# From Transfac to HOCOMOCO: using cross-validation and human curation to take most from the high throughput data compiling a complete collection of transcription factor binding motifs

Vsevolod J. Makeev

Vavilov Institute of General Genetics, Russian Academy of Sciences, Moscow

March 8, 2018

- We started to work with regulatory genomics in 1998
- Dima Papatsenko studied *Drosophila* enhancers
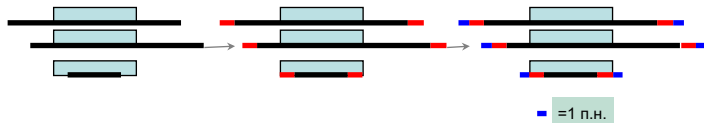- he was interested in TF binding sites

**Table 1.** Comparison between the Refined and Consistent Maps

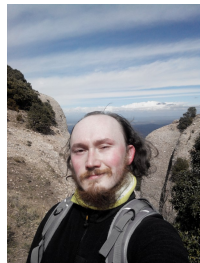| POSITION | SITE | REFINED MAP | SCORE | CONSISTENT MAP |
|----------|------|-------------|-------|----------------|
| 5-21c | Giant | | 10.46 | ATTATTGGGTTATATTG |
| 10-18 | Krüppel | TAACCCAAT | 5.94 | TAACCCAAT |
| 143-151 | Bicoid | GTTAATCCG | 7.93 | GTTAATCCG |
| 145-153 | Krüppel | TAATCCGTT | 7.11 | TAATCCGTT |
| 164-172c | Bicoid | AATAATCTC | 5.06 | |
| 167-183 | Giant | ATTATTAGTCAATTGCA | 9.11 | ATTATTAGTCAATTGCA |
| 229-245 | Giant | TTTATTGCAGCATCTTG | 9.36 | TTTATTGCAGCATCTTG |
| 314-322 | Bicoid | TATAATCGC | 4.70 | |
| 331-339c | Krüppel | CAACCCGGT | 5.47 | CAACCCGGT |
| 407-415c | Bicoid | GCTAATCCC | 8.09 | GCTAATCCC |
| 472-480 | Krüppel | | 5.90 | CAATCCCTT |
| 500-507c | Hunchback | TTTTTATG | 8.58 | TTTTTATG |
| 502-518c | Giant | ATTATTATGTGTTTTTA | 9.32 | ATTATTATGTGTTTTTA |
| 526-534c | Krüppel | | 6.59 | TAATCCCTT |
| 528-536c | Bicoid | CCTAATCCC | 8.17 | CCTAATCCC |
| 576-584c | Krüppel | | 5.94 | TAACCCAGT |
| 585-592 | Hunchback | TTTTTTTG | 8.77 | TTTTTTTG |
| 618-626 | Bicoid | | 5.71 | CTTAACCCG |
| 620-628 | Krüppel | TAACCCGTT | 7.55 | TAACCCGTT |
| 668-675 | Hunchback | | 8.77 | TTTTTTTG |

Distribution of sites shown for the *even-skipped* strip 2 region. Most of the experimentally verified binding sites shown are shared between the two maps (hits, shown in red). Two known Bicoid sites false-negatives in blue) are missing in the consistent map due to their low positional weight matrix score. In vitro binding assays support the suggestion of low affinity for these two Bicoid sites (Wilson et al. 1996). High-scoring matches (false-positives) to Bicoid, Krüppel, and Giant are shown in green.

