# Анализ экспериментальных данных для определения участков ДНК, связывающих регуляторные белки

Макеев В.Ю.

Институт общей генетики РАН им. Н.И. Вавилова
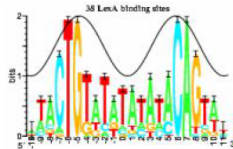
4 октября 2018 г.

# Experimentally verified TF binding regions often contain similar words related to protein binding



(a) LacI with DNA

```
TACTGTATATATATACAGTA Site1
TACTGGTTACGTACACAGTA Site2
TAATGTATATATATACATTA Site3
TACTGTACTTAAGTACAGTA Site4
TACTGGGAGCGCGACCAGTA Site5
```

(b) LacI sites



38 LexA binding sites

From Tom Schneider website

(c) LacI Logo

# D.melanogaster enhancers



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

# Our first collection of TFBS



**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 |
| 528-534c | 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)
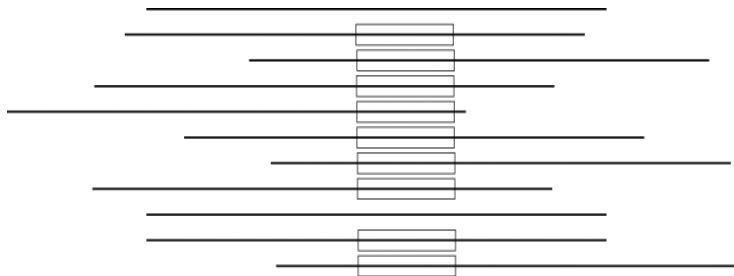- Aligned with CLUSTALW and manually and cut the flanks

| 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 | $\sim$30 | 20-50 | in bacteria |
| SELEX, RSS | *in vitro* | synthetic | $\sim$20 | 20-50 | saturation |
| HT-SELEX | *in vitro* | synthetic | $\sim$50 | 5000 | saturation |
| PBA | *in vitro* | synthetic | $\sim$50 | 10000 | overlapping |
| Footprints | either | native | $\sim$100 | 20 - 10000 | indirect |

Таблица: Experimental methods of TF binding identification

# SeSiMCMC

=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

629 sites total
sequences lengths from 5 to 98
(22 average)

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

233 sites total
sequences lengths from 9 to 60
(25 average)

small-BiSMark

database engine

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* → or *read tag*

*protein of interest*

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
100
80
60
40
20
~3000bp

# ChIPmunk page

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

,and supplies us with a new version of SITE database (for free)
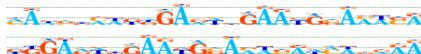
STAT3

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
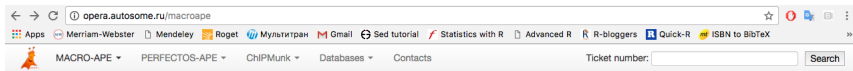(known to have two distinct binding motifs)

XRCC4 - no significant motif
(long and unstructured)

# 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

- 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
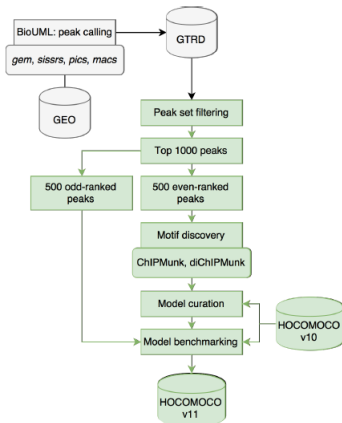- or at least with a clearly exhibited consensus (based on LOGO representation, manually assessed).

Gather as many datasets as possible

Motif discovery in all datasets
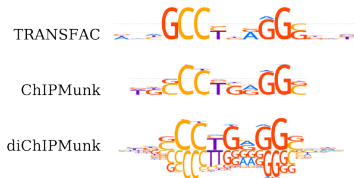
Benchmarking and conservative filtering

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

TFBS recognition quality
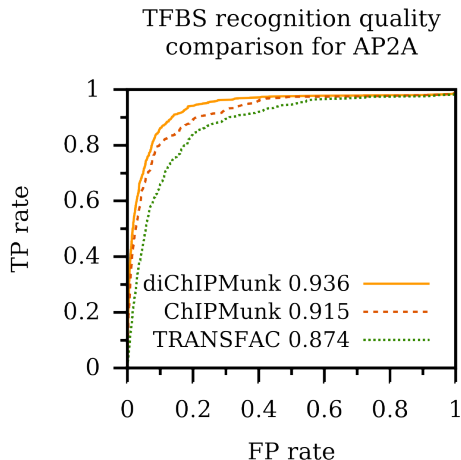comparison for AP2A

TRANSFAC

ChIPMunk

diChIPMunk

GATA:
G'A'T'A or GA'AT'TA

mononucleotide    dinucleotide
alphabet          superalphabet
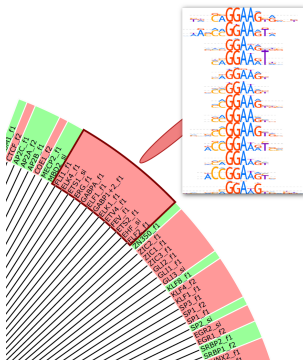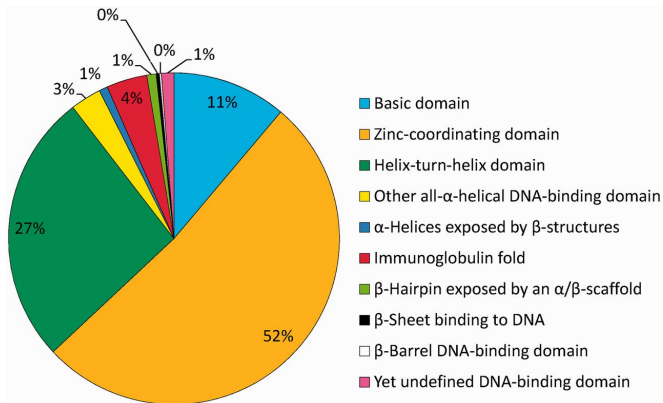{A,C,G,T}         {AA,AC,...TT}

diChIPMunk 0.936
ChIPMunk 0.915
TRANSFAC 0.874

TP rate

FP rate

Рис.: 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*

*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

- 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