

# Discriminative de novo motif discovery from high-throughput data

Jan Grau<sup>1</sup>, Stefan Posch<sup>1</sup>, Ivo Grosse<sup>1</sup>, and Jens Keilwagen<sup>2</sup>

<sup>1</sup>Institute of Computer Science, Martin Luther University Halle–Wittenberg, Germany  
<sup>2</sup>Julius Kühn-Institut (JKI) - Federal Research Centre for Cultivated Plants, Quedlinburg, Germany



CSC ChIP- and DNase-seq data analysis workshop

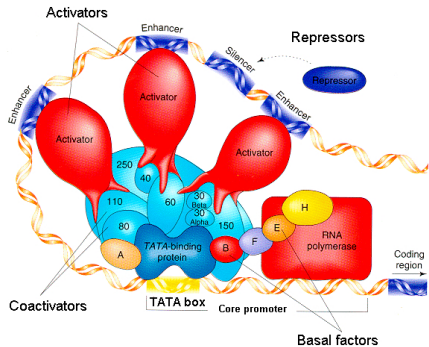
# Transcriptional regulation by transcription factors

## Biological question

- reasons for phenotypic observations
  - regulation of gene expression
  - first step: transcriptional regulation
- ⇒ transcription factor binding sites

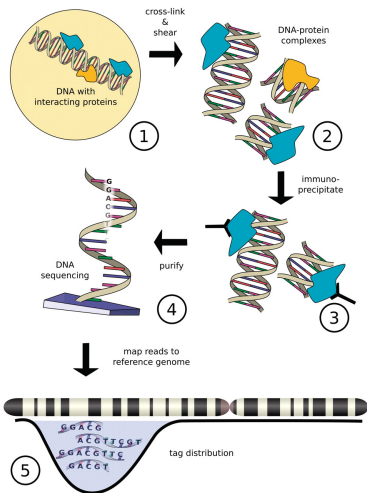
## De-novo motif discovery without knowledge of

- motif
  - exact location of sites
- from set of input sequences



[Based on Robert Tjian, "Molecular Machines that Control Genes"]

# Experimental techniques - ChIP-seq

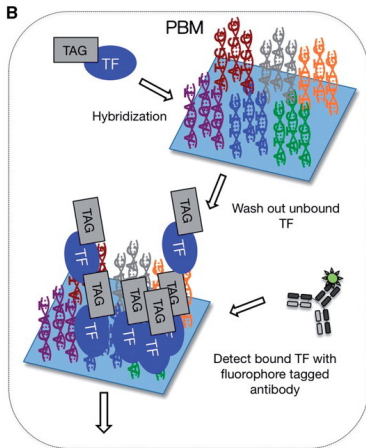


[Szalkowski & Schmid, Brief Bioinform, 2010]

## Data

- ChIP-seq peaks:
  - approximate binding regions
- ⇒ extract sequences under peaks
- ChIP-seq peak statistics:
  - information about TF abundance at binding region

# Experimental techniques - PBM



[Geertz & Maerkl, Brief Func Genom, 2010]

## Data

- PBM probes:  
contain all possible DNA  
10-mers
- ⇒ probe sequences  
(length 35 bp + linker)
- Probe intensities:  
information about TF binding  
frequency

# Requirements for a novel approach

## Use all sequences

thresholding to extract top peaks/probes arbitrary  
⇒ use all peaks and probe sequences, respectively

## Use all information present in the data

ChIP-seq sequence under peak

peak statistics

binding more likely around peak center

PBM probe sequence (including part of linker)

probe intensities

## Use discriminative learning principle

which often yield better results than generative principles

## Allow for flexible choice of motif models

e.g., position weight matrices, weight array matrices,...

## Retain acceptable runtime

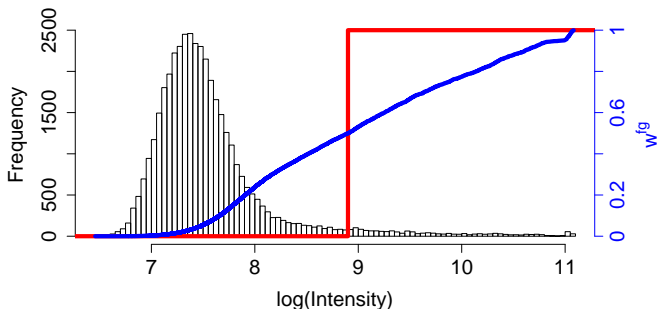
below 1h for majority of data sets

# Weighting schema for integrating ChIP and PBM data

allows for using ChIP peak statistics and PBM probe intensities in a common approach

