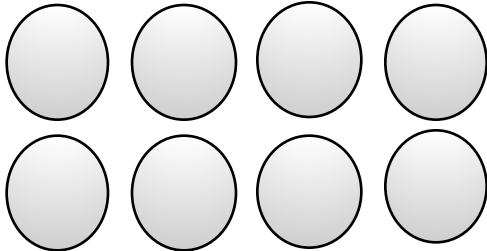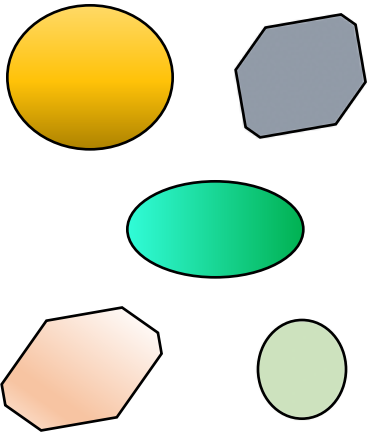# Computational prediction of cis-regulatory elements in eukaryotic genomes.

Dmitry Svetlichnyy

# Single input multiple outputs
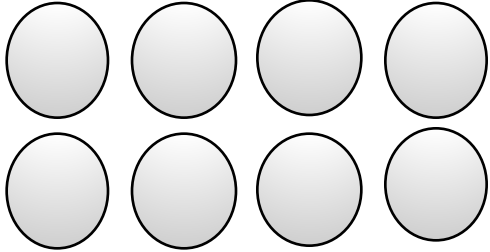


Undifferentiated cells
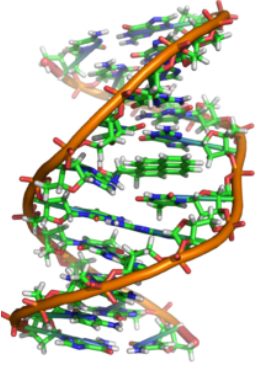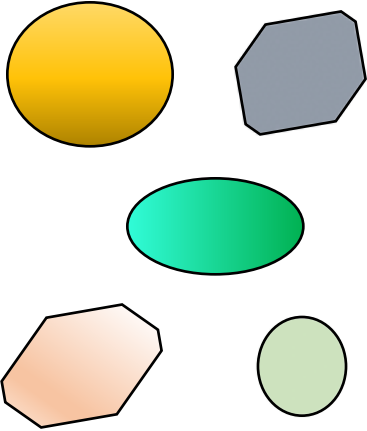
Differentiation

Differentiated cells

# Single input multiple outputs



Undifferentiated cells

Differentiation

Differentiated cells

Single input multiple outputs

Undifferentiated cells

Differentiation

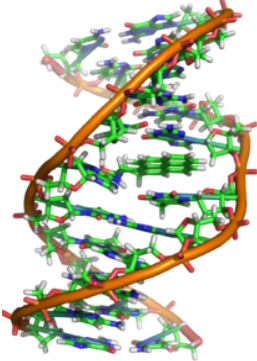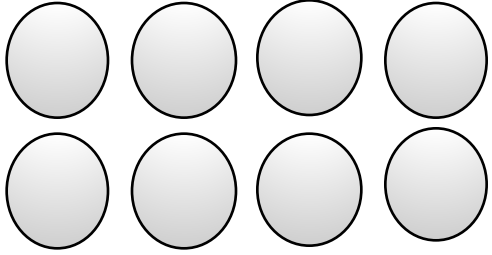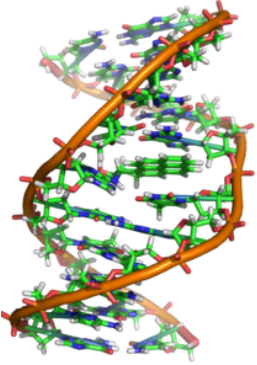Control of gene expression

Differentiated cells

# Single input multiple outputs
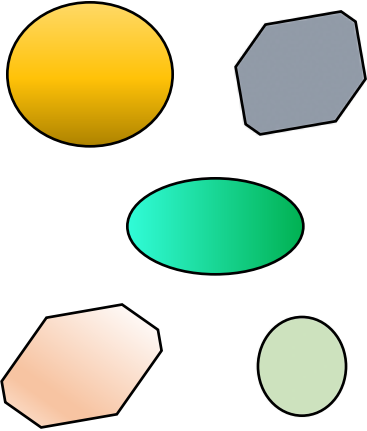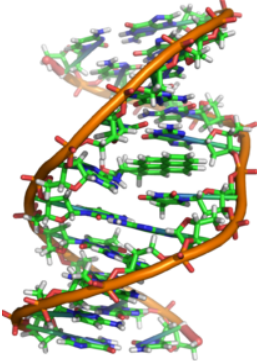


Undifferentiated cells

Differentiation

Differentiated cells

Control of gene expression

**coding**
2%

98%

**non-coding**

**Human genome**

2%

98%

Control of gene expression

exon1 exon2 exon3

Human genome

2%

98%

Control of gene expression

exon1   exon2   exon3

Non-coding functional RNA

Introns

Telomeres

Pseudogenes

Repeat sequences, transposons

Cis-regulatory elements

Control of gene expression

**Cis-regulatory elements**

Control of gene expression

Cis-regulatory elements

Walter Schaffner Biol. Chem. 2015

Cis-regulatory elements

Enhancer

Enhancer

Enhancer

P

Exon    Intron    Exon

Silencer

Silencer

Insulator/
boundary element

Transcription unit

insulator/
boundary element

Walter Schaffner Biol. Chem. 2015

CACGTGG    CACgTG    AACAAT

Cis-regulatory elements

Enhancer

Enhancer

Enhancer

Silencer

P

Exon — Intron — Exon

Transcription unit

Insulator/ boundary element

insulator/ boundary element

Silencer

Walter Schaffner Biol. Chem. 2015

cohesin

CBP/p300

Mediator

GTFs

RNA Pol II

H3K4me3

Problems

**a Motif orientation**

Cooperative binding — Site 2 is non-palindromic
Site 1    Site 2
5'–CACTTGA–3'

No cooperative binding
Site 1    Site 2
5'–TCAAGTG–3'

**b Motif relative position**

Cooperative binding
Site 3    Site 1    Site 2
5'–3'

No cooperative binding
Site 2    Site 1    Site 3
5'–3'

**c Motif distance: helical phasing**

Cooperative binding — Sites are separated by complete helical turns
10 bp (or 20 bp, 30 bp etc.)
Site 1    Site 2

No cooperative binding — TFs are on opposite faces of the helix
5 bp (or 15 bp, 25 bp etc.)
Site 2
Site 1

Spitz et al.

# Problems

**a** Enhanceosome

**b** Billboard

Enhancer 1

Enhancer 1

Enhancer 1

**c** TF collective

Enhancer 1

Enhancer 2

Enhancer 3

Spitz et al.

# Problems



**a Enhanceosome**

**b Billboard**

Enhancer 1

Enhancer 1

**c TF collective**

Enhancer 1

Enhancer 2

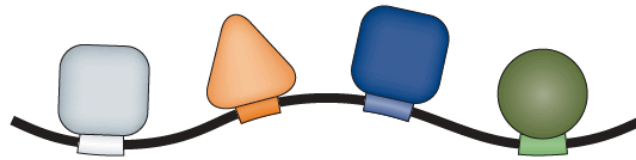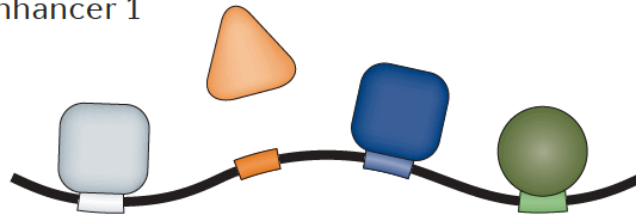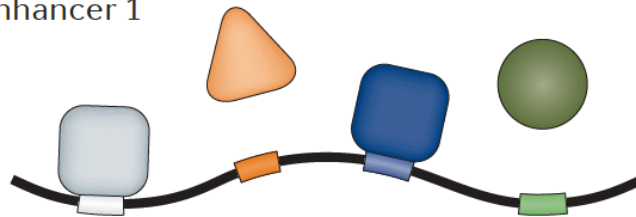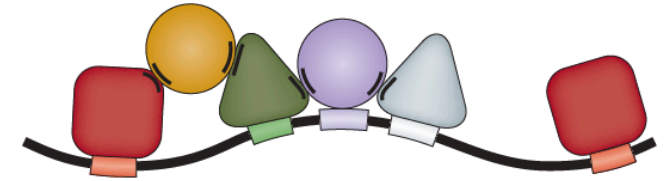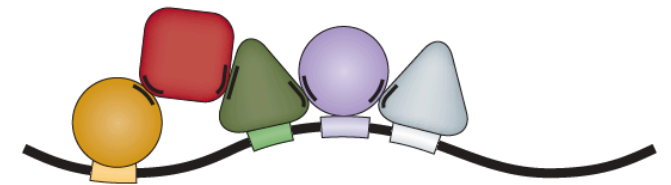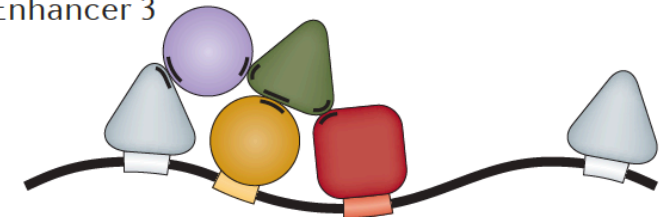| | Enhanceosome | Billboard | TF collective |
|---|---|---|---|
| **Protein DNA** | • Highly cooperative DNA binding<br>• DNA sequence acts as scaffold | • Cooperative and additive binding | • Cooperative DNA binding<br>• Both DNA and protein may act as scaffold |
| **Protein interface** | • Fixed (formed from a higher-order TF–DNA complex) | • Variable | • Variable |
| **Motif** | • Fixed motif composition (sites for all factors must be there)<br>• Fixed motif positioning (grammar) | • Fixed motif composition<br>• Flexible motif grammar | • Flexible motif composition (as different TFs directly bind to DNA)<br>• Flexible motif grammar |
| **Output** | • Unitary (requires the integrated activity of all TFs) | • Only requires a subset of factors to be active | • Collective (requires most TFs but not clear if it requires all TFs) |

Spitz et al.

# Chromatin modifications and various classes of CRMs



(a) sheared, digested, or PCR-isolated genomic DNA

(b) ATAC - seq
FAIRE - seq
DNase - seq

(c) ChIP-seq

sequence

sequence

Suryamohan et al.