- A site verified by at least two methods from footprints, mutant, or highly conserved blocks
- Bicoid (34 sites), Caudal (15), Ftz (25), Hunchback (43), Knirps (47), Kruppel (21), and Tramtrak (7)
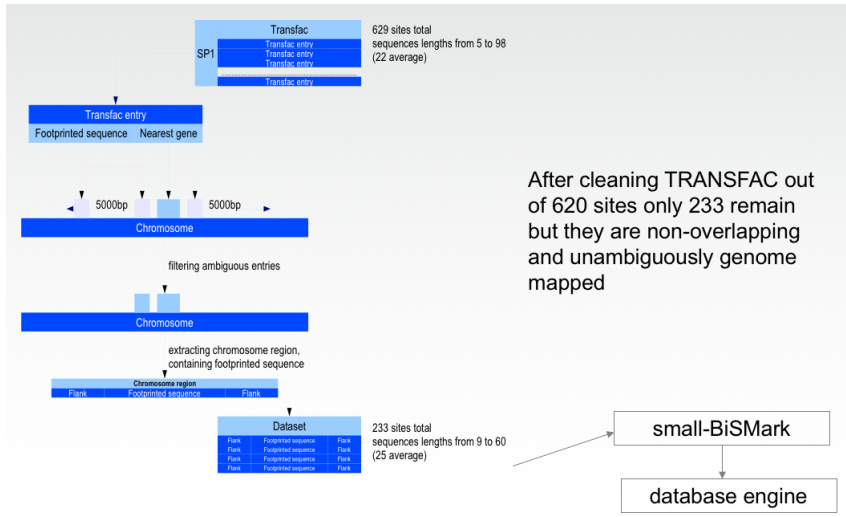- Aligned with CLUSTALW and manually and cut the flanks

■ =1 п.н.

- 2008
- Mapping footprints on the genome allows recovering up to 40
- Usually it is enough to add only two letters
- Genome data may be very useful for interpretation *in vitro* results
- http://autosome.ru/dmmpmm/ DMMPMM collection



Ivan Kulakovskiy

After cleaning TRANSFAC out of 620 sites only 233 remain but they are non-overlapping and unambiguously genome mapped

Sp1 JASPAR 2007
(SELEX data)



Sp1 Remapped and realigned
TRANSFAC 2008

- Chip-on-chip yielded long regions (up to 20K)
- Wasn't suitable for motif discovery
- But perhaps could be helped with *in vitro* data

Subsampling on many sets of sequences then optimization on total set of weighted sequencies

# Background

The task of identification of transcription factor binding motifs in a limited number of short DNA sequences has a long history.

Recently upcoming ChIP-Seq data provided a new challenge for motif discovery. Such data consist of thousands of sequences where a short overrepresented motif is to be found.

*peak*

*ChIP-Sequencing* ⟶

*protein of interest*

or *read* *tag*

Fortunately, in the case of a ChIP-Seq data one has additional information, which helps to select the correct signal. This information is the coverage profile constructed for DNA fragments obtained from ChIP-Seq experiments.

typical *ideal* peak
~100bp
~1000bp

...and its *real* brother
~3000bp

# ChIPmunk page

Peak shape and motif shape prior (like double box)
available at http://autosome.ru/ChIPMunk/

# TRANSFAC comes into view again

... and supplies us with a new version of SITE database (for free)

STAT3

CTCF

From a set of (**f1**,**f2**,**si**,**do**) motifs we **manually** select
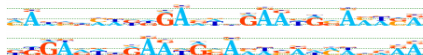reasonable ones according to the following criteria:
- select similar motifs for the TFs from a particular family;
- select motifs having higher weight / number of aligned sequences;
- for huge sequence sets: trust flanking regions;
- for small sequence sets: take motif cores;
- take >1 motifs for one TF when the motifs have completely different consensi;
- use information from other sources (compare to known existing motifs).



KAISO - both motifs are significant
(known to have two distinct binding motifs)
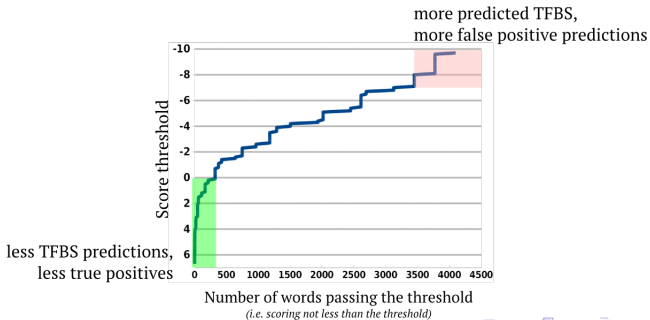