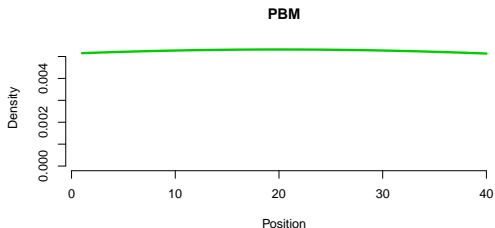
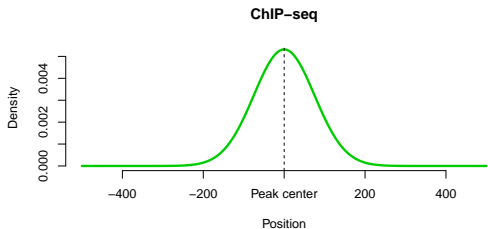
$$w_n^{fg} := \frac{1}{1 + \frac{h_n}{1-h_n} \cdot \frac{1-q}{q}}, \quad w_n^{bg} := 1 - w_n^{fg}$$

$h_n$ : relative rank of sequence  $\mathbf{x}_n$  based on peak statistic or probe intensity,  
 $q$ : weighting factor, i.e., a-priori fraction of foreground sequences



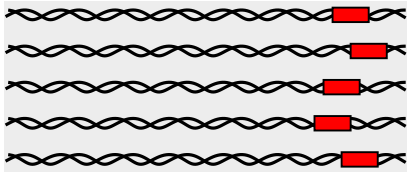
# A-priori position distribution

represents that binding occurs close to peak center



# Discriminative learning - Motivation

## ChIP-seq positives

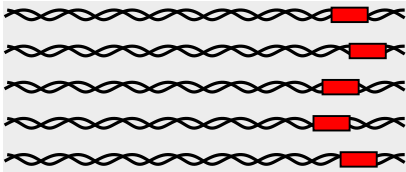


**over-represented**

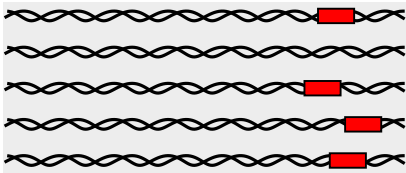


# Discriminative learning - Motivation

## ChIP-seq positives



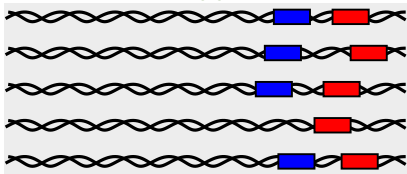
## ChIP-seq negatives



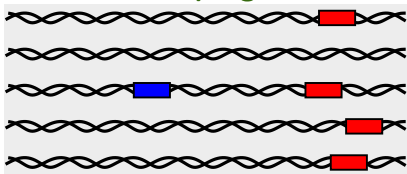
**over-represented**

# Discriminative learning - Motivation

## ChIP-seq positives



## ChIP-seq negatives



**over-represented**  
**differentially abundant**

⇒ discriminative learning

Discriminative weighted maximum supervised posterior principle

$$\hat{\lambda} = \operatorname{argmax}_{\lambda} \underbrace{\sum_{n=1}^N \sum_{c \in \mathcal{C}} w_n^c \log \left( \frac{P(c|\lambda)P_c(\mathbf{x}_n|\lambda)}{\sum_{\tilde{c} \in \mathcal{C}} P(\tilde{c}|\lambda)P_{\tilde{c}}(\mathbf{x}_n|\lambda)} \right)}_{\text{Weighted conditional likelihood}} + \underbrace{Q(\lambda|\alpha)}_{\text{Prior}},$$

where  $\mathcal{C} = \{fg, bg\}$ : set of classes,

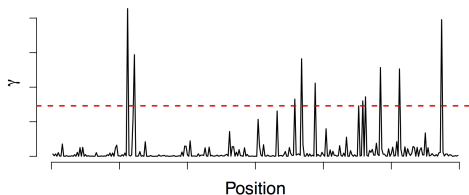
$Q(\lambda|\alpha)$ : prior on the parameters  $\lambda$  given hyper-parameters  $\alpha$ ,

$P(c|\lambda)$ : a-priori class probability, and

$P_c(\mathbf{x}_n|\lambda)$ : class-conditional likelihood, “model”

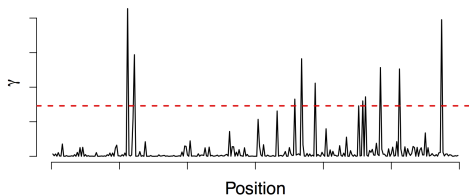
$$P_{fg}(\mathbf{x}|\boldsymbol{\lambda}) = P(\text{motif}|\boldsymbol{\lambda}) \cdot \frac{1}{|\Sigma|^{L-w}} \sum_{\ell \in \mathcal{L}} P(\ell) P_{\text{motif}}(x_{\ell}, \dots, x_{\ell+w-1}|\boldsymbol{\lambda}) \\ + (1 - P(\text{motif}|\boldsymbol{\lambda})) \cdot \frac{1}{|\Sigma|^L}$$

- Dimont uses standard ZOOPS model ( $P_{fg}(\mathbf{x}|\boldsymbol{\lambda})$ )
- sequence flanking the motif: uniform, i.e., all nucleotides with equal probability
- motif model: strand model enclosing
  - position weight matrix (PWM): assumes nucleotide independence or
  - weight array matrix (WAM): allows dependencies between neighboring nucleotides or
  - higher-order Markov models
- background model: uniform or Markov model ( $P_{bg}(\mathbf{x}|\boldsymbol{\lambda})$ )



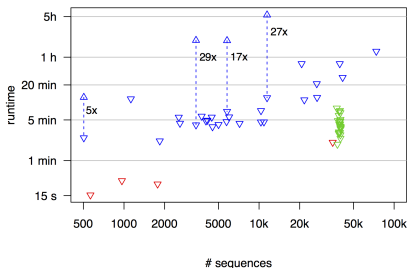
## Idea:

- pre-optimization on reduced data set
- evaluation of only highest-scoring motif occurrences



## Idea:

- pre-optimization on reduced data set
- evaluation of only highest-scoring motif occurrences



- |                                  |                                         |
|----------------------------------|-----------------------------------------|
| △ ChIP-seq, no speed-up, 1000 bp | ▽ ChIP-exo, speed-up, 100/200 bp (CTCF) |
| ▽ ChIP-seq, speed-up, 1000 bp    | ▽ PBM, speed-up, 40 bp                  |

## 66 PBM data sets (Weirauch *et al.*)

- protein binding microarray data for 66 TFs
- two different array designs (HK/ME) with different probes

### Task:

Learn motif on one design, predict binding intensities for other design

Algorithm	Pearson corr.	AUC-ROC	Final
<b>Dimont</b>	0.695	<b>0.951</b>	<b>1.002</b>
FeatureREDUCE	0.693	0.949	0.997
Team_D	0.691	0.938	0.984
Team_E	<b>0.696</b>	0.906	0.952

# Benchmark on ChIP-seq data

## 26 ChIP-seq data sets (Ma *et al.*)

- 26 ChIP-seq data sets for TFs with known motifs
- human, mouse, fly

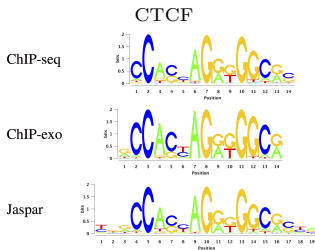
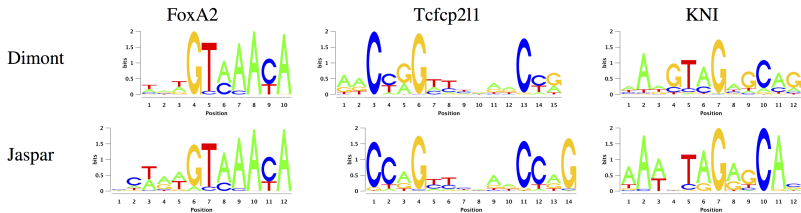
### Task:

Discover motif consistent with literature

Algorithm	Total successes	Average rank
<b>Dimont</b>	<b>26</b>	1.23
POSMO	23	1.00
ChIPMunk	23	1.00
MEME	22	1.32
DME	22	1.45
DREME	22	1.45
HMS	12	1.00

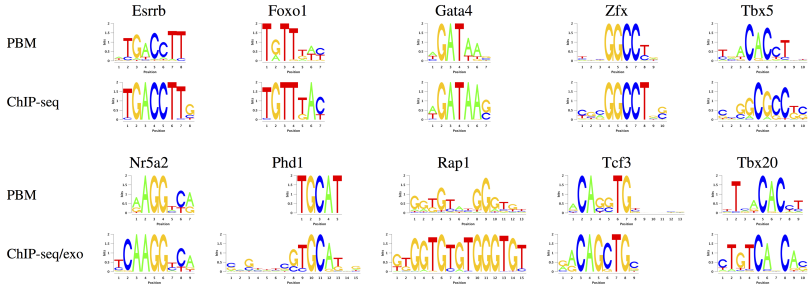


# Example motifs



⇒ most motifs fit the literature well

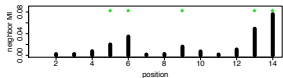
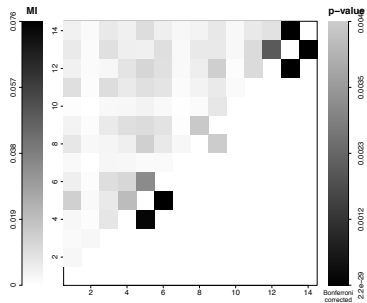
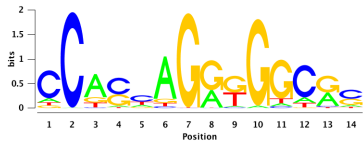
# In-vivo vs in-vitro binding



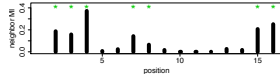
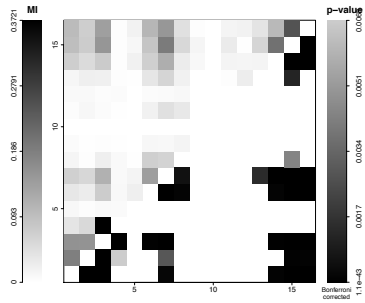
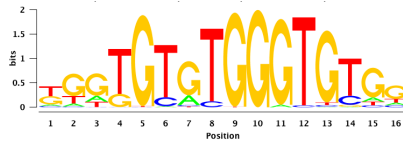
⇒ good accordance between *in-vivo* and *in-vitro* binding, but notable exceptions

# Dependencies between neighboring positions

## CTCF



## Rap1



Chipster 3.0.0 (build 1440)

File Edit View Workflow Help

**Datasets**

- dimont-logo-rc-1.png
- dimont.log
- extracted.fasta
- dimont-predictor-predictions.tsv
- extracted.fasta
- dimont-logo-1.png
- dimont-model-1.xml
- dimont-logo-rc-1.png**
- dimont-model-2.xml
- dimont-predictions-2.tsv
- dimont.log
- dimont-logo-rc-2.png
- dimont-predictions-1.tsv

**Analysis tools – ChIP-seq and DNase-seq – Find motifs with Dimont**

Position tag: peak  Hide parameters Run

Value tag: signal

Standard deviation: 75.0

Weighting factor: 0.2

Starts: 20

Initial motif width: 15

Markov order of motif model: 0

Markov order of background model: -1

Equivalent sample size: 4.0

Delete BSs from profile: yes

Dimont is a universal tool for de-novo motif discovery. Dimont has successfully been applied to ChIP-seq, ChIP-exo and protein-binding microarray (PBM) data.

More help Show tool sourcecode

**Workflow**

tsv tsv fasta tsv png

fasta fasta fasta

png xml txt png log png xml

png log png tsv

**Visualisation**

Show image Maximise Detach Close

bits

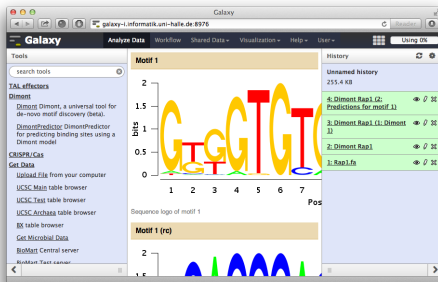
Position

Connected to chipster.csc.fi

Ready 271M / 800M

## Galaxy application

- public server
- convenient user interface
- also available in Galaxy Tool-Shed

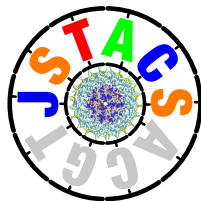


`galaxy.informatik.uni-halle.de`

- Galaxy-Server: 45 registered users, 500 runs (est.)
- Galaxy Tool-Shed: 60 clones

## Command line application

- <key>=<value> interface
- easily scriptable
- multi-threaded



```
java -jar Dimont.jar data=myseqs.fa infix=myresult position=peak  
value=signal threads=8
```

- available from [www.jstacs.de/index.php/Dimont](http://www.jstacs.de/index.php/Dimont)
- 290 downloads of command line program

## **Dimont, a general approach for motif discovery**

- reliably discovers motifs from ChIP-seq and PBM data
- achieves an acceptable runtime

## **In-vitro and in-vivo binding**

- often in good accordance
- but notable exceptions

## **Availability**

- Chipster since version 2.11
- public Galaxy at [galaxy.informatik.uni-halle.de](http://galaxy.informatik.uni-halle.de)  
and Galaxy Tool-Shed
- command line application: [www.jstacs.de/index.php/Dimont](http://www.jstacs.de/index.php/Dimont)

## Dimont, a general approach for motif discovery

- reliably discovers motifs from ChIP-seq and PBM data
- achieves an acceptable runtime

## In-vitro and in-vivo binding

- often in good accordance
- but notable exceptions

## Availability

- Chipster since version 2.11
- public Galaxy at [galaxy.informatik.uni-halle.de](http://galaxy.informatik.uni-halle.de)  
and Galaxy Tool-Shed
- command line application: [www.jstacs.de/index.php/Dimont](http://www.jstacs.de/index.php/Dimont)

**Thank you for your attention!**