# Chromatin modifications and various classes of CRMs



| Regulatory element | DHS | H3K4Me1 | H3K4Me3 | H3K27Ac |
|---|---|---|---|---|
| Promoter | + | - | + | ++ |
| Enhancer | + | + | - | ++ |
| Insulator | + | - | +/- | - |



Suryamohan et al.

# Chromatin modifications and various classes of CRMs



Suryamohan et al.

# Approaches for computational prediction of CRMs

**Comparative genomics**

- Aligned sequence

# Approaches for computational prediction of CRMs



**Comparative genomics**
- Aligned sequence

**Motif based**

**Motif blind**

```
...CGGAATCACCACTGGATGCGGATACTGGGGAATCAC...
```

Li et al. Biosystems 2015

**Interpolated Markov Model**

$0^{th}$ MM, $1^{th}$ MM, ..., $n^{th}$ MM → $\omega_0$, $\omega_1$, $\omega_n$ → $n^{th}$ IMM

Kazemian et al. NAR 2011

**String kernel**

```
AGTAGGGTAGG  CAGTGATAGAT   AAATTTTCGCG
     TAGGTCAGTGA           TTTCGCGCTAT
          TAGATAGAAAT    GCGCTATCGAT
AGTAGGGTAGGTCAGTGATAGATAGAAATTTTCGCGCTATCGAT
```

# Computational prediction of CRMs using training data

**Training CRMs**

**Background CRMs**

**Model**

# Computational prediction of CRMs using training data

Training CRMs

**Chromatin profiling**  **TF ChIP-seq**  **In vivo CRM activity**

Training
CRMs

Background
CRMs

Model

DNase signal

H3K27Ac signal

MED1

TCF3

TCF4

TLE

SOX9

https://www.encodeproject.org/data/annotations/

# Computational prediction of CRMs using training data

Training CRMs

Background CRMs

Model

Cross Validation

# Computational prediction of CRMs using training data

**Training CRMs**

**Background CRMs**

**Model**

**CRM**

**Cross Validation**

M3

M1

M0

True Positive Rate

False Positive Rate

1.0

0.8

0.6

0.4

0.2

0.0

0.2    0.4    0.6    0.8    1.0

Computational prediction of CRMs using training data

# CRM prediction using machine learning methods

# CRM prediction using machine learning methods

HMM

# CRM prediction using machine learning methods

HMM

SVM

# CRM prediction using machine learning methods

**HMM**

**SVM**

**RF**

CRM prediction using machine learning methods

HMM

SVM

RF

Deep Learning

ACTAGCGACGTGTGCGACG
ACGAGCTTCGAGAGCGACG
ATAGGTTACGTTCGAGACG

Convolution

Pooling

Features

Prediction

Output

F1>3.4
F1>3.4
F15>3.6
F5<3.2
F2<5.4
F2>8.1

# Conservation to find CRMs



Jessica Alfoldi et al. Genome Research 2013

# Models for predicting transcription-factor binding sites



**a**

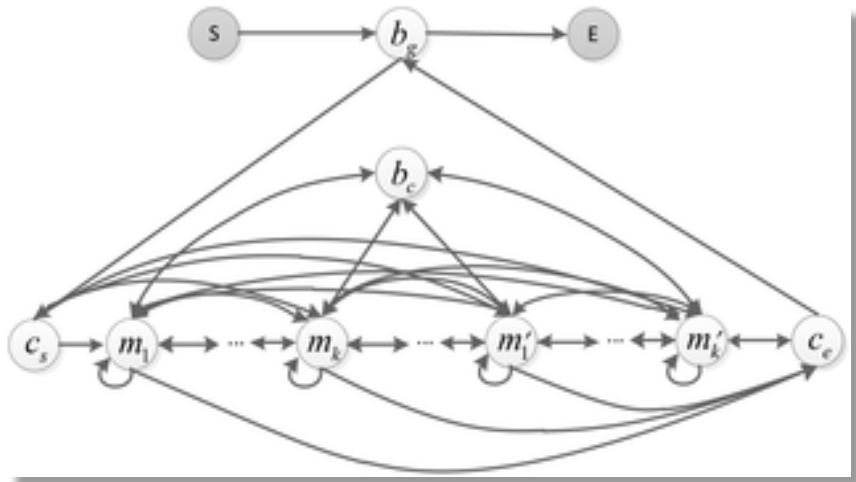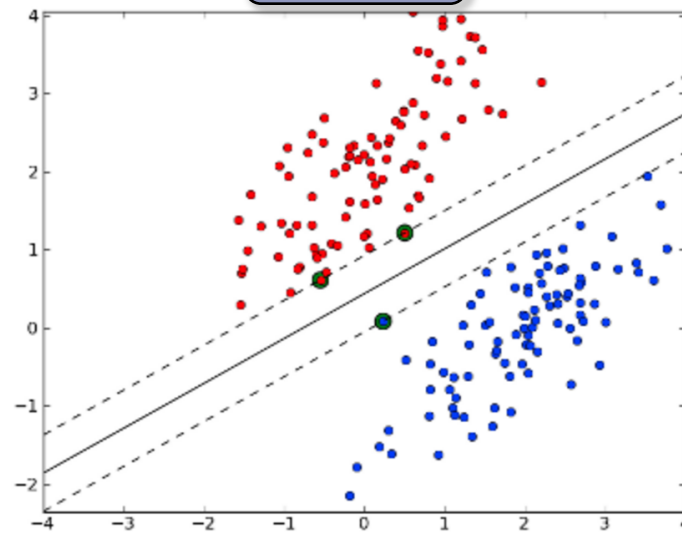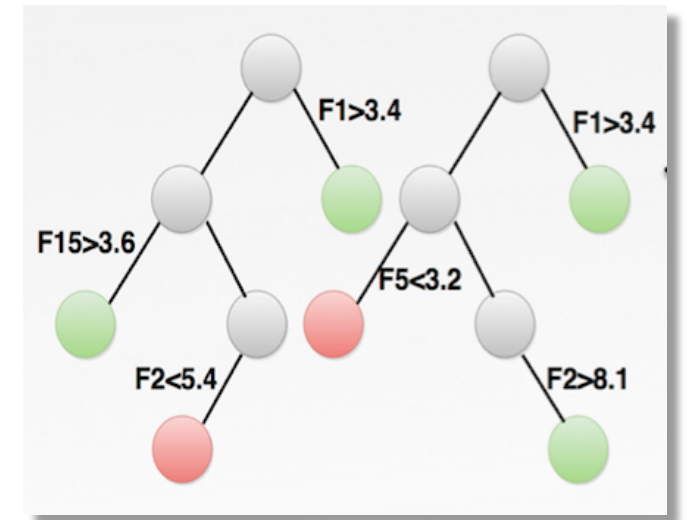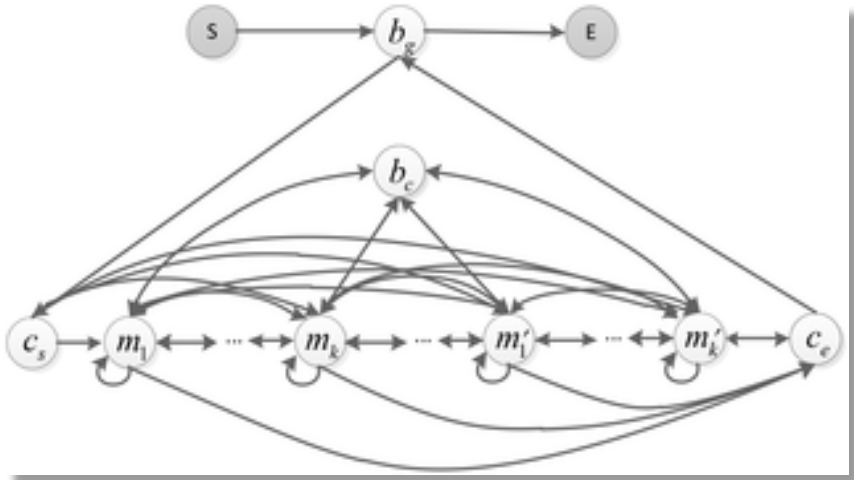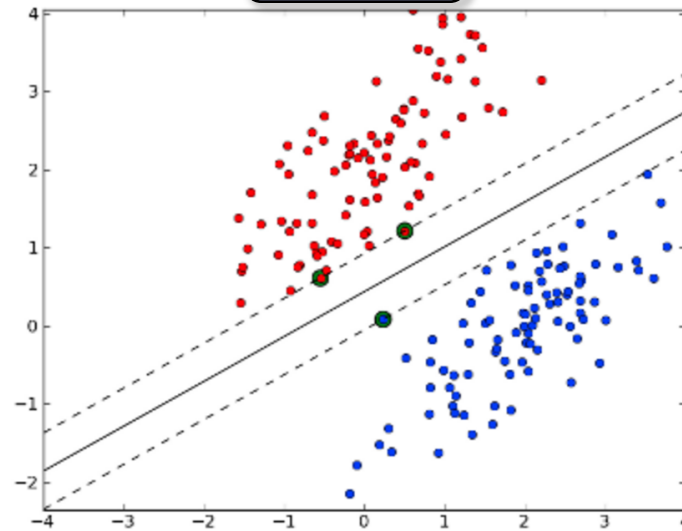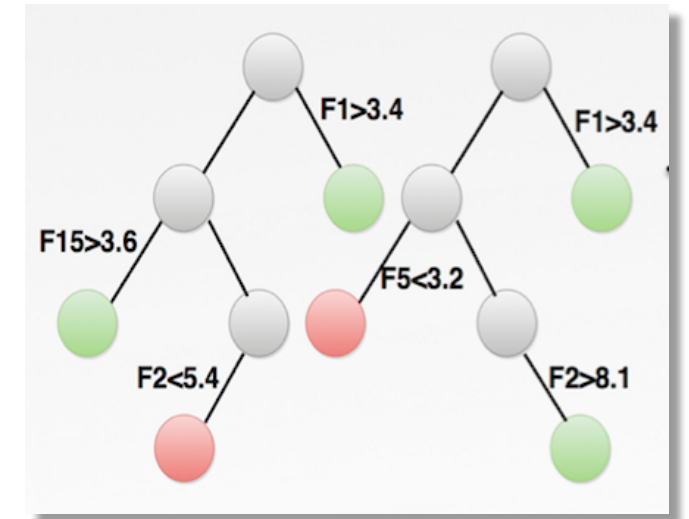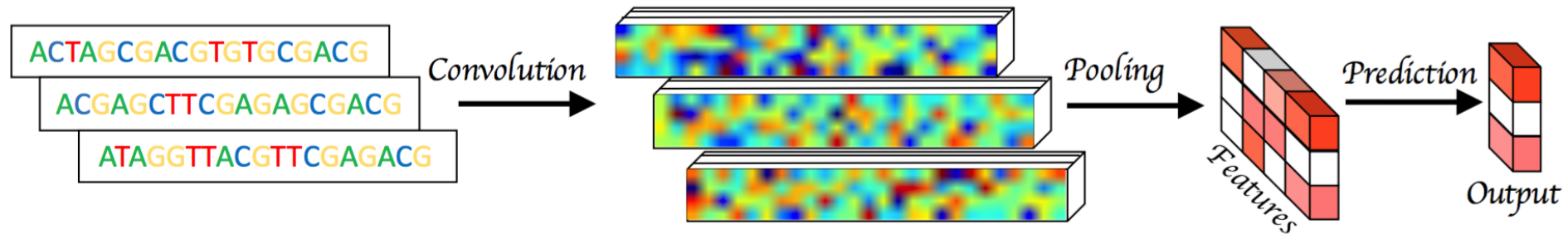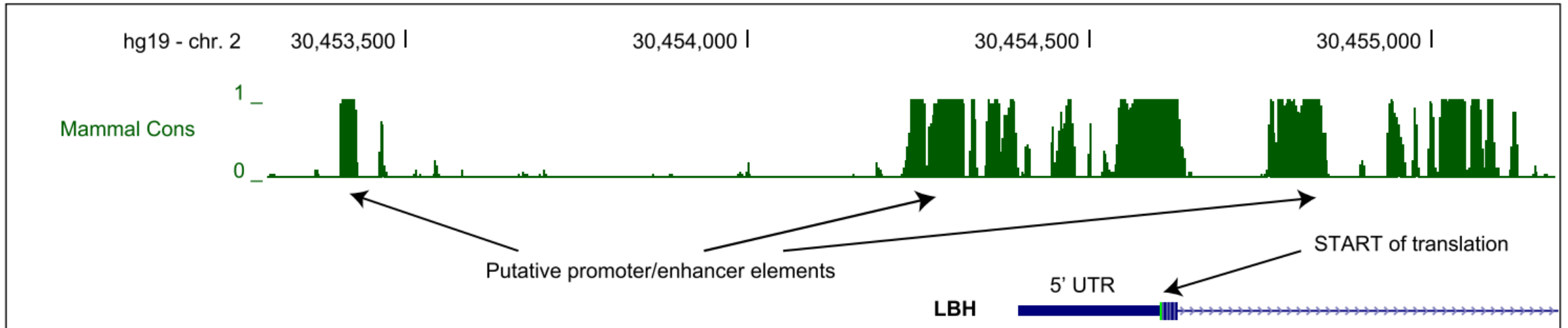|        | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|--------|---|---|---|---|---|---|---|---|---|----|----|----|----|----|
| Site 1 | G | A | C | C | A | A | A | T | A | A  | G  | G  | C  | A  |
| Site 2 | G | A | C | C | A | A | A | T | A | A  | G  | G  | C  | A  |
| Site 3 | T | G | A | C | T | A | T | A | A | A  | A  | G  | G  | A  |
| Site 4 | T | G | A | C | T | A | T | A | A | A  | A  | G  | G  | A  |
| Site 5 | T | G | C | C | A | A | A | A | G | T  | G  | G  | T  | C  |
| Site 6 | C | A | A | C | T | A | T | C | T | T  | G  | G  | G  | C  |
| Site 7 | C | A | A | C | T | A | T | C | T | T  | G  | G  | G  | C  |
| Site 8 | C | T | C | C | T | T | A | C | A | T  | G  | G  | G  | C  |

Source binding sites

**b**

| B | R | M | C | W | A | W | H | R | W | G | G | B | M |

Consensus sequence

**c** Position frequency matrix (PFM)

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|
| A | 0 | 4 | 4 | 0 | 3 | 7 | 4 | 3 | 5 | 4  | 2  | 0  | 0  | 4  |
| C | 3 | 0 | 4 | 8 | 0 | 0 | 0 | 3 | 0 | 0  | 0  | 0  | 2  | 4  |
| G | 2 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0  | 6  | 8  | 5  | 0  |
| T | 3 | 1 | 0 | 0 | 5 | 1 | 4 | 2 | 2 | 4  | 0  | 0  | 1  | 0  |

Wasserman et al. Nature Reviews 2007

# Models for predicting transcription-factor binding sites



Wasserman et al. Nature Reviews 2007

# Models for predicting transcription-factor by Coupling binding-site prediction with phylogenetic footprinting



Wasserman et al. Nature Reviews 2007

# Models using clusters of binding sites



**Window Clustering**

Significant clustering of high densities of binding sites within a sequence window
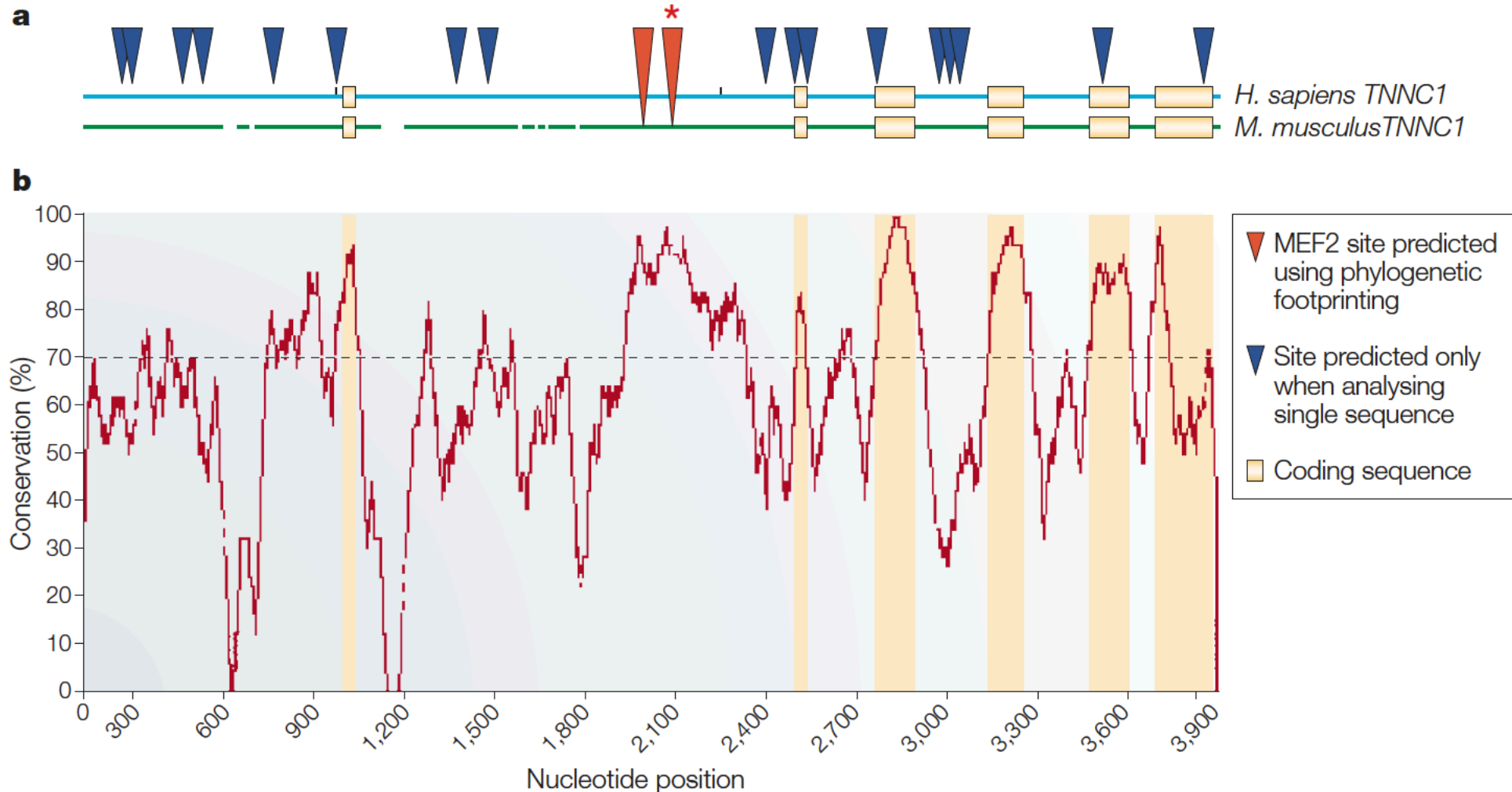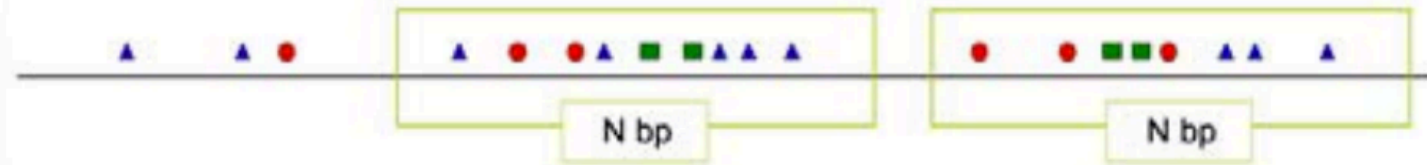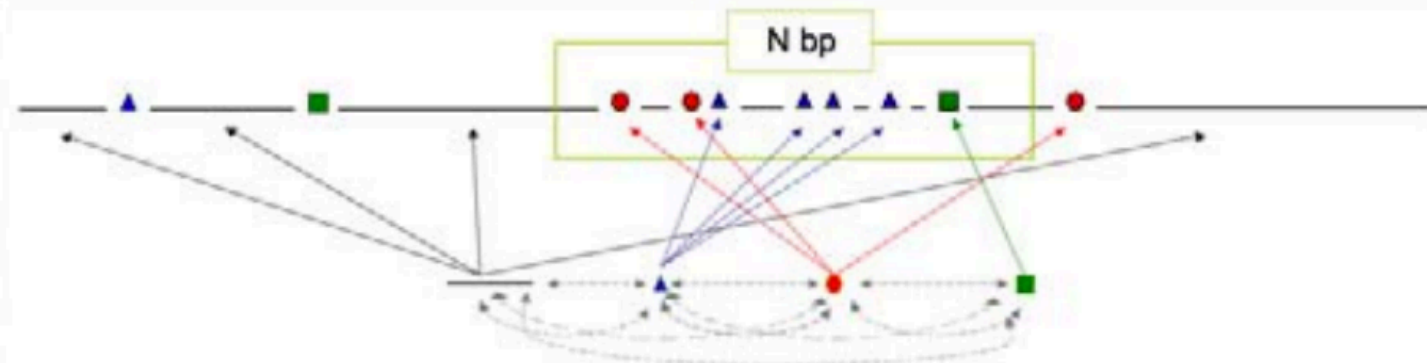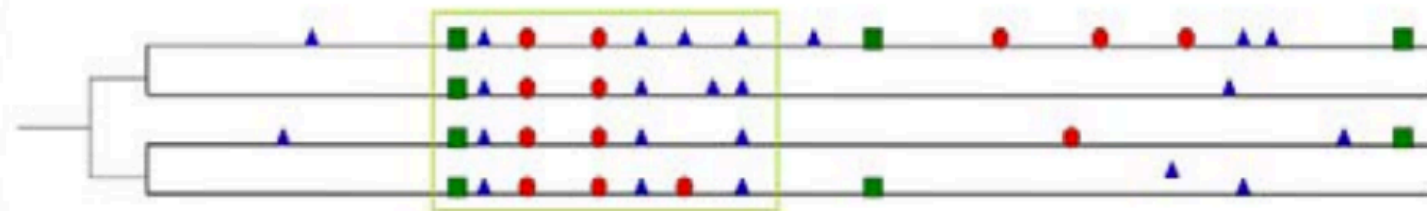
**Probabilistic Modelling**

Region of the sequence that resembles a statistical model of a binding site cluster more than a model of background DNA
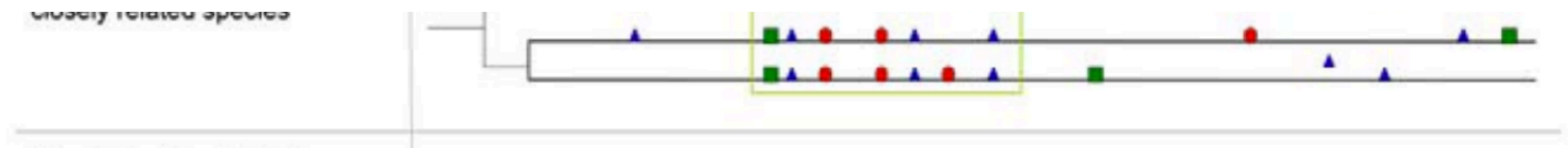
**Phylogenetic Footprinting**

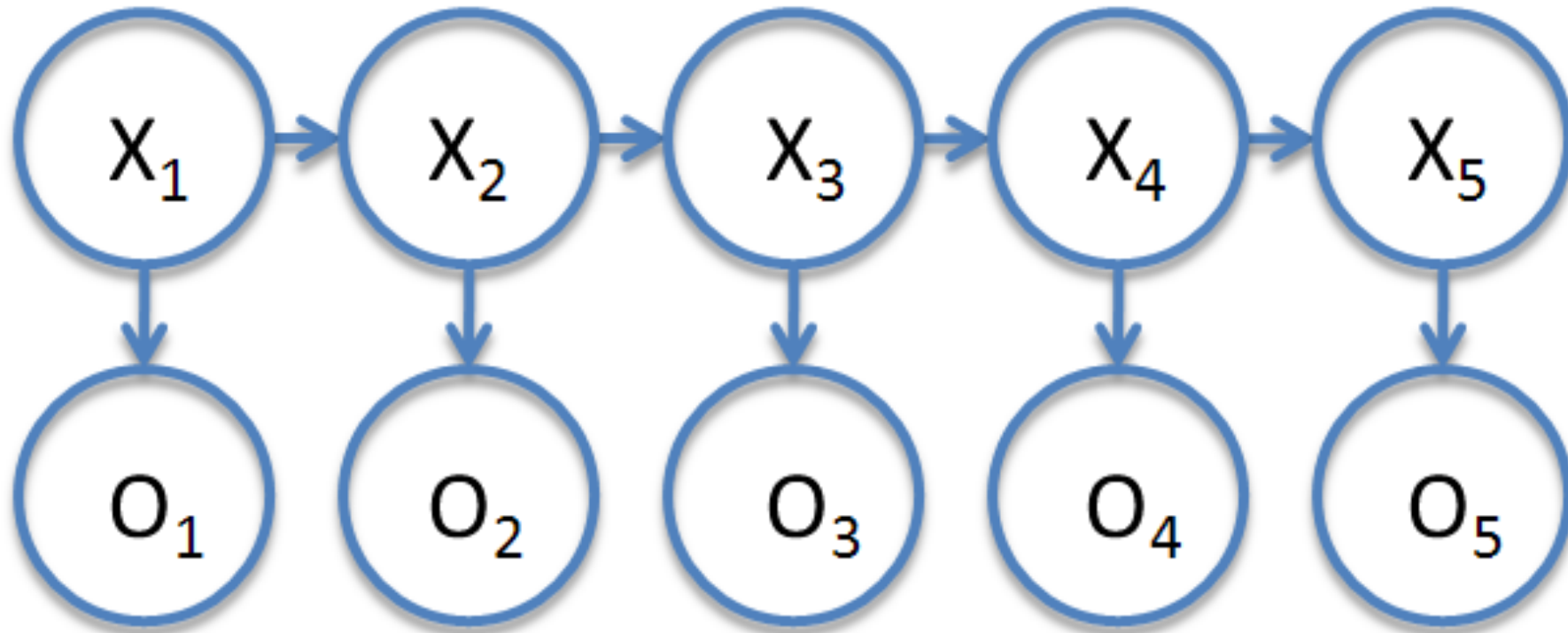High density region of binding sites conserved between closely related species

Su et al. PLoS Computational Biology 2010

# Models using clusters of binding sites

| Method | Search Strategy | | | | Input Data | | | |
|---|---|---|---|---|---|---|---|---|
| | Window Clustering | Probabilistic Modelling | Phylogenetic Footprinting | Discriminative Modelling | Single Genome | Multiple Alignment | Motif Library | CRM Annotation |
| MSCAN | ✓ | | | | ✓ | | ✓ | |
| MCAST | ✓ | | | | ✓ | | ✓ | |
| ClusterBuster | | ✓ | | | ✓ | | ✓ | |
| Stubb | | ✓ | | | ✓ | | ✓ | |
| StubbMS | | ✓ | ✓ | | | ✓ | ✓ | |
| MorphMS | | ✓ | ✓ | | | ✓ | ✓ | |
| CisModule | | ✓ | | | ✓ | | ✓ | |
| MultiModule | | ✓ | ✓ | | | ✓ | ✓ | |
| CisPlusFinder | ✓ | | ✓ | | | ✓ | | |
| EEL | ✓ | | ✓ | | | ✓ | ✓ | |
| RP | | ✓ | ✓ | ✓ | | ✓ | | ✓ |
| HexDiff | | | | ✓ | ✓ | | | ✓ |
| PhylCRM | | ✓ | ✓ | | | ✓ | ✓ | |
| EMMA | | ✓ | ✓ | | | ✓ | | ✓ |

Su et al. PLoS Computational Biology 2010

# Models using clusters of binding sites



Guo et al. PLoS One 2016

# Computational prediction of CRMs using training data

**Training CRMs**



**Training CRMs**

**Background CRMs**

**Model**

**Chromatin profiling**

DNase signal

H3K27Ac signal

**TF ChIP-seq**

MED1
TCF3
TCF4
TLE
SOX9

**In vivo CRM activit**

https://www.encodeproject.org/data/annotations/
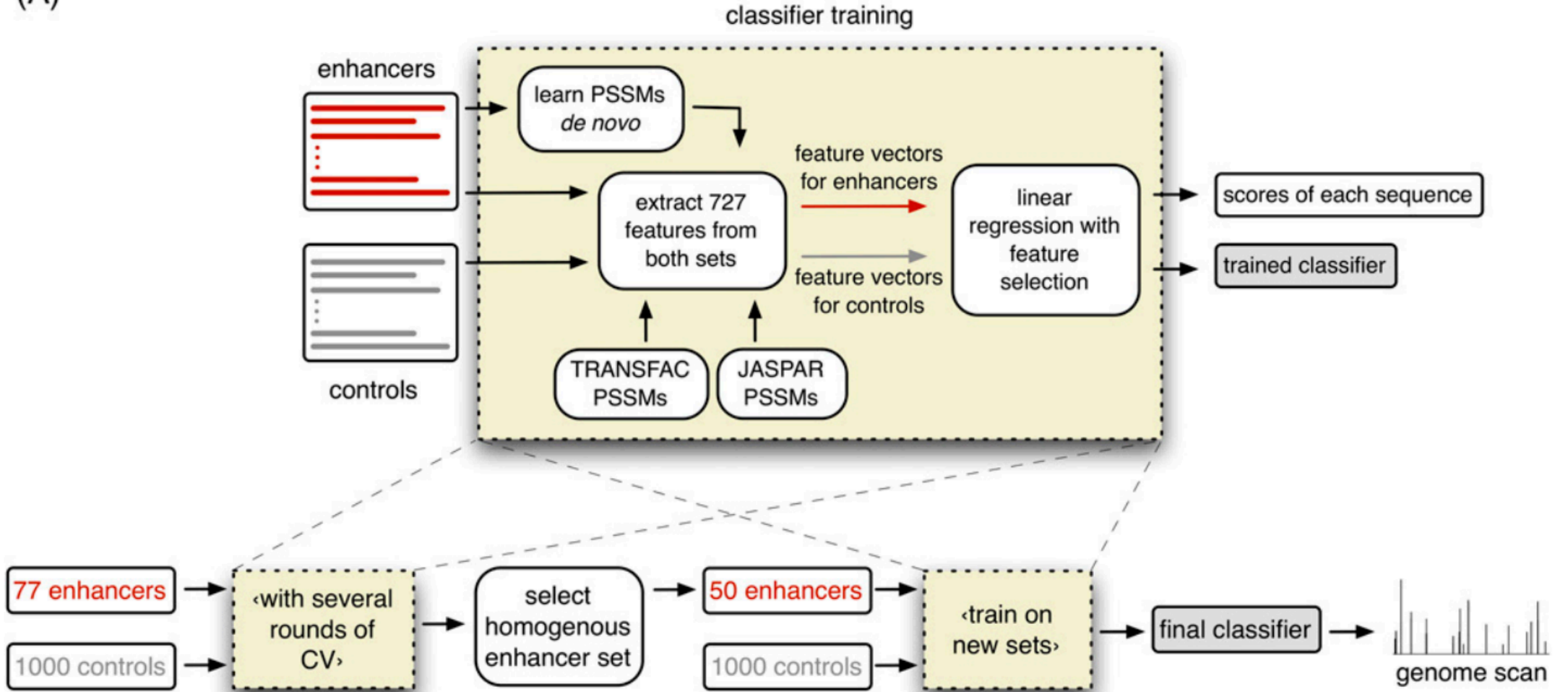
# Genome-wide discovery of human heart enhancers

Leelavati Narlikar,[1] Noboru J. Sakabe,[2] Alexander A. Blanski,[2] Fabio E. Arimura,[2] John M. Westlund,[2] Marcelo A. Nobrega,[2,3] and Ivan Ovcharenko[1,3]

[1]Computational Biology Branch, National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health (NIH), Bethesda, Maryland 20894, USA; [2]Department of Human Genetics, The University of Chicago, Chicago, Illinois 60637, USA
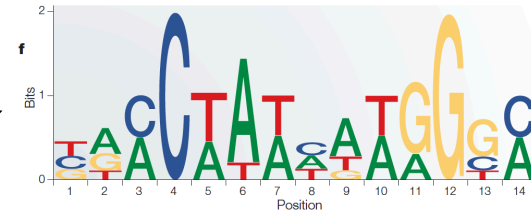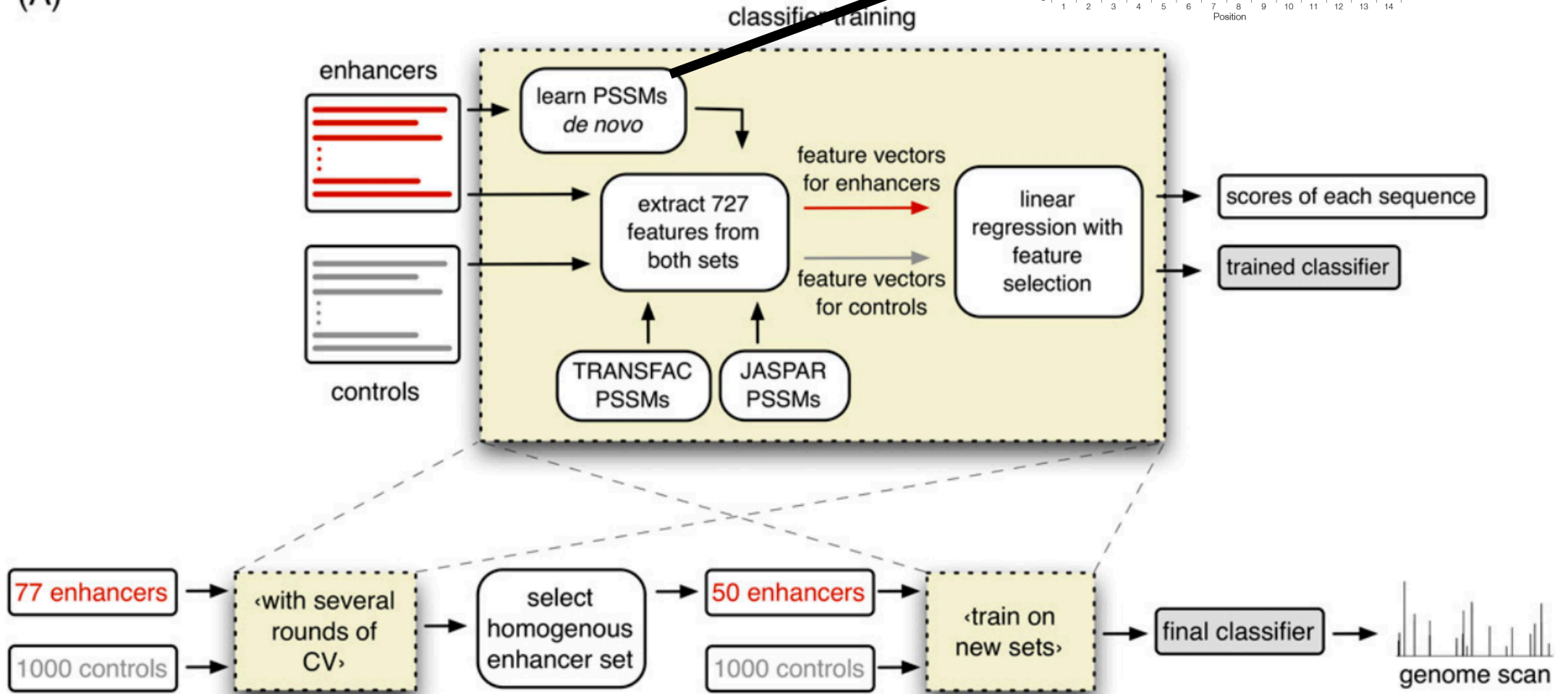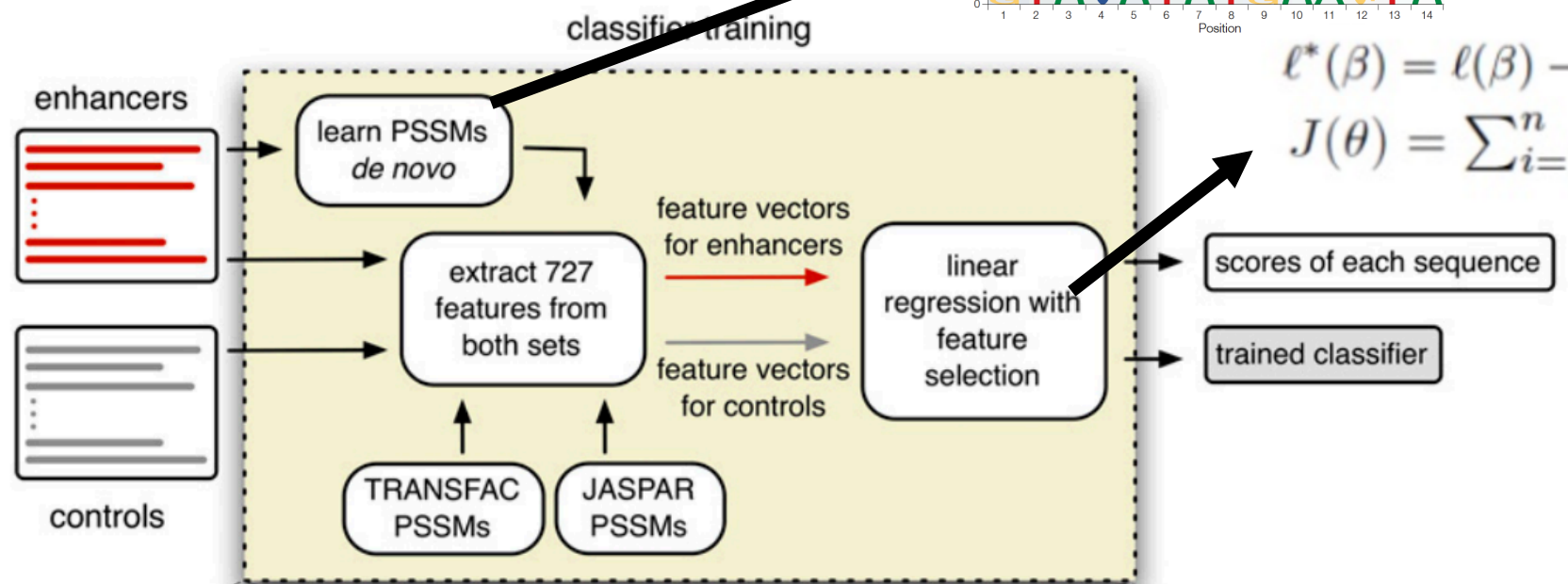
**Case 1**

(A)

classifier training

enhancers → learn PSSMs *de novo*

extract 727 features from both sets

feature vectors for enhancers → linear regression with feature selection → scores of each sequence

feature vectors for controls → trained classifier

TRANSFAC PSSMs    JASPAR PSSMs

controls

77 enhancers → ‹with several rounds of CV› → select homogenous enhancer set → 50 enhancers → ‹train on new sets› → final classifier → genome scan

1000 controls    1000 controls

Narlikar et al.

Case 1

(A)

classifier training

enhancers

learn PSSMs *de novo*

feature vectors for enhancers

extract 727 features from both sets

feature vectors for controls

linear regression with feature selection

scores of each sequence

trained classifier

TRANSFAC PSSMs

JASPAR PSSMs

controls

77 enhancers → ‹with several rounds of CV› → select homogenous enhancer set → 50 enhancers → ‹train on new sets› → final classifier → genome scan

1000 controls →

1000 controls →

Narlikar et al.

Case 1

(A)

classifier training

enhancers

learn PSSMs de novo

extract 727 features from both sets

feature vectors for enhancers

feature vectors for controls

TRANSFAC PSSMs

JASPAR PSSMs

controls

linear regression with feature selection

scores of each sequence

trained classifier

$$\ell^*(\beta) = \ell(\beta) - \lambda J(\theta),$$

$$J(\theta) = \sum_{i=1}^{n} \|\beta_i\|$$

77 enhancers → ‹with several rounds of CV› → select homogenous enhancer set → 50 enhancers → ‹train on new sets› → final classifier → genome scan

1000 controls

1000 controls

Narlikar et al.

Case 1

Narlikar et al.

Case 1

weight

chr5:172,109,161-172,109,994

chr10:29,208,862-29,210,340

chr15:32,876,216-32,876,751

chr8:30,389,108-30,389,687

MEF2A SRF LMO2 D-A TATA D-B Markov-5 ETS GATA1 D-C GATA2 XBP1 OCT API D-D UF1H3B TP53 NFAT MEF3 E2F1 MTATA CAC-BP NKX61 MYOGNF1 GC box D-E GLI SMAD4 PBX1 MYC-MAX WT1 E2A IRF Markov-1 ZF5

Narlikar et al.

# Discriminative prediction of mammalian enhancers from DNA sequence

Dongwon Lee,[1] Rachel Karchin,[1,2] and Michael A. Beer[1,3,4]

[1]*Department of Biomedical Engineering, Johns Hopkins University, Baltimore, Maryland 21205, USA;* [2]*Institute for Computational Medicine, Johns Hopkins University, Baltimore, Maryland 21218, USA;* [3]*McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University, Baltimore, Maryland 21205, USA*

# Discriminative prediction of mammalian enhancers from DNA sequence

Dongwon Lee,[1] Rachel Karchin,[1,2] and Michael A. Beer[1,3,4]

[1]Department of Biomedical Engineering, Johns Hopkins University, Baltimore, Maryland 21205, USA; [2]Institute for Computational Medicine, Johns Hopkins University, Baltimore, Maryland 21218, USA; [3]McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University, Baltimore, Maryland 21205, USA

**String kernel**

# Case 2



Lee et al.

Case 2

Lee et al.

Case 2

Lee et al.

## Case 2

**Table 1.** Predictive 6-mers of EP300 forebrain

### (A) Fifteen 6-mers with the largest positive SVM weights

| 6-mers | Reverse complement | SVM weight | Database family match | Top matched transcription factors (q-val < 0.1) |
|---|---|---|---|---|
| AATGAG | CTCATT | 3.94 | Homeodomain | POU6F1 |
| AATTAG | CTAATT | 3.85 | Homeodomain | VSX2, PRRX2, EVX2, PDX1, GBX2 |
| AGCTGC | GCAGCT | 3.65 | HLH | NHLH1, HEN1, ASCL2, REPIN1, TCF3 |
| CAATTA | TAATTG | 3.62 | Homeodomain | BARHL2, PRRX2, NKX2-5, NKX6-1, BARHL1 |
| CAGCTG | CAGCTG | 3.32 | HLH | NHLH1, HEN1, REPIN1, ASCL2, MYOD1, TCF3 |
| ACAAAG | CTTTGT | 3.29 | SOX | SOX4, SOX11, SOX10, HNF4A |
| TAATTA | TAATTA | 3.24 | Homeodomain | OTP, PROP1, HOXA, ALX1, LHX3 |
| CAGATG | CATCTG | 3.15 | HLH | ZFP238, TAL1:TCF3, TAL1:TCF4, TCF3 |
| TAATGA | TCATTA | 3.03 | Homeodomain | POU6F1, POU4F3, LHX3, HOXC9, NKX6-3 |
| AATTAA | TTAATT | 2.94 | Homeodomain | LHX3, OTP, PRRX2, PROP1, LHX5 |
| ATTAGC | GCTAAT | 2.90 | Homeodomain | VSX2, POU3F2, EVX2, PITX3, LHX8 |
| GGCAAC | GTTGCC | 2.86 | — | — |
| ACAATG | CATTGT | 2.63 | SOX | SOX17, SOX9, SOX5, SOX10, SOX30 |
| CATTCA | TGAATG | 2.45 | SOX | HBP1 |
| AATTAC | GTAATT | 2.18 | Homeodomain | PRRX2, HOXA6, HOXA1, HOXC8, DLX1 |

### (B) Five 6-mers with the largest negative SVM weights

| 6-mers | Reverse complement | SVM weight | Database family match | Top matched transcription factors (q-val < 0.1) |
|---|---|---|---|---|
| AGGTAG | CTACCT | −1.79 | – | – |
| AAGTCA | TGACTT | −1.89 | – | – |
| AGGTGA | TCACCT | −1.97 | Zinc-finger | ZEB1 |
| ACCTGG | CCAGGT | −2.03 | Zinc-finger | ZEB1, TCF3 |
| CAGGTA | TACCTG | −2.06 | Zinc-finger | ZEB1 |

Lee et al.

Case 3

# Enhanced Regulatory Sequence Prediction Using Gapped *k*-mer Features

Mahmoud Ghandi[1][☯][¤], Dongwon Lee[1][☯], Morteza Mohammad-Noori[2,3], Michael A. Beer[1,4]*

1 Department of Biomedical Engineering, Johns Hopkins University, Baltimore, Maryland, United States of America, 2 School of Mathematics, Statistics and Computer Science, University of Tehran, Tehran, Iran, 3 School of Computer Science, Institute for Research in Fundamental Sciences (IPM), Tehran, Iran, 4 McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University, Baltimore, Maryland, United States of America

Ghandi et al.

Case 4

Deep Learning

ACTAGCGACGTGTGCGACG
ACGAGCTTCGAGAGCGACG
ATAGGTTACGTTCGAGACG

Convolution

Pooling

Features

Prediction

Output

# Case 4

## Label
● / ✗

## Layer 2
High order features

## Layer 1

## Raw data

### One Hot Code Sequence
ATTCCCGTAATCTACGATTAAGTCACAACCAAACCATGGATTACGGTCTGCGTTGGAATCAGGGCCGTGC

### Convolution Layers
Convolve filters
ReLU
Max pool

### Fully Connected Layer
Linear transformation
ReLU

### Multi-task Prediction
Linear transformation
Sigmoid
Prediction
Actual

Kelley et al. 2016

Case 4

A
DNase I hypersensitivity sites
Cells

B
mean AUC 0.895
Basset AUC
0.95
0.90
0.85
0.80
0.75
0.70
mean AUC 0.780
0.70  0.75  0.80  0.85  0.90  0.95
gkm-SVM AUC

C
True positive rate
1.0
0.8
0.6
0.4
0.2
0.0
0.0  0.2  0.4  0.6  0.8  1.0
False positive rate

K562     AUC: 0.837
H7-hESC  AUC: 0.886
Globla   AUC: 0.901
AoSMC    AUC: 0.914
H9ES     AUC: 0.927

Kelley et al. 2016

# Models based on chromatin features

# Mapping and analysis of chromatin state dynamics in nine human cell types

Jason Ernst[1,2], Pouya Kheradpour[1,2], Tarjei S. Mikkelsen[1], Noam Shoresh[1], Lucas D. Ward[1,2], Charles B. Epstein[1], Xiaolan Zhang[1], Li Wang[1], Robbyn Issner[1], Michael Coyne[1], Manching Ku[1,3,4], Timothy Durham[1], Manolis Kellis[1,2]* & Bradley E. Bernstein[1,3,4]*

Case 5
Models based on chromatin features

Ernst et al.

# Integrating Diverse Datasets Improves Developmental Enhancer Prediction

**Genevieve D. Erwin[1,2], Nir Oksenberg[2,3], Rebecca M. Truty[1], Dennis Kostka[4], Karl K. Murphy[2,3], Nadav Ahituv[2,3], Katherine S. Pollard[1,2,5]\*, John A. Capra[6]\***

**1** Gladstone Institute of Cardiovascular Disease, San Francisco, California, United States of America, **2** Institute for Human Genetics, University of California San Francisco, San Francisco, California, United States of America, **3** Department of Bioengineering and Therapeutic Sciences, University of California San Francisco, San Francisco, California, United States of America, **4** Department of Developmental Biology and Department of Computational and Systems Biology, University of Pittsburgh, Pittsburgh, Pennsylvania, United States of America, **5** Department of Epidemiology and Biostatistics, University of California San Francisco, San Francisco, California, United States of America, **6** Center for Human Genetics Research and Department of Biomedical Informatics, Vanderbilt University, Nashville, Tennessee, United States of America

**Case 6**

**MKL training:**

Genome

ATTC
ATAG

ATAG

AATC
ATCC

MKL algorithm

Feature 2

Feature 1

10-fold cross validation

ROC evaluation

Trained classifier

Feature 2

Feature 1

$$f(\mathbf{x}) = \sum_{i=1}^{N} \alpha_i \sum_{j=1}^{M} \beta_j k_j(\mathbf{x}, \mathbf{x_i}) + b$$

Case 6

# Case 7
## Prediction of functional TFBS using bag of motifs

**Select functional CRMs for training:**

Target gene 1

Target gene 2

10 kB

Feature selection for training data

**Input**
Functional CRMs

**Output**
DNA motifs

10 most enriched motifs of the TF

10 most enriched co-regulatory TF motifs

Data tracks

Top 5 enriched DHS tracks

Top 5 active chromatin marks tracks

Top 5 co-regulatory TF tracks

Cross-validation

| | |
|---|---|
| M0 | Single PWM |
| Mk | kmer-SVM |
| M1 | RF with multiple PWMs |
| M2 | RF with ChIP/DHS |
| M3 | RF with multiple PWMs and ChIP/DHS |

Cross-validation

| | |
|---|---|
| M0 | Single PWM |
| Mk | kmer-SVM |
| M1 | RF with multiple PWMs |

# CHEQ-seq shows distinct sets of active enhancers



MAplot barcode expression

CHEQ-seq positives show marks of functional enhancers

p53 ENHANCERS

| NS | Nutlin | NS | Nutlin | NS | estradiol |
|----|--------|-----|--------|-----|-----------|

**p53 ChIP-seq**   **H3K27ac**   **DHS-seq**

CHEQ-seq positives show marks of functional enhancers

**p53 ENHANCERS**

**NEGATIVES**

| NS | Nutlin | NS | Nutlin | NS | estradiol |
|---|---|---|---|---|---|
| **p53 ChIP-seq** | | **H3K27ac** | | **DHS-seq** | |

| NS | Nutlin | NS | Nutlin | NS | estradiol |
|---|---|---|---|---|---|
| **p53 ChIP-seq** | | **H3K27ac** | | **DHS-seq** | |

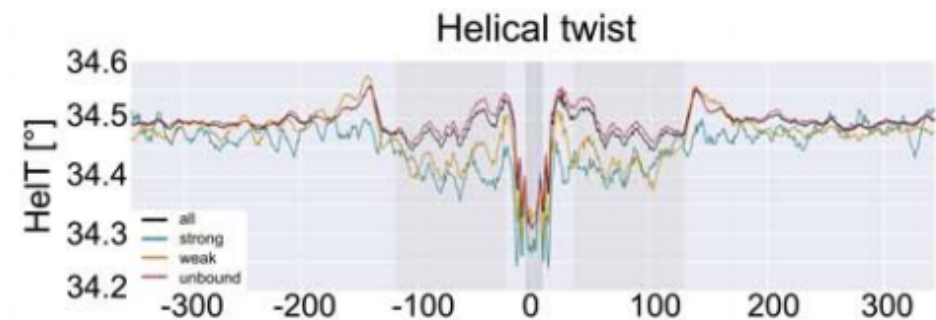Enhancer classification using motifs

# Enhancer classification using motifs

# Enhancer classification using motifs

# Enhancer classification using motifs



HOMER
- denovo len20
- denovo len19
- known

RSAT_peak motif
- 6nt_m4
- 7nt_m1
- halfsite
- dyad

i-cisTarget
- M01655
- M01656
- M00034

Sequence 1 $= \begin{bmatrix} X_{1,1} & X_{1,2} & \cdots\cdots\cdots\cdots & X_{1,p} \end{bmatrix}$

Sequence 2 $= \begin{bmatrix} X_{2,1} & X_{2,2} & \cdots\cdots\cdots\cdots & X_{2,p} \end{bmatrix}$

Sequence 3 $= \begin{bmatrix} X_{n,1} & X_{n,2} & \cdots\cdots\cdots\cdots & X_{n,p} \end{bmatrix}$

F1>3.4
F1>3.4
F15>3.6
F5<3.2
F2<5.4
F2>8.1

# Enhancer classification using motifs

# Enhancer classification using motifs

Enhancer classification using motifs
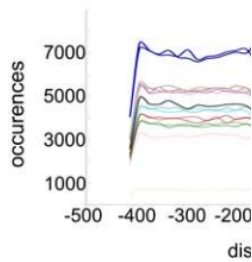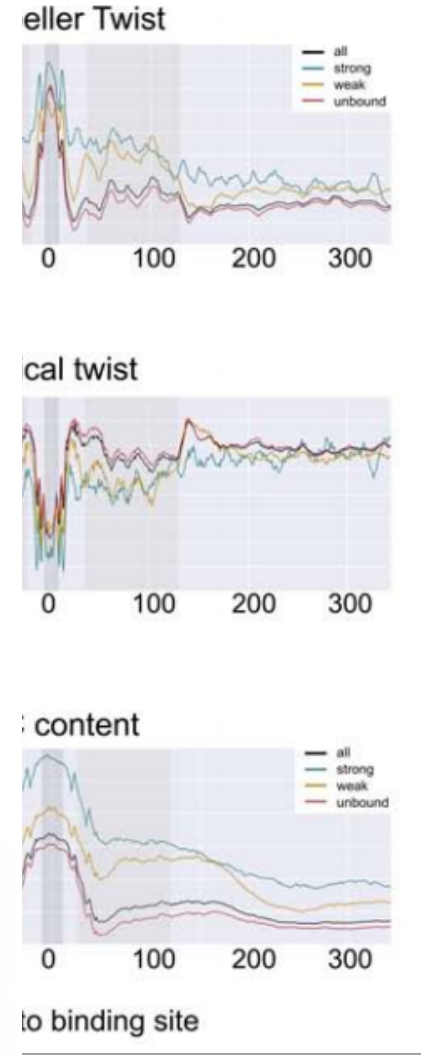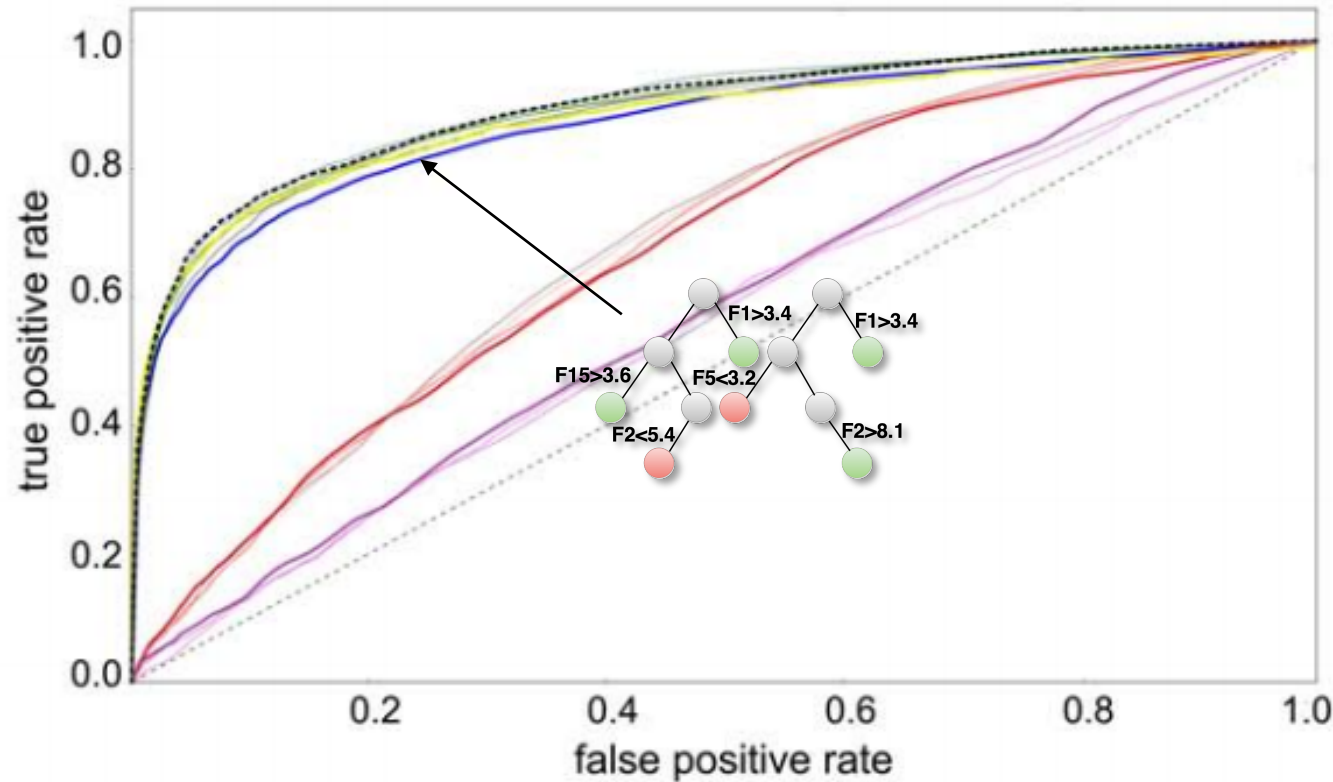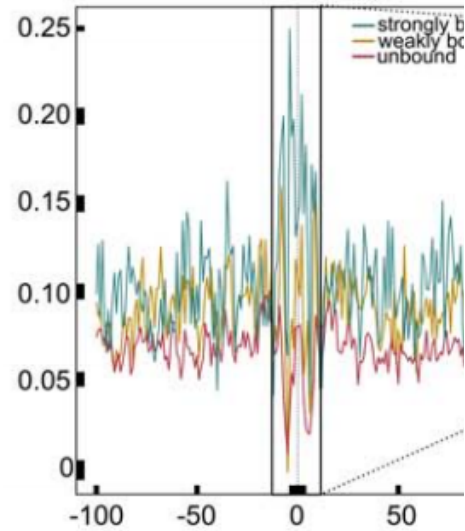
# Discriminate bound versus unbound sites in the genome

motif (0.87)
shape (0.67)
CpG islands (0.56)
deep learning (0.88)
motif + shape (0.89)
CpG-TATA-TSS (0.56)
motif + CpG island (0.89)

shape + CpG island (0.68)
shape + CpG-TATA-TSS (0.69)
motif + CpG-TATA-TSS (0.89)
motif + +shape CpG island (0.90)
CpG-TATA-TSS + CpG islands(0.57)
motif + shape + CpG-TATA-TSS (0.90)
motif + CpG-TATA-TSS +CpG island (0.89)
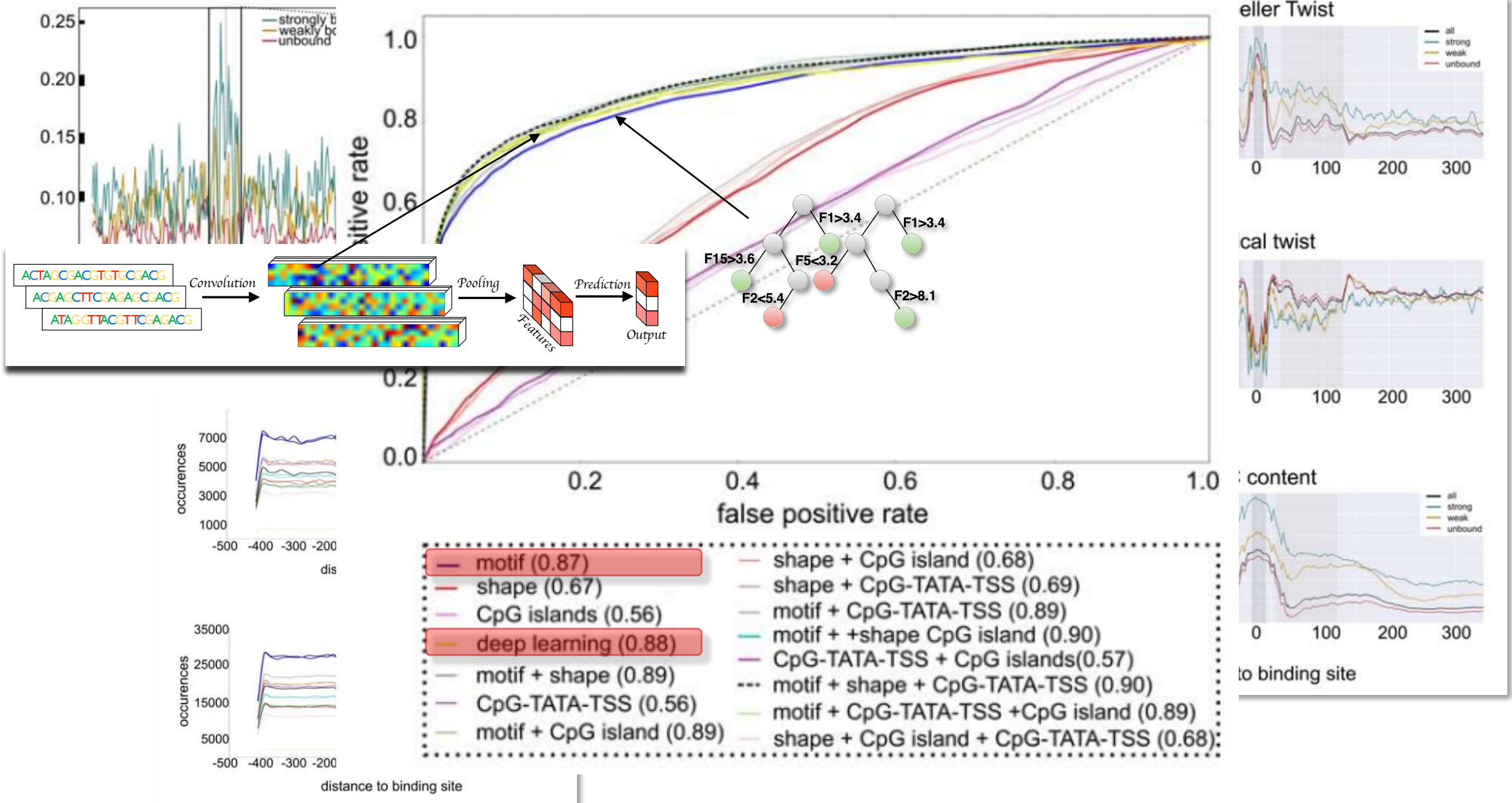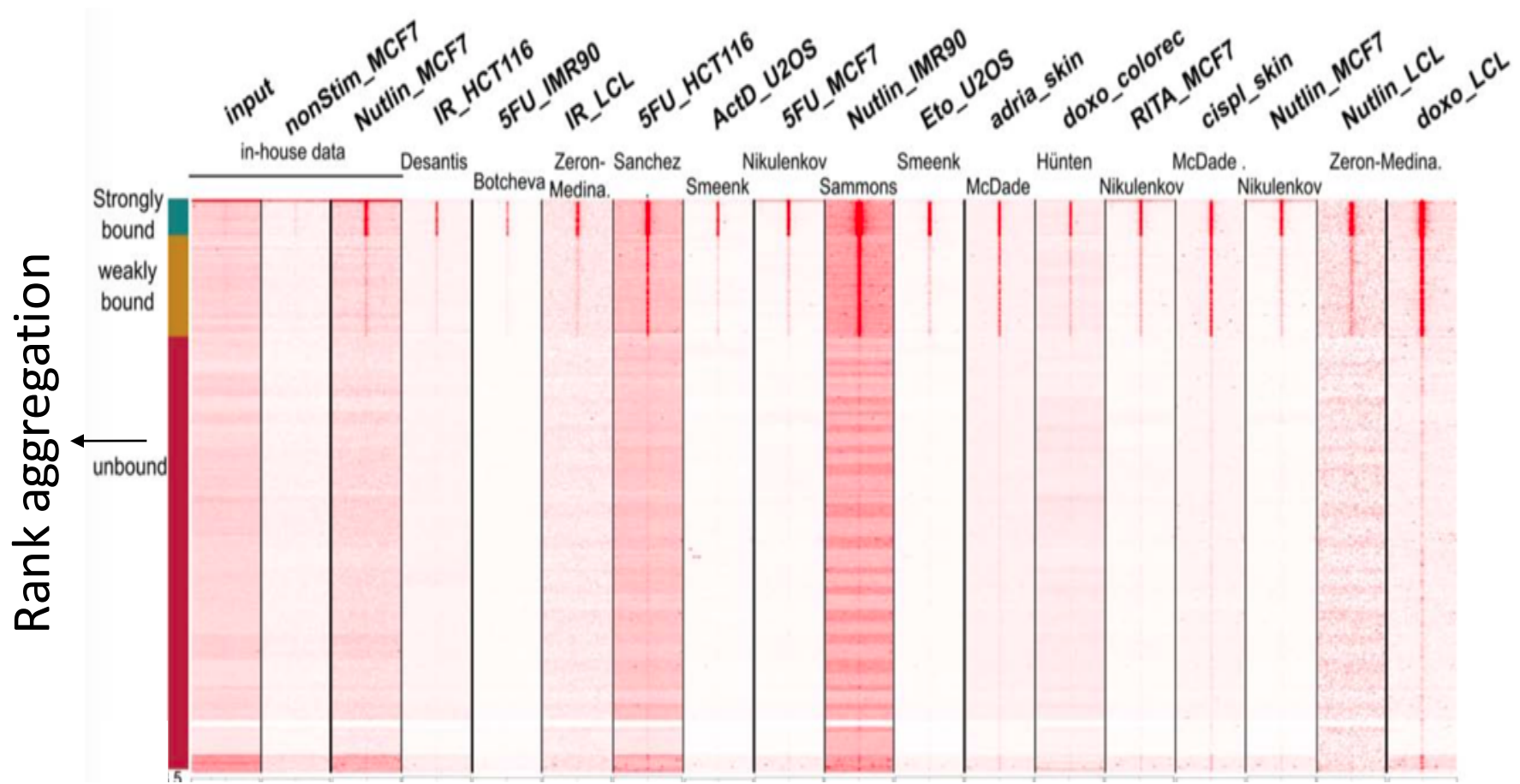shape + CpG island + CpG-TATA-TSS (0.68)

Strength of binding site predicts quantitative TP53 binding

# Strength of binding site predicts quantitative TP53 binding
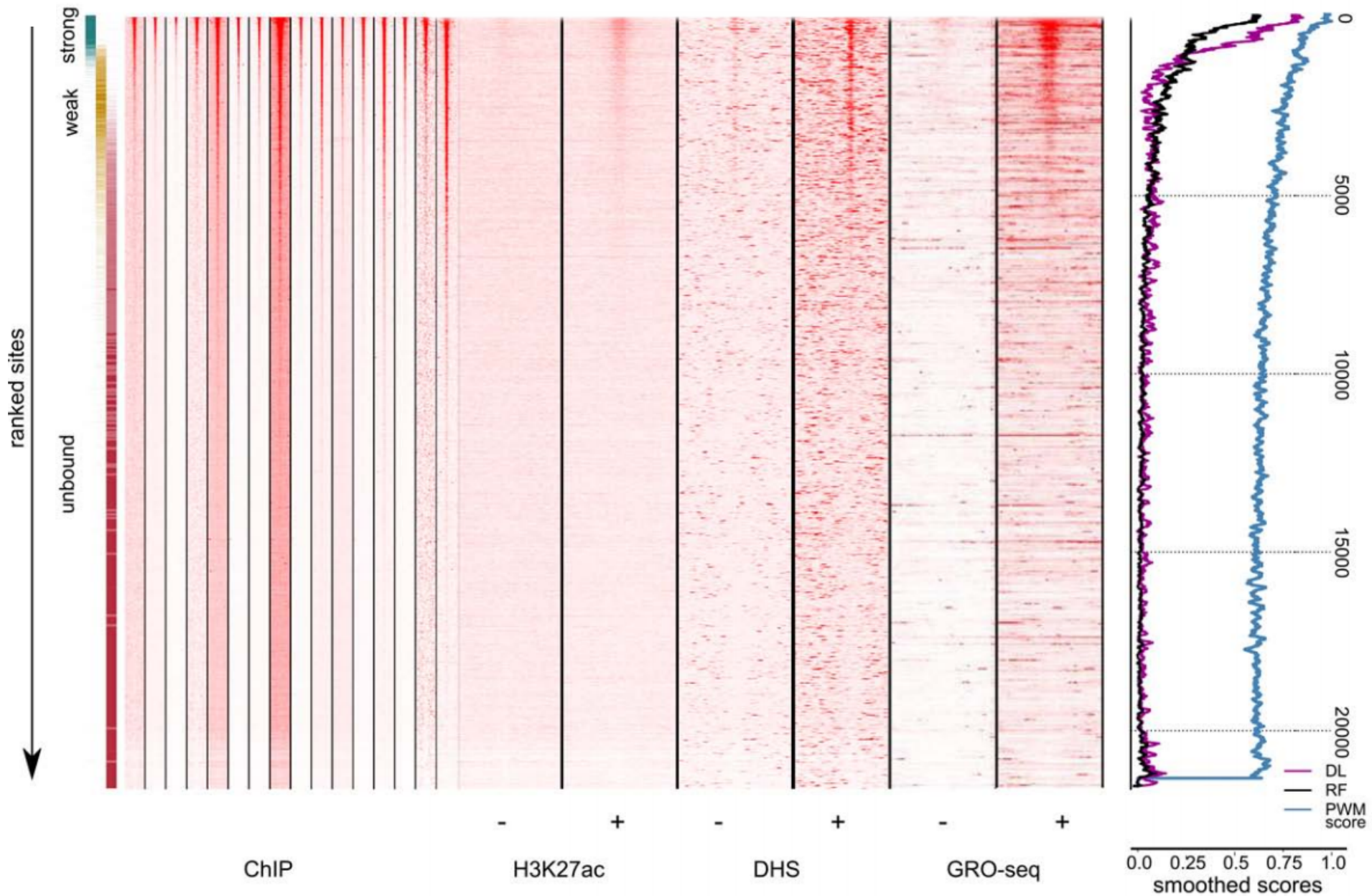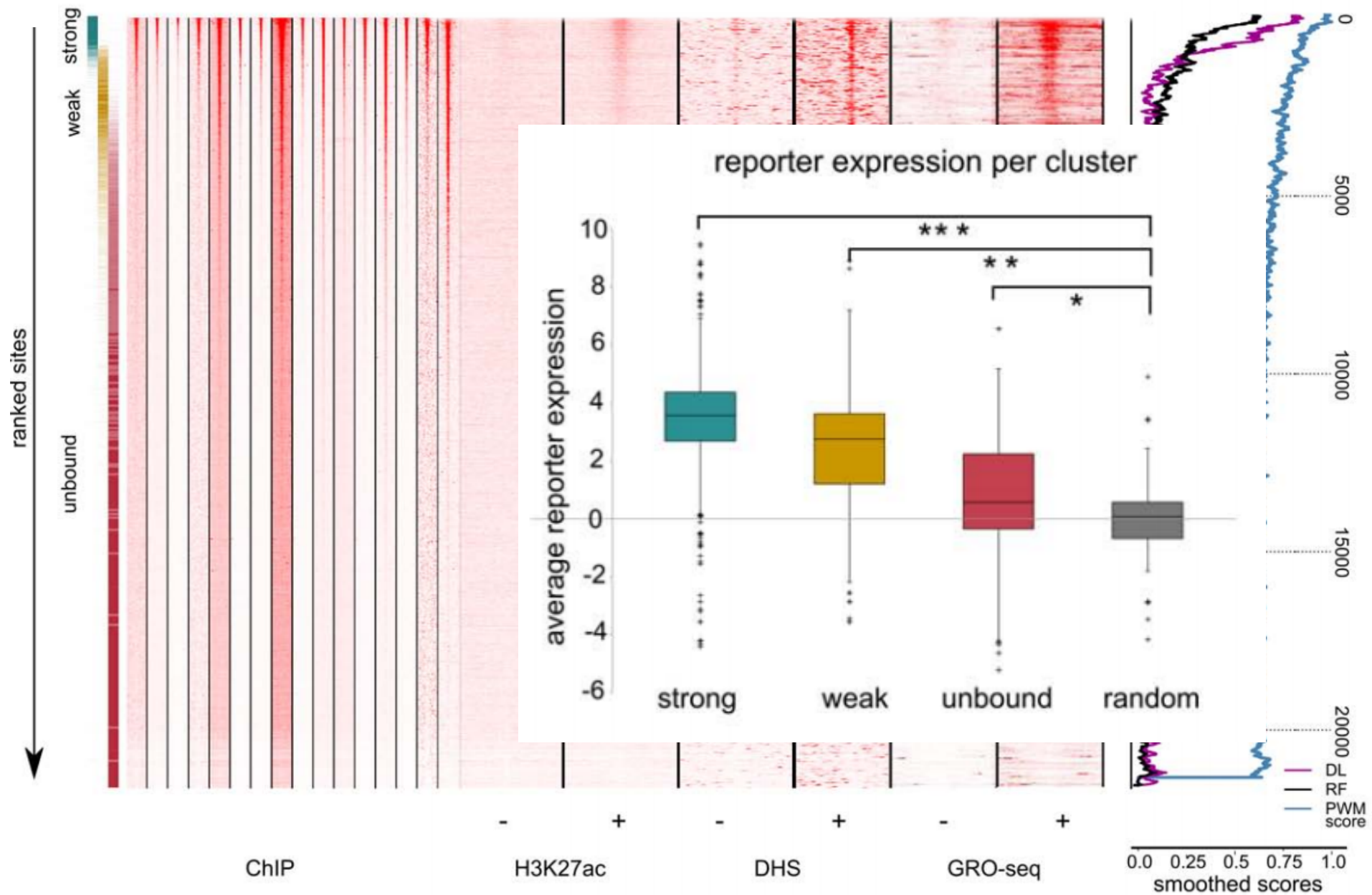
Strength of binding site predicts quantitative TP53 binding

Strength of binding site predicts quantitative TP53 binding

# Conclusions

1) Only a small subset of experimentally determined binding events represents TP53 responsive elements

2) TP53 binds the DNA strictly as a tetramer, to a duplicate of the consensus palindromic responsive element

3) Strength of binding site predicts quantitative TP53 binding

4) TP53 acts on its own, without co-regulatory transcription factors that bind to the same enhancer

5) TP53 is activator but not a repressor