping to a target gene

• But much remains to be explained…

- Functional variants outside enhancers
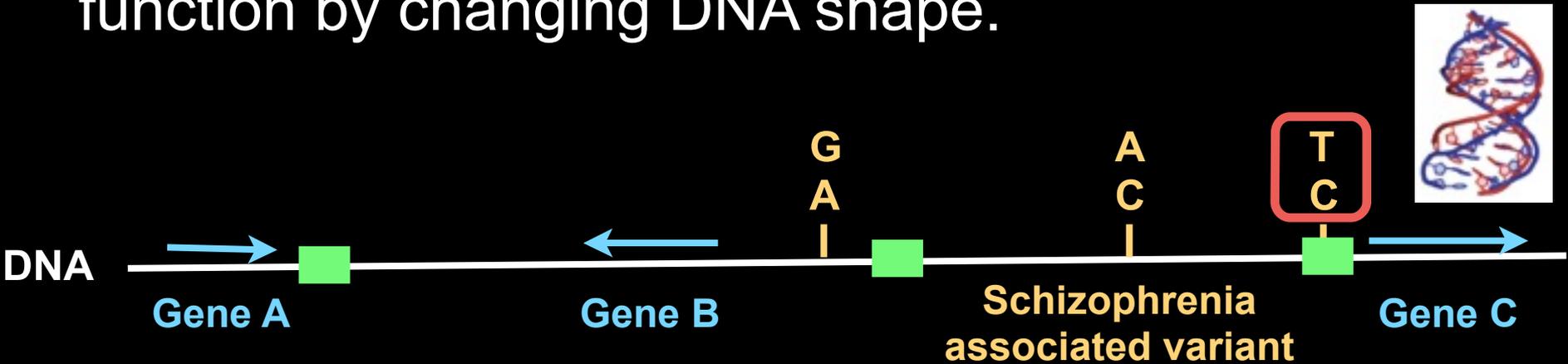- Enhancer variants outside sequence motifs

# Summary and Challenges

• Machine-learning on biologically validated enhancers identifies non-coding variants most likely to affect gene regulation and the targeted genes.

  – Massive integration of functional genomics data enables cell type specific predictions

  – Many enhancer-like regions are minimally active and not consistently looping to a target gene

• But much remains to be explained…

  – Functional variants outside enhancers

  – Enhancer variants outside sequence motifs

For a typical ENCODE TF 23% of the top 2000 ChIP-seq peaks have no sequence motif (range = 1%-63%)
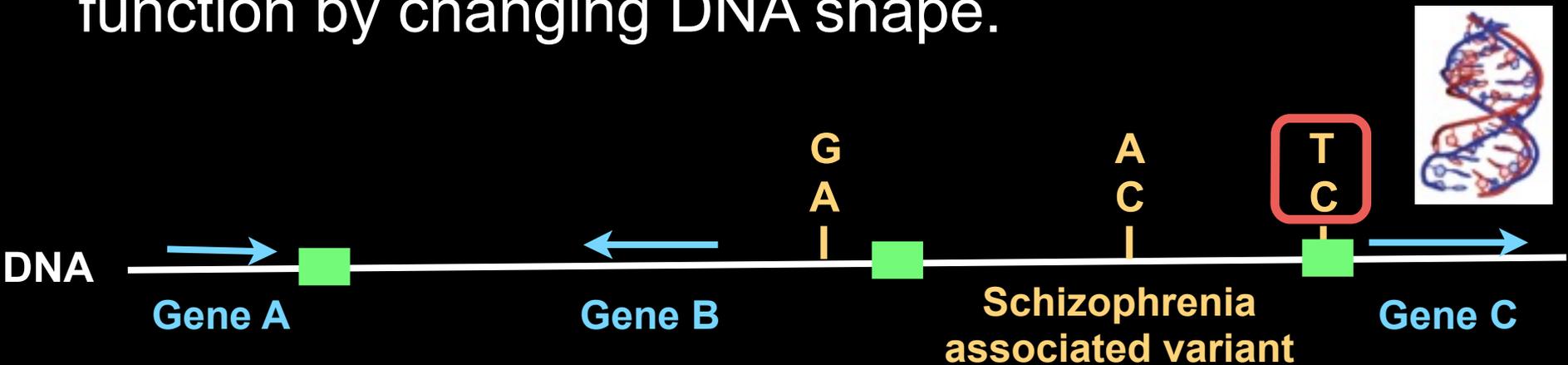
# Many enhancer mutations are outside known or de novo sequence motifs

Hypothesis 3: Non-coding variants alter enhancer function by changing DNA shape.



