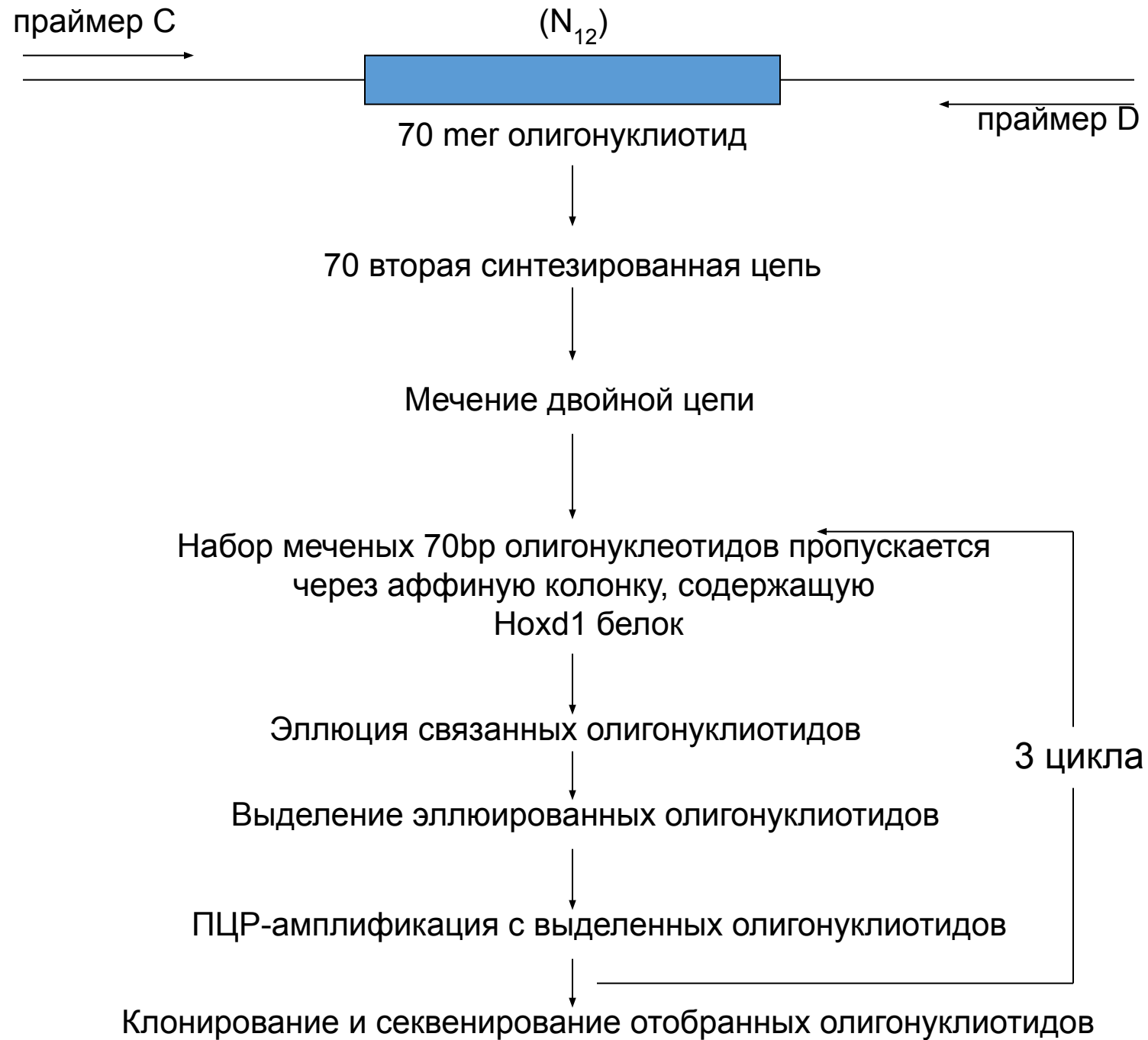


**Методы изучения регуляторных
районов генов**
Лекция III

Меркулова Татьяна Ивановна

Институт цитологии и генетики СО РАН



Method

Multiplexed massively parallel SELEX for characterization of human transcription factor binding specificities

Arttu Jolma,^{1,2} Teemu Kivioja,^{1,3} Jarkko Toivonen,³ Lu Cheng,³ Gonghong Wei,¹
Martin Enge,² Mikko Taipale,¹ Juan M. Vaquerizas,⁴ Jian Yan,¹ Mikko J. Sillanpää,⁵
Martin Bonke,¹ Kimmo Palin,³ Shaheynoor Talukder,⁶ Timothy R. Hughes,⁶
Nicholas M. Luscombe,⁴ Esko Ukkonen,³ and Jussi Taipale^{1,2,7}

¹Department of Molecular Medicine, National Public Health Institute (KTL) and Genome-Scale Biology Program, Institute of Biomedicine and High Throughput Center, University of Helsinki, Biomedicum, Helsinki, Finland; ²Department of Biosciences and Nutrition, Karolinska Institutet, Stockholm, Sweden; ³Department of Computer Science, FI-00014 University of Helsinki, Helsinki, Finland; ⁴EMBL–European Bioinformatics Institute, Cambridge CB10 1SD, United Kingdom; ⁵Department of Mathematics and Statistics, FI-00014 University of Helsinki, Helsinki, Finland; ⁶Department of Molecular Genetics and Banting and Best Department of Medical Research, University of Toronto, Toronto, ON M4T 2J4, Canada

Genome Research
www.genome.org

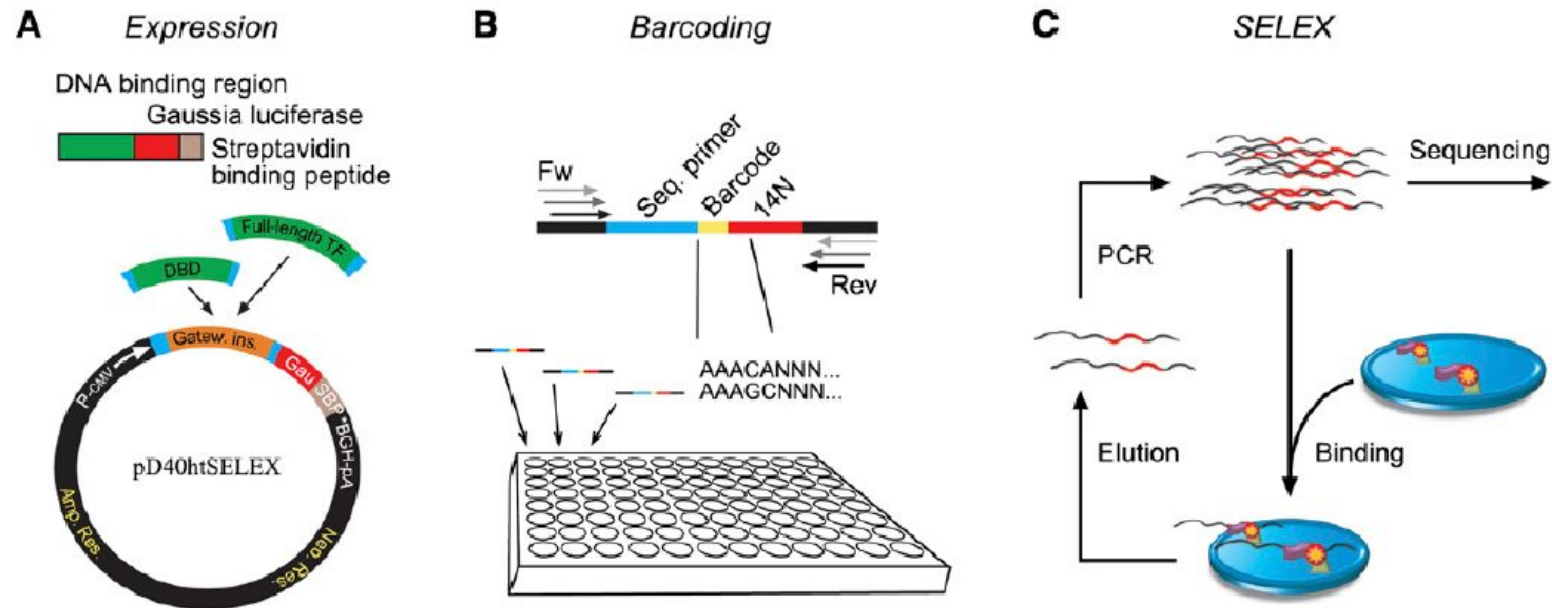


Figure 1. Schematic description of the high-throughput SELEX process. (A) Protein expression. (Top) Proteins are expressed as fusion proteins with SBP-tagged *Gaussia*-luciferase. (Bottom) The GATEWAY recombination cloning system is used to transfer DNA sequences encoding DBDs or TFs from donor-vectors to the pD40htSELEX expression vector. (B) Ligand design that accommodates multiplexing of samples using barcodes. Each DNA ligand contains a 14-bp randomized region (14N), and a 5-bp barcode (Barcode) that uniquely identifies the individual SELEX sample. To increase specificity, each barcode differs from all other barcodes by at least 2 bp. These variable sequences are flanked by constant sequences that include an Illumina Genome Analyzer sequencing primer site (Seq. primer) and bridge amplification primer binding regions (Fw, Rev; arrows), which are extended in their 5' regions to accommodate partially nested primers (used in successive SELEX rounds). (C) Basic principle of high-throughput SELEX. A double-stranded DNA mixture containing all possible 14-bp sequences (from B) is incubated with a DNA-binding protein immobilized into a well of a 96-well plate, resulting in binding of DNA to the protein. After washing and elution, the resulting population of more specific sequences is amplified by PCR and subjected to high-throughput single-molecule sequencing. The specificity of the TF can then be constructed by iterating the process and calculating the abundance of distinct sequences after different numbers of cycles. In each cycle, multiple reactions are mixed into a single sequencing lane, and the TFs are identified using the barcode sequences.

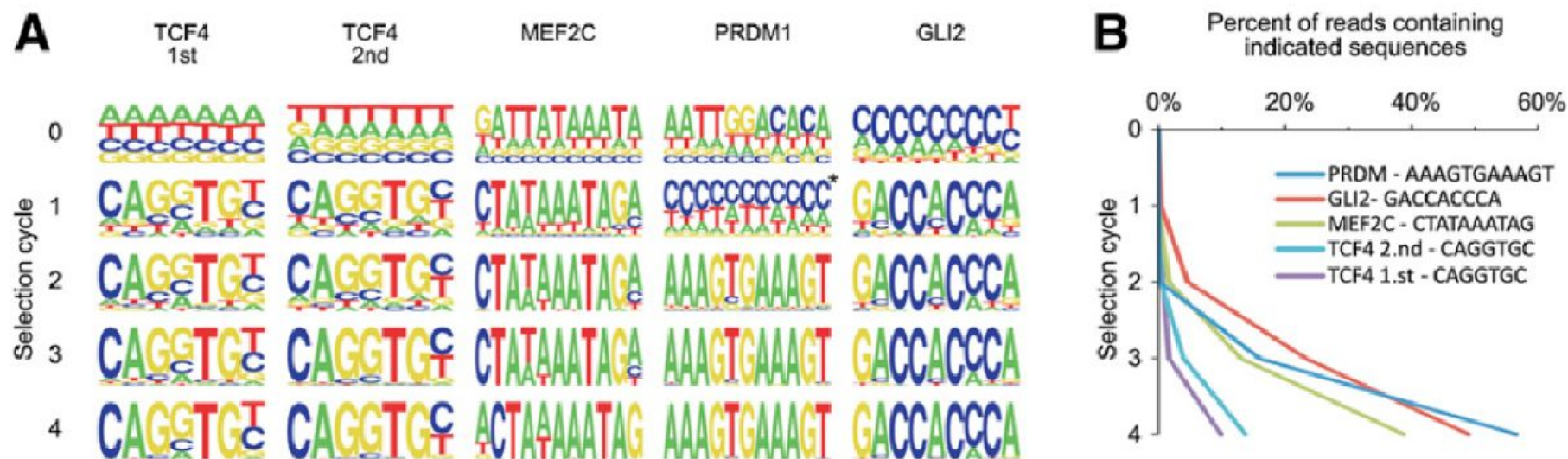


Figure 2. Enrichment of specific sequences during the SELEX process. (A) Position weight matrices built around the most enriched sequence for four different TFs (see Methods for details). The height of the letter at each position is directly proportional to the incidence of the indicated base in sequences where all other bases exactly match the most enriched sequence. Note that clear enrichment of sequences is observed after one or two SELEX rounds, and that two separate experiments for TCF4 result in a very similar enrichment pattern. In the first cycle, the algorithm used here detects incorrect binding profile for PRDM1 (asterisk) due to a low number of the relatively long consensus sequences. The enrichment of high-affinity sequences can, however, be detected by seeding the algorithm with consensus from the later cycles (see Supplemental Fig. S4A). (B) The fraction of all fragments containing the most enriched sequence from the third SELEX cycle plotted as a function of the SELEX cycle.

A E-Box (CANNTG)

```

A1      ggggtg CTGGTGT CACGTG TTT ggaga
A2      ggggtg TAT CACGTG AGCGTTC ggaga
A3      ggggtg GGTTCAT CACGTG AGCA ggaga
A4      tctccc ATATAGTAT CACGTG A cacgc
A5      ggggtg TCAT CATGTG ATGCCC ggaga
A6      tctccc ACAT CATGTG ATTATG cacgc
A7      ggggtg GTAT CACCTG ACTAGC ggaga
A8      tctccc CAACGGAT CACCTG AT cacgc
A9      tctccc TGAT CACATG CCAAGTA cacgc
A10     ggggtg TTGATTCAT CACGTG C ggaga
A11     tctccc TGGCACAAT CACGTG A cacgc
A12     tctccc TTATAGT CACGTG ACT cacgc
A13     tctccc ATCTAGT CACGTG ACA cacgc
A14     tctccc GCACT CACCTG ATTTG cacgc
A15     tctccc ACAGACGGT CACGTG A cacgc
A16     tctccc AACACAC CACGTG ACC cacgc
A17     ggggtg GCGAT CACGTG TATTA ggaga
A18     tctccc ACAGATGAT CACGTG A cacgc
A19     ggggtg TTCAT CACCTG AGTCA ggaga
A20     ggggtg GCGAT CAGGTG ATGCA ggaga
A21     ggggtg GGT CACGTG GTGTCCA ggaga
A22     tctccc CAAGCGGAT CACGTG A cacgc
A23     tctccc TGCACACTA CACGTG A cacgc
A24     tctccc TGAT CACATG CCATGT cacgc
A26     ggggtg GGTCAT CACGTG AGCC ggaga
A27     ggggtg CAACGGAT CACCTG AT ggaga
A28     tctccc TCTGTTCAT CACGTG A cacgc
A29     tctccc TCATGCAT CACGTG AT cacgc
A30     ggggtg TGCACAAT CACGTG AT ggaga
    
```

Position	-6	-5	-4	-3	-2	-1	1	2	3	4	5
G	11	5				1	22		30	1	5
A	4	23	1		30		2			21	2
T	5	1	28			2		30		2	11
C	10	1	1	30	27	6				3	12

G/C A T C A C G T G A Py

ADD1/SREBP1 E-Box Consensus Binding Sequence:

5'-ATCACGTGA-3'

B Non-E-Box

```

S1      tctccc GATCTTGAG ATCACCCcac gc
S2      tctccc AGAGCTAGT GTCACCCcac gc
S3      tctccc TTGTGTCTG ATCACCCcac gc
S4      tctccc ATCACCCCAC CTGTCCGcacgc
S5      ggggtg TCACAG ATCACCCCAC ggaga
S6      ggggtg CTAGCTCGT GTCACCCcac gc
S7      tctccc ACGTATTGT ATCACCCcac gc
    
```

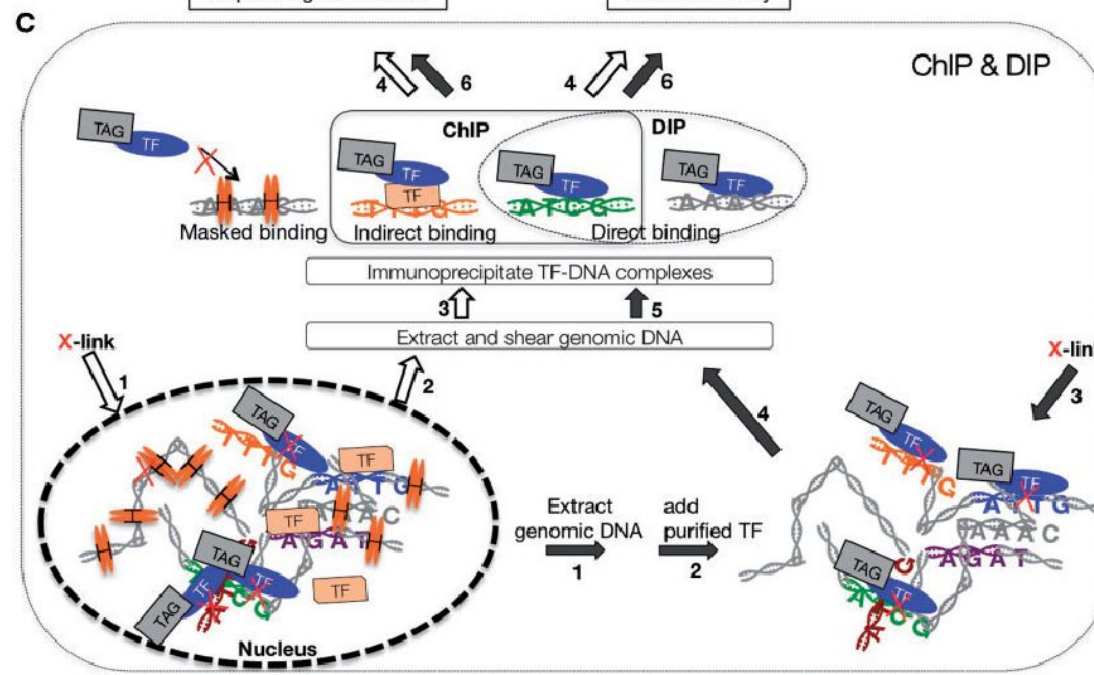
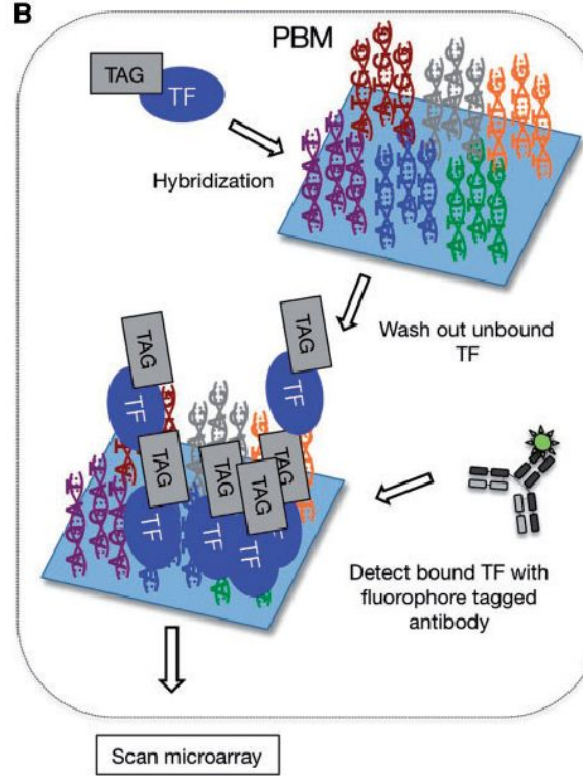
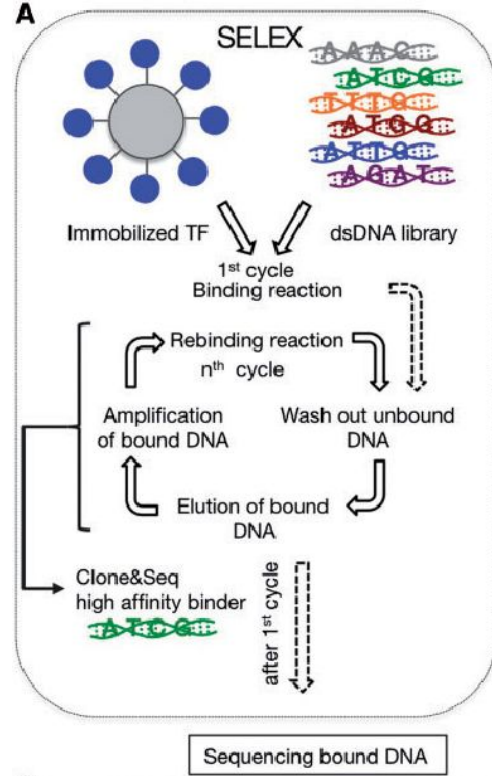
Position	-1	1	2	3	4	5	6	7	8	9	10
G	3	2									
A		5			7					7	
T	3		7								
C	1			7		7	7	7	7		7

N A T C A C C C C A C

ADD1/SREBP1 Non-E-Box Consensus Binding Sequence:

SAAB Consensus sequence: 5'-ATCACCCAC-3'

Published SRE-1 sequence: 5'-ATCACCCAC-3'



A

МТдт 5' -AGCACTATA **GGGACATGATGTTCC**ACACGTCACATGGTCGTCC-3'

МТ27 5' -GATGTCGCG **GGAACA**CAG**TGTTCC**GTGTACTGTGCAACTACTT-3'

25 5' -GATCCCCCGGGCATCACCGTGCAGGGGGGA **GGTACA**GAG**TGTTCT**
GCGAGGATGCG-3' G/C богатое окружение (28% A/T п.н.)

111 5' -GATCCTGTTACCATAGTGТААСТТССТАТТСА **GGTACA**АТА**TGTTCT**
АТАСТТССТАТТТ-3' A/T богатое окружение (72% A/T п.н.)

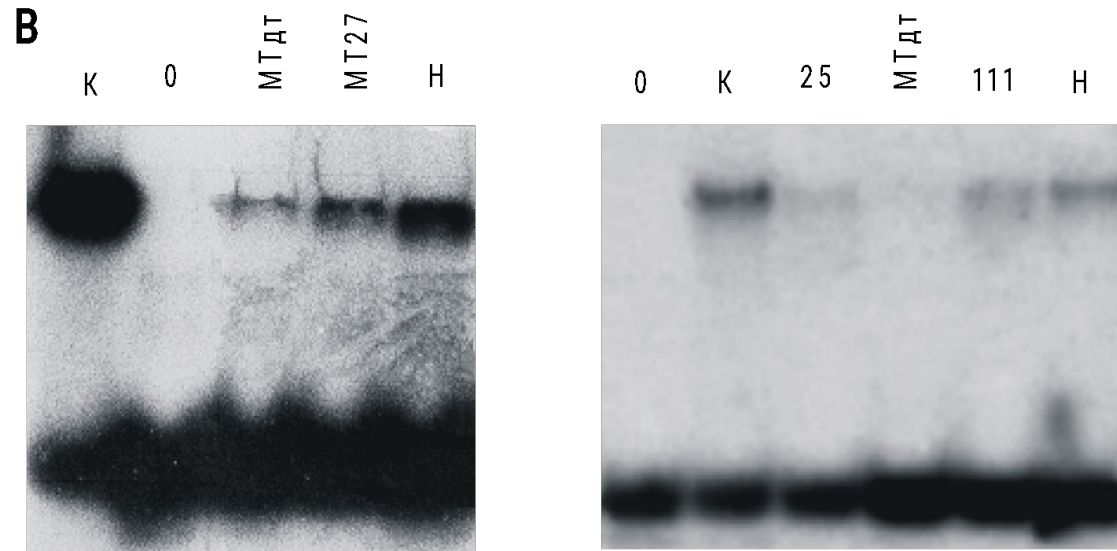
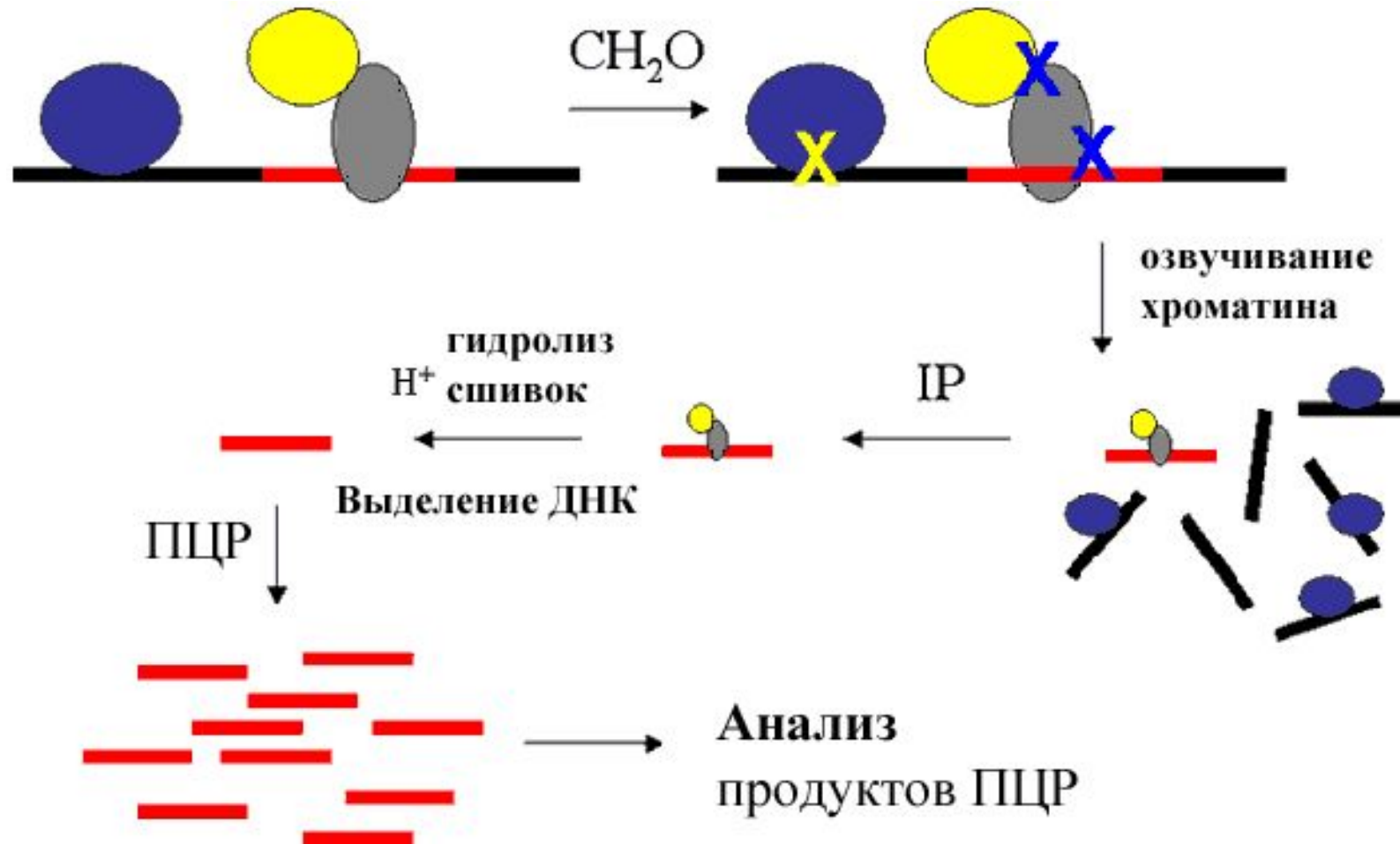


Рис. 2. Влияние первичной структуры ДНК в окружении консенсуса сайтов связывания рецептора глюкокортикоидов на эффективность связывания ДНК с рецептором.

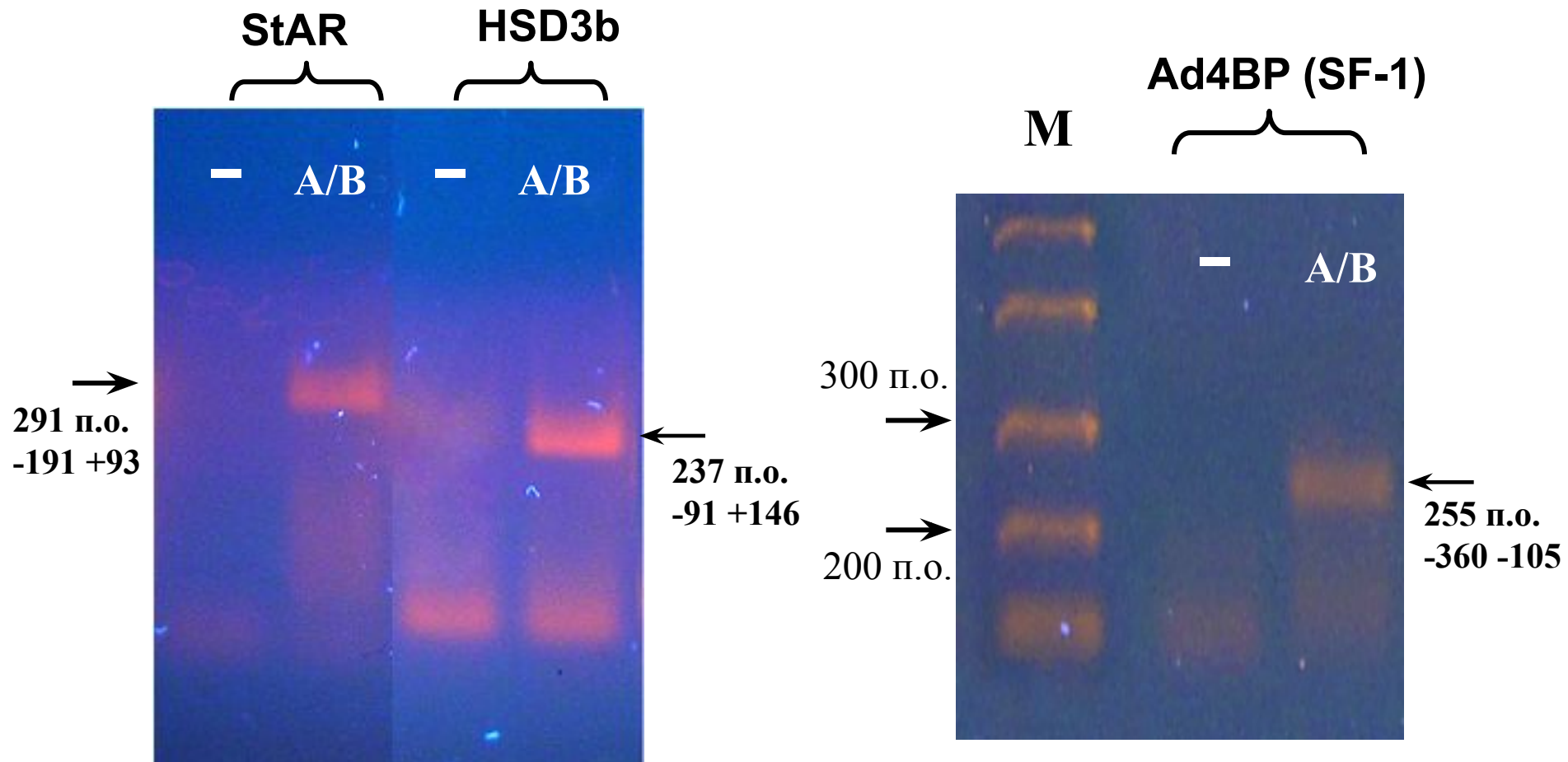
К - ДНК-проба в комплексе с рецептором глюкокортикоидов; 0 - свободная ДНК-проба; 3-5 и 8-11 - вытеснение меченой пробы из коплекса с ГР различными фрагментами немеченой конкурентной ДНК: МТдт - *EcoRI-HindIII* фрагмент рМТдт, МТ27 - *EcoRI-HindIII* фрагмент рМТ27, 25 - *EcoRI-HindIII* фрагмент р25, 111 - *EcoRI-HindIII* фрагмент р111, Н - *PvuII-EcoRI* фрагмент рUC19 (неспецифическая ДНК).

Принцип метода иммунопреципитации *in vivo*

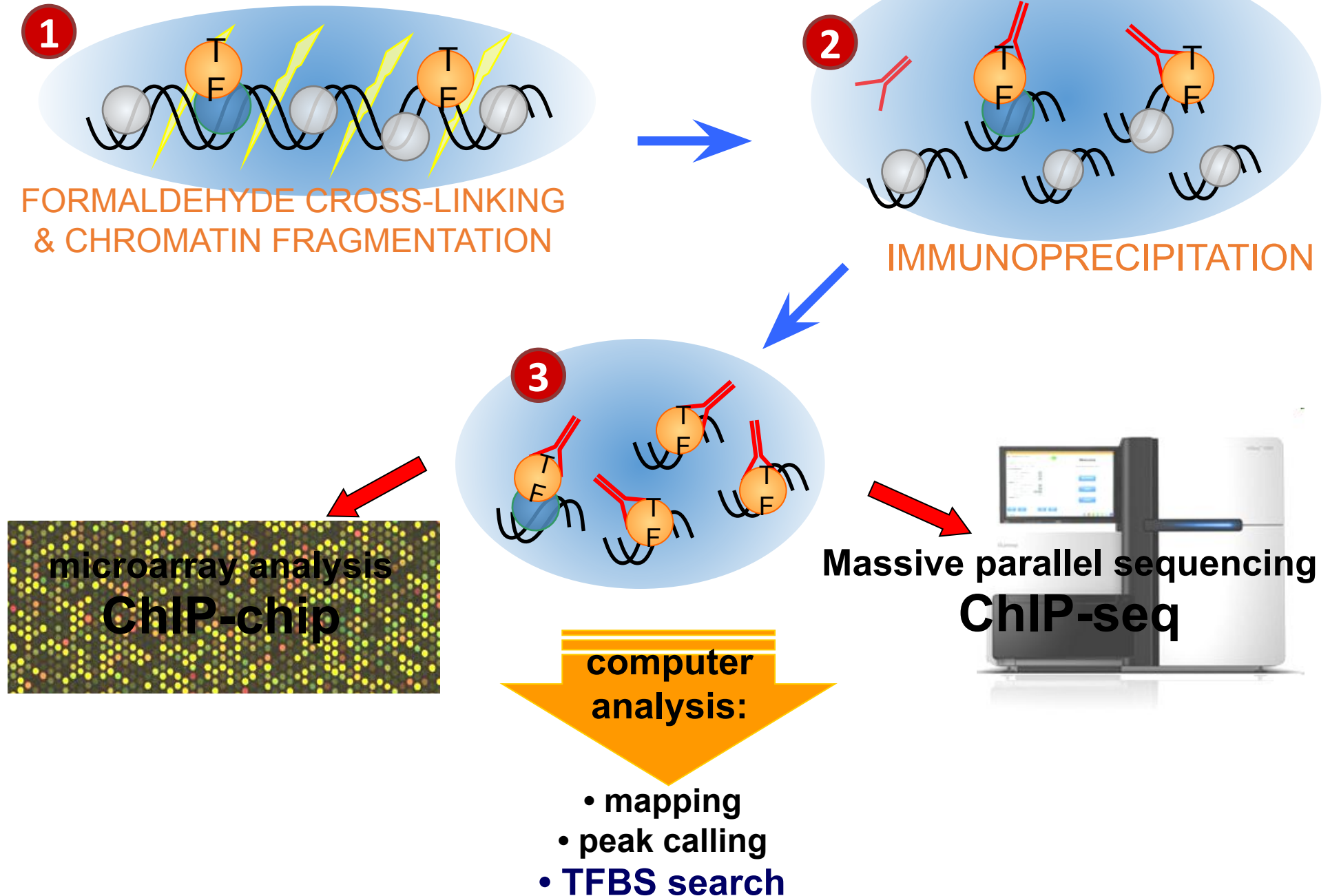


Иммунопреципитация хроматина *in vivo* с использованием антител против SF-1 для трёх генов мыши

Семенники мышей линии C57Br возраста 7-10 дней



Scheme of chromatin immunoprecipitation workflow



Features of p65-binding sites on chromosome 22

(A) p65 binding relative to Sanger-annotated genes and hybridizing regions on chromosome 22 is illustrated.

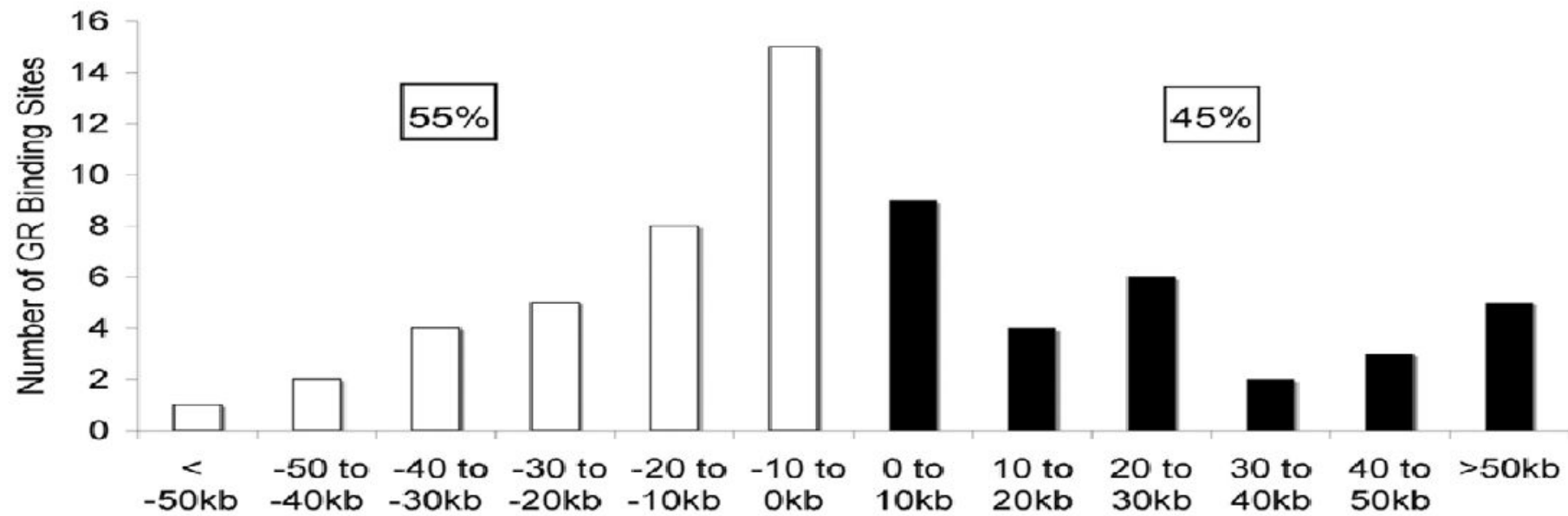
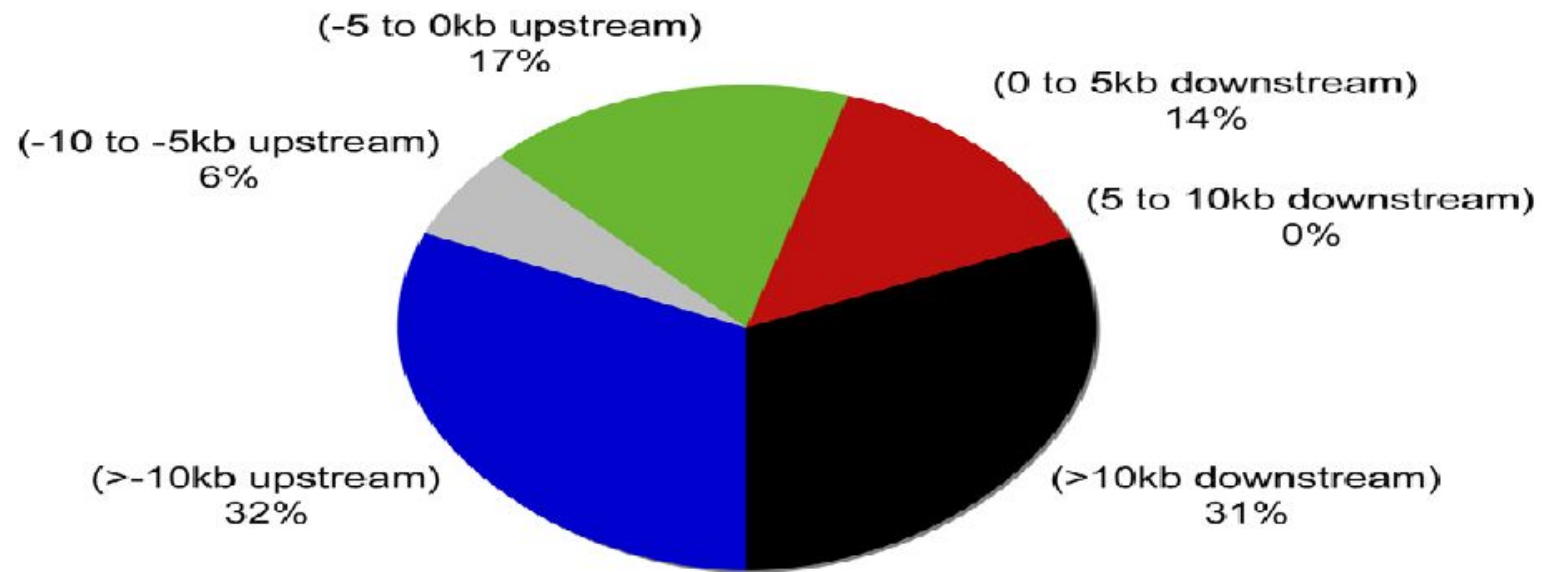
(B) Distribution of NF- κ B consensus sequences in p65-binding sites.

The sequences of p65-bound fragments on the microarray were searched for NF- κ B consensus sites by using both an in-house chromosome annotation system and the TFSEARCH database.

www.cbrc.jp/researchdb/TFSEARCH.html

A

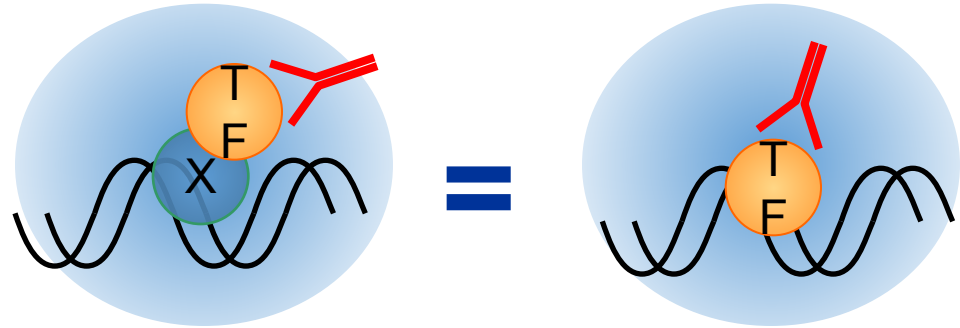
Location and Position of GREs

**B**

Weak points of TFBS identification by ChIP-seq technology

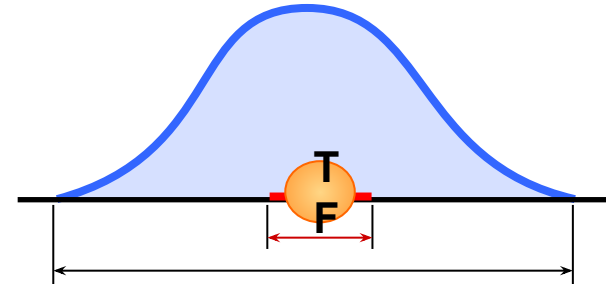
1 Indirect TF binding to DNA

The cases of direct and indirect interaction of TF with DNA are indistinguishable



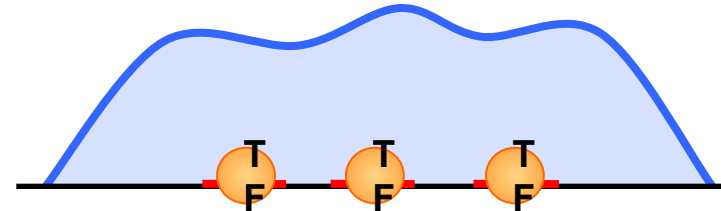
2 Insufficient resolution

Typical ChIP-seq peak – 100-500 bp
Typical TF footprint on DNA – 5-25 bp



3 TFBS clusters

A group of neighboring TFBSs is usually processed by peak-callers as one peak



4 False positive peaks (antibody cross-specificity, overrepresented seqs)

5 The spectrum of detected binding loci depends on cell type and state



Используемые для интерпретации данных ChIP-seq биоинформационные методы

Методы предсказания ССТФ

De novo методы поиска ССТФ

- методы поиска наиболее представленных мотивов в участках связывания ТФ с хроматином (пиках ChIP-seq). (без предварительной информации о ССТФ)

monoChIPMunk

Kulakovskiy I.V., et al.,
Bioinformatics 2010, 26(20):2622-3

diChIPMunk

Методы, основанные на использовании обучающих выборок ССТФ

- Консенсус, весовые матрицы (PWM) и пр.

FoxA oPWM
SiteGA

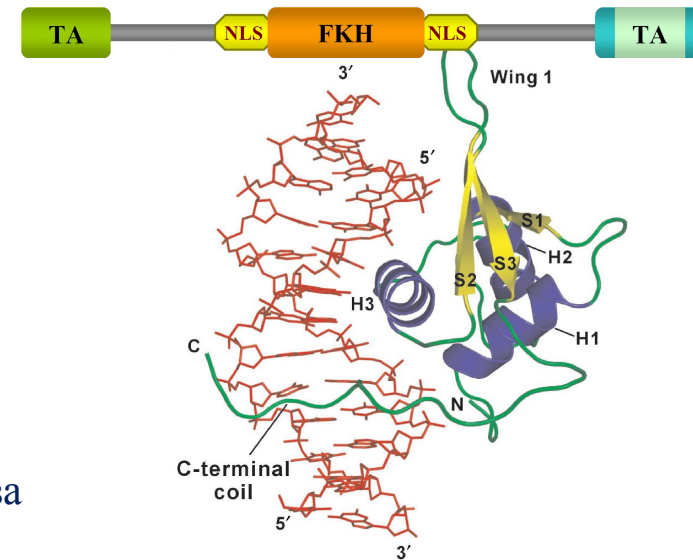
[Levitsky V.G., et al., BMC Bioinformatics 2007, 8:481]
[Levitsky V.G., et al., Dokl Biochem Biophys. 2011, 436:12-5]



Объект исследования: сайты связывания транскрипционных факторов FoxA

ТФ FoxA : FoxA1, FoxA2, FoxA3

- FoxAs имеют высококонсервативный ДСД типа winged-helix (FKH).
- FoxA обладают способностью вытеснять линкерные гистоны и открывать компактизованный хроматин, т.е. могут служить как **pioneer transcription factors**.
- FoxAs играют важную роль в раннем эмбриогенезе, органогенезе и поддержании метаболического гомеостаза организма.
- Все три ТФ FoxA остаются активными в печени взрослого, играя активную роль в регуляции метаболизма

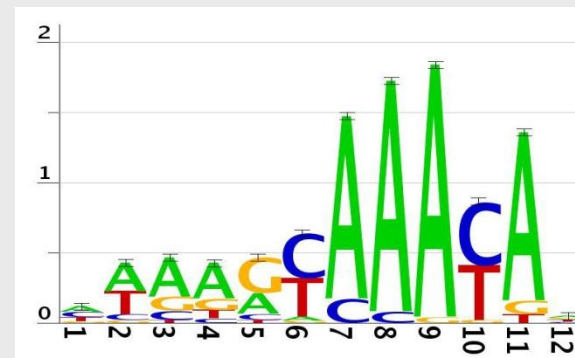


Консенсус FoxA :
5'-RYMAAYA-3'

[Overdier et al., 1994; Roux et al., 1995]

FoxA2 ChIP-seq в печени мыши [Wederell et al., 2008]

Число пиков— **11 475**
Средняя высота пика— **16**
Средняя длина пика— **727 bp**





PWM и SiteGA

Для создания обучающей выборки из 81 известных ССТФ FoxA использовали данные из базы TRRD и литературных источников.

Критерием для отбора сайтов служило наличие доказательств взаимодействия FoxA с соответствующим районом ДНК, полученных одним из следующих методов:

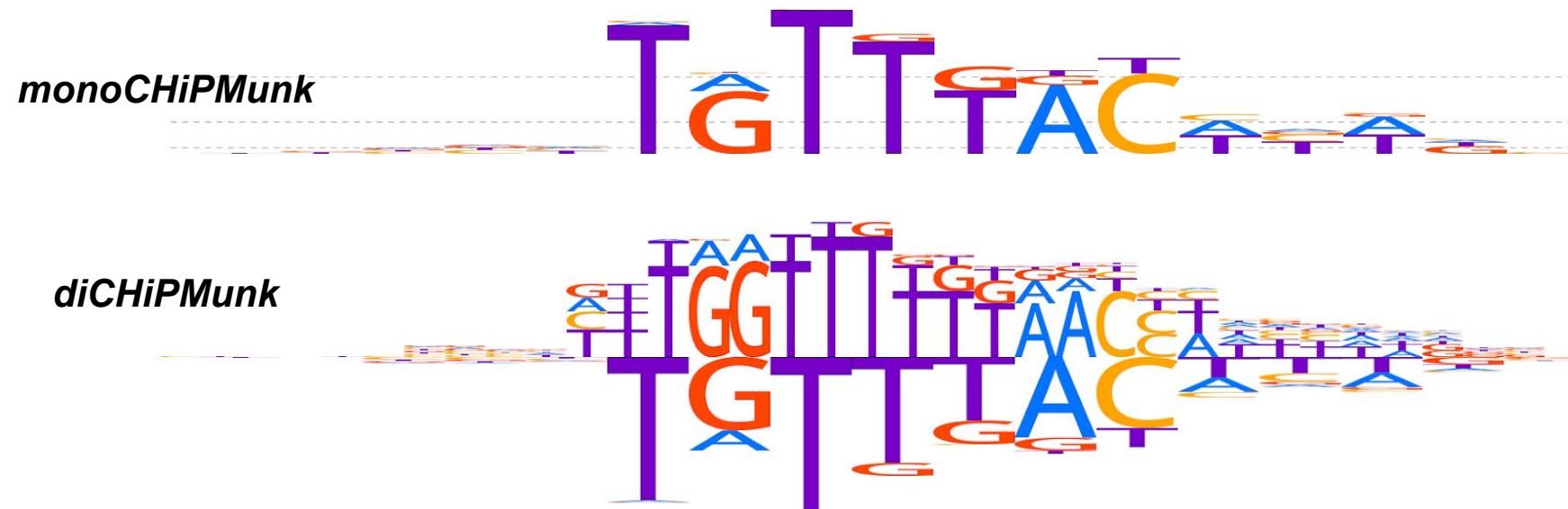
- футпринтинг ДНКазой I с использованием очищенного белка;
- EMSA с использованием очищенного белка;
- EMSA с использованием ядерного экстракта белков и специфических антител (EMSA supershift).

Выравнивание последовательностей проводилось относительно наиболее представленного мотива (в кодировке IUPAC) [Levitsky et al., 2011]. В итоге, для построения методов PWM и SiteGA использовалось выравнивание 53 последовательностей, содержащих мотив TRTTTRYH (R=A/G, Y=C/T, H=A/C/G).

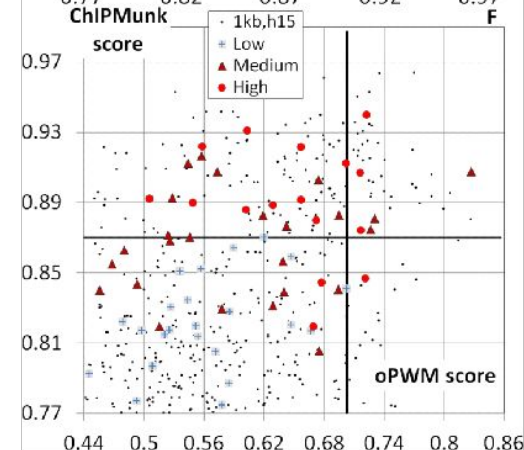
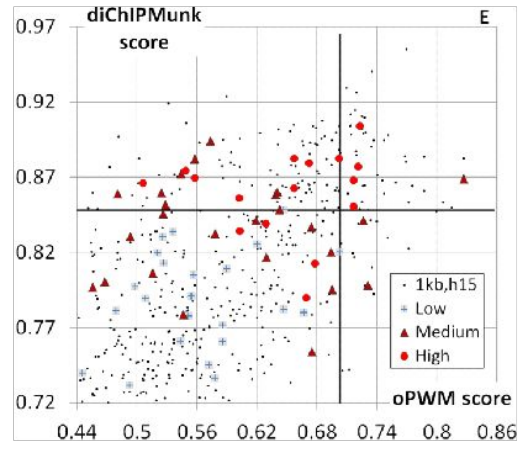
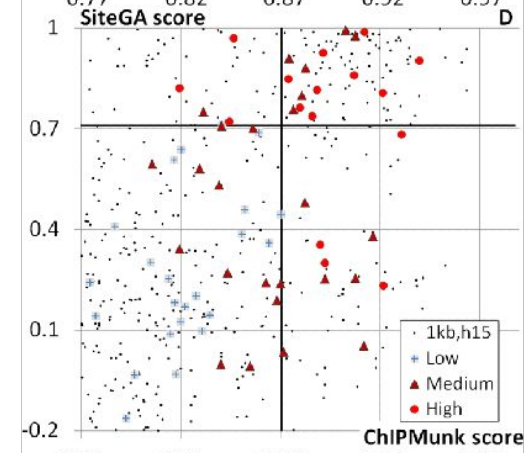
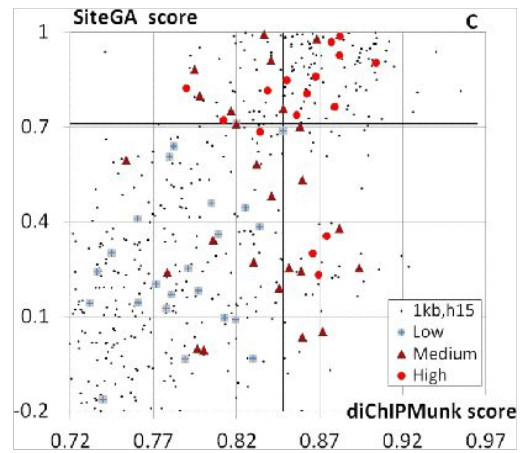
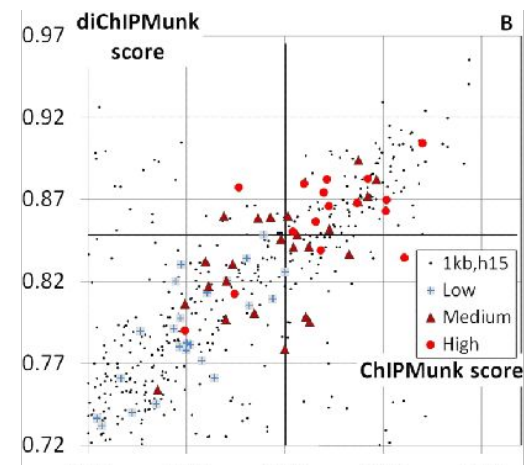
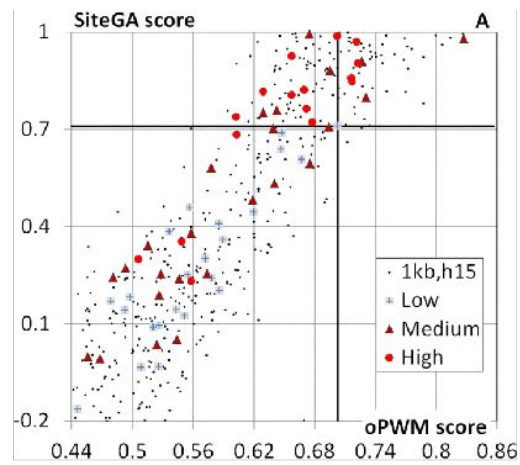


CHiPMunk

Для обучения CHiPMunk были использованы 4455 локуса связывания FoxA2 (CHiP-seq; высота пика >15). Длины мотива для mono- & di-версий Chiptunk установлены как 20 нт и 28 нт, соответственно. Выявленные мотивы имеют высокий уровень сходства как с мотивом TRTTTRYH, полученным на выборке 81 CC FoxA, так и с известным консенсусом сайтов FoxA2 [Overdier et al., 1994; Roux et al., 1995].



Лого наиболее представленных мотивов, полученных с помощью mono- and diChiptunk для CHiP-seq FoxA2



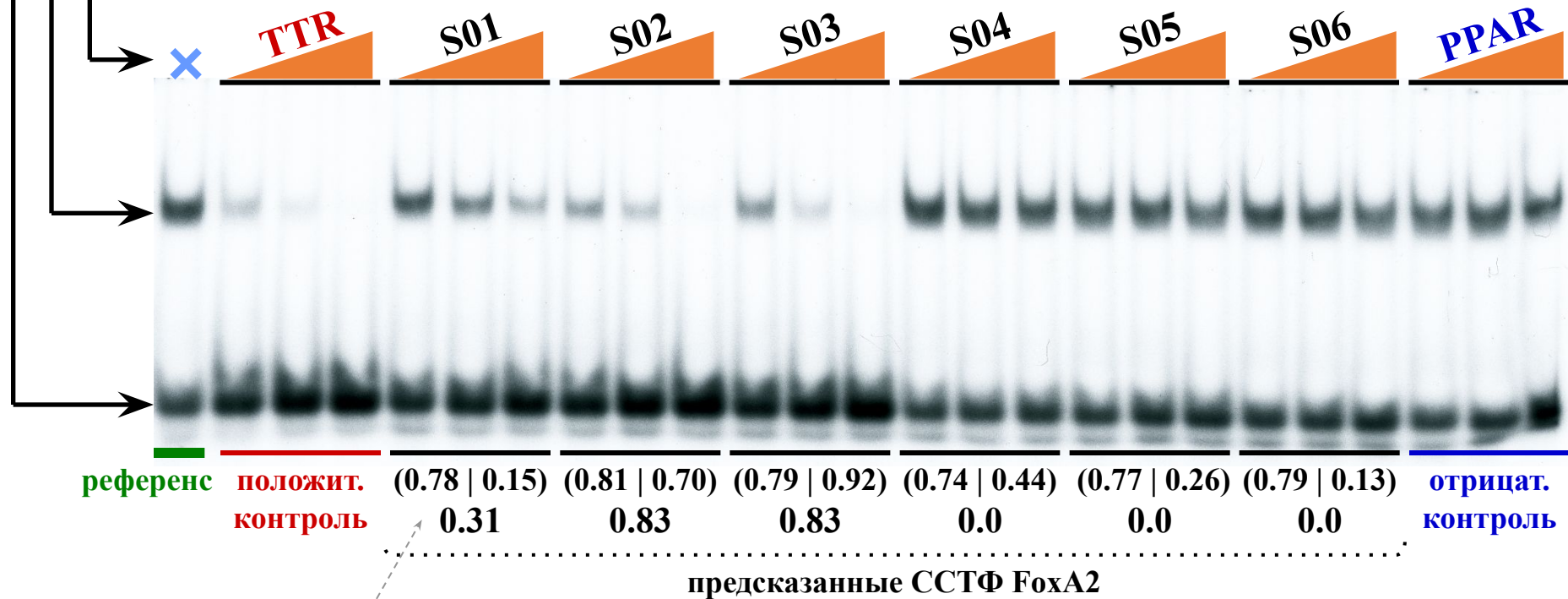


Экспериментальная верификация 64 предсказанных ССТФ FoxA методом задержки в геле (EMSA), конкурентный анализ

[³²P]-TTR: меченый олигонуклеотид, соответствующий сайту связывания FoxA из промотора гена *ttr* мыши

GST-DBD-FoxA2: Задержка в геле GST-слитым белком, соответствующим ДНК-связывающему домену TF FoxA2

немеченый конкурент: ×2, ×5 или ×20 молярный избыток конкурентного олигонуклеотида

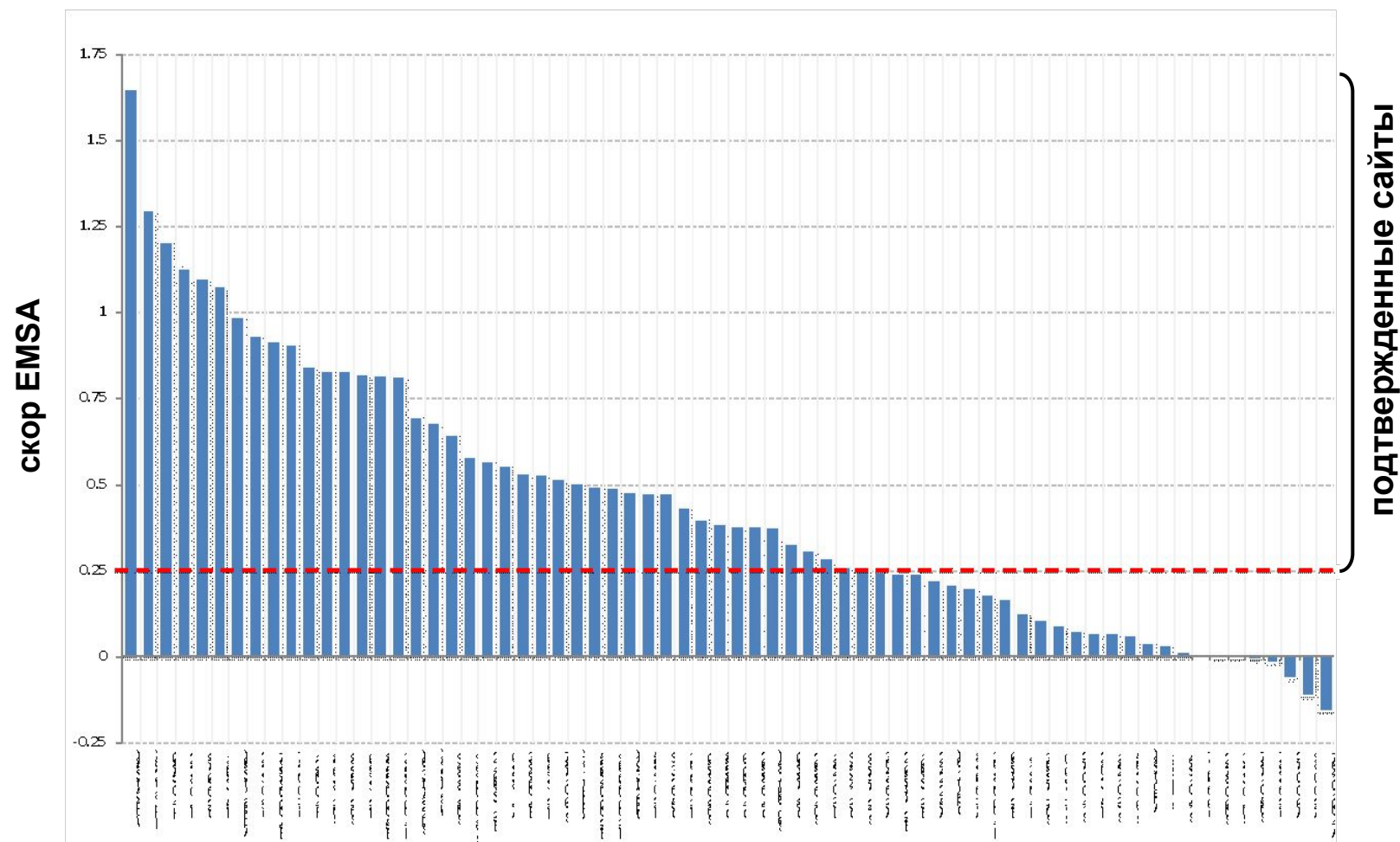


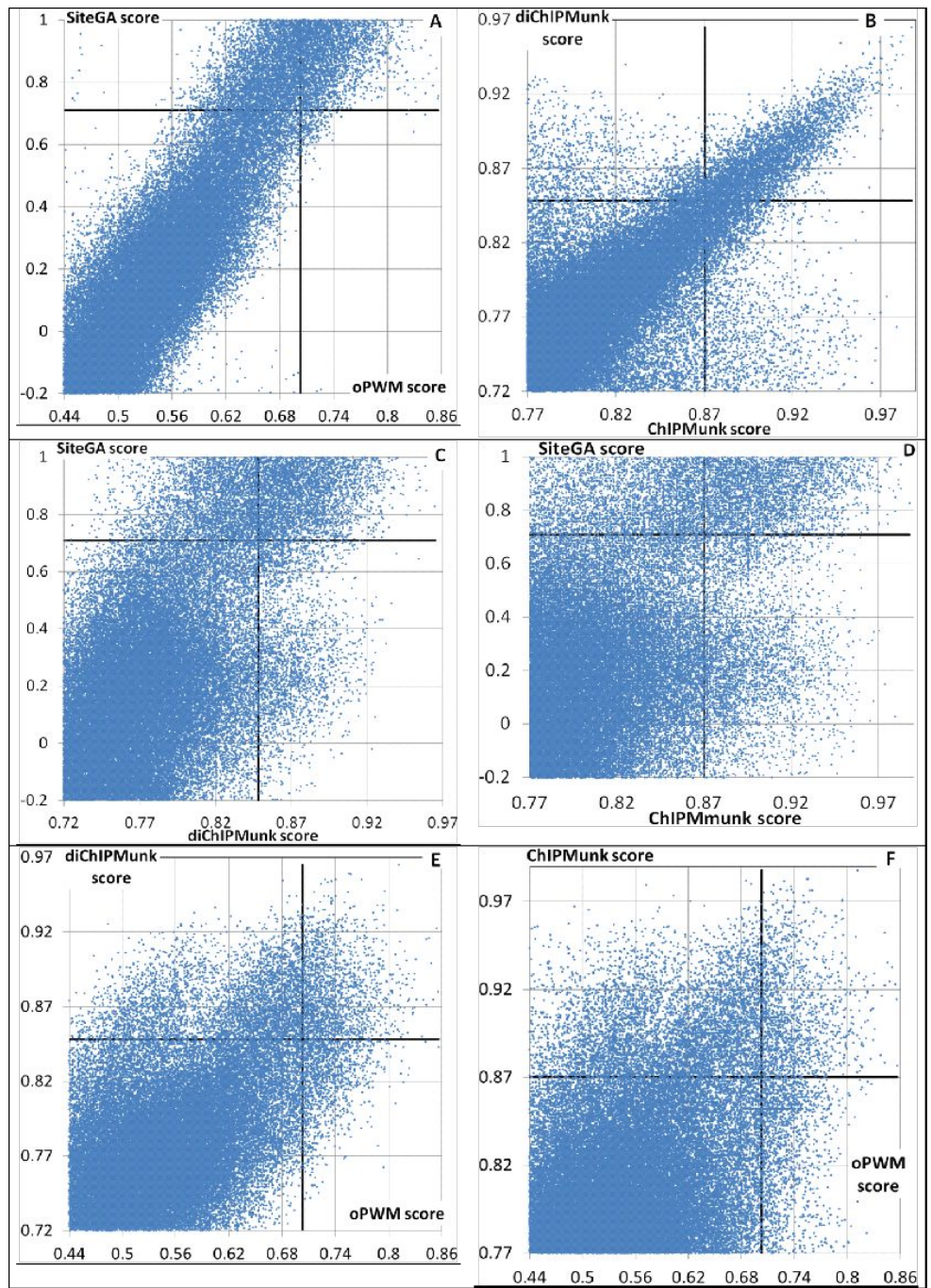
(скор ChIPMunk | скор SiteGA)
скор EMSA

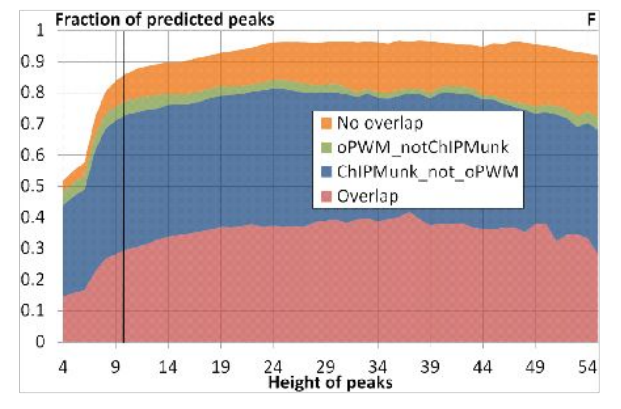
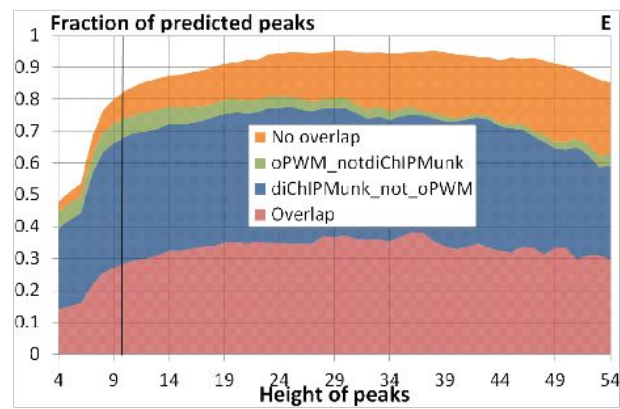
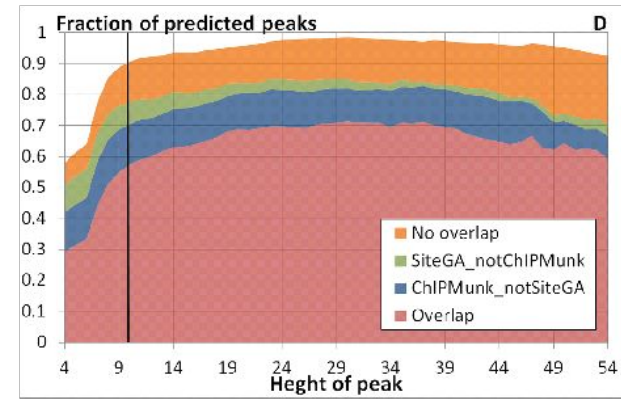
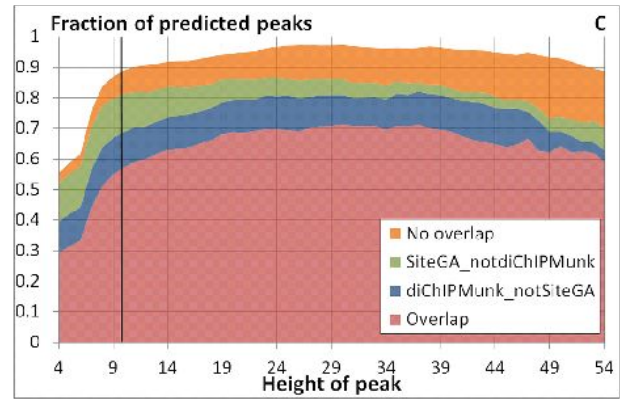
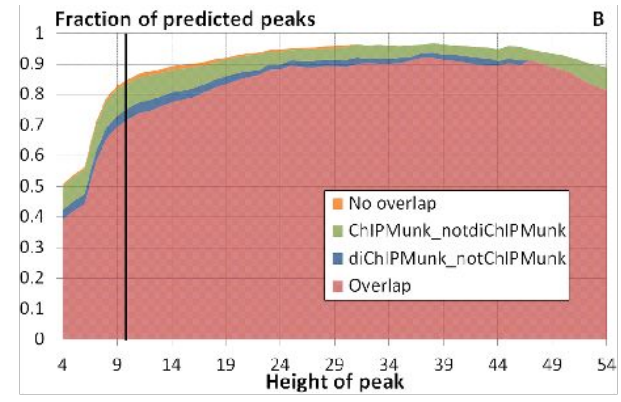
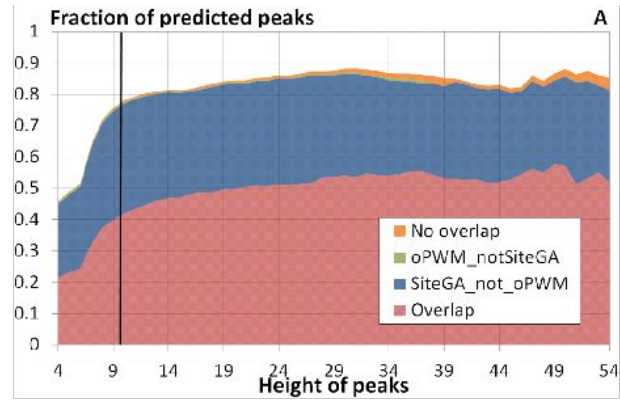


Экспериментальная верификация 64 предсказанных ССТФ ФохА методом задержки в геле (EMSA), конкурентный анализ

Скор EMSA рассчитывался как отношение угла наклона линейной регрессии log-кривой зависимости от концентрации тестируемого олигонуклеотида к положительному контролю (TTR) и использовался в качестве оценки относительной аффинности этого олигонуклеотида. Поскольку скор EMSA распределен достаточно равномерно, порог отсеечения неподтвержденных сайтов был выбран вручную (0.25).







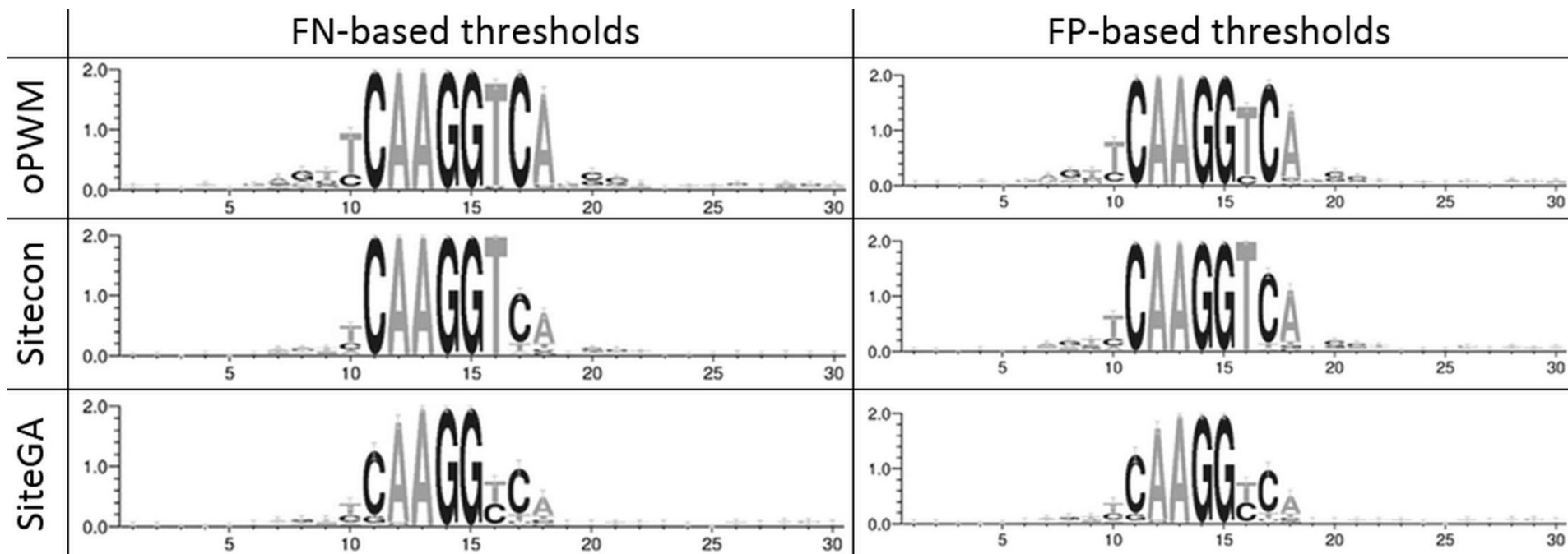
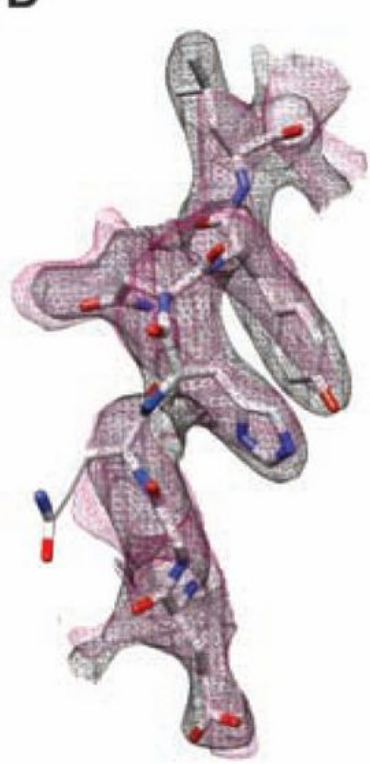


Figure 7. Context patterns of the SF-1 BSs recognized by sequence analysis of the top-scoring ChIP-seq peaks (peak height of 15 or higher). The frequency matrices were visualized by the WebLogo tool (Crooks et al., 2004) for the samples of sites predicted by oPWM, Sitecon, and SiteGA. The *X* axis shows nucleotide positions and *Y* axis, bits.

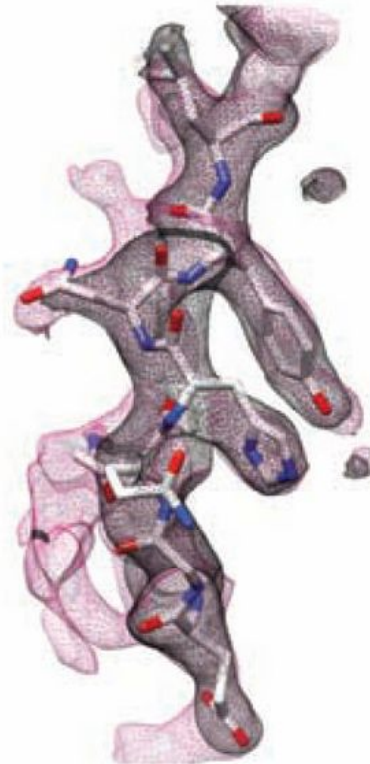


GBS	Sequence	Fold activation dex/etoh	K_D (μM)
Random/-	-	0.8 ± 0.1	>10
GilZ	AGAACAAttgGGTTCC	2.3 ± 0.4	0.89 ± 0.20
Pal	AGAACAaaaTGTTCT	2.0 ± 0.1	0.08 ± 0.01
Sgk	AGAACAAtttTGTC CG	5.4 ± 0.8	0.15 ± 0.03
Tat	AGAACAAtcccTGTACA	22.3 ± 4.6	1.39 ± 0.35
Cgt	AGAACAAtttTG TACG	8.2 ± 1.0	0.24 ± 0.09
Cons	AGAACAaaaTGTACC	5.9 ± 0.7	0.19 ± 0.06
FKBP5	AGAACAagggTGTTCT	5.9 ± 0.4	0.44 ± 0.12

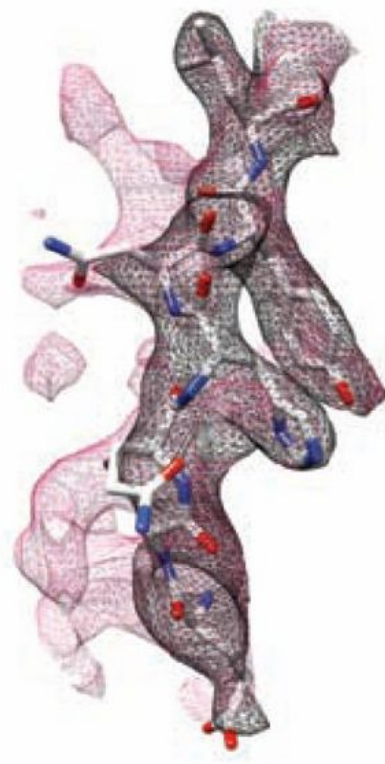
2



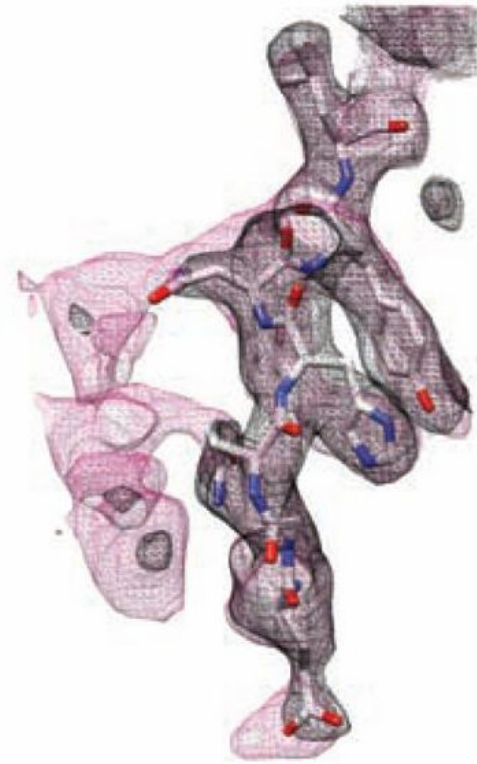
GR-DBD:GilZ



GR-DBD:Sgk

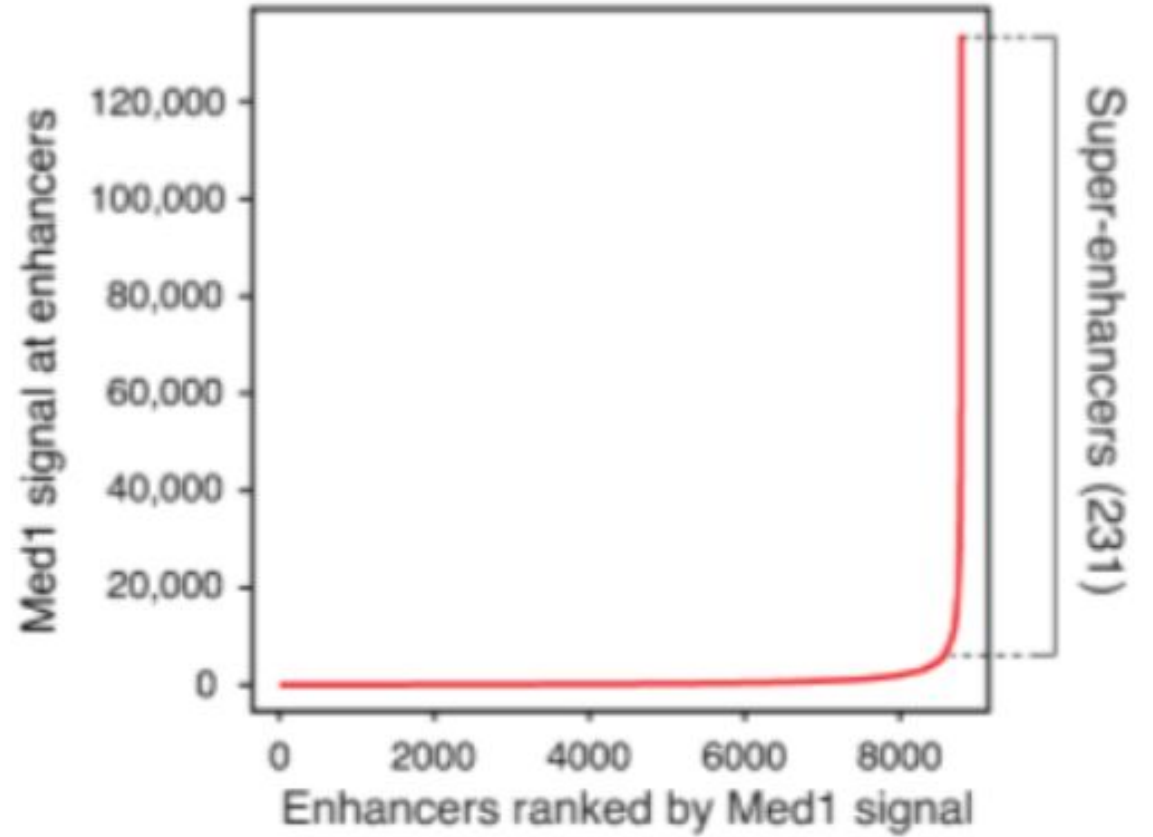
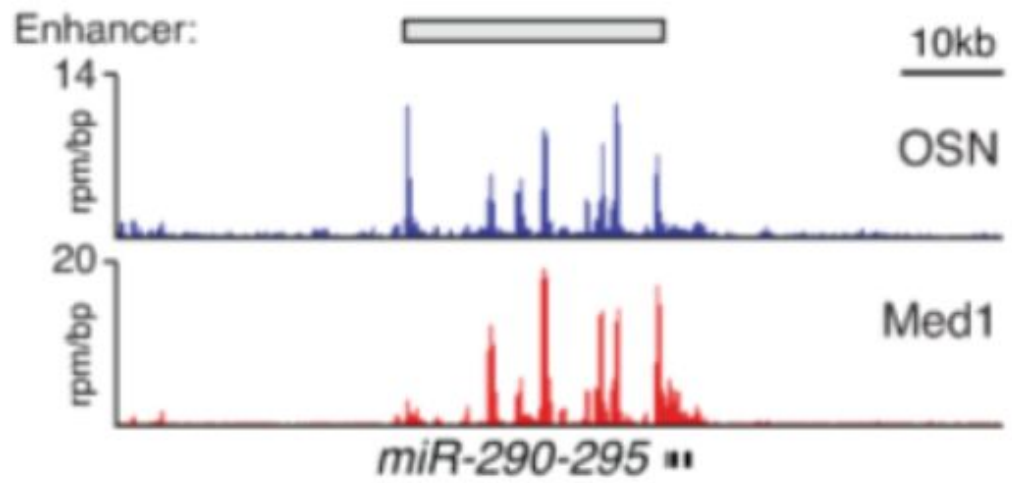
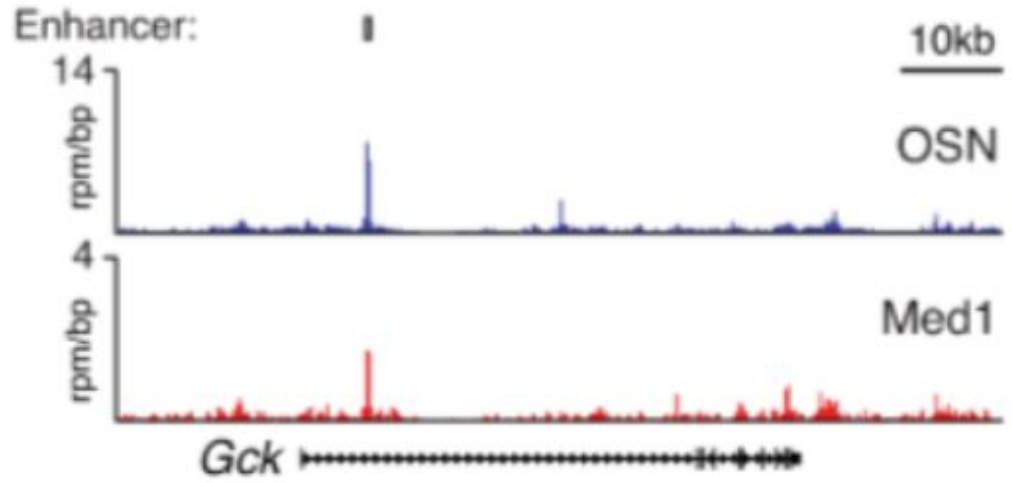


GR-DBD:FKBP5

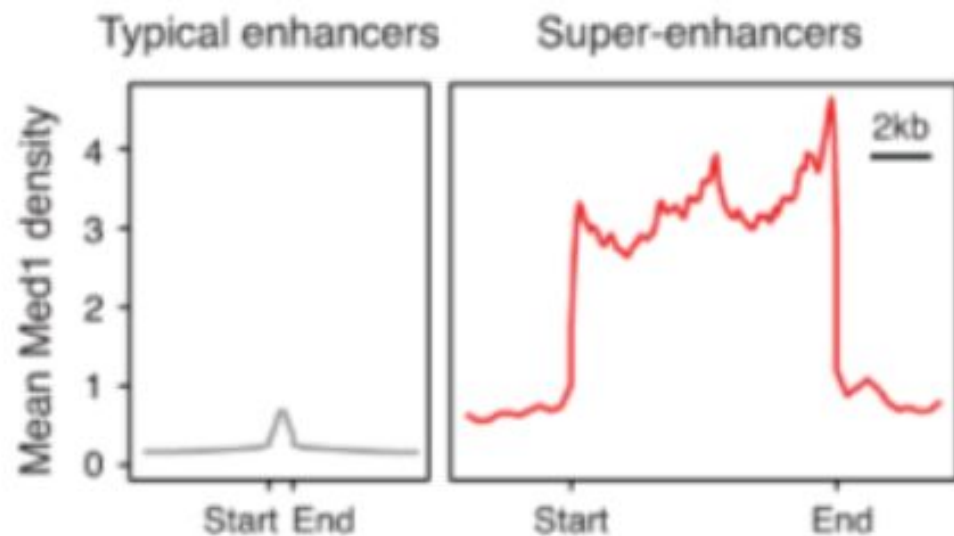


GR-DBD:Pal

СУПЕРЭНХАНСЕРЫ

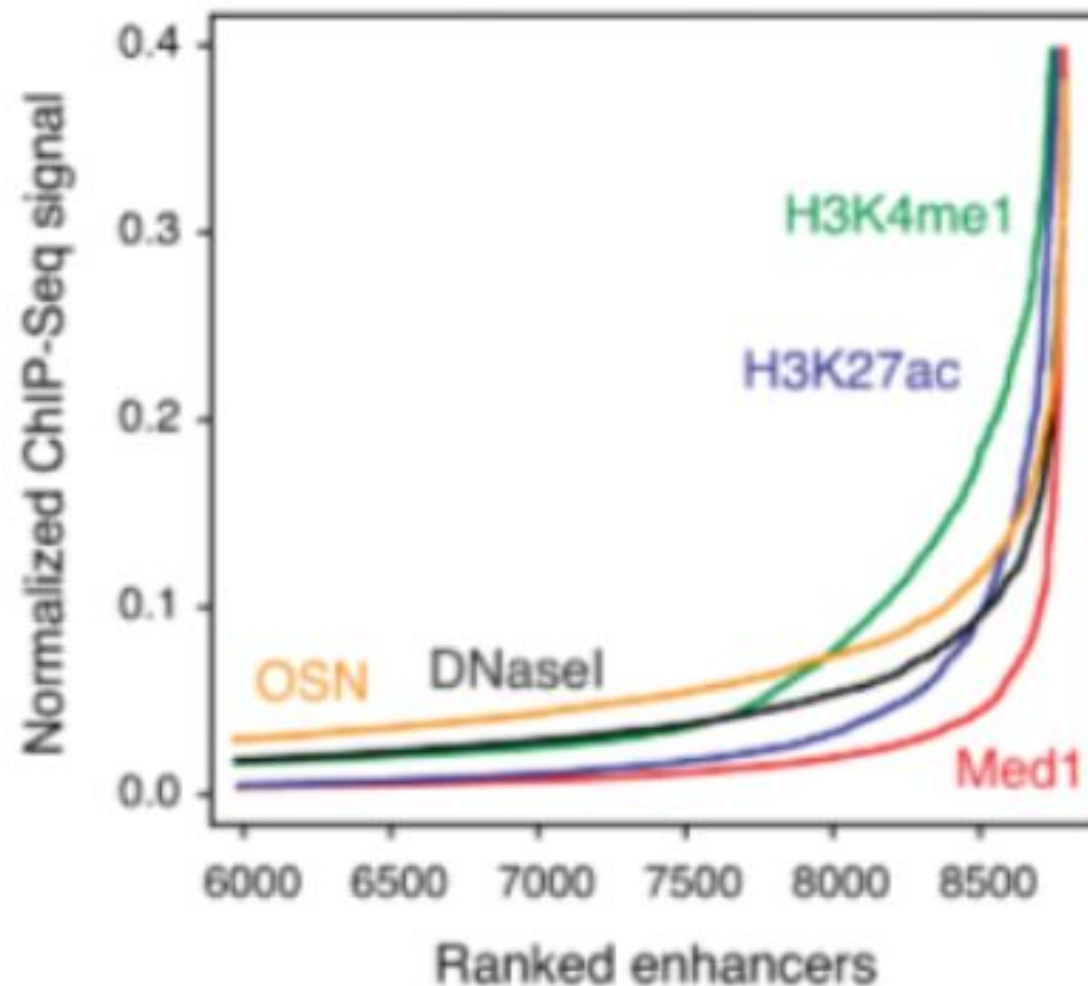


СУПЕРЭНХАНСЕРЫ



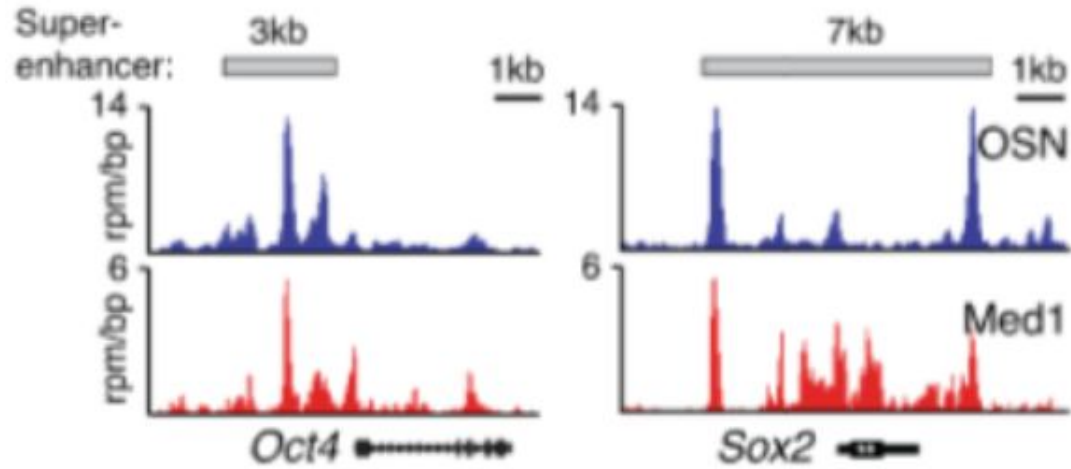
Number: 8563 231
 Median size: 703 bp 8667 bp

	Total signal	Density at constituents	Total signal	Density at constituents
Med1:	1x	1x	28x	8.1x
H3K27ac:	1x	1x	26x	4.8x
H3K4me1:	1x	1x	10x	1.3x
DNaseI:	1x	1x	8x	2.2x



СУПЕРЭНХАНС

ERLI



Selected genes associated with super-enhancers

Transcription factors

Oct4, Sox2, Nanog, Klf4, Esrrb, n-Myc, Utf1, Sall4, Prdm14, Dppa5a, Tbx3, Zfp42

Chromatin modifying enzymes

Tet1, Tet2

Micro RNAs

miR-290-295

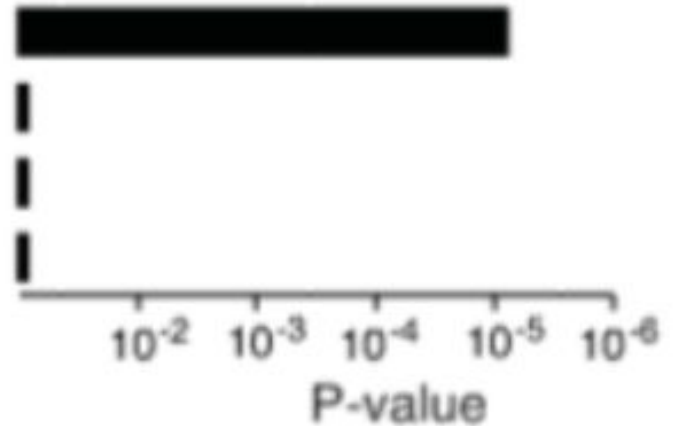
Function

Transcription factor activity

DNA synthesis

Protein synthesis

Metabolism



СУПЕРЭНХАНСЕРЫ