XRCC4 - no significant motif
(long and unstructured)

- PWM can be used to calculate a score for any sequence

- $Score[j] = \sum\limits_{j}^{j+L-1} PWM[j, s(j)]$

- $s(j)$ is the letter in the position $j$ of the alignment of PWM with the sequence

- $L$ is the PWM length

Each pair **( PWM , threshold )** classifies any word as a motif hit (YES/NO)



|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| A | -1.6 | -1.6 | 0.96 | -1.6 | -1.6 | 0.96 |
| C | -1.6 | -1.6 | 0.00 | -1.6 | -1.6 | -1.6 |
| G | 1.22 | 1.22 | -1.6 | -1.6 | -1.6 | -1.6 |
| T | -1.6 | -1.6 | -1.6 | 1.22 | 1.22 | 0.00 |

PWM
GGATTA → $S_{GGATTA}$=1.22+1.22+0.96+1.22+1.22+0.96=**6.8** the best score

$S_{GGGGGG}$=2.44-6.4=**-3.96**

$S$=**-9.6** the worst score

more predicted TFBS, more false positive predictions



less TFBS predictions, less true positives

Number of words passing the threshold
(i.e. scoring not less than the threshold)

# Fast exact calculation of motif P-vlaue

- Suppose there is a probability distribution upon the *l*-words
- Motif *P*-value is the sum of probabilities of all words scoring above the threshold
- In 2007 Hélèn Touzet and Jean Stéfan Varré designed nice precise algorithm



AMB Algorithms for Molecular Biology

| Home | About | Articles | Submission Guidelines |

What do you think about BMC? Take part in our survey our short survey

Abstract
Background
Complexity of the P-value problem
Algorithms for the P-value problems
Experimental Results
Discussion and Conclusion
Declarations
References

Research    Open Access

## Efficient and accurate P-value computation for Position Weight Matrices

Hélène Touzet ✉ and Jean-Stéphane Varré ✉

*Algorithms for Molecular Biology* 2007 **2**:15
https://doi.org/10.1186/1748-7188-2-15   © Touzet and Varré; licensee BioMed Central Ltd. 2007
Received: 06 July 2007 | Accepted: 11 December 2007 | Published: 11 December 2007

Download PDF
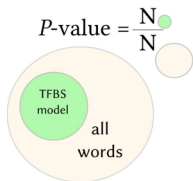
Export citations ▾

Metrics

Article accesses: 11491

Citations: 41 more information

Altmetric Attention Score: 1

- One needs to set both thresholds
- ... but after that it is possible to calculate the percentage of common words recognized by both motfs
- and compare it with a larger set of words recognized by any of them
- Matrices of different origine (or even PWM and PCM) can be compared without additional normalization

$P\text{-value} = \dfrac{N}{N}$

TFBS model

all words

A B  $\mathcal{J} = \dfrac{N}{N}$

...or in case of probabilities:

$$\mathcal{J}1(\Omega_1, \Omega_2) = \frac{P(\{\omega\} : \omega \in \Omega_1 \cap \Omega_2)}{P(\{\omega\} : \omega \in \Omega_1 \cup \Omega_2)}$$

# MacroApe to compare motifs

We modified Touzet - Varré algorithm to compare PMWs
Available at http://opera.autosome.ru/macroape
Can be used to extract motifs from various motif databases

We can use theoretically calculated P-values for a false-positive rate

This allows us to compare performance of different motifs on the same benchmark datasets

# Hocomoco database log

- 2011 first website published
- 2012, first publication, v.9, *Nucleic acids research, database 2013*
- 2015, second publication, v.10, *Nucleic acids research, database, 2016*
- 2017, third publication, v.11, *Nucleic acids research, database 2018*
- http://hocomoco11.autosome.ru/
- http://www.cbrc.kaust.edu.sa/hocomoco11

- large number of HT-SELEX data and new ChIP-seq data allowed us to extend the core base only by benchmarking and curation

- similar to known models (0.05 Jaccard similarity)
- consistent within a TF family, TFclass families are taken
- or at least with a clearly exhibited consensus (based on LOGO representation, manually assessed).

