

A new era in transcription factor research

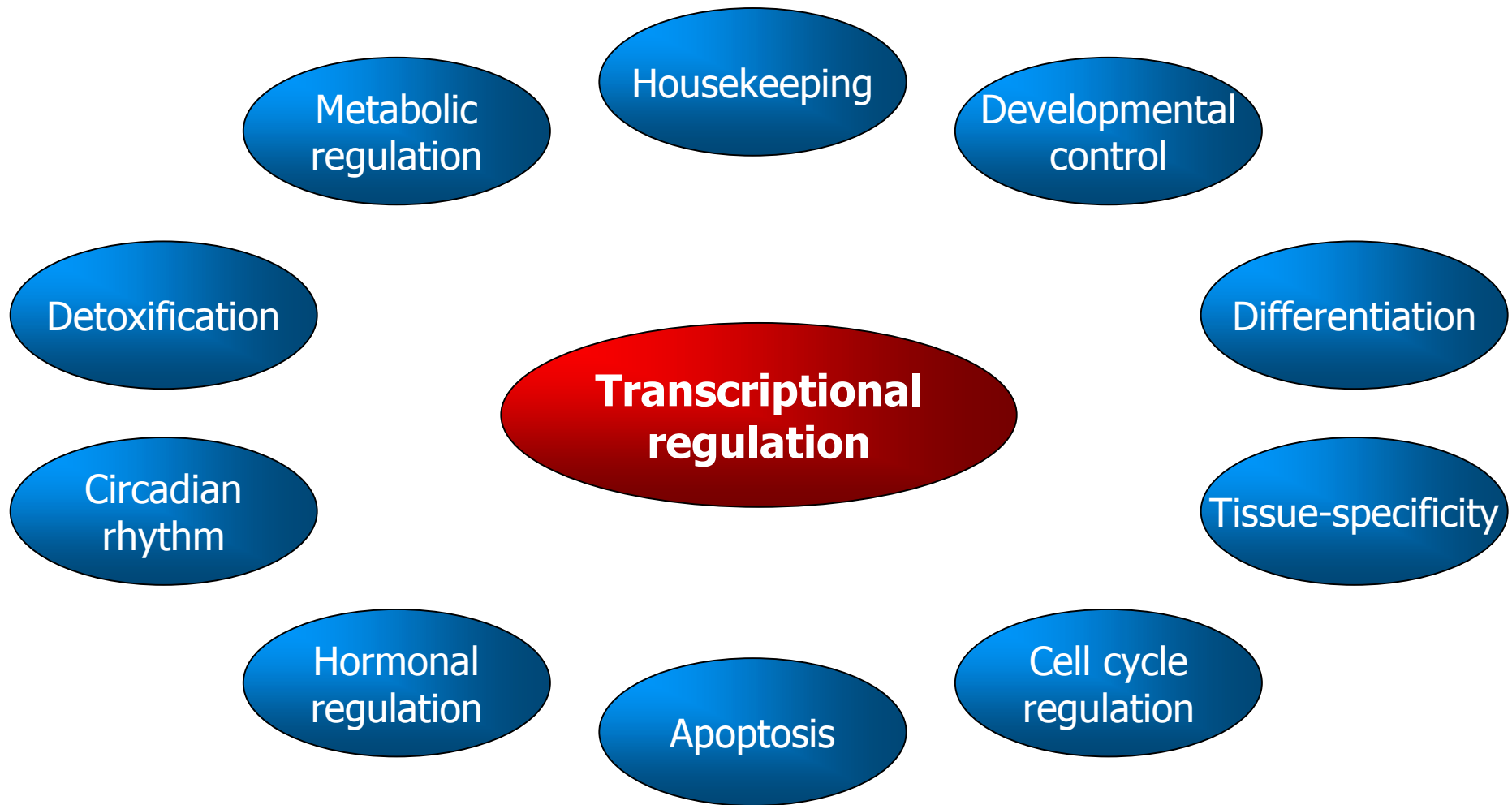
Edgar Wingender

UNIVERSITÄTSMEDIZIN :UMG
GÖTTINGEN

Dept. Bioinformatics,
Medical School
Georg August University
Göttingen, Germany

BIOBASE
BIOLOGICAL DATABASES

Wolfenbüttel, Germany
Beverly, MA / USA
Bangalore, India
Yokohama, Japan



The goals of transcription factor research:

- (1) To understand how the gene-specificity of transcriptional regulation is achieved
- (2) To develop genome-wide maps of transcription factor binding sites (TFBSs)
- (3) To enable prediction of new TFBSs
- (4) To comprehend the complex structure of regulatory genome regions (promoters, enhancers, etc.)
- (5) To predict the DNA-binding specificity of new transcription factors (TFs)
- (6) To construct system-wide transcription networks
- (7) To understand transcriptional dysregulation under disease conditions
- (8) To render transcriptional regulation amenable for targeted alterations

The goals of transcription factor research:

(1) To understand how the gene-specificity of transcriptional regulation is achieved

Biology

(2) To develop genome-wide maps of transcription factor binding sites (TFBSs)

(3) To enable prediction of new TFBSs

Bioinformatics

(4) To comprehend the complex structure of regulatory genome regions (promoters, enhancers, etc.)

(5) To predict the DNA-binding specificity of new transcription factors (TFs)

(6) To construct system-wide transcription networks

Systems biology

(7) To understand transcriptional dysregulation under disease conditions

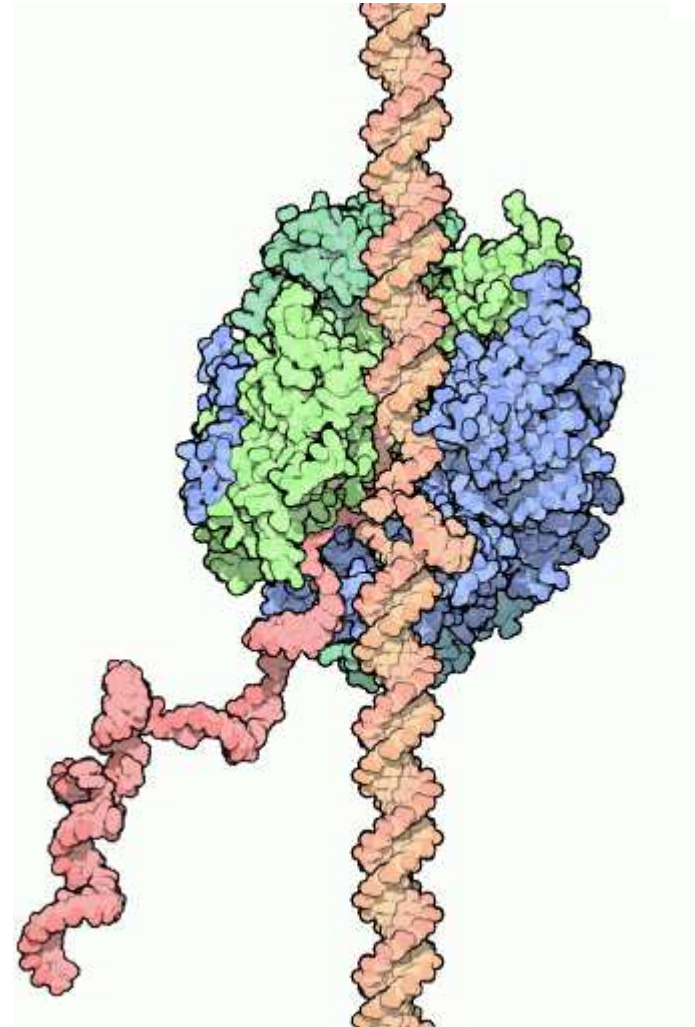
(8) To render transcriptional regulation amenable for targeted alterations

Synthetic biology

The DNA-dependent RNA-polymerases

RNA polymerase II
from *S. cerevisiae*

PDB entry
1I6H



The DNA-dependent RNA-polymerases

In *E. coli*:

RNA polymerase (RNAP):	all genes	5 subunits
------------------------	-----------	------------

In Eukaryotes:

RNA polymerase I (A):	45S-rRNA	7-14 subunits
-----------------------	----------	---------------

RNA polymerase II (B):	mRNA	12 subunits
------------------------	------	-------------

RNA polymerase III (C):	tRNA, 5S-rRNA	10 subunits
-------------------------	---------------	-------------

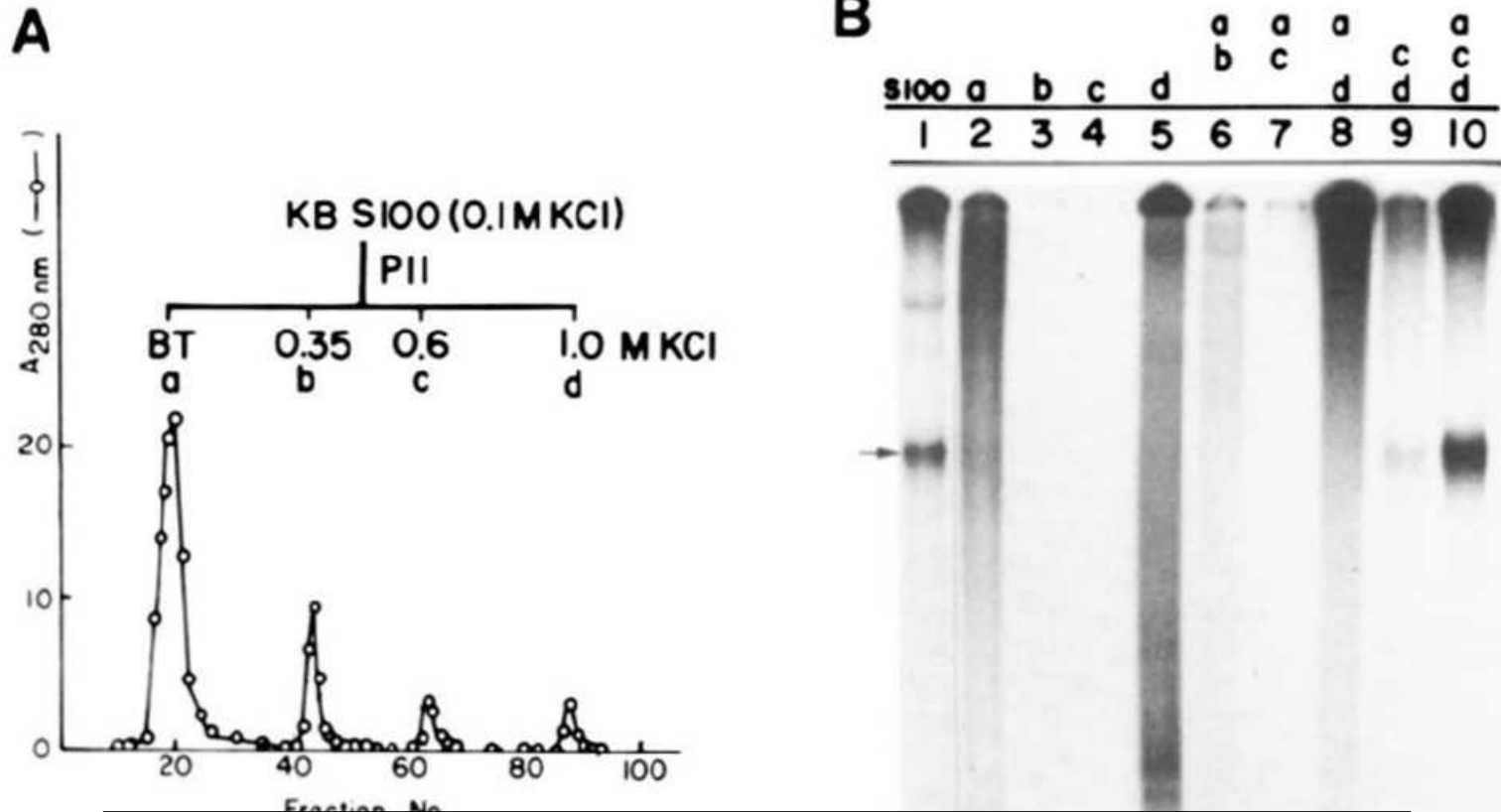
RNA polymerase IV:	siRNA in plants	
--------------------	-----------------	--

The DNA-dependent RNA-polymerases

Goal #1:

Understand how the specificity of eukaryotic transcriptional regulation is achieved?

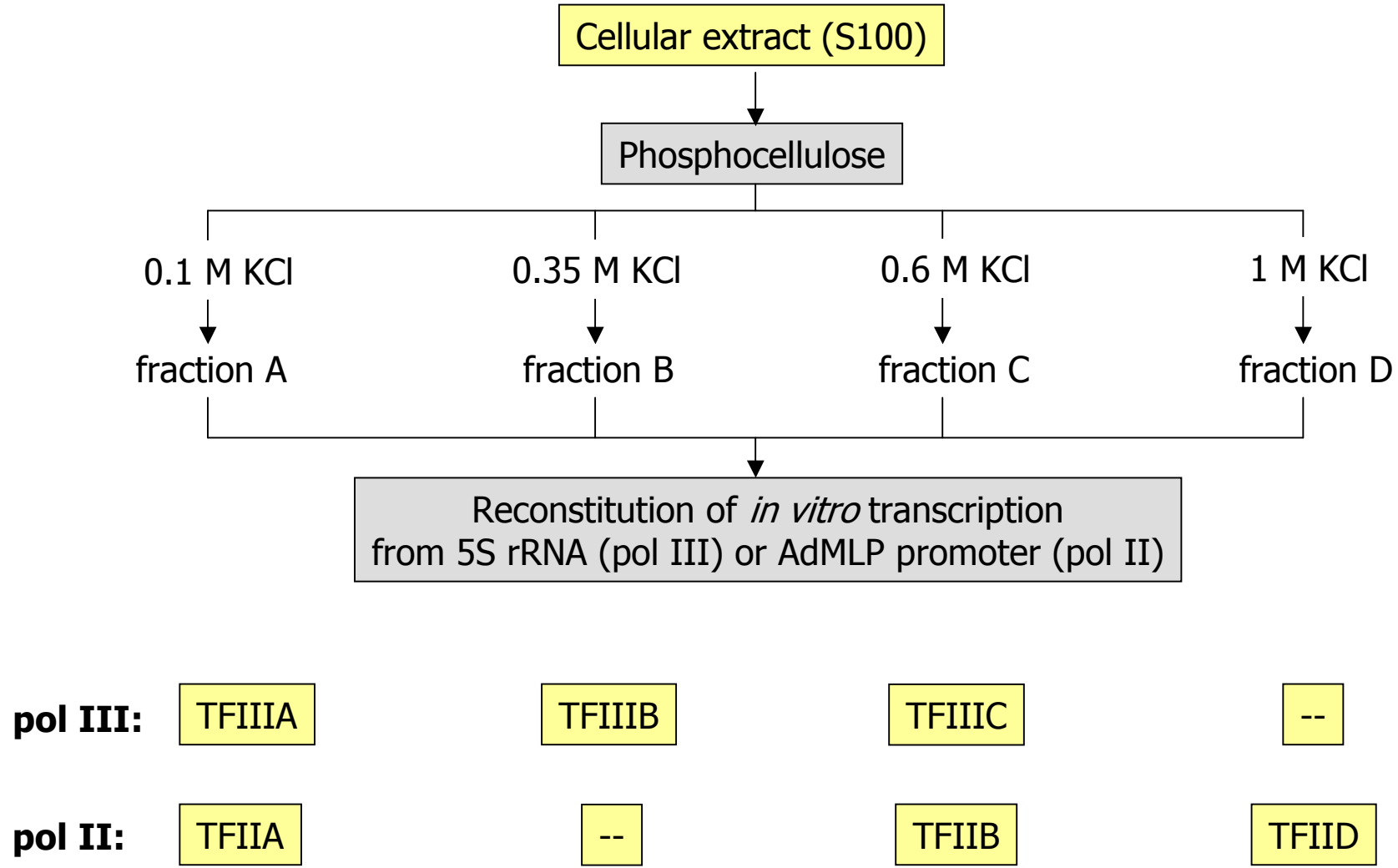
General transcription factors



RNA polymerase requires additional factors for specific promoter recognition.

Multiple factors required for accurate initiation of transcription by purified RNA polymerase II. Matsui T, Segall J, Weil PA, Roeder RG. J Biol Chem. 1980 Dec 25;255(24):11992-6. PMID: 7440580

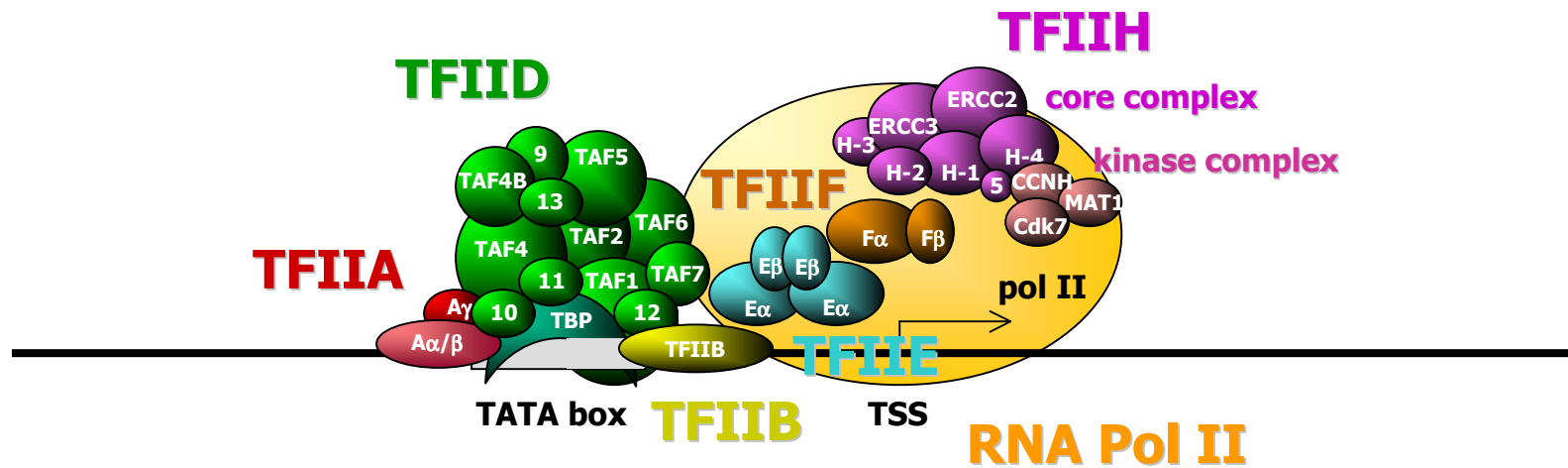
General transcription factors



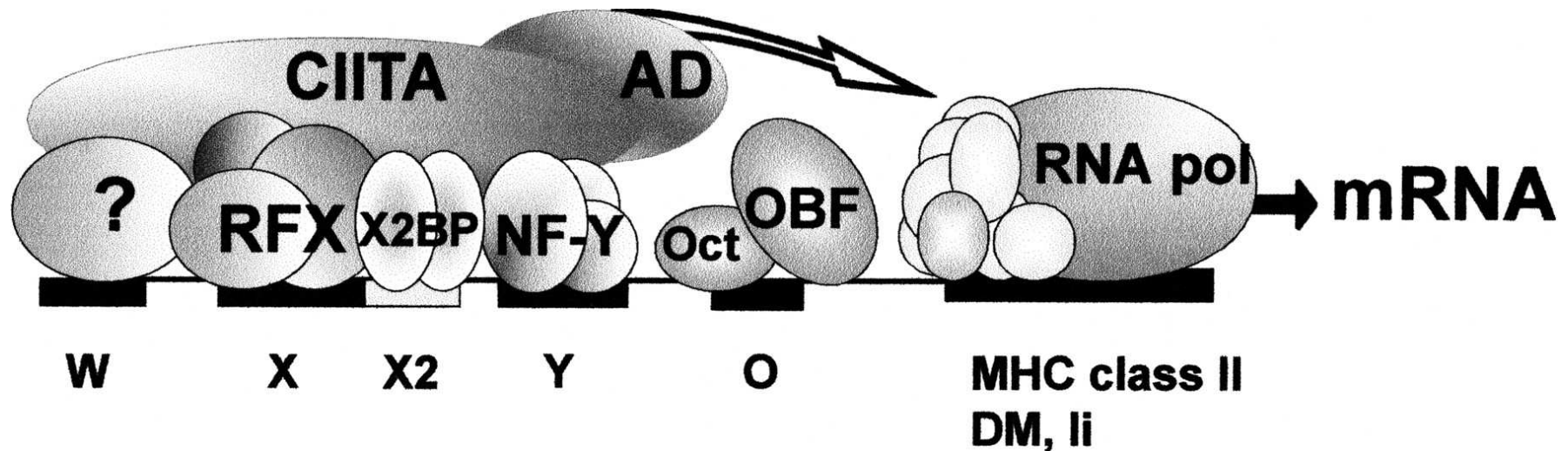
Multiple factors required for accurate initiation of transcription by purified RNA polymerase II. Matsui T, Segall J, Weil PA, Roeder RG. J Biol Chem. 1980 Dec 25;255(24):11992-6. PMID: 7440580

General transcription factors

The pre-initiation transcription complex



Enhanceosome



Master

In addition to the assembly of general transcription factors, „upstream“ factors are required. Many of them are sequence-specific DNA-binding proteins.