**DNA**

**Gene A**

**Gene B**

G
A
I

A
C
I

T
C

**Schizophrenia associated variant**

**Gene C**

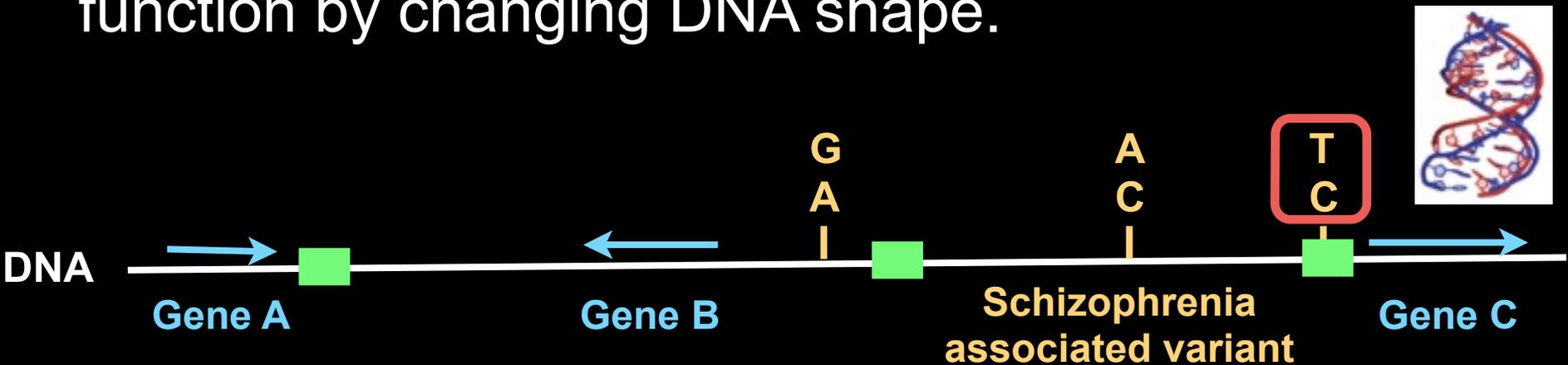# Many enhancer mutations are outside known or de novo sequence motifs

Hypothesis 3: Non-coding variants alter enhancer function by changing DNA shape.



- TFs can recognize shape in addition to sequence.
- DNA shape differentiates similar sequence motifs.
- Distinct sequences can encode same shape.

# Many enhancer mutations are outside known or de novo sequence motifs

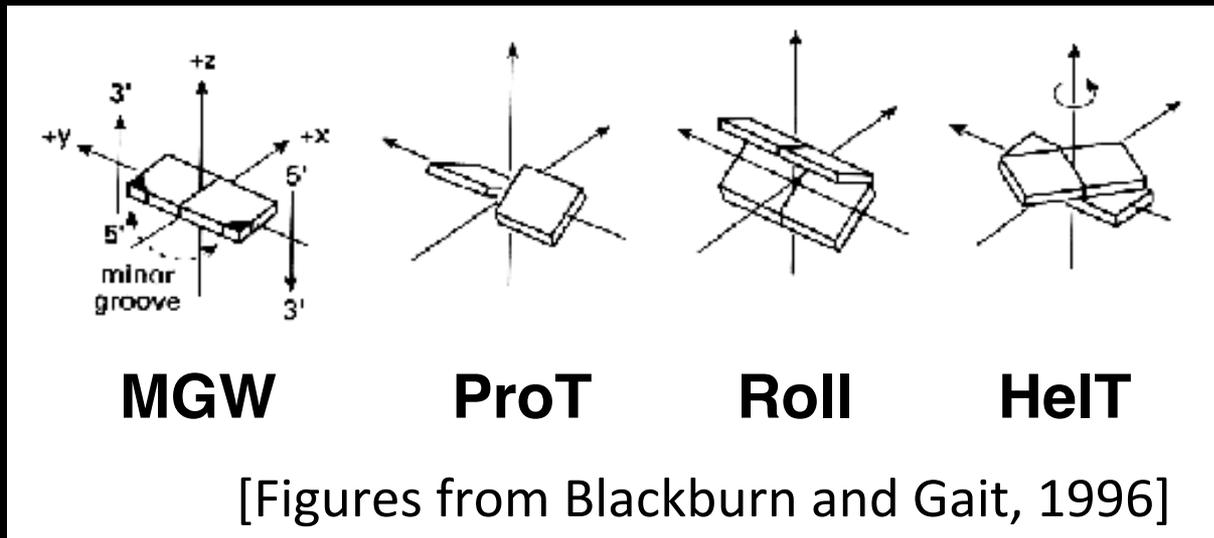Hypothesis 3: Non-coding variants alter enhancer function by changing DNA shape.



DNA

Gene A

Gene B

G
A
I

Schizophrenia
associated variant

A
C
I

T
C

Gene C

- TFs can recognize shape in addition to sequence.
- DNA shape differentiates similar sequence motifs.
- Distinct sequences can encode same shape.

Approach: Algorithm to learn **shape motifs** de novo for all ENCODE TFs, predict shape motif hits in ChIP-seq peaks, compare to sequence motifs

# De novo shape motif discovery

1. <u>Estimate DNA structure</u>: `DNAshape` (Zhou et al. 2013)
   - Maps 5-mer sequences to structural features.
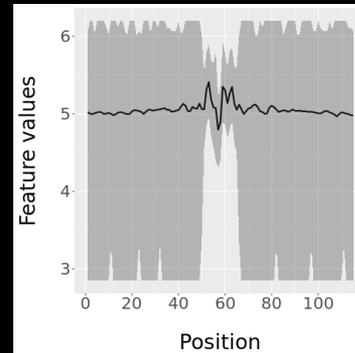   - Based on molecular dynamics simulations.



**MGW**      **ProT**      **Roll**      **HelT**

[Figures from Blackburn and Gait, 1996]

# De novo shape motif discovery

1. <u>Estimate DNA structure</u>: `DNAshape` (Zhou et al. 2013)
   - Maps 5-mer sequences to structural features
   - Based on molecular dynamics simulations
2. <u>Learn TF shape motifs</u>: Search ChIP-seq peaks for windows with similar values of a shape feature.
   - Gibbs sampling with scores ~ $\exp(\sum D_{ij})$
   - Vary window size 5-25bp
   - Train on 1000 of top 2000 peaks for ~250 TFs