Gather as many datasets as possible

Motif discovery in all datasets

Benchmarking and conservative filtering

- Cross-validation based dataset filtering
- If known motif performs better than the genuine dataset motif the entire dataset is discarded

# Dinucleotide models



TFBS recognition quality comparison for AP2A

Figure: ETC family

Difficulties for MARA style analysis. SwissRegulon contains small number of "isolated" motifs

*Adapted from TFclass database, Wingender et al., 2015*

- models for 453 mouse and 680 human transcription factors
- contains 1302 mononucleotide and 576 dinucleotide PWMs
- build from more than 3000 ChIP-seq tracks and four peak callers

A:A brown eye colour, 80%
A:G brown eye colour
G:G blue eye colour, 99%

Found in the intron of HERC2, the non-pigment gene
21kb upstream of OCA2, the non-pigment gene



*Mike Visser et al. Genome Research, 2012; 22:446-455*

# No experimental location of TFBS

| method | *in vitro* *in vivo* | native or synthetic | segment length | # segments | comment |
|---|---|---|---|---|---|
| ChIP | *in vivo* | native | 40 (exo) 5000 | 150 - 50000 | indirect binding |
| One-hybrid | *in vivo* | synthetic | ∼30 | 20-50 | in bacteria |
| SELEX, RSS | *in vitro* | synthetic | ∼20 | 20-50 | saturation |
| HT-SELEX | *in vitro* | synthetic | ∼50 | 5000 | saturation |
| PBA | *in vitro* | synthetic | ∼50 | 10000 | overlapping |
| Footprints | either | native | ∼100 | 20 - 10000 | indirect |

Table: Experimental methods of TF binding identification

*From Levo and Segal, 2014,*
*Nat Rev Genet*

Because many other processes
(mostly chromatin related)
contribute to the protein
positioning at the genome

| | |
|---|---|
| Functional genomics (genome structure, annotation, etc) | 15 |
| Genetics: annotation of loci and rSNP | 13 |
| Systems biology (regulatory networks from DE data) | 10 |
| Algorithms and Machine learning assisted genome annotation | 7 |
| "Stories" about particular promoters etc | 7 |
| DNA - protein interaction studies | 6 |
| TF studies - databases, structure of DNA recognition motifs etc | 4 |
| Genetic engineering - prediction of genemics manipulation | 2 |
| General Molecular biology (transctiption initiation etc) | 1 |

An advertisment slot: autosome.ru software

Integrative motif discovery with ChIPMunk (for CHromatin ImmunoPrecipitation)



Motif comparison by Jaccard Similarity with MACRO-APE (for Approximate P-value Estimation)



Efficient motif finding with SPRY-SARUS (for Super Alphabet Representation)



Functional annotation of genetic variants with PERFECTOS-APE

# Who contributed this?

- VIGG RAS:
- Artem Kasianov
- Ivan Kulakovskiy
- Ilya Vorontsov
- Seva Makeev
- KAUST:
- Haitham Ashoor
- Wail Ba-alawi
- Arturo Magana-Mora
- Ulf Schaefer
- Vlad Bajic

- CB RAS:
- Julya Medvedeva
- ISB Ltd:
- Ruslan Shapirov
- Ivan Yevshin
- Fedor Kolpakov
- Skolkovo Tech:
- Dima Papatsenko
- students
- Alla Fedorova, MSU FBB
- Eugen Rumynskiy, MIPT
- Nastya Soboleva, MIPT

# Thank you!

- Russian Fund of Basics Research
- Russian Scientific Fund
- Ministry of Science and Education of Russian Federation
- Biobase and personally Edgar Wingender and Alexander Kel
- RIKEN Fantom Project
- Ecole Polytechnique and personally Mireille Regnier