Their task is to provide a favorable chromatin structure and/or to facilitate the assembly of the general transcription factors.

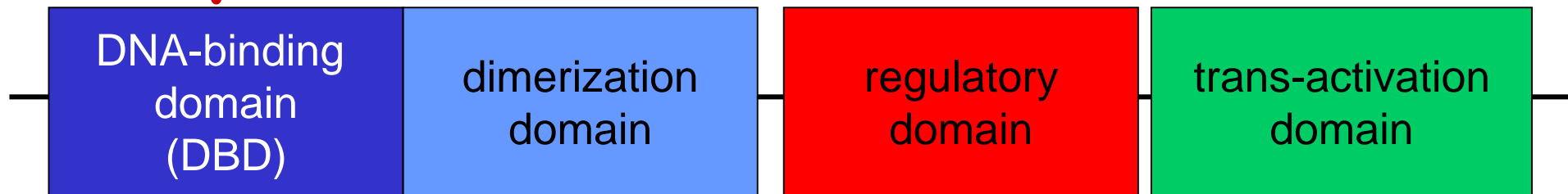
Definition

What is a transcription factor?

A transcription factor is a protein that regulates transcription by specific interaction with DNA or, after nuclear translocation, by stoichiometric interaction with a protein that can be assembled into a sequence-specific DNA-protein complex.

Modular structure of a transcription factor

Primary function:
To recognize *cis*-regulatory elements



Ultimate function:
To activate transcription

Goal #2:

To generate a comprehensive map
of genome-bound proteins.

The first compilation

Volume 16 Number 5 1988

Nucleic Acids Research

Compilation of transcription regulating proteins

Edgar Wingender

Gesellschaft für Biotechnologische Forschung mbH, Mascheroder Weg 1, D-3300 Braunschweig, FRG

Received November 28, 1987; Revised and Accepted January 28, 1988

The genome-wide map: TRANSFAC Site Table



The first compilation

<u>gene</u> <u>gene product</u>	<u>species/tissue</u> ^(a)	<u>protein interacting region</u> ^(b)	<u>method</u> ^(c)	<u>sequence motif</u>	<u>factor bound</u>	<u>ref.</u>
α-actin	chicken/rat myocytes /rat non-myocytes	-83 to -78	3, 4b		myoc.-spec.f. distinct factor	1
			3, 4b			1
β-actin	rat/HeLa		3 (compet. with c-fos)	small, dyad symmetry	SRF ?	2
actin 5C	Drosophila	-38 to +34	1a	TATAAAA	B factor	3
actin (cytoskeletal)	Xenopus laevis /HeLa	-94 to -75 (SRE)	1a, 4a, 4b	AAGATgcCCATATtTGGcgATCTT	SRF (?)	4
Ad2MLP (adenovirus 2 major late prom.)	adenovirus/HeLa	-68 to -49	1a, 3, 4a, 4b	TGTAGGCCACGTGACCCGG	UEF, USF, MLTF	5-10
			1b	GGCCACGTGACC	USF	9
			1a	TATAAAA	?	8
			1a	TATAAAA	TFIID	9
			1b	TATAAAA	TFIID	9
Adh (alcohol dehydrogenase) -distal prom. -proxim.prom.	Drosophila	- 85 to - 47 -269 to -229 -151 to -105 - 98 to - 77	1a	BCTGctGCTGcatcCBTCGacGTcG	Adf-1	11
			1a	TACTAA (4x)		11
			1a	GCAGcGCTGcCCTcGccggctgaGCAGC	Adf-1 ?	11
			1a	GAGATcGCcTAACcGTATGATAA		11
Adh1 (alcohol dehydrogenase)	maize	-190 to -186 -145 to -138 (after ind.) -120 to -117 (after ind.) -108 to -100	4a (in vivo)	CCACG		12
			4a (in vivo)	CCCCGG		12
			4a (in vivo)	CGTGG		12
			4a (in vivo)	CCCACAGGC		12
ADH2 (alcohol dehydrogenase)	yeast	-257 to -216	deletions	22 bp dyad symmetry	ADR1	13
albumin	rat/liver	-156 to -141 -126 to -107 -105 to - 89 - 72 to - 35	1a	GCAAGGGATTTAGTTA	NF-1	14
			1a	TTTTTGGCAAGGAT		14
			1a	ATTTTGTAAAT		14
			1a	TGGTTAATGATCTACAGTTATTGGTTA		14
A-MuLV (amphotrop. murine leukemia virus)	/F9, PCC4	-87 to -59	1a, 3	CCAAT	EPBF	15
aP2 (adipocyte P2)	mouse/adipocytes	-124 to -108	1a, 3	AACATGACTCAGAGGAA	c-fos ?	16

Wingender E.
Compilation of
transcription regulating
proteins. Nucleic Acids
Res. 1988 Mar
25;16(5):1879-902.
PMID: 3282223

The genome-wide map: TRANSFAC Site Table



The present status

A BKL Locus Report
(rat albumin)

Locus Report

Rat Alb (Serum albumin) Albumin, a drug and fatty acid transporter that acts in protein nitrosylation and complex formation, MAPK3/6 cascade, and TGF-beta1 production; human ALB is associated with familial form of hyperthyroxinemia and hypoalbuminemia

Navigator [expand]

ALB [orthogene]
 ALB [cattle, Bos taurus] [gene]
 ALB [Human] [gene]
 Alb [Mouse] [gene]
Alb [Rat] [gene]
 ALB [taxonomic class Mammalia] [gene]
 Albumin [orthologous]

Gene

Gene Regulation [hide]

Visualization Click icon to view in the BKL Pathfinder.

Chromosome 14p22

Promoter 19141231..19152231 PM000039510; albumin

Regulatory Elements

* Note: Only binding sites whose location is relative to the TSS are graphically displayed.

Binding Sites (18 interactions)	Identifier	Location	Binding Factor(s)	DNA Binding Reaction	Effect
	RAT\$ALBU_24	-177 to -154	HNF-3alpha(t)	HNF-3alpha(t) --> ALB(t)	DNA binding
	RAT\$ALBU_24	-177 to -154	HNF-3alpha(t)	HNF-3alpha(t) --> ALB(t)	DNA binding
	RAT\$ALBU_24	-177 to -154	HNF-3beta(t)	HNF-3beta(t) --> ALB(t)	transactivation
	RAT\$ALBU_18	-156 to -141			
	RAT\$ALBU_25	-152 to -132 *	HNF-3alpha(t)	HNF-3alpha(t) --> ALB(t)	DNA binding
	RAT\$ALBU_25	-152 to -132 *	HNF-3beta(t)	HNF-3beta(t) --> ALB(t)	transactivation
	RAT\$ALBU_19	-126 to -107	NF-1(t)	NF-1(t) --> ALB(t)	DNA binding
	RAT\$ALBU_20	-106 to -88			
	RAT\$ALBU_21	-105 to -89	C/EBPalpha(m.s.)	C/EBPalpha(m.s.) --> ALB(t)	transactivation
	RAT\$ALBU_21	-105 to -89	C/EBPalpha(t)	C/EBPalpha(t) --> ALB(t)	transactivation
	RAT\$ALBU_21	-105 to -89	C/EBPbeta(LAP)(t)	C/EBPbeta(LAP)(t) --> ALB(t)	DNA binding
	RAT\$ALBU_21	-105 to -89	C/EBPbeta(t)	C/EBPbeta(t) --> ALB(t)	transactivation
	RAT\$ALBU_21	-105 to -89	C/EBPdelta(m)	C/EBPdelta(m) --> ALB(t)	DNA binding
	RAT\$ALBU_21	-105 to -89	C/EBPepsilon(t)	C/EBPepsilon(t) --> ALB(t)	DNA binding
	RAT\$ALBU_22	-89 to -64	ACF(t)	ACF(t) --> ALB(t)	DNA binding
	RAT\$ALBU_23	-72 to -35	HNF-1alpha(t)	HNF-1alpha(t) --> ALB(t)	DNA binding
	RAT\$ALBU_23	-72 to -35	HNF-1beta(t)	HNF-1beta(t) --> ALB(t)	DNA binding
	RAT\$ALBU_23	-72 to -35	HNF-3alpha(t)	HNF-3alpha(t) --> ALB(t)	DNA binding

[less]

The genome-wide map: TRANSFAC Site Table



The present

Site Report



R26428

Site: RAT\$ALBU_24 - ALB (albumin)

Binding Site Information [hide]



Identifier	RAT\$ALBU_24
Gene	Rat Alb (albumin)
Region	promoter
Sequence	agcttcagaTGGCAAACATACgca
Export	FASTA
Sequence Type	DNA
Reference Point for Sequence Start	
Element	Site X
Element Range	from -177 to -154
External Identifiers	EMBL/GenBank/DDBJ: S82890 ; (327:350)
Element Mapped to Promoter(s)	
Binding Factors	HNF-3beta(r) Quality:1 HNF-3alpha(r) Quality:2 HNF-3alpha(h) Quality:5
Factor Source	rec(rat-COS-7); Rat; recombinant expression: rat factor has been expressed in COS-7 cells HepG2; Human;
Method	functional analysis, direct gel shift, gel shift competition
Comments	Positive regulatory element [1];

A Site entry

References [hide] (1 entry)



[1] *PMID 10024498*. Hsiang, C. H., Marten, N. W., Straus, D. S. Upstream region of rat serum albumin gene promoter contributes to promoter activity: presence of functional binding site for hepatocyte nuclear factor-3. *Biochem J* 338 241-9. (1999)

The present status

Number of Entries

	Compilation 1988	TRANSFAC 2009.1 ^a
Transcription factors	145	12,183 ^b
Binding sites	464	24,745 ^c
Factor-site links	361	33,513
Genes	122	36,317
ChIP-chip fragments	--	155,306
Matrices	--	885
References	209	20,072

New high-throughput technologies quickly populate the genome-wide map.

Additional 130,000 ChIP-Seq fragments coming up with the Summer release.

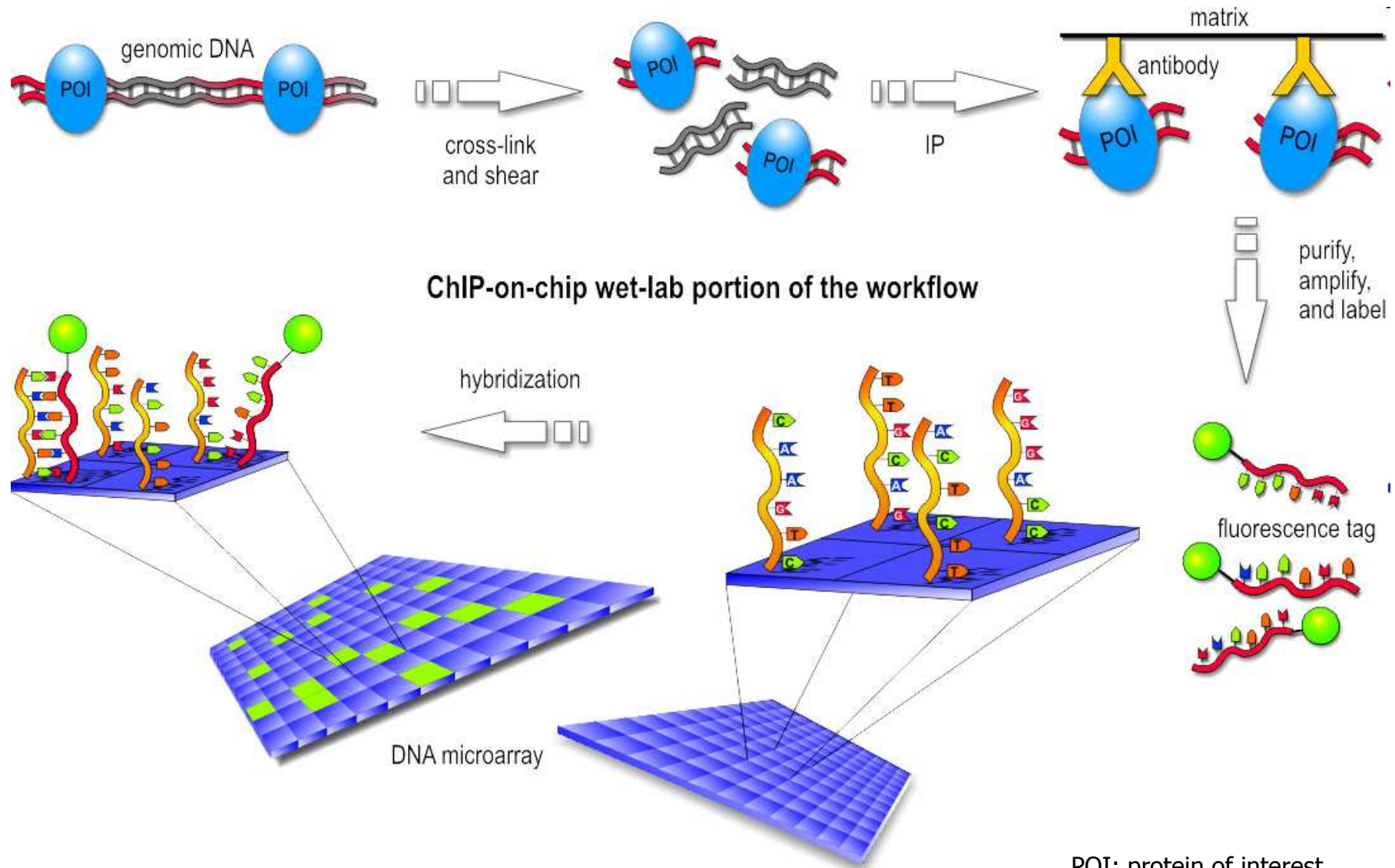
^a as of March 31, 2009

^b including 237 miRNAs

^c including 577 miRNA target sites

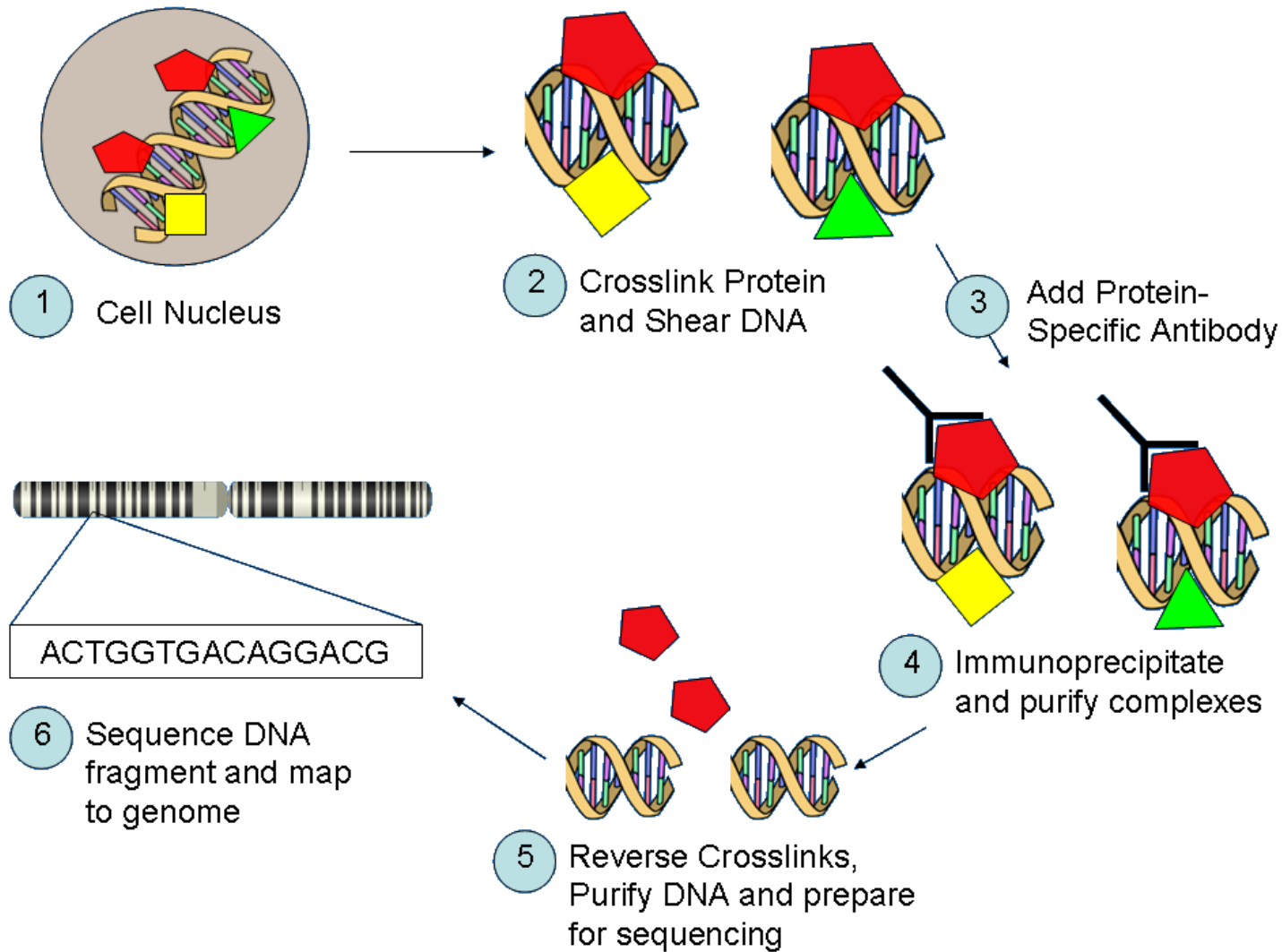
The genome-wide map: TRANSFAC

High throughput approaches: ChIP-chip



The genome-wide map: TRANSFAC

High throughput approaches: ChIP-seq



The present status

Additional 130,000 ChIP-Seq fragments coming up with the Summer release ...

Current data sets (TRANSFAC 2009.1)

T00140	c-Myc(h)	751
T00163	CREB(h)	215
T02284	CTCF(h)	13711
T00368	HNF-1alpha-A(h)	288
T03828	HNF-4alpha(h)	2085
T03286	HNF-6(h)	309
T08969	Nanog(m)	2932
T00671	p53(h)	586
T09363	POU5F1(m)	1049
T00594	RelA-p65(h)	208
T00759	Sp1(h)	125062
T00781	TAF(II)250-isoform2(h)	8523
	sum	155719

ChIP-chip

Available with release of Sep09

T10194	(GABP-alpha(h))2:(GABP-beta(h))2	16489
T02284	CTCF(h)	27119
T02512	HNF-3A(h)	14761
T06124	REST(h)	9607
T00764	SRF(h)	27822
T04759	STAT1(h)	34334
	sum	130132

ChIP-seq

total sum 285851 interacting fragments for 17 TFs

The present status

... and more coming up until the end of 2009, totaling ~790,000 genome fragments interacting with 40 different TFs

Most of the fragments come from the international ENCODE project.