# De novo shape motif discovery

1. <u>Estimate DNA structure</u>: `DNAshape` (Zhou et al. 2013)
   - Maps 5-mer sequences to structural features
   - Based on molecular dynamics simulations
2. <u>Learn TF shape motifs</u>: Search ChIP-seq peaks for windows with similar values of a shape feature.
   - Gibbs sampling with scores ~ $\exp(\sum D_{ij})$
   - Vary window size 5-25bp
   - Train on 1000 of top 2000 peaks for ~250 TFs
3. <u>Call hits</u>: Scan ChIP-seq peaks with shape motifs.
   - Null distribution on distance from mean shape feature value at each position
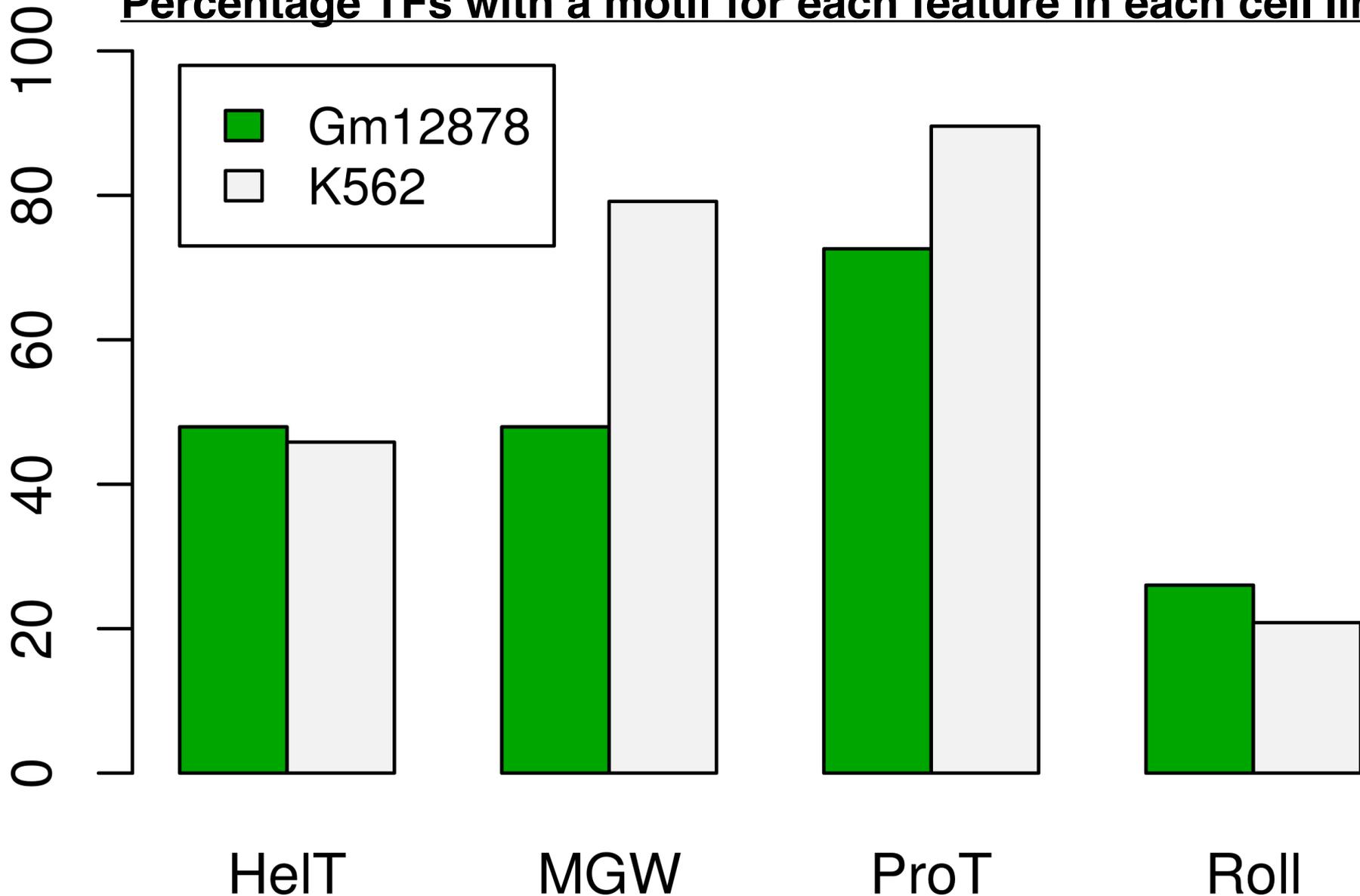
# De novo shape motif discovery

1. <u>Estimate DNA structure</u>: `DNAshape` (Zhou et al. 2013)
   - Maps 5-mer sequences to structural features
   - Based on molecular dynamics simulations
2. <u>Learn TF shape motifs</u>: Search ChIP-seq peaks for windows with similar values of a shape feature.
   - Gibbs sampling with scores ~ $\exp(\sum D_{ij})$
   - Vary window size 5-25bp
   - Train on 1000 of top 2000 peaks for ~250 TFs
3. <u>Call hits</u>: Scan ChIP-seq peaks with shape motifs.
   - Null distribution on distance from mean shape feature value at each position
   - Apply to remaining 1000 of top 2000 peaks and flanking non-peak regions for each TF

# De novo shape motif discovery

1. <u>Estimate DNA structure</u>: `DNAshape` (Zhou et al. 2013)
   - Maps 5-mer sequences to structural features
   - Based on molecular dynamics simulations
2. <u>Learn TF shape motifs</u>: Search ChIP-seq peaks for windows with similar values of a shape feature.
   - Gibbs sampling with scores ~ $\exp(\sum D_{ij})$
   - Vary window size 5-25bp
   - Train on 1000 of top 2000 peaks for ~250 TFs
3. <u>Call hits</u>: Scan ChIP-seq peaks with shape motifs.
   - Null distribution on distance from mean shape feature value at each position
   - Apply to remaining 1000 of top 2000 peaks and flanking non-peak regions for each TF
4. <u>Enrichment test</u>: Hypergeometric p-value.

# Shape motifs are common



Percentage TFs with a motif for each feature in each cell line

Legend:
- Gm12878 (green)
- K562 (light gray)

X-axis categories: HelT, MGW, ProT, Roll

# Shape complements sequence motifs

• Most peaks without sequence motifs have at least one shape motif. It is typically at the peak center.

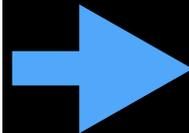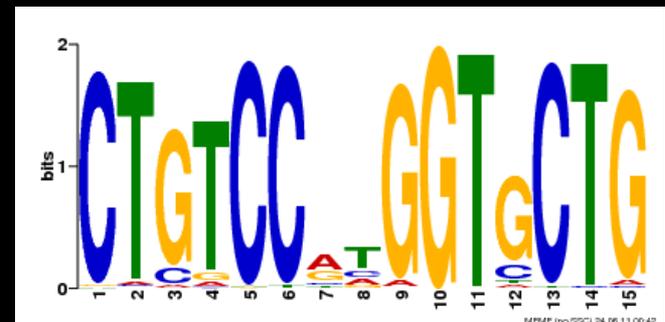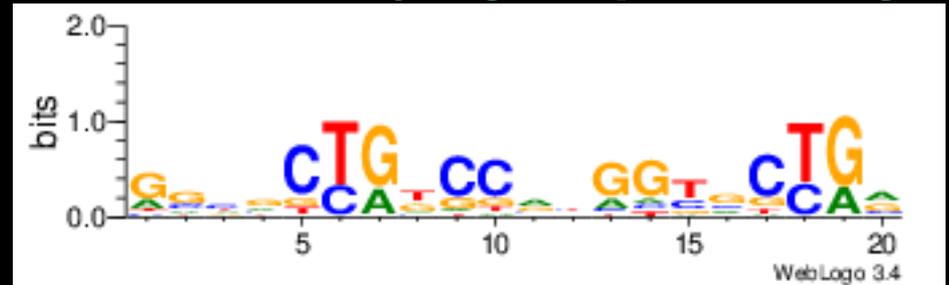# Shape complements sequence motifs

• Most peaks without sequence motifs have at least one shape motif. It is typically at the peak center.

• ~25% of peaks have sequence and shape motifs.

# Shape complements sequence motifs

• Most peaks without sequence motifs have at least one shape motif. It is typically at the peak center.

• ~25% of peaks have sequence and shape motifs.
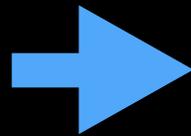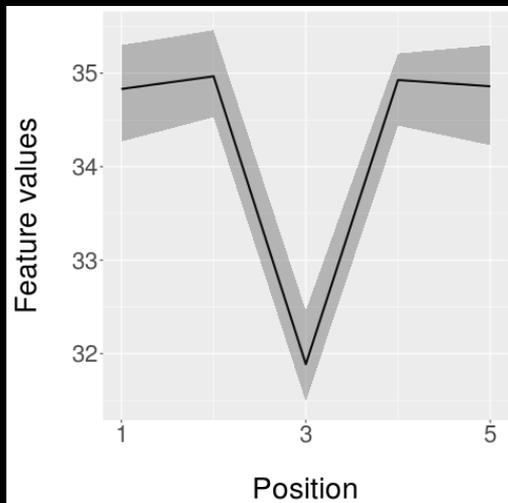  – These can be similar,

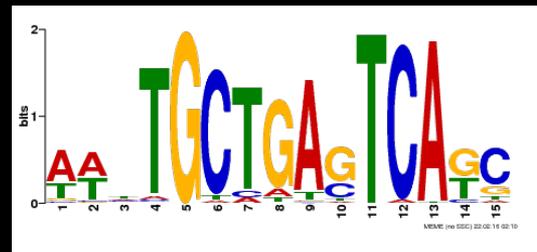**Underlying sequence logo**





**Nrsf Roll motif in K562**

**Nrsf FactorBook sequence motif**

# Shape motifs are complementary

• Most peaks without sequence motifs have at least one shape motif. It is typically at the peak center.

• Many peaks have sequence and shape motifs.
  – These can be similar,
  – Extensions or refinements of one another,



**Cfos ProT motif in K562**

**Underlying sequence**

**FactorBook sequence motif**

# Shape motifs are complementary

• Most peaks without sequence motifs have at least one shape motif. It is typically at the peak center.

• Many peaks have sequence and shape motifs.
- These can be similar,
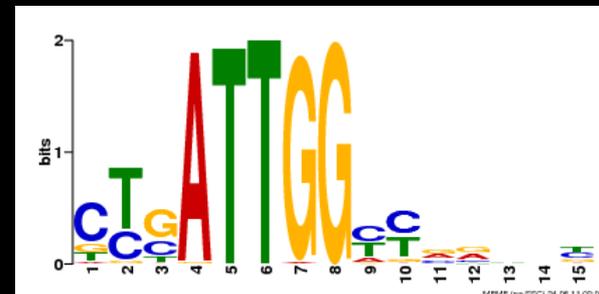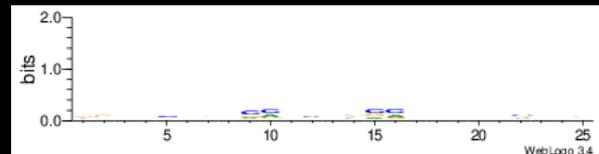- Extensions or refinements of one another,
- Or very different



**Maff HelT motif in K562**
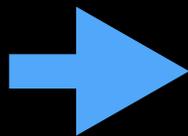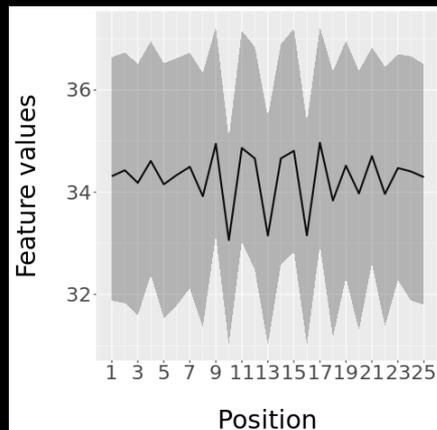


**Underlying sequence**



**FactorBook sequence motif**

# Shape motifs are complementary

• Most peaks without sequence motifs have at least one shape motif. It is typically at the peak center.

• Many peaks have sequence and shape motifs.
  - These can be similar,
  - Extensions or refinements of one another,
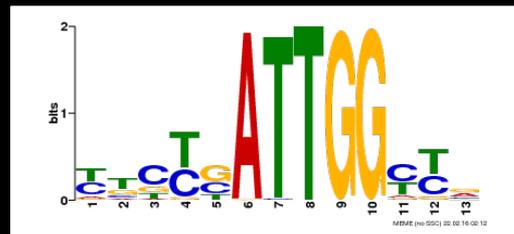  - Or very different

• Shape motifs can flank sequence motifs



**Nfya HelT motif in K562**

**Underlying sequence**

**FactorBook sequence motif is 3bp upstream**

# Shape motifs are complementary

• Most peaks without sequence motifs have at least one shape motif. It is typically at the peak center.

• Many peaks have sequence and shape motifs.
  – These can be similar,
  – Extensions or refinements of one another,
  – Or very different

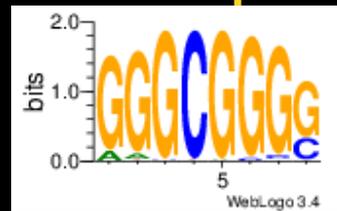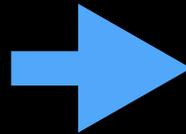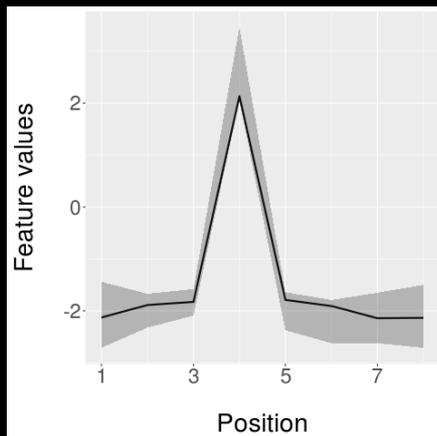• Shape motifs can flank sequence motifs



**Cfos Roll motif in K562**





**Underlying sequence**

**FactorBook sequence motif is <u>30bp</u> away**

# Shape motifs are complementary

• Most peaks without sequence motifs have at least one shape motif. It is typically at the peak center.

• Many peaks have sequence and shape motifs.
   – These can be similar,
   – Extensions or refinements of one another,
   – Or very different

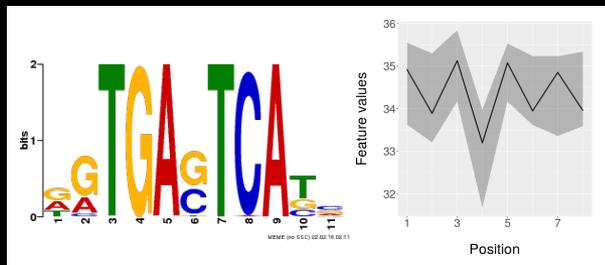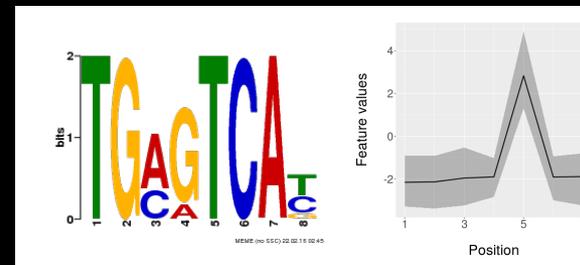• Shape motifs can flank sequence motifs

• Shape motifs can differ between TFs with similar sequence motifs and/or the same protein fold.



**FosI1 has a HelT motif**          **Atf3 has a Roll motif**

# Ongoing Work

• Hierarchical or mixture model of TF binding with sequence and shape motifs
- Decompose sequence motifs by shape types
- Spectrum of recognition modes

# Ongoing Work

- Hierarchical or mixture model of TF binding with sequence and shape motifs
  - Decompose sequence motifs by shape types
  - Spectrum of recognition modes
- Shape motifs in different contexts
  - Co-factors and complexes
  - Weak ChIP-seq peaks

# Ongoing Work

- Hierarchical or mixture model of TF binding with sequence and shape motifs
  - Decompose sequence motifs by shape types
  - Spectrum of recognition modes
- Shape motifs in different contexts
  - Co-factors and complexes
  - Weak ChIP-seq peaks
- Role of shape in ectopic binding of TFs when co-factors are absent [Luna-Zurita et al. 2016]

# Ongoing Work

- Hierarchical or mixture model of TF binding with sequence and shape motifs
  - Decompose sequence motifs by shape types
  - Spectrum of recognition modes
- Shape motifs in different contexts
  - Co-factors and complexes
  - Weak ChIP-seq peaks
- Role of shape in ectopic binding of TFs when co-factors are absent [Luna-Zurita et al. 2016]
- Evolutionary modeling of DNA shape
  - Conservation of shape without sequence
  - Scoring SNPs for effects on shape motifs

# Collaborators

**EnhancerFinder**
Tony Capra
Gen Haliburton
**DNA Shape**
**Hassan Samee**

**TargetFinder**
Rebecca Truty
**Sean Whalen**
**MotifDiverge**
Dennis Kostka

**Functional Assays**
Hane Ryu
Alex Pollen
**Nadav Ahituv**
**Arnold Kriegstein**