Further upcoming datasets for human TFs:

c-Fos	29391
c-Myc	27921
JunD	6352
NF-E2	13325
SREBP-1a	1147
Tcf4	52293
c-Jun	60574
GATA-1	11097
Max	27084
Pol2	56968
SREBP2	1145
ZNF263	32101
Sum	319398

Further upcoming datasets for mouse TFs:

c-Myc	3422
CTCF	39601
E2F1	20696
ESRRB	21644
Klf4	10873
NANOG	10343
n-Myc	7182
Oct4	376
P300	524
SmaD1	1126
Sox2	4526
Stat3	2546
Suz12	4215
TCFCP2l1	26908
Zfx	10336
FoxA2	11462
P300	2453
Sum	181617

The ENCODE Project



- **ENCyclopedia Of DNA Elements**
- **Goal:** To compile a comprehensive parts list of functional elements in the human genome
- **Pilot Project:**
 - 2003-2007
 - Experiments focused on a limited set of genomic regions comprising about 1% of the human genome
- In September 2007, the ENCODE project was scaled up from pilot to productive phase, aiming at the coverage of the entire human genome

The ENCODE Project

- Protein-coding and non-protein coding genes
 - Full-length coding sequence and variants
 - Transcriptional regulatory elements
 - All pseudogenes
- Global sequence features:
 - Methylation/CpG islands, sequence variation, evolutionary history of sequence blocks, and repetitive elements
 - Non-coding chromosomal elements:
 - Origins of replication, nuclease hypersensitive sites, matrix attachment sites, and histone modifications

The ENCODE Project: experimental protocols applied

Feature Class	Experimental Technique(s)	Abbreviations	No. Exp. Data Points
Transcription	Tiling array, Integrated annotation	TxFrag, RxFrag, GENCODE	63,348,656
5' Ends of transcripts	Tag sequencing	GIS-PET, CAGE	864,964
Histone modifications	Tiling array	Histone nomenclature, RFBR	4,401,291
Chromatin structure	QT-PCR, Tiling array	DHS, FAIRE	15,318,324
Sequence-specific factors	Tiling array, tag sequencing, Promoter assays	STAGE, ChIP-Chip, ChIP-PET, RFBR	324,846,018
Replication	Tiling array	TR50	14,735,740
Computational analysis	Computational methods	CCI, RFBR Cluster	NA
Comparative sequence analysis	Genomic sequencing, multi-sequence alignments, computational analyses	CS	NA
Polymorphisms	Resequencing, copy no. variation	CNV	NA

The ENCODE Project: Status

- Pilot project results (Nature **447**, 799-816, 2007):
 - Transcription more complex than expected
 - Transcription Start Sites (TSS) more numerous than protein-coding genes
 - Regulatory information is distributed:
 - clustered across the genome, distribution near TSSs is symmetrical
 - Many distal DNaseI Hypersensitive Sites (DHSs)
 - Replication correlated with histone structure in a more detailed manner than previously known
- 9 Dec 2008 - First ENCODE whole-genome data freeze completed

However, even with the new HTP technologies, we will not be able to investigate the complete „promotome“ for:

- all TFs (~ 1800 TF genes in human) in
- all cell types (~ 300 human cell types),
- all organs (~ 14000 morphologically distinguishable structures) at
- all developmental stages (≥ 24 in human) and
- under all environmental conditions ($\rightarrow \infty$).

But what we can do is to learn about the rules behind the observations and to develop predictive models.

Goal #3:

To enable prediction of new
transcription factor binding sites.

Prediction-relevant contents in TRANSFAC

Number of Entries

	Compilation 1988	TRANSFAC 2009.1 ^a
Transcription factors	145	12,183 ^b
Binding sites	464	24,745 ^c
Factor-site links	361	33,513
Genes	122	36,317
ChIP-chip fragments	--	155,306
Matrices	--	885
References	209	20,072

Sites and matrices can be used to predict TFBS.

^a as of March 31, 2009

^b including 237 miRNAs

^c including 577 miRNA target sites

Matrix construction

TRANSFAC


A	3	3	1	11	0	0	0	0	2	2	0	0
C	6	2	1	0	12	0	0	0	7	3	2	4
G	1	7	3	1	0	12	0	12	3	6	8	4
T	2	0	7	0	0	0	12	0	0	1	2	4

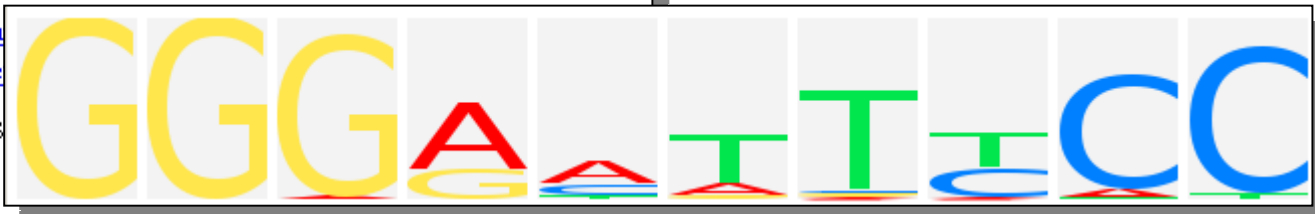


GCCCTACGTGCTGTCTCA
CAGGCAACGTGCAGCCGGA G
CAGTGCATACGTGGGCTCCA
CTTTGTGTGTACGTGCAGGAA
GAAATACGTGCGCTTTGTGTG
CGCGAGCGTACGTGCCTCAGG
CCCCCTCGGACGTGACTCGGACCAC
AGGGCCGGACGTGGGGCCCC
GGAGTACGTGACGGAGCCCC
ACGCTGAGTGCGTGCGGGAC
CCCAGCCTACACGTGGGGTTC
GGAGCCCAGCGGACGTGCGGGAA

Prediction of TFBSs



<u>Accession Number</u>	M00054																																																																		
<u>Identifier</u>	V\$NFKAPPAB_01																																																																		
<u>Created</u>	13.04.1995 by h																																																																		
<u>Updated</u>	30.07.1996 by e																																																																		
Copyright	Copyright (C), B																																																																		
<u>Name</u>	NF-kappaB																																																																		
<u>Factor Description</u>	NF-kappaB																																																																		
<u>Binding factors</u>	T00587 ; NF-kappaB; Species: rat, Rattus norvegicus. T00588 ; NF-kappaB; Species: mouse, Mus musculus. T00590 ; NF-kappaB; Species: human, Homo sapiens. T00593 ; p50; Species: human, Homo sapiens. T00594 ; RelA-p65; Species: human, Homo sapiens.																																																																		
<u>Binding Matrix</u>	<table border="1"> <thead> <tr> <th></th> <th>A</th> <th>C</th> <th>G</th> <th>T</th> <th>Consensus</th> </tr> </thead> <tbody> <tr> <td></td> <td>0</td> <td>0</td> <td>40</td> <td>0</td> <td>G</td> </tr> <tr> <td></td> <td>0</td> <td>0</td> <td>40</td> <td>0</td> <td>G</td> </tr> <tr> <td></td> <td>1</td> <td>0</td> <td>39</td> <td>0</td> <td>G</td> </tr> <tr> <td></td> <td>27</td> <td>0</td> <td>13</td> <td>0</td> <td>A</td> </tr> <tr> <td></td> <td>21</td> <td>13</td> <td>1</td> <td>5</td> <td>M</td> </tr> <tr> <td></td> <td>8</td> <td>1</td> <td>3</td> <td>28</td> <td>T</td> </tr> <tr> <td></td> <td>1</td> <td>2</td> <td>2</td> <td>35</td> <td>T</td> </tr> <tr> <td></td> <td>2</td> <td>18</td> <td>0</td> <td>20</td> <td>Y</td> </tr> <tr> <td></td> <td>3</td> <td>36</td> <td>0</td> <td>1</td> <td>C</td> </tr> <tr> <td></td> <td>0</td> <td>38</td> <td>0</td> <td>2</td> <td>C</td> </tr> </tbody> </table> <p>  Click for fullsize view! </p>		A	C	G	T	Consensus		0	0	40	0	G		0	0	40	0	G		1	0	39	0	G		27	0	13	0	A		21	13	1	5	M		8	1	3	28	T		1	2	2	35	T		2	18	0	20	Y		3	36	0	1	C		0	38	0	2	C
	A	C	G	T	Consensus																																																														
	0	0	40	0	G																																																														
	0	0	40	0	G																																																														
	1	0	39	0	G																																																														
	27	0	13	0	A																																																														
	21	13	1	5	M																																																														
	8	1	3	28	T																																																														
	1	2	2	35	T																																																														
	2	18	0	20	Y																																																														
	3	36	0	1	C																																																														
	0	38	0	2	C																																																														
<u>Basis</u>	40 binding sites from 30 genes (26 cellular genes and 4 viral genomes)																																																																		



TRANSFAC Matrix Table

TFBS detection with TRANSFAC matrices and the Match™ algorithm:

$$q = \left(\sum_{i=1}^L I(i) f_{i,b_i} - \sum_{i=1}^L I(i) f_i^{\min} \right) / \sum_{i=1}^L I(i) f_i^{\max}$$

with: b_i , nucleotide b found in the i -th position of test sequence,
 f_{bi} , frequency of nucleotide b in the i -th position of the aligned training sequences,
 f_i^{\min} , minimum frequency in position i ,
 f_i^{\max} , maximum frequency in position i ,
and

$$I(i) = \sum_{B \in \{A, T, G, C\}} f_{i,B} \ln(4 f_{i,B}), \quad i = 1, 2, \dots, L$$

Validation of potential TFBSs by comparative genome analysis:

General idea:

Genomic sites that are functionally important are under evolutionary pressure and, thus, are more conserved among related genomes than the genomic background („phylogenetic footprinting“).

This could add an independent criteria to the pattern-based prediction of TFBSs.

Validation of potential TFBSs by comparative genome analysis:

However:

What has to be conserved, the sequence or the pattern?

Validation of potential TFBSs by comparative genome analysis:

Sequence-only conservation:

core_sim/matrix_sim

GGGGAATTTCC (NF- κ B consensus):

1.000 / 0.997

** *****

GGAGAATTTCC (10/11 match, 91%):

0.713 / 0.834

** * *****

GGAGTATTTCC (9/11 match, 82%):

0.426 / 0.672

Matrix V\$NFKB_Q6, M0194

Pattern-only conservation:

AATGCCTGAGGCGCT (AP-2 α):

0.991 / 0.983

*** ** (6/15 match, 40%)

TTCGCCCCAGGGCGC (AP-2 α):

0.957 / 0.951

Matrix V\$AP2ALPHA_02, M01045

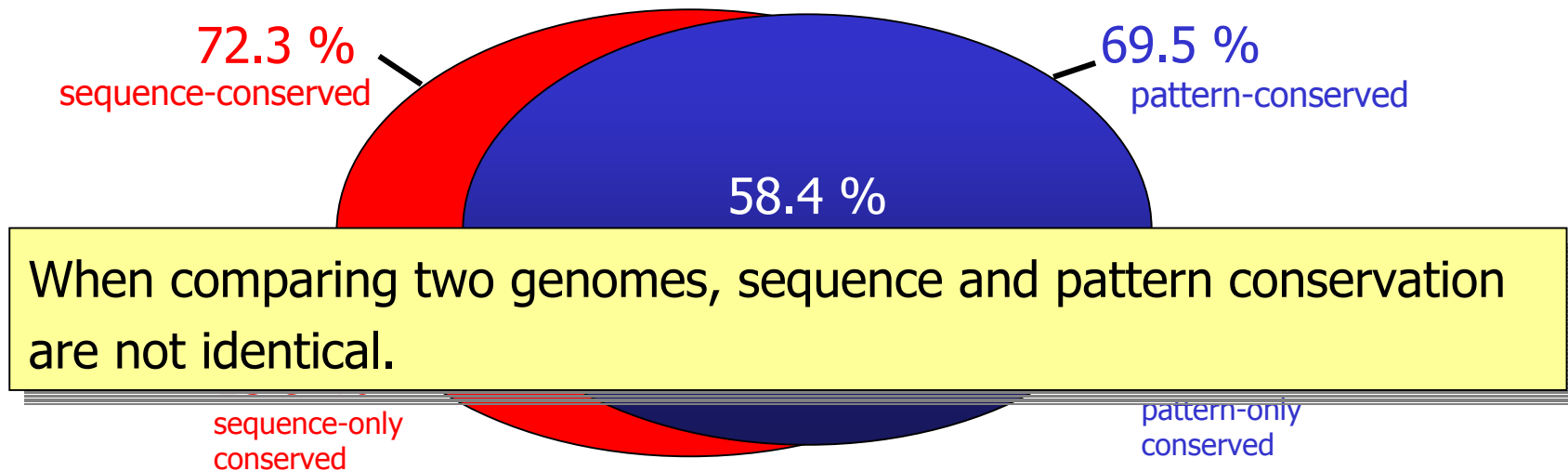
Validation of potential TFBSs by comparative genome analysis:

Human / rodent comparison for all corresponding TRANSFAC sites:

General sequence conservation of TFBS: 72.3 %

Background conservation in upstream sequences: 35.2 %

Pattern-conserved TFBS: 69.5 %



Validation of potential TFBSs by comparative genome analysis:

human *c-jun* gene

```

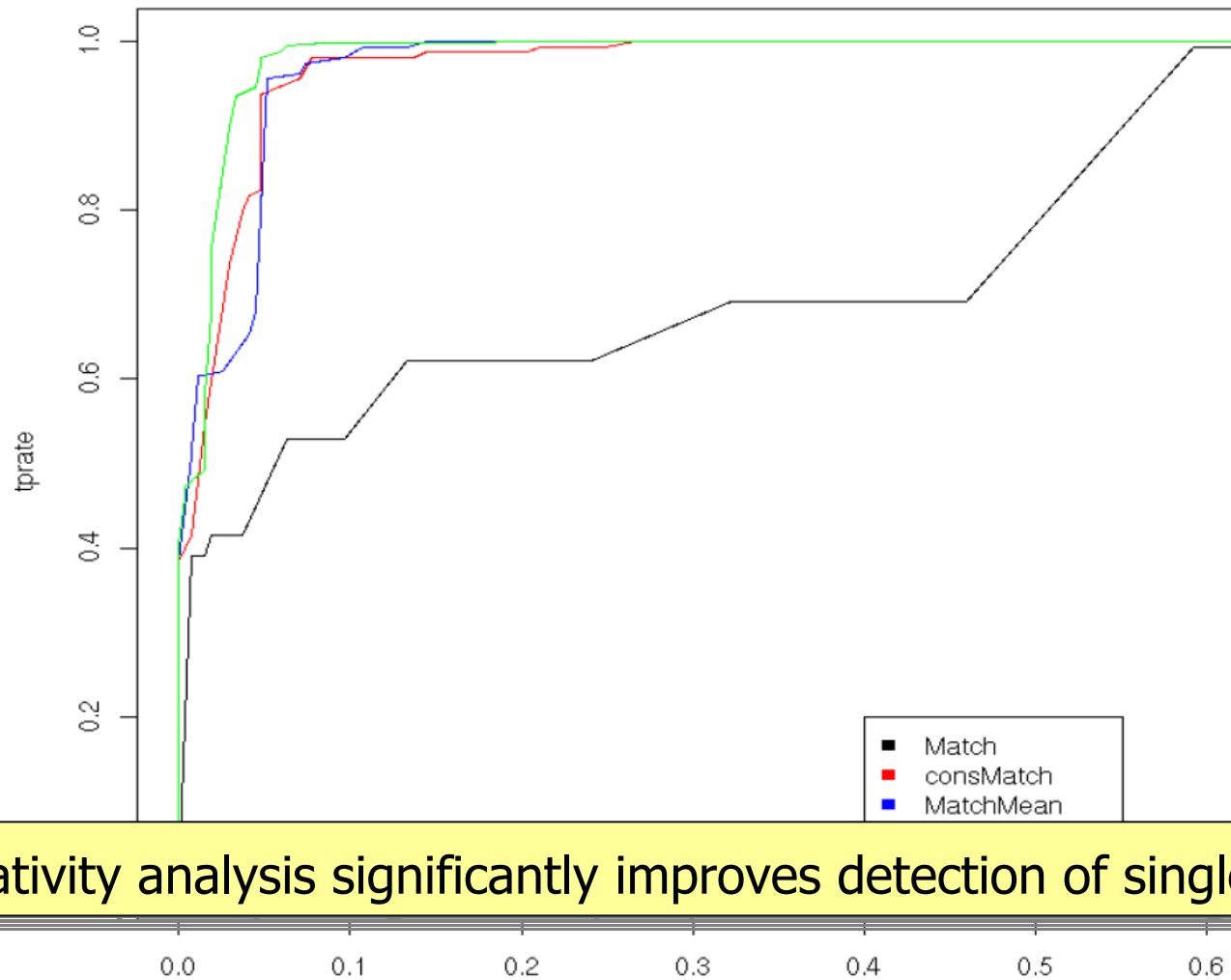
                                R00960 (AP-1)
cow   : GCTGCG-----CTCGAGAGAGCTCCGTGAGTGACCGCGACTTT @ 953/1171
dog   : GCCGCGCCCCGAGAGCGCCCCGAGCGCGccccgtgggtgaCCGCGACTTT @ 945/1153
human : GCTGCG-----CACGAAGAGCGCTCAGTGAGTGACCGCGACTTT @ 1001/1217
mouse : GCTGCT-----CCCCGAGAGCGCTCCGTGAGTGACCGCGACTTT @ 964/1173
rat   : GCTGCT-----TCCCGAGAGCGCTCCGTGAGTGACCGCGACTTT @ 957/1165
      = 1051      1061      1071      1081      1091

                                M00925 (AP-1)
cow   : T--CAAAGCCGGGCGGGCGCGCGCG--AGCCGACAAGTAAGAGCGCGGGC @ 998/1171
dog   : TCCCGAGGCGGGCAGCGCGCGCGCCAGCTGACCCAGGAGGAGCGCGG-- @ 993/1153
human : T--CAAAGCCGGGTAGCGCGCGCG--AGTCGACAAGTAAGAGTGCGGGA @ 1046/1217
mouse : T--CAAAGCTCGGCATCGCGCGGG--AGCCTACCAACGTGAGTGCTAGC @ 1009/1173
rat   : T--CAAAGCTCCGGATCGCGCGGG--AGCCAACCAACGTGAGTGCAAGC @ 1002/1165
      = 1101      1111      1121      1131      1141
```

Multiple genome comparison may better reveal conserved sites, and may pinpoint to „surrogate“ sites.

Validation of potential TFBSs by comparative genome analysis:

ROC-curve NFK(m774)

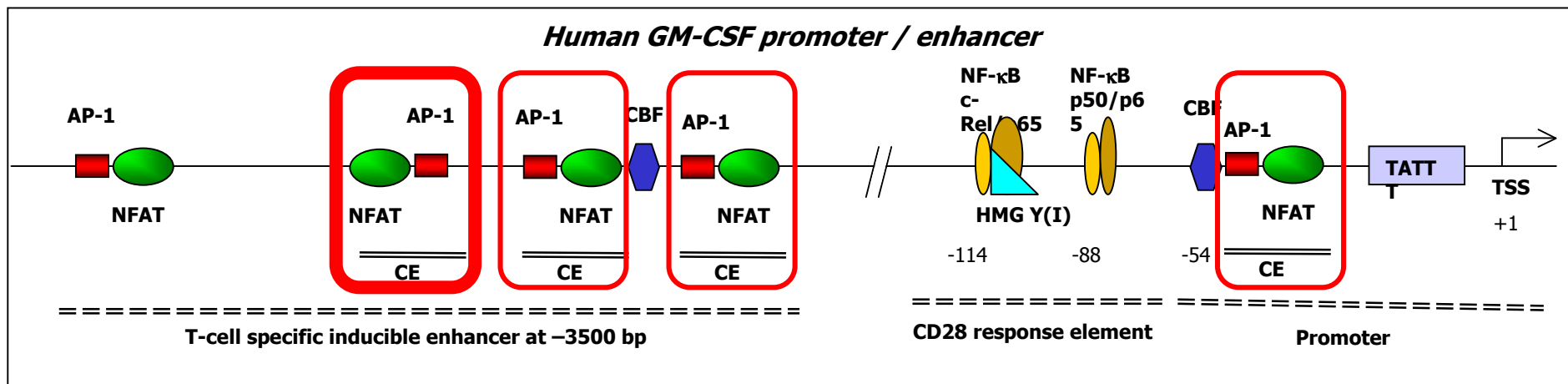
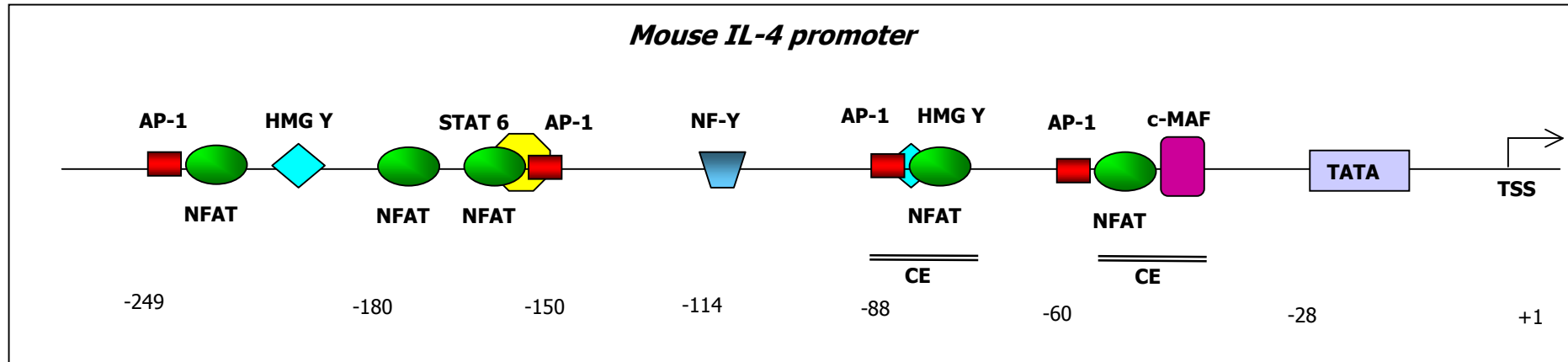


Conservativity analysis significantly improves detection of single TFBSs.

Goal #4:

To reveal the complex structure of regulatory genome regions (promoters, enhancers, etc.).

Promoters and enhancers: examples

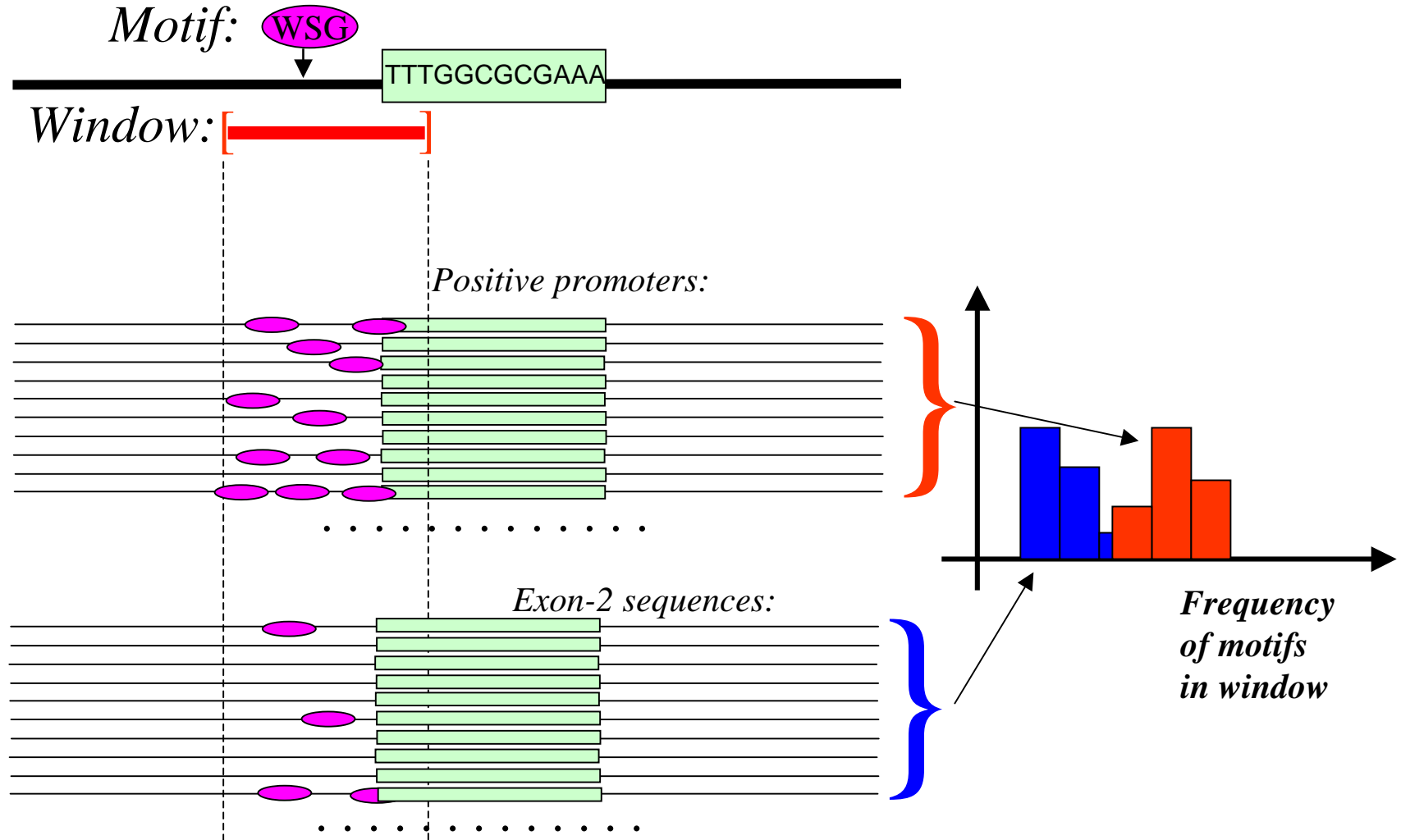


Local context analysis

Task:

Finding nucleotide patterns around specific TFBSs which may facilitate or impede TF binding.

Local context analysis



Local context analysis

$$f(\lambda, w, S) = \frac{N(\lambda, w, S)}{t_2 - t_1}$$

$N(\lambda, w, s)$: Number of motifs λ in window $w=[t_1, t_2]$ of sequence S

Context Score:

$$d(S) = \beta + \sum_{i=0}^k \alpha_i \times f(\lambda_i, w_i, S)$$

	N	Motif (λ)	Window (w) ¹⁾	$\hat{f}^Y / \hat{f}^{N \ 2)}$	Utility	α_i
Positive characteristics	1	MGCG	[27,34]	0.0048 / 0.0041 = 1.179	0.80	0.394
	2	TTT	[39,41]	0.0112 / 0.0032 = 3.536	0.75	0.9618
	3	CGSK	[17,38]	0.0851 / 0.0341 = 2.499	0.90	0.5353
	4	HKCG	[13,16]	0.0675 / 0.0095 = 7.071	0.79	0.5904
	5	VDWW	[17,46]	0.1233 / 0.0536 = 2.299	0.72	0.223
	6	DWTT	[21,26]	0.0337 / 0.0000	0.80	0.5036
	7	GSDM	[3,69]	0.0980 / 0.0559 = 1.754	0.82	0.595
Negative characteristics	8	VWS	[7,66]	0.1258 / 0.1932 = 0.651	0.91	-0.095
	9	HSWY	[26,65]	0.0413 / 0.0813 = 0.508	0.79	-0.2297
	10	VTV	[19,34]	0.0427 / 0.1354 = 0.315	0.71	-0.261
	11	BAY	[7,65]	0.0274 / 0.0614 = 0.447	0.78	-0.566
						$\beta = -5.6767$

E2F binding sites start at pos. 31; Y and N: positive and negative sequence set; "utility": $-1 < U < +1$, calculated from the average difference, distribution overlap, normality of the distribution, *bootstrap* test, etc.

Global context analysis

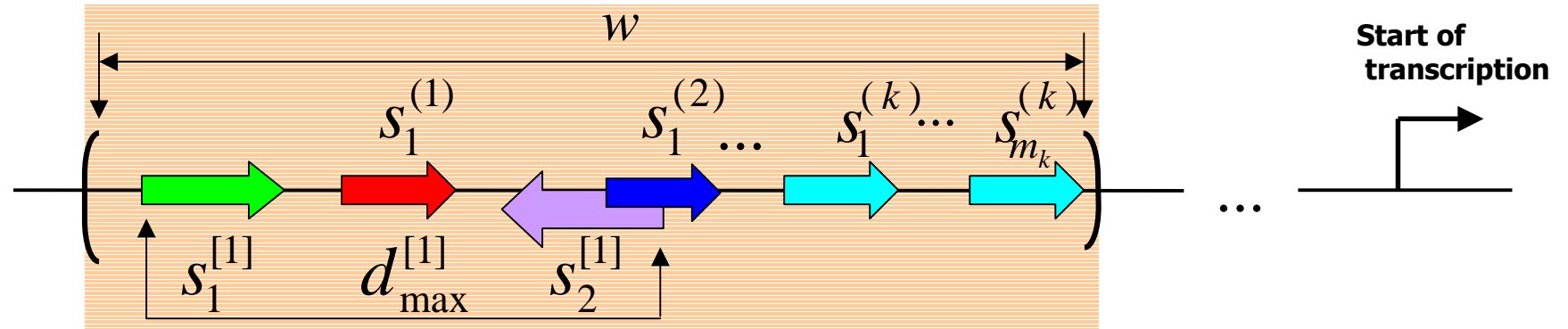
Task:

Finding (a) combination(s) of TFBS („Composite Modules“) that is/are characteristic for a given set of promoters.

Composite module analyzer (CMA)

Kel et al., Bioinformatics 22, 1190-1197 (2006).

Mathematical model (v.2)

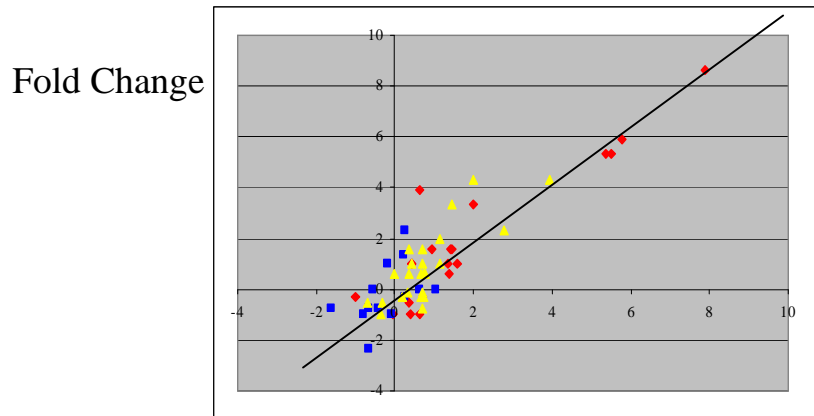
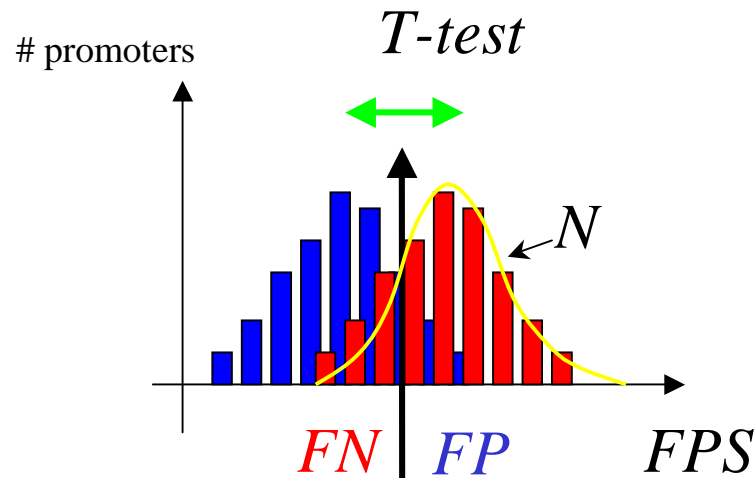


A CM contains single elements as well as composite elements (site pairs)

$d_{\max}^{[1]}$	$d_{\max}^{[1]}$...	$d_{\max}^{[R]}$	} Parameters of the model to be estimated by a Genetic Algorithm
$q_{cut-off}^{(1)}$	$q_{cut-off}^{(2)}$...	$q_{cut-off}^{(k)}$	
$\phi^{(1)}$	$\phi^{(2)}$...	$\phi^{(k)}$	

Fitness function of the Genetic Regression Algorithm (GRA)

$$F = \alpha \cdot R + \beta \cdot (1 - FN) + (1 - \beta) \cdot (1 - FP) + \gamma \cdot T + \delta \cdot N - \mu \cdot k$$



R – linear regression

FN – false negatives

FP – false positives

T – T-test (difference between mean values)

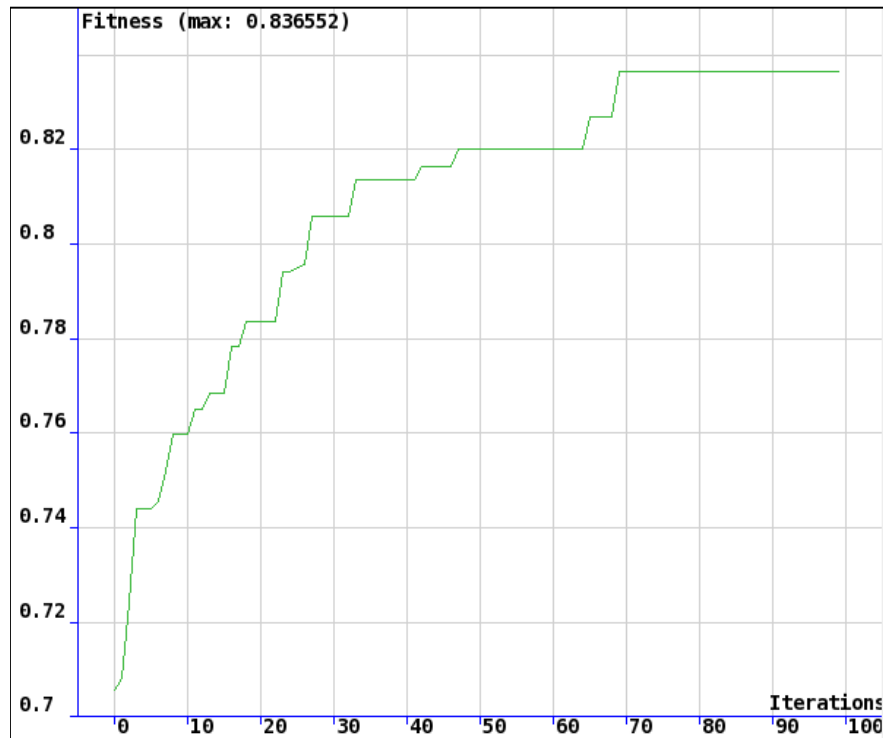
N – normal likeness

k – number of free parameters

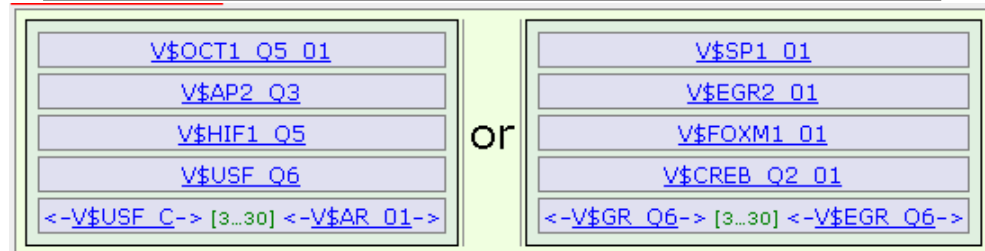
Composite module analyzer (CMA)

Promoter model generation

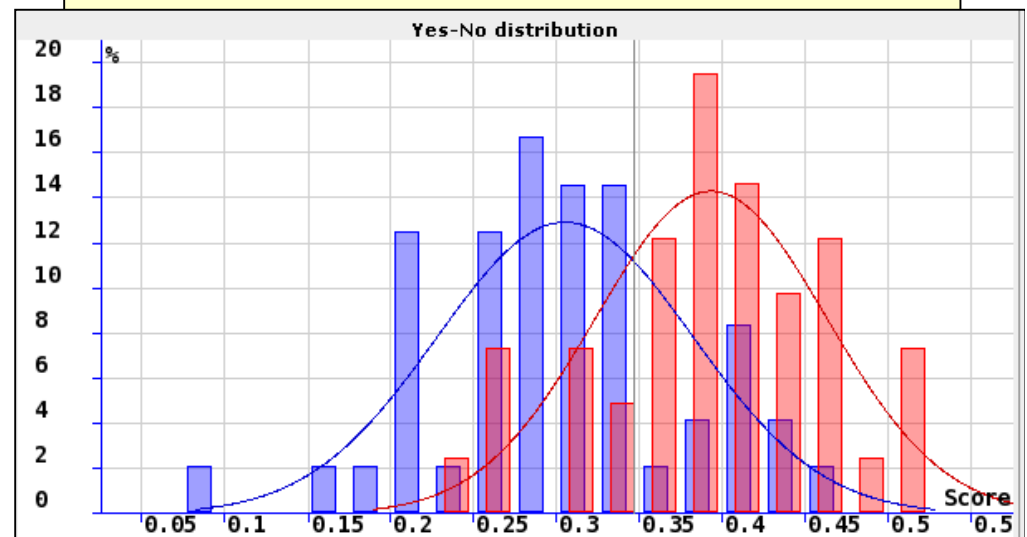
Increase of Fitness function with number of iterations



Composition of the promoter model

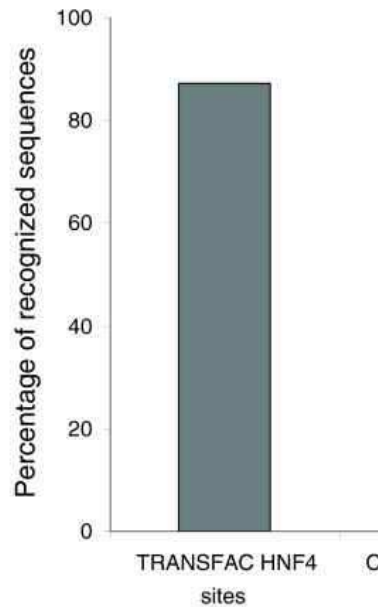


Sequences of YES and NO sets are well separated by the selected promoter model



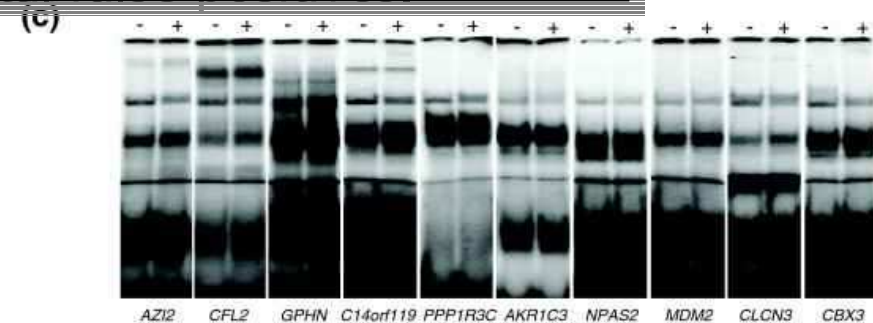
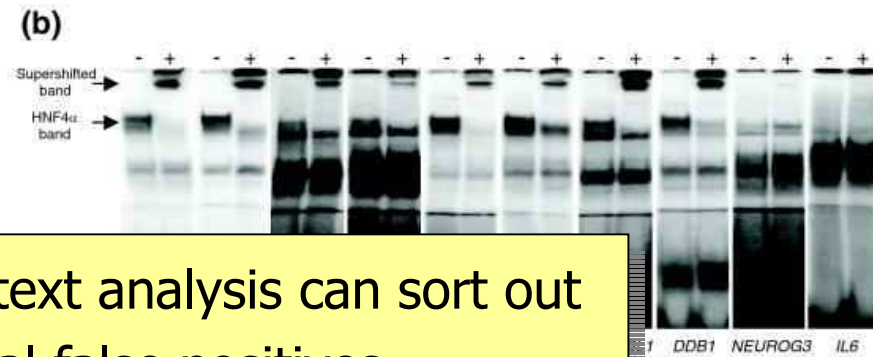
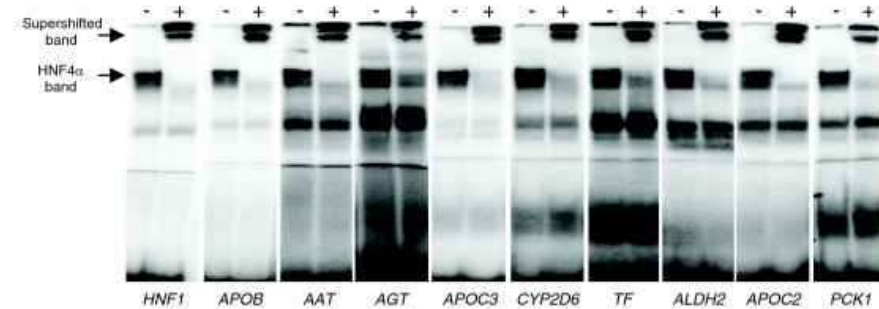
Composite module analyzer (CMA)

HNF-4 α sites of ChIP-chip fragments* revisited:

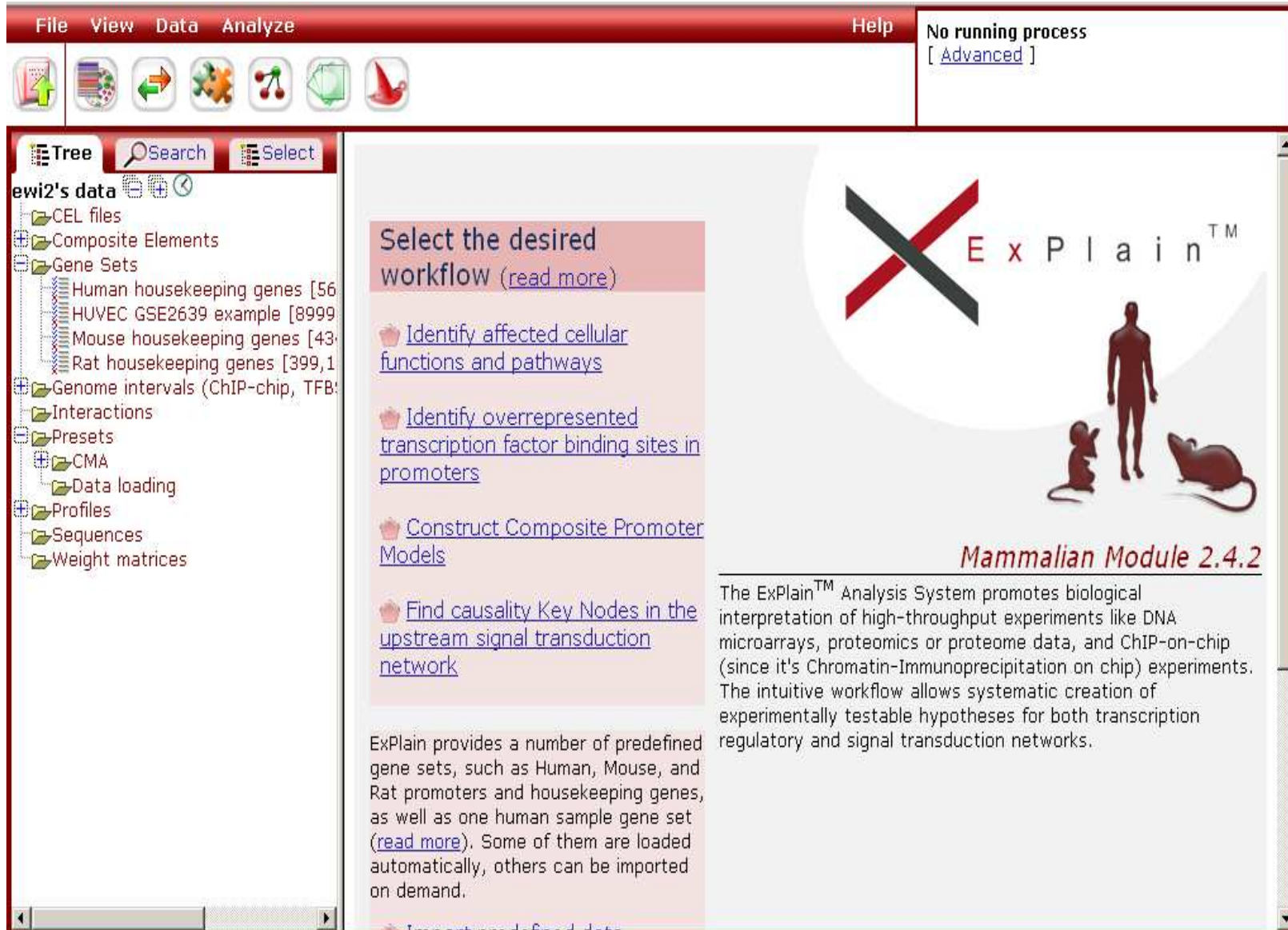


Computational context analysis can sort out experimental false positives.

(a) * Odom et al., *Science* 2004, **303**:1378-1381.



The Explain™ System



The screenshot shows the Explain™ System software interface. At the top is a menu bar with 'File', 'View', 'Data', 'Analyze', and 'Help'. Below the menu bar is a toolbar with several icons. The main window is divided into a left sidebar and a main content area. The sidebar, titled 'ewi2's data', contains a tree view with folders for 'CEL files', 'Composite Elements', 'Gene Sets', 'Genome intervals (ChIP-chip, TFB...', 'Interactions', 'Presets', 'CMA', 'Data loading', 'Profiles', 'Sequences', and 'Weight matrices'. The main content area features a large graphic with a red 'X' over the 'Explain™' logo and silhouettes of a human, a mouse, and a rat. Below the graphic, the text reads 'Mammalian Module 2.4.2'. The main content area also contains a list of workflow options: 'Select the desired workflow (read more)', 'Identify affected cellular functions and pathways', 'Identify overrepresented transcription factor binding sites in promoters', 'Construct Composite Promoter Models', and 'Find causality Key Nodes in the upstream signal transduction network'. A paragraph at the bottom explains that Explain provides predefined gene sets and workflows for biological interpretation of high-throughput experiments.

File View Data Analyze Help

No running process
[[Advanced](#)]

Tree Search Select

ewi2's data

- CEL files
- Composite Elements
- Gene Sets
 - Human housekeeping genes [56]
 - HUVEC GSE2639 example [8999]
 - Mouse housekeeping genes [43]
 - Rat housekeeping genes [399,1]
- Genome intervals (ChIP-chip, TFB...
- Interactions
- Presets
- CMA
- Data loading
- Profiles
- Sequences
- Weight matrices

Select the desired workflow ([read more](#))

- [Identify affected cellular functions and pathways](#)
- [Identify overrepresented transcription factor binding sites in promoters](#)
- [Construct Composite Promoter Models](#)
- [Find causality Key Nodes in the upstream signal transduction network](#)

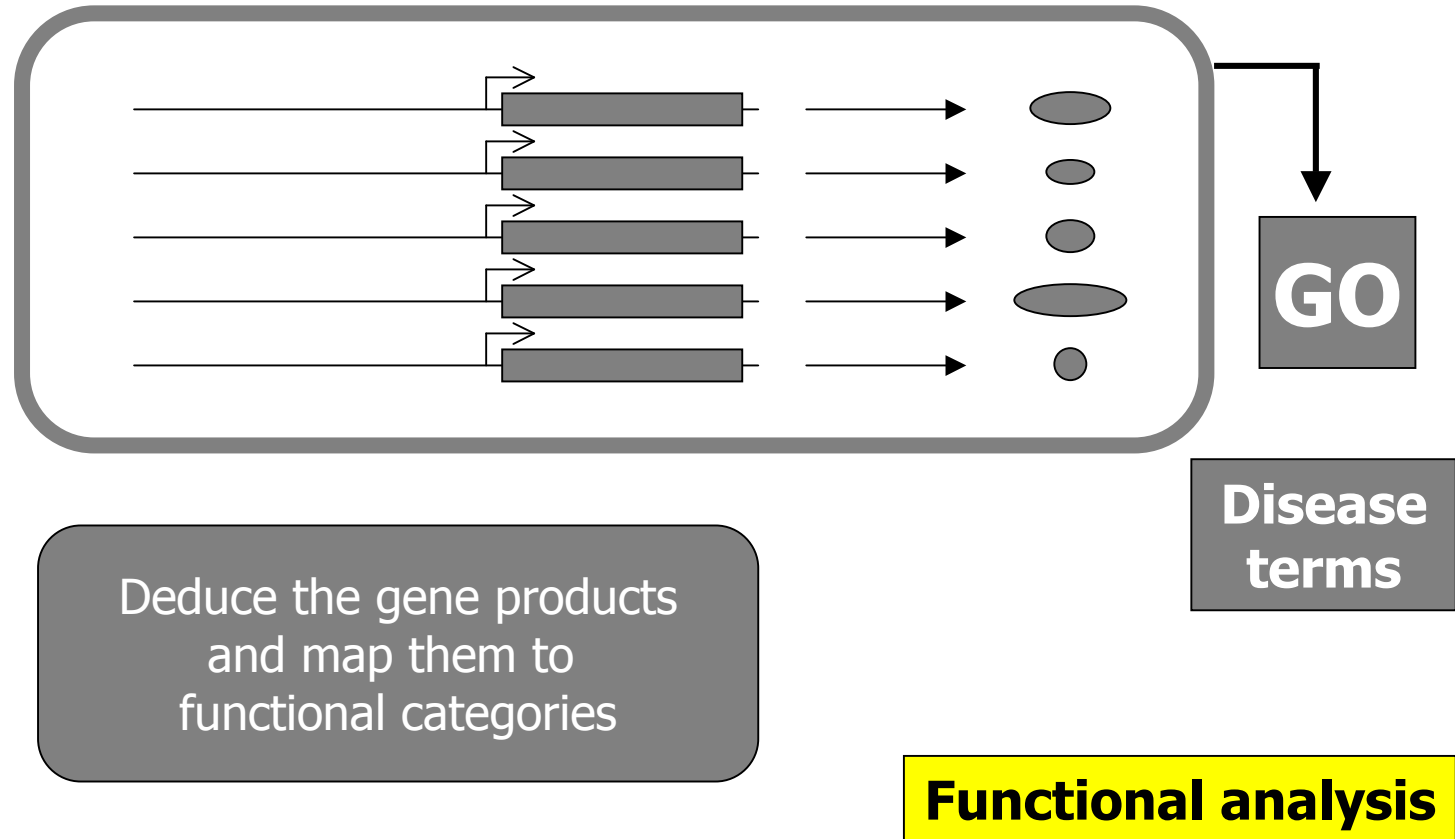
Explain provides a number of predefined gene sets, such as Human, Mouse, and Rat promoters and housekeeping genes, as well as one human sample gene set ([read more](#)). Some of them are loaded automatically, others can be imported on demand.

Explain™

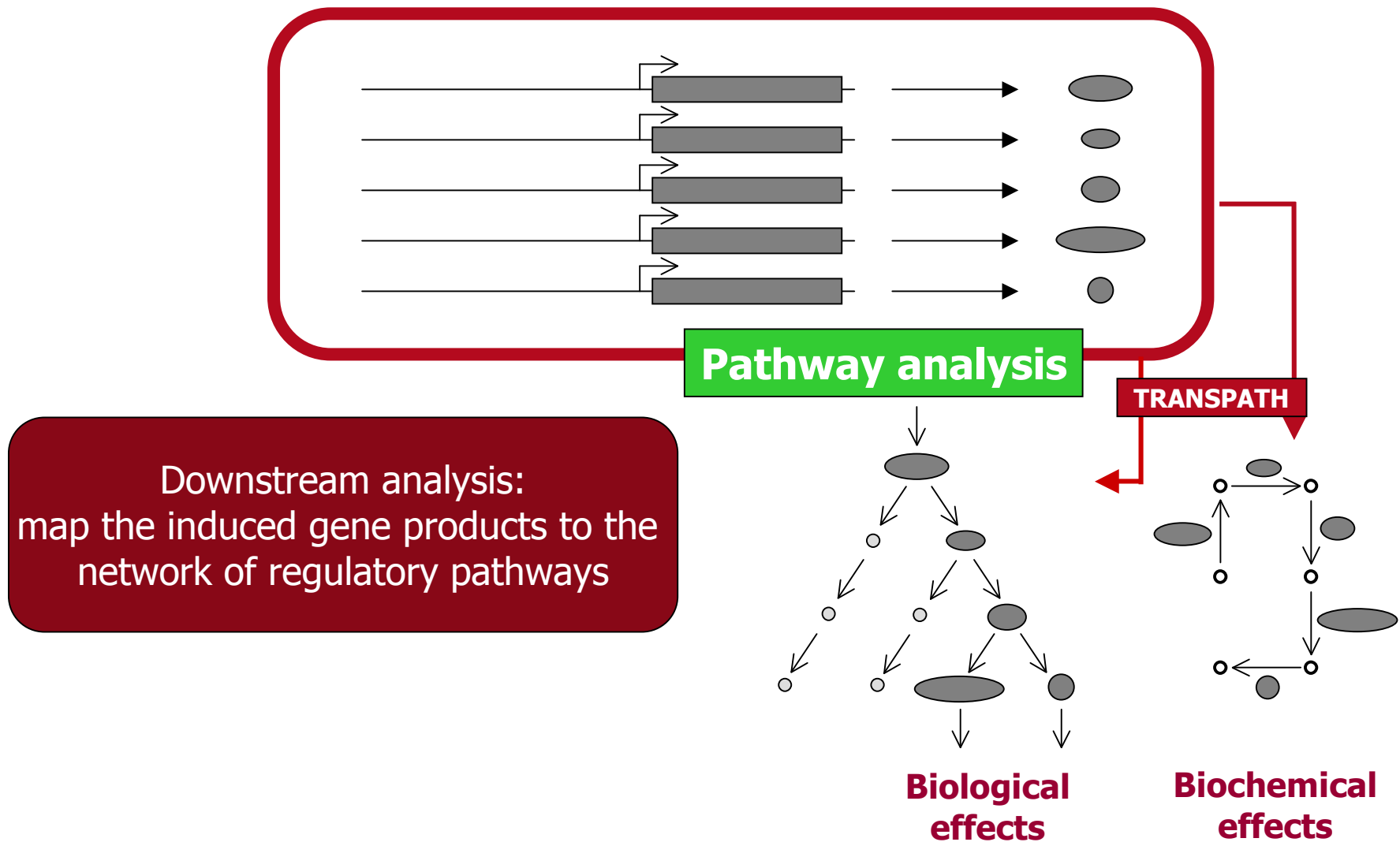
Mammalian Module 2.4.2

The Explain™ Analysis System promotes biological interpretation of high-throughput experiments like DNA microarrays, proteomics or proteome data, and ChIP-on-chip (since it's Chromatin-Immunoprecipitation on chip) experiments. The intuitive workflow allows systematic creation of experimentally testable hypotheses for both transcription regulatory and signal transduction networks.

The Explain™ System

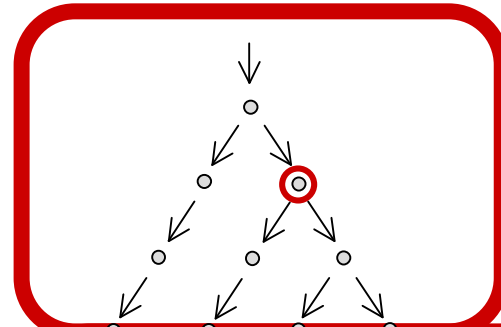


The Explain™ System



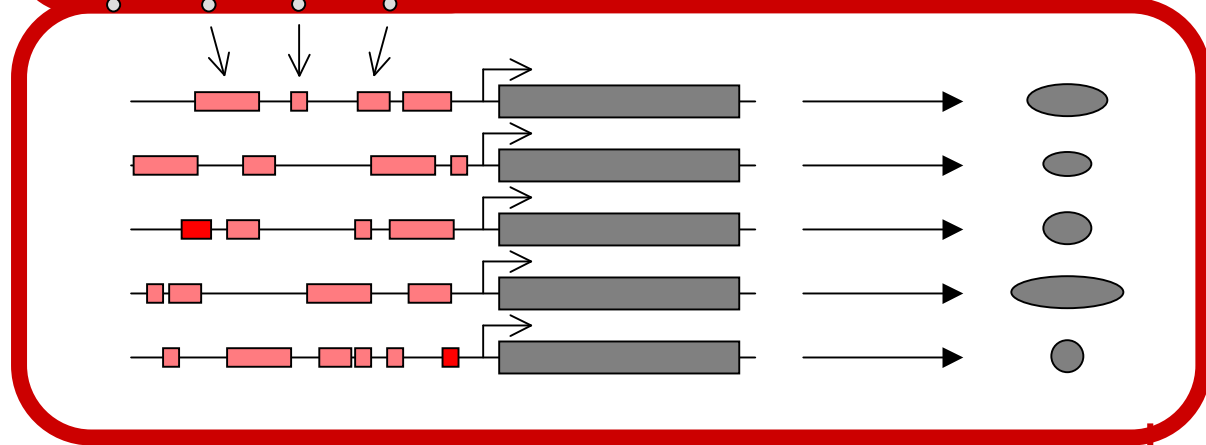
The Explain™ System

Pathway analysis



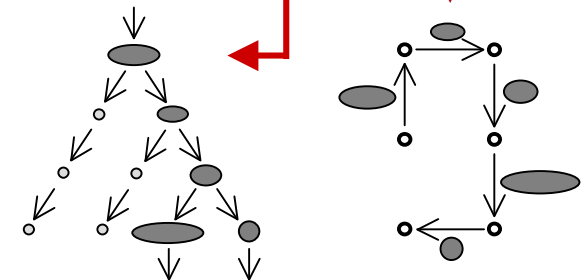
Identification of new targets

Promoter analysis



Reasoning of experimental findings: promoter analysis of induced genes connected to network mapping

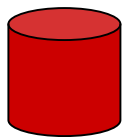
TRANSPATH



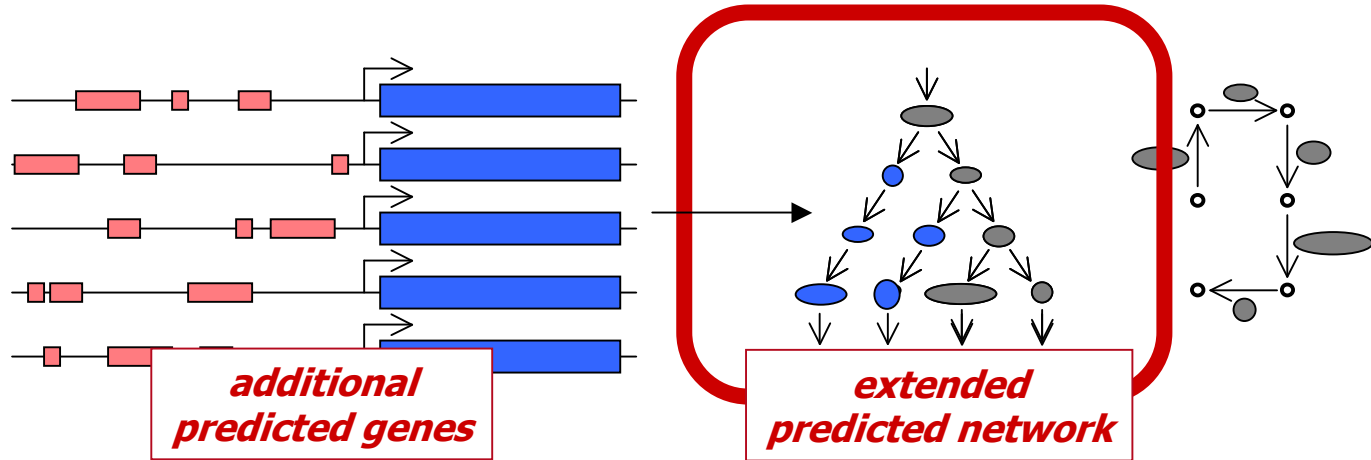
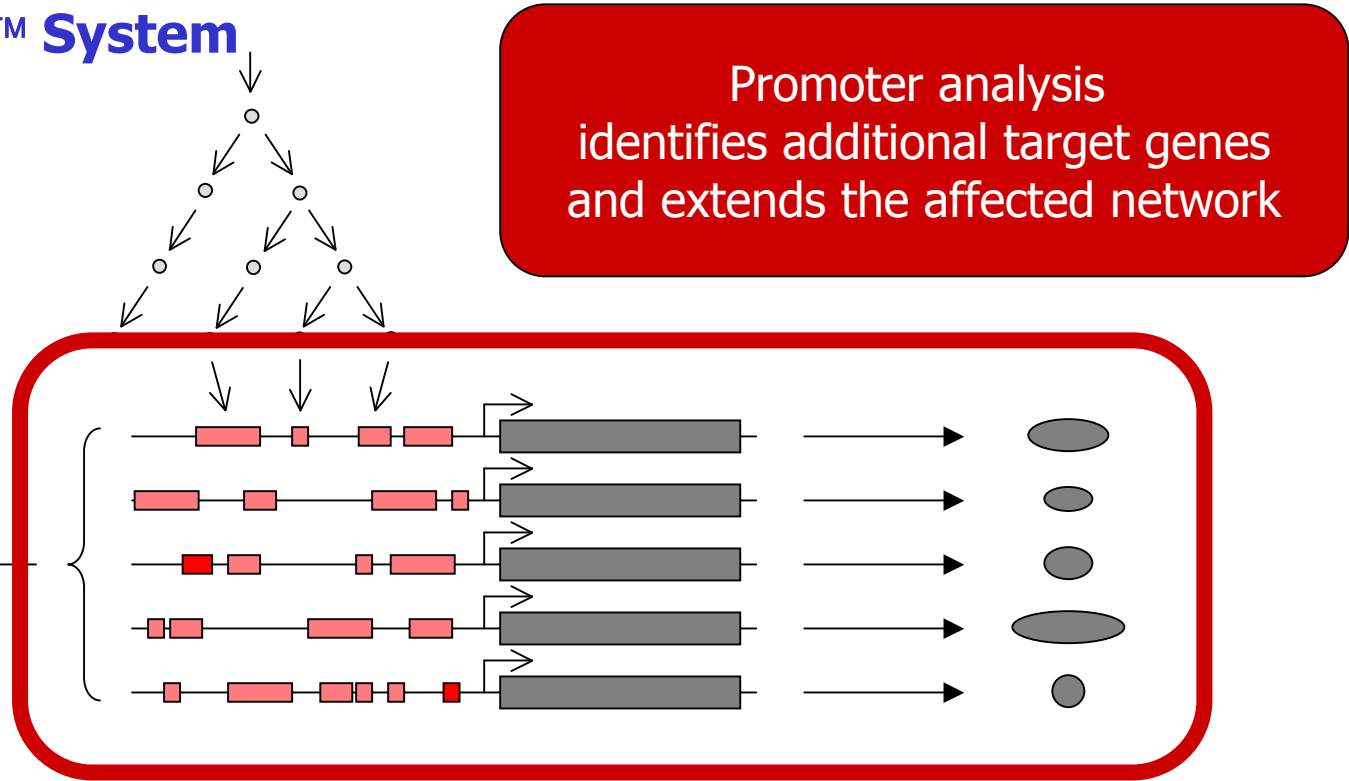
The Explain™ System

Promoter screening

promoter model



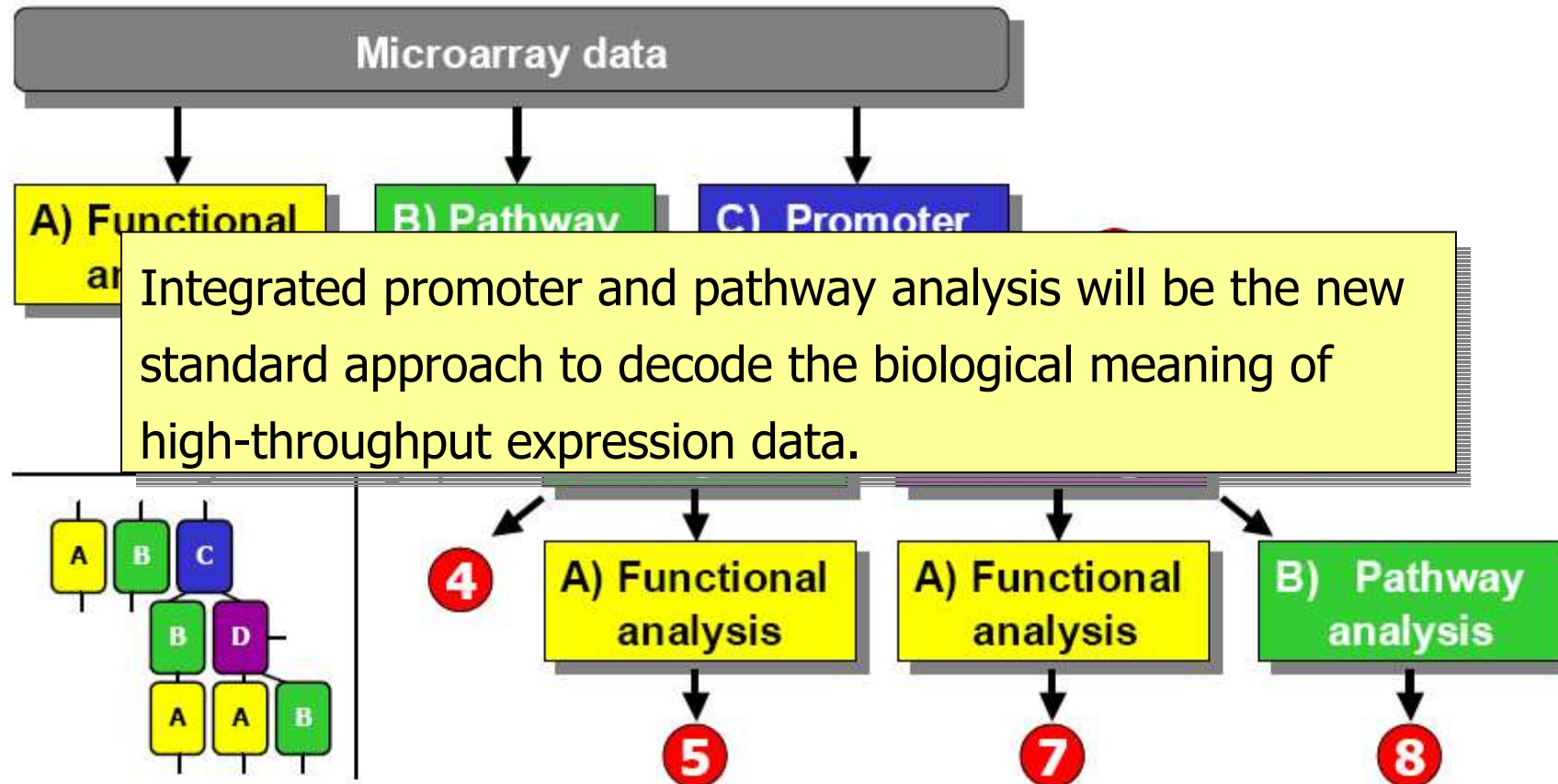
TRANSGENOME database



additional predicted genes

extended predicted network

The Explain™ System: workflow example



Goal #5:

Deciphering the Protein-DNA interaction code in order to predict the DNA-binding specificity of new transcription factors.

Protein-DNA interaction code



The first compilation: The FACTOR table

The zinc-finger structure:
The eukaryotic alternative
to the bacterial helix-turn-
helix motif?

<u>factor</u>	<u>source</u>	<u>gene</u>	<u>M. W. (kDa)/ finger str.</u>	<u>synonyms, equivalents</u>	<u>ref.</u>
60K protein	soybean	lectin	60		104
α 4 protein	human(HeLa)	α 0 (HSV) α 4 *	2x 170	ICP4,Vmw175	84, 176 84
Adf-1	Drosophila	adh			11
ADR1	yeast	ADH2	151 / 2 CH		13, 177
AP-1	human	collagenase IL-2 (?) metallothionein IIA polyoma virus (?) rat stromelysin (?) SV40 enhancer	47		21, 110, 178 21
AP-2	HeLa	BPV GH metallothionein IIA c-myc MHC class I H-2K ^b SV40 enhancer	50; 52	KBF1 (mouse) ?	18, 110 18 18 18 18 18, 110, 162
AP-3	HeLa	enhancer (SV40)			110, 162
B factor	Drosophila " "	actin 5C histone H3 histone H4		TFIID (HeLa) ?	3
CBF	sea urchin	histone H2B-1			64
CBP	mouse human rat rat	α -globin hsp70 ? LTR (HSV) tk (HSV)		CTF ?	51 79 125 125
CBF	sea urchin	histone H2B-1			64
CCAAT-binding factor	human sea urchin mouse	hsp70 histone H2B-1 α 2(I) collagen		CTF	78 64 20
CCAAT box bind.protein	mouse, rat, human(HeLa)	α globin	64 ?	CCAAT-bind.f.; CTF;CBP ?	51

Wingender E. Compilation of
transcription regulating proteins.
Nucleic Acids Res. 1988 Mar
25;16(5):1879-902.
PMID: 3282223

The first compilation: Zinc finger alignments

<u>gene, gene product^(*)</u> (ref.)	<u>position</u>	<u>finger sequences</u>
TFIIIA (210, 211)	1	H G E K A L
	7	P V V Y K R Y I C S F A D C G A A Y N K N W K L Q A H L C K H
	38	T G E K P F P C K E E G C E K G F T S L H H L T R H S L T H
	68	T G E K N F T C D S D G C D L R F T T K A N M K K H F N R F H
	99	N I K I C V Y V C H F E N C G K A F K K H N Q L K V H Q F S H
	130	T Q Q L P Y E C P H E G C D K R F S L P S R L K R H E K V H
	160	A G Y P C K K D D S C S F V G K T W T L Y L K H V A E C H
	189	Q D L A V C D V C N R K F R H K D Y L R D H Q K T H
	215	E K E R T V Y L C P R D G C D R S Y T T A F N L R S H I Q S F H
	247	E E Q R P F V C E H A G C G K C F A M K K S L E R H S V V H
Xfin (212)	103	S A K K S H I C S H T G K L F S C T A A V V R H Q R M H
	131	Q L Q K S H H C P H C K K S F V Q R S D F I K H Q R T H
	159	T G E R P Y Q C V E C Q K K F T E R S A L V N H Q R T H
	187	T G E R P Y T C L D C Q K T F N Q R S A L T K H R R T H
	215	T G E R P Y R C S V C S K S F I Q N S D L V K H L R T H
	243	T G E K P Y E C P L C V K R F A E S S A L M K H K R T H
	271	S T H R P F R C S E C S R S F T H N S D L T A H M R K H
	299	T E F R ...
...	320	N V A S S P Y S C S K C R K T F K R W K S F L N H Q Q T H
	349	S R E K P Y L C S H C N K G F I Q N S D L V K H F R T H
	377	T G E R P Y Q C A E C H K G F I Q K S D L V K H L R T H
	405	T G E K P F K C S H C D K K F T E R S A L A K H Q R T H
	433	T G E K P Y K C S D C G K E F T Q R S N L I L H Q R I H

Wingender E. Compilation of transcription regulating proteins. Nucleic Acids Res. 1988 Mar 25;16(5):1879-902. PMID: 3282223

The helix-turn-helix motif as paradigm of prokaryotic DNA-binding domains

Sauer RT, Yocum RR, Doolittle RF, Lewis M, Pabo CO (1982)
Nature 298, 447-451:

Homology among DNA-binding proteins suggests use of a conserved super-secondary structure.

The amino acid sequences of the repressor and cro proteins of phages lambda, 434 and P22 are homologous, especially in a region in which repressor and lambda cro have a similar **alpha-helix-turn-alpha-helix secondary structure**. Model-building studies indicate that this structure is important in DNA binding, and we suggest it may be a common feature of many DNA-binding proteins.

The helix-turn-helix motif as paradigm of prokaryotic DNA-binding domains

PDB entry
1CGP



Schultz, S.C., Shields, G.C., Steitz, T.A. (1991) Crystal structure of a CAP-DNA complex: the DNA is bent by 90 degrees. *Science* **253**: 1001-1007

dating back to:

McKay, D.B., Steitz, T.A. (1981) Structure of catabolite gene activator protein at 2.9 Å resolution suggests binding to left-handed B-DNA. *Nature* **290**:744-749.

The zinc finger motif as general paradigm for eukaryotic DNA-binding domains?

PDB entry
1TF3



Foster, M.P., Wuttke, D.S., Radhakrishnan, I., Case, D.A., Gottesfeld, J.M., Wright, P.E. (1997) Domain packing and dynamics in the DNA complex of the N-terminal zinc fingers of TFIID. *Nat.Struct.Biol.* **4**: 605-608

The zinc finger motif as general paradigm for eukaryotic DNA-binding domains?

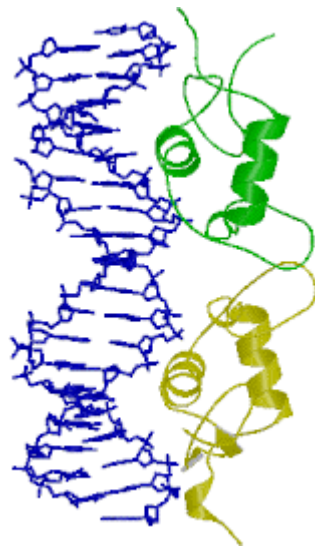
GR zinc finger

Similarity with TFIIIA zinc finger structure suggested for the first time.

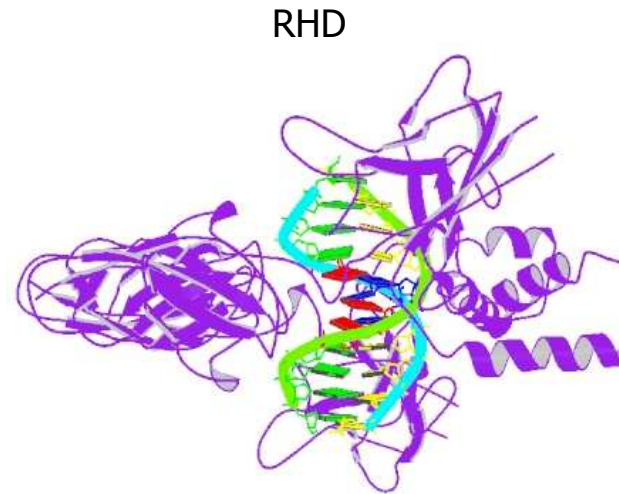
hGR	S	N	Q	Y	S	S	P	S	M	R	P	D	V	S	S	P	P	S	S	S	S	T	A	T	T	G	P	P	P	K	L	-	420	
v-erbA	E	D	T	R	W	L	D	G	K	H	K	R	K	S	S	Q	C	L	V	E	S	S	M	S	G	Y	I	P	S	C	L	D	32	
hGR	-	-	-	-	C	L	V	C	S	D	E	A	S	G	C	H	Y	G	V	L	I	C	Q	S	C	K	V	F	F	K	R	-	447	
v-erbA	K	D	E	Q	C	V	V	C	G	D	K	A	T	Q	Y	H	R	C	I	T	C	E	G	C	K	S	F	F	R	R	T	64		
hGR	-	A	V	E	Q	Q	H	N	Y	L	C	A	G	R	N	D	C	I	I	D	R	I	R	R	K	I	C	P	A	C	R	Y	478	
v-erbA	I	Q	K	N	L	H	P	T	Y	S	C	T	Y	D	G	C	C	V	I	D	K	I	T	R	N	C	Q	L	C	R	F	96		
hGR	R	K	C	L	Q	A	G	M	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	N	L	E	A	490		
v-erbA	K	K	C	I	S	V	G	M	A	M	D	L	V	L	D	D	S	K	R	V	A	K	R	K	L	I	E	E	N	R	E	R	128	
hGR	R	K	T	K	K	K	I	K	Q	I	Q	Q	A	T	T	G	V	S	Q	E	T	S	E	N	P	G	N	K	T	I	V	P	522	
v-erbA	R	R	K	E	E	M	I	K	S	L	Q	H	R	P	S	P	S	A	E	E	W	E	L	I	H	V	V	T	E	A	H	R	160	
hGR	A	T	L	P	Q	-	-	-	L	T	P	T	L	V	S	L	L	E	V	I	E	P	E	V	L	Y	A	G	Y	D	S	S	551	
v-erbA	S	T	N	A	Q	G	S	H	W	K	Q	R	R	K	F	L	L	E	D	I	G	Q	S	P	M	A	S	M	L	D	G	D	192	
hGR	V	P	D	S	T	W	R	I	M	T	T	L	N	M	L	G	G	R	Q	V	I	A	A	V	K	W	A	K	A	I	P	G	583	
v-erbA	K	V	D	L	E	A	P	S	E	F	T	-	-	K	I	I	T	P	A	I	T	R	V	V	D	F	A	K	N	L	P	M	222	
hGR	F	R	N	L	H	L	D	D	Q	M	T	L	L	Q	Y	S	W	M	F	L	M	A	F	A	L	G	W	R	S	Y	R	Q	615	
v-erbA	F	S	E	L	P	C	E	D	Q	I	I	L	L	K	Q	C	C	M	E	I	M	S	L	R	A	A	V	R	Y	D	P	E	254	
hGR	S	S	A	N	L	L	C	P	A	P	D	L	I	N	E	Q	R	M	T	L	P	C	M	Y	D	Q	C	K	H	M	L	647		
v-erbA	S	E	T	L	T	L	S	O	E	M	A	V	K	R	R	Q	L	K	N	G	O	L	G	V	V	S	D	A	I	F	D	L	286	
hGR	Y	V	S	S	E	L	H	R	L	Q	V	S	Y	E	E	Y	L	C	M	K	T	L	L	L	S	S	Y	P	K	D	Q	679		
v-erbA	O	K	S	L	S	A	F	N	L	D	D	T	R	V	A	L	L	-	-	-	Q	A	V	L	L	M	S	S	D	R	T	Q	315	
hGR	L	K	S	Q	E	L	F	D	E	I	R	M	T	Y	I	K	E	L	G	K	A	I	V	K	R	E	G	M	S	S	Q	N	711	
v-erbA	L	I	C	V	D	K	I	E	K	C	Q	E	S	Y	L	L	A	F	E	H	Y	I	N	Y	R	K	H	N	I	P	H	F	347	
hGR	W	Q	R	F	Y	Q	L	T	K	L	L	D	S	M	H	R	V	V	E	N	L	L	N	Y	C	F	Q	T	F	L	D	K	743	
v-erbA	W	-	-	-	-	-	-	-	-	S	K	L	L	M	K	V	A	D	L	R	M	I	Q	A	Y	H	A	S	R	F	L	H	M	372
hGR	T	M	S	I	E	F	P	E	M	L	A	E	I	I	T	N	Q	I	P	K	Y	S	N	G	N	I	K	K	L	L	P	H	775	
v-erbA	K	V	E	C	P	T	E	L	P	P	R	R	C	R	A	L	Q	I	L	G	S	I	L	P	P	V	-	-	-	-	-	-	398	

Weinberger C, Hollenberg SM, Ong ES, H hGR Q K 777
 Identification of human glucocorticoid receptor complementary DNA clones by epitope selection. Science. 1985 May 10;228(4700):740-2.
 PMID: 2581314

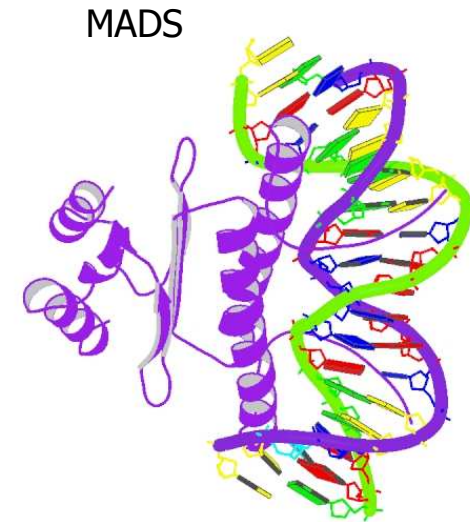
Many more structural classes of eukaryotic DNA-binding domains found later on



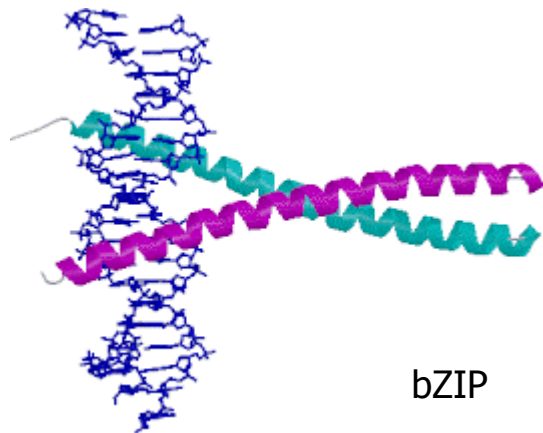
HTH



RHD



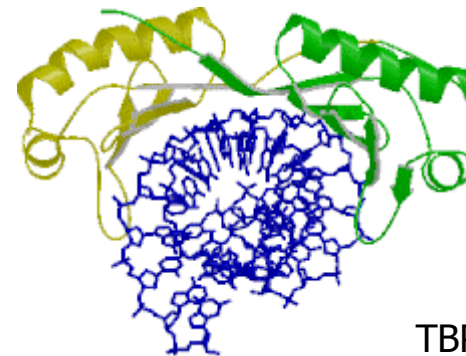
MADS



bZIP

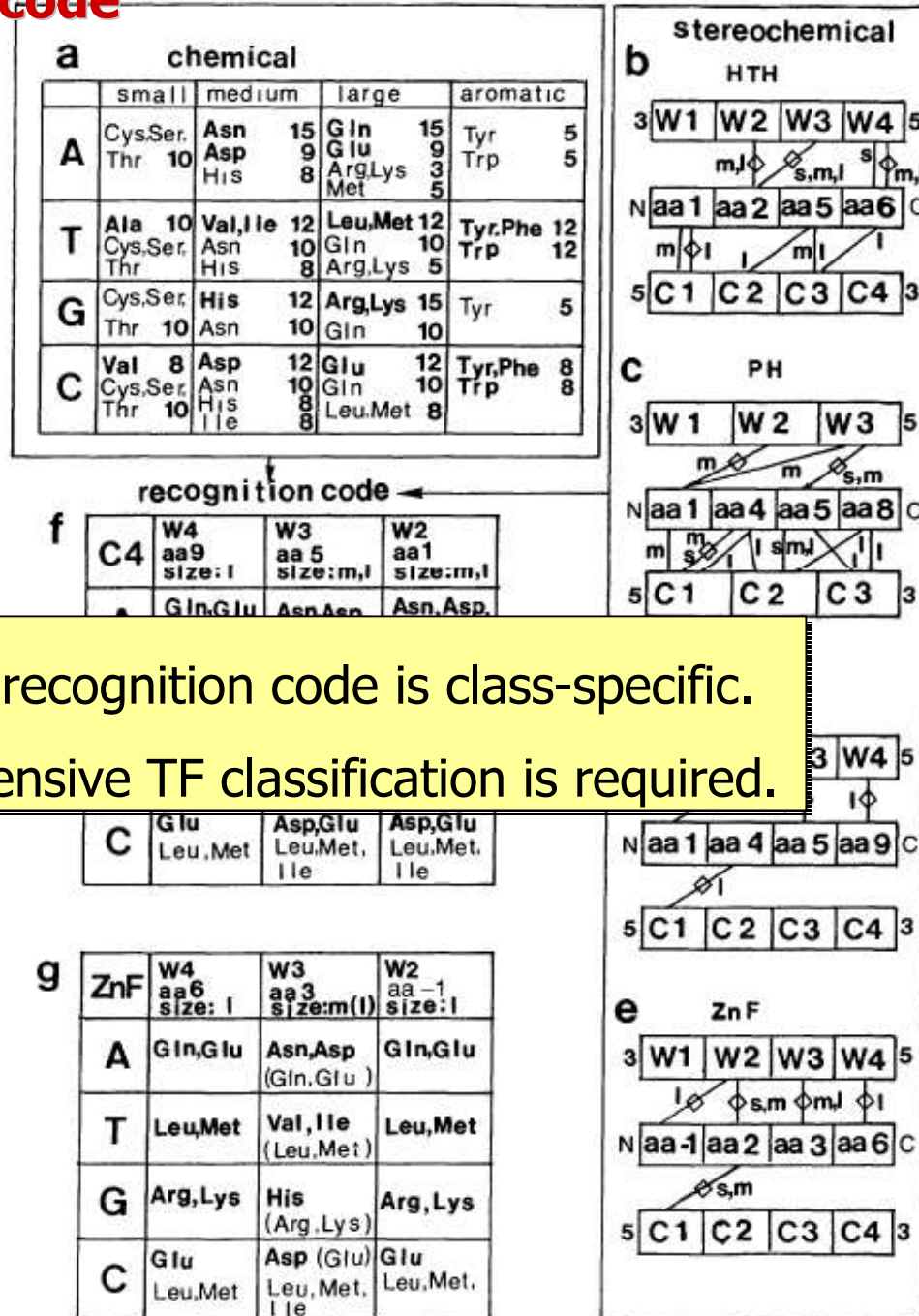


TBP



Protein-DNA interaction code

Protein-DNA recognition code?



The protein-DNA recognition code is class-specific.
 Thus, a comprehensive TF classification is required.

Suzuki & Yagi, PNAS 91:12357-12361, 1994

Transcription factor classification

Table 1. Classification scheme for eukaryotic transcription factors

Level	Pattern of classification no.	group designation	description	example
1	N	superclass	general topology of DBD	zinc-coordinating domains
2	N.N	class	structural blueprint of the DBD	zinc-finger nuclear receptors
3	N.N.N	family	functional criteria such as protein-DNA-complex formation (DNA-binding specificity, multimerization behaviour) or biological effect	T ₃ R/RAR (in contrast to steroid hormone receptors)
4	N.N.N.N	subfamily	mainly according to sequence similarity of the DBDs	RAR (retinoic acid receptor)
5	N.N.N.N.N	genus	according to factor gene	RAR- α , RAR- β
6	N.N.N.N.N.N	factor "species"	initiation/ splice/ processing variants	RAR- α 1, RAR- α 2

Transcription factors

- 1 Basic Domains
 - 1.1 Leucine zipper factors (bZIP)
 - 1.2 Helix-loop-helix factors (bHLH)
 - 1.3 Helix-loop-helix / leucine zipper factors (bHLH-ZIP)
 - 1.6 bHSH
- 2 Zinc-coordinating DNA-binding domains
 - 2.1 Cys4 zinc finger of nuclear receptor type
 - 2.2 diverse Cys4 zinc fingers
 - 2.3 Cys2His2 zinc finger domain
 - 2.4 Cys6 cysteine-zinc cluster
 - 2.5 Zinc fingers of alternating composition

DNA-binding domains of eukaryotic TFs can be classified in (at least) 4 superclasses and 30 classes.

- 4 beta-Scaffold Factors with Minor Groove Contacts
 - 4.1 RHR (Rel homology region)
 - 4.2 STAT

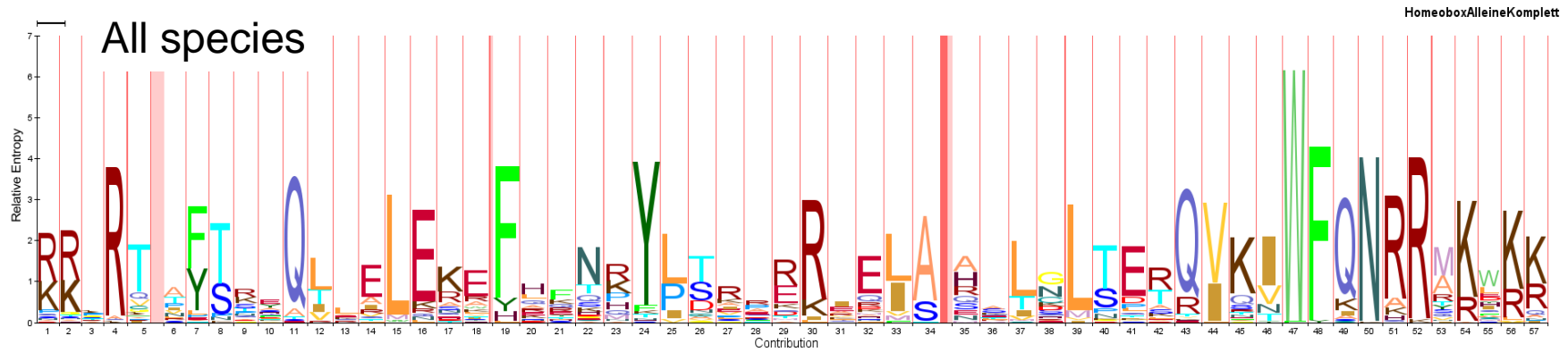
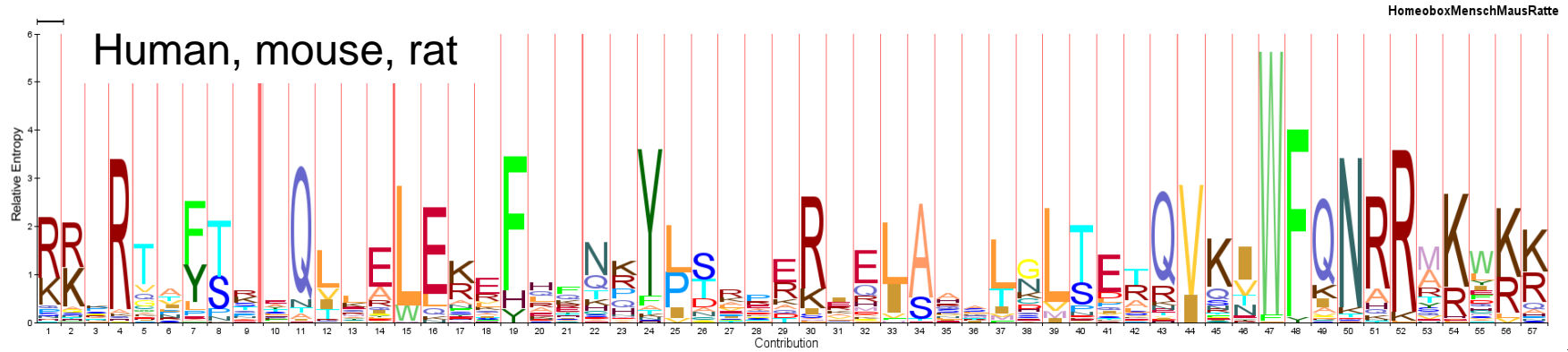
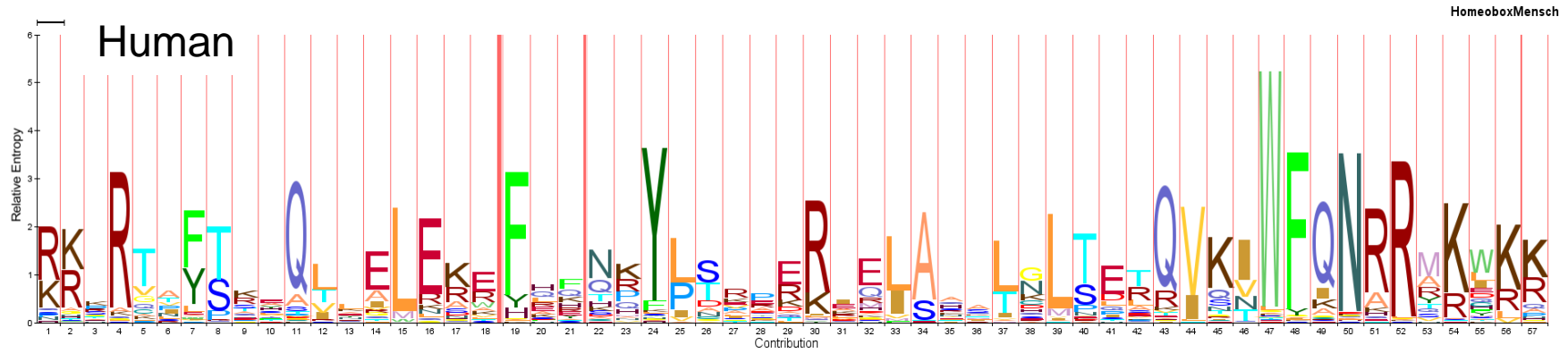
More recent data, however, suggest that there may be up to 13 superclasses (mostly spreading out from previous superclass 4).

- 0 Other Transcription Factors
 - 0.1 Copper fist proteins
 - 0.2 HMGI(Y)
 - 0.3 Pocket domain
 - 0.4 E1A-like factors
 - 0.5 AP2/EREBP-related factors

Transcription factor classification

- 1 [Basic Domains](#)
- 2 [Zinc-coordinating DNA-binding domains](#)
 - 2.1 [Cys4 zinc finger of nuclear receptor type](#)
 - 2.1.1 Steroid hormone receptors (NR3)
 - 2.1.1.1 Corticoid receptors (NR3C)
 - 2.1.1.1.1 **GR (NR3C1):** [GR](#) [GR](#) [GR](#) [GR](#) [GR](#) [GR](#) [GR](#)
 - 2.1.1.1.1.1 **GR-alpha :** [GR-alpha](#)
 - 2.1.1.1.1.2 **GR-beta:** [GR-beta](#)
 - 2.1.1.1.2 **MR (NR3C2):** [MR](#) [MR](#)
 - 2.1.1.2 Progesterone receptor (NR3C)
 - 2.1.1.3 Androgen receptor (NR3C)
 - 2.1.1.4 Estrogen receptor (NR3A)
 - 2.1.1.5 Estrogen related receptor (NR3B)
 - 2.1.2 Thyroid hormone receptor-like factors (NR0, NR1, NR2, NR4, NR5)
 - 2.2 [diverse Cys4 zinc fingers](#)
 - 2.3 [Cys2His2 zinc finger domain](#)
 - 2.4 [Cys6 cysteine-zinc cluster](#)
 - 2.5 Zinc fingers of alternating composition
- 3 [Helix-turn-helix](#)
- 4 [beta-Scaffold Factors with Minor Groove Contacts](#)
- 0 Other Transcription Factors

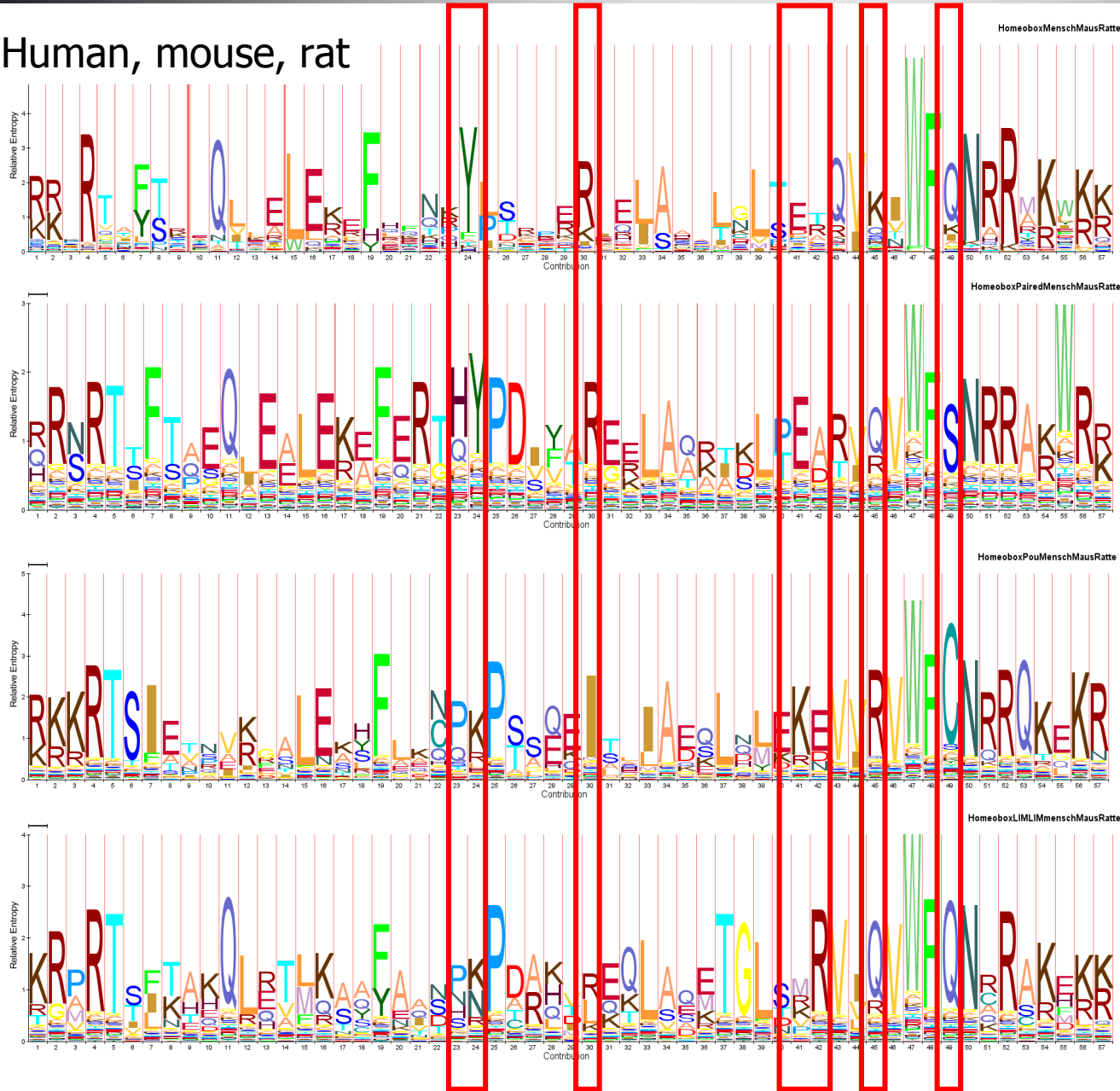
Homeo domains



Protein-DNA interaction code



Human, mouse, rat



Homeo domain only

Homeo domain in hom-paired proteins

Homeo domain in hom-POU proteins

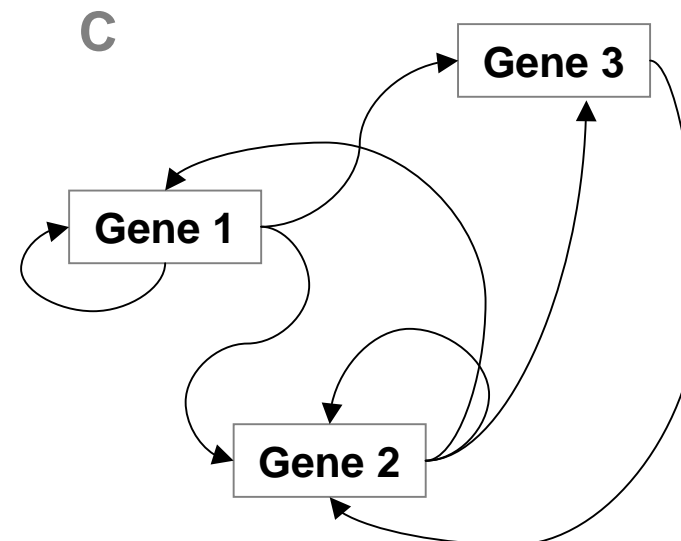
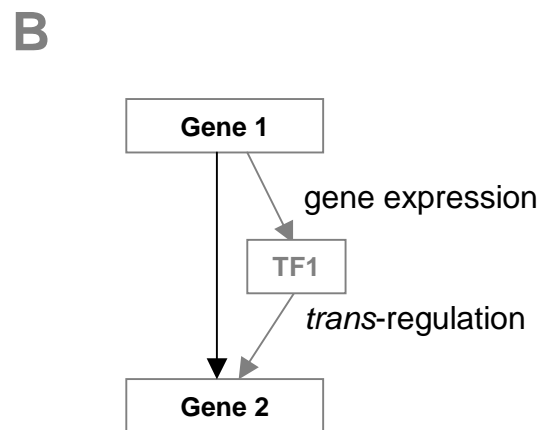
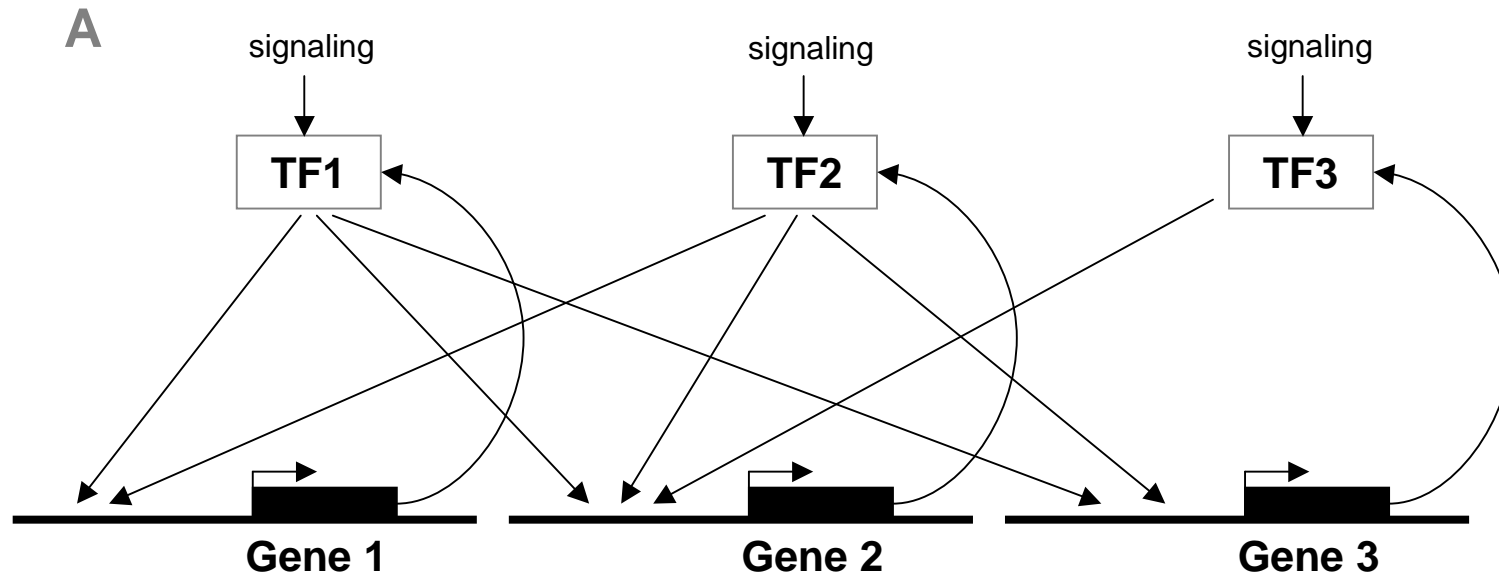
Homeo domain in hom-lim proteins

In a next step, it will be investigated whether whole DBDs or just a few residues are required to reveal mutual dependencies between DBDs and their target DNA sequences.

Goal #6:

Re-engineering and analyzing the complete transcription network of a cell / an organism.

Definition

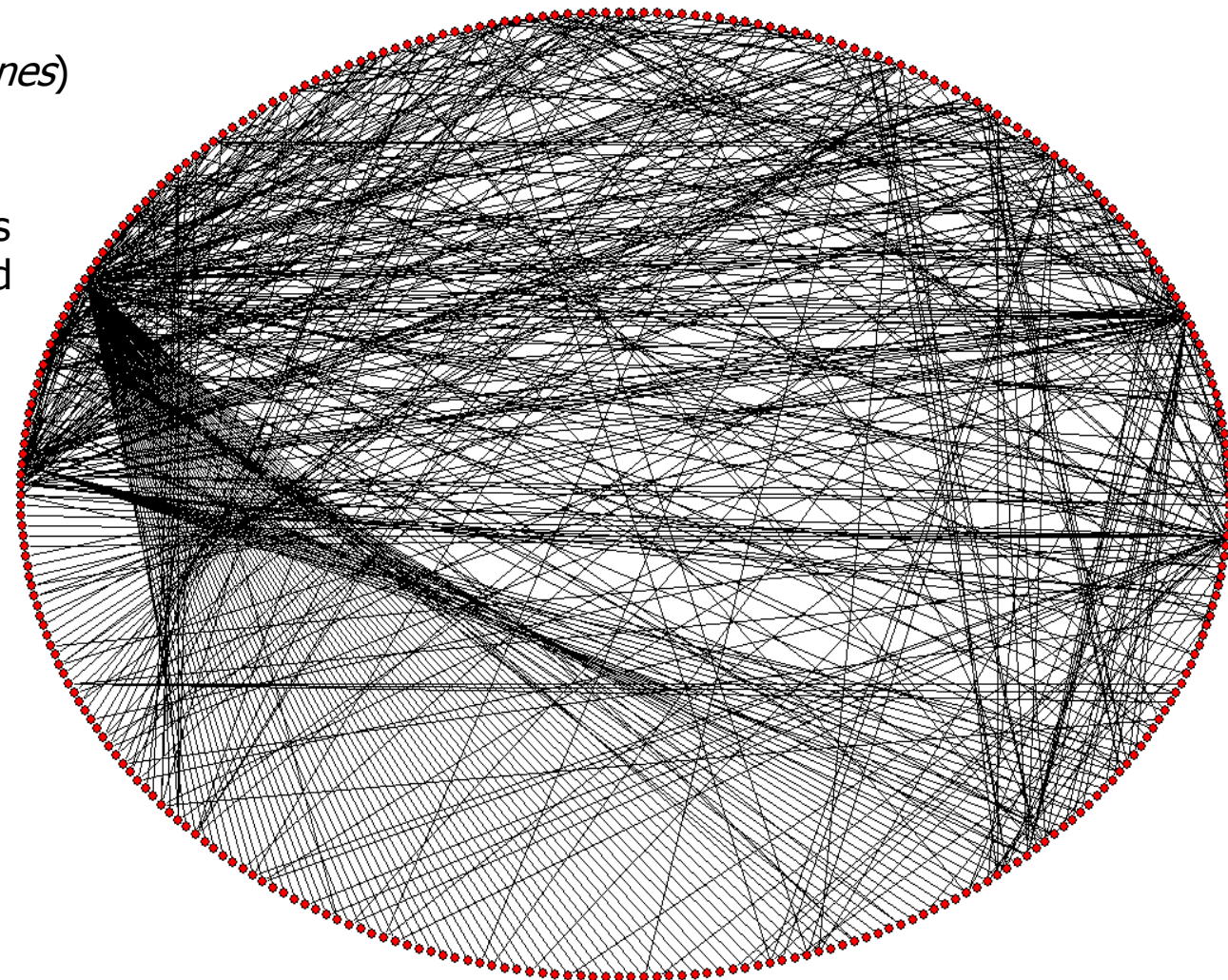
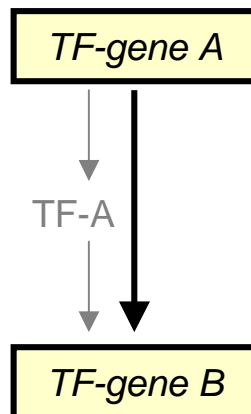


Mammalian network of transcription factor genes (TFG_RN)

298 vertices (*TF-genes*)

638 edges

Each edge combines
gene expression and
trans-regulation



Disconnectivity network analysis

Disconnectivity: Effect of removing a vertex

$$Dis(v) = \frac{N_0 - N_{-v}}{N_0} = 1 - \frac{N_{-v}}{N_0}$$

N_0 - number of all-pairs shortest paths

N_{-v} - number of all-pairs shortest paths after removing vertex v from the graph.

$$Dis(v) = \frac{\sum_{s \neq v \in V} \delta_{sv} + \sum_{s \neq v \neq t \in V} \delta_{st}(v) + \sum_{t \neq v \in V} \delta_{vt}}{\sum_{s \neq t \in V} \delta_{st}}$$

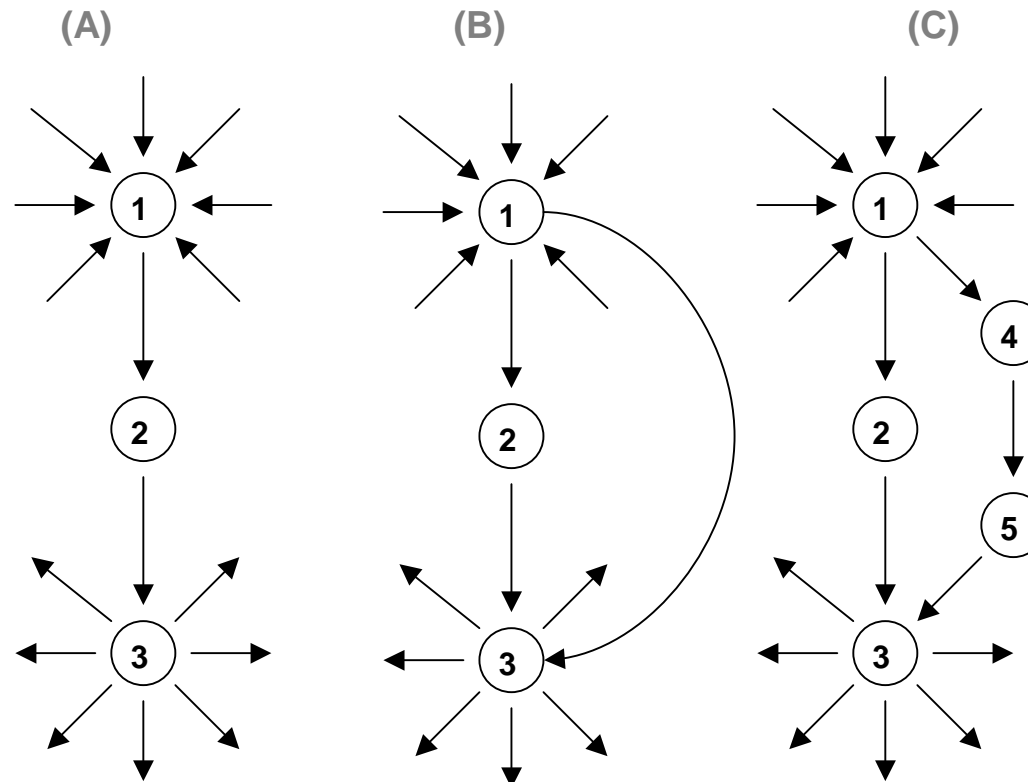


Fraction of all-pairs connections for which vertex v is crucial: there is no other parallel path(s) between vertices of such a pair and these vertices will not be connected anymore if vertex v is deleted.

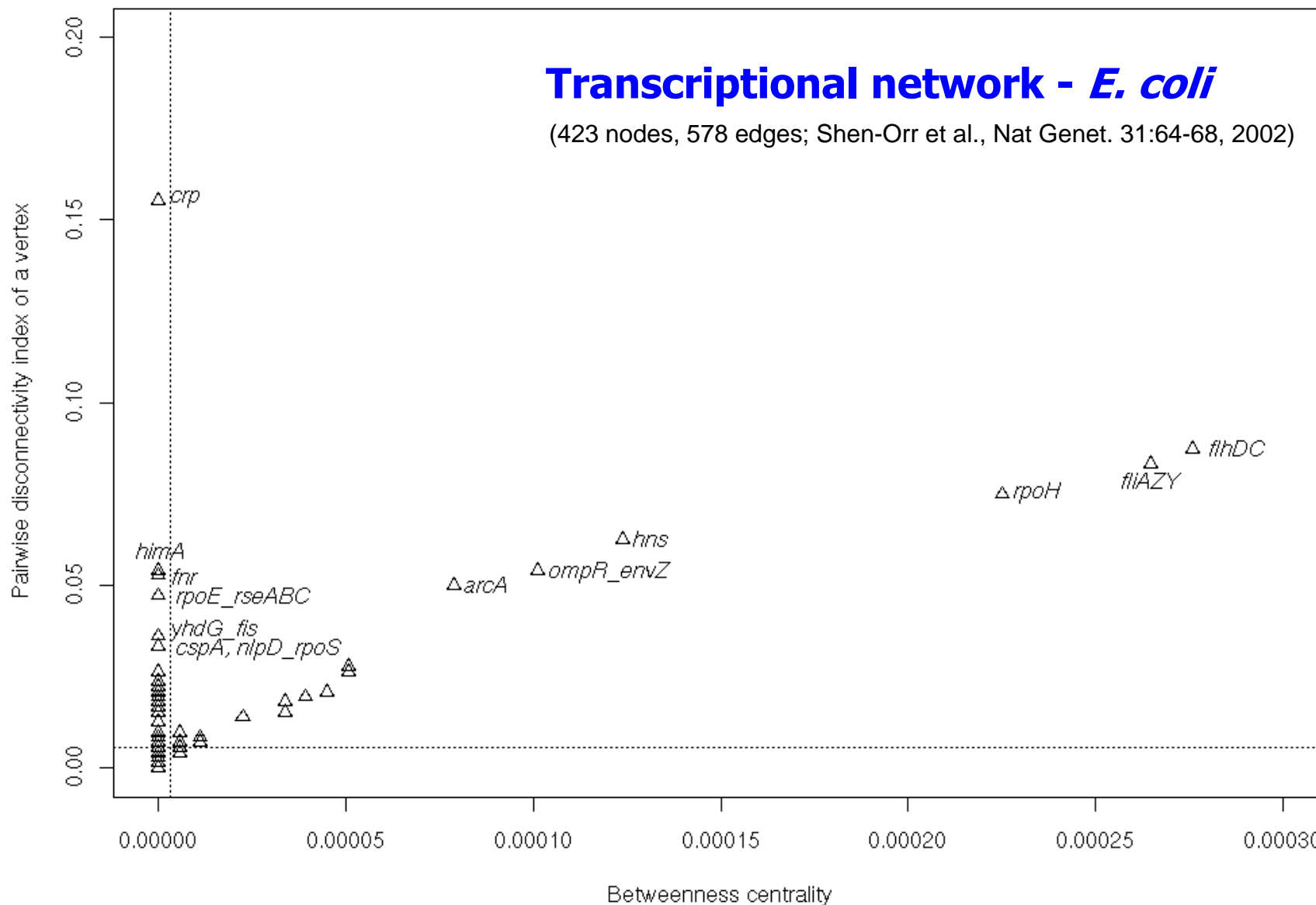
The betweenness centrality and the disconnectivity index D_i display different aspects of individual vertex topology in a whole network.

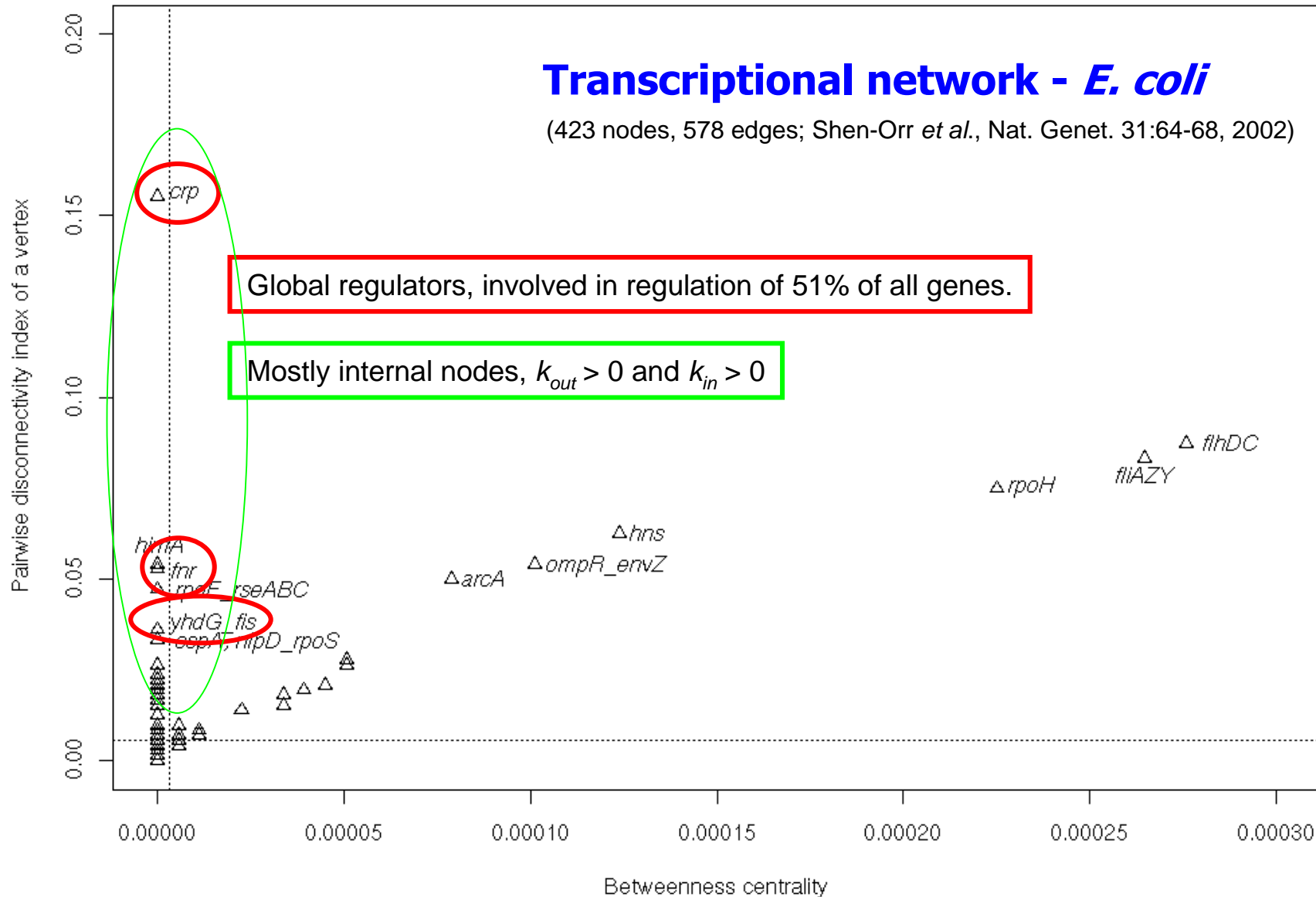
Disconnectivity network analysis

Disconnectivity sensitive towards bypasses

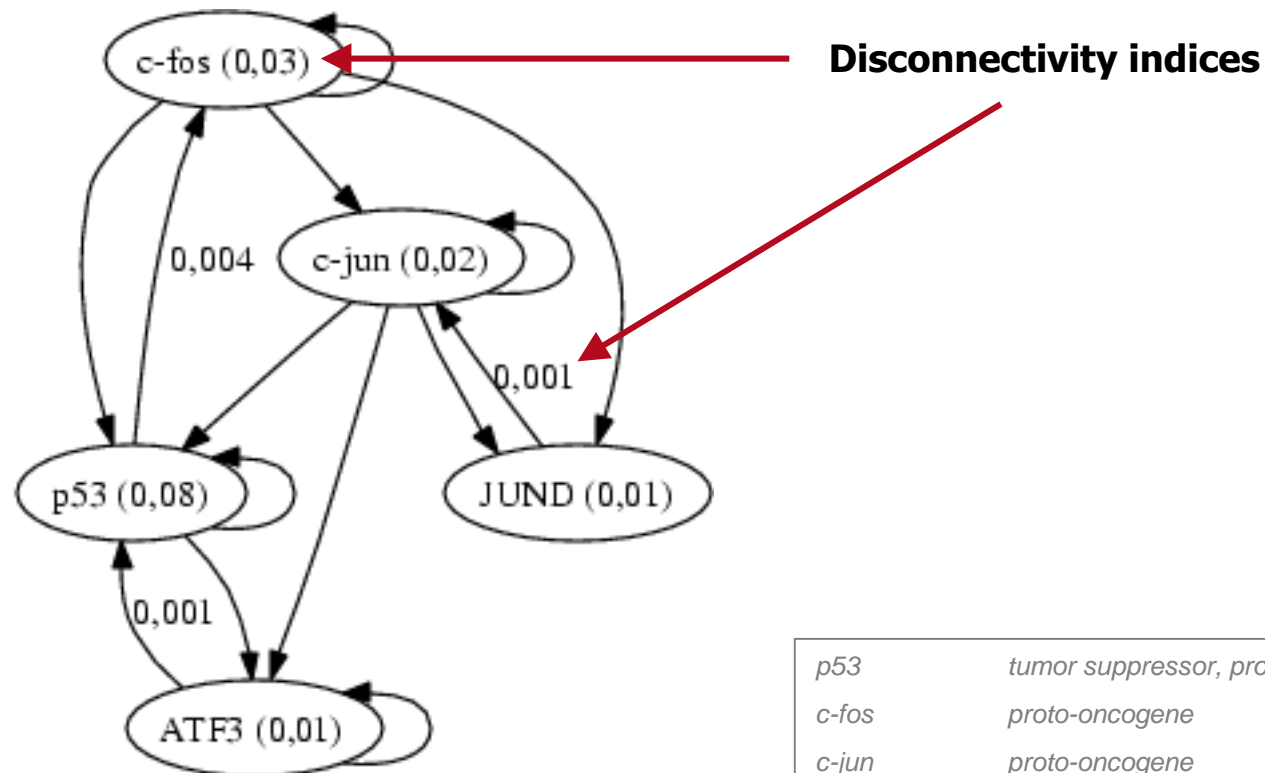


$B(v_2)$:	high	low	high
$Dis(v_2)$:	high	low	low





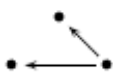
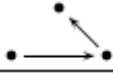
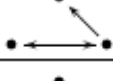
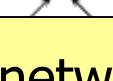
'p53' strongly connected sub-graph in TFG_RN



p53	tumor suppressor, pro-apoptotic
c-fos	proto-oncogene
c-jun	proto-oncogene
JUND	proto-oncogene
ATF3	tumor suppressor, pro-apoptotic




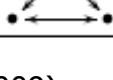
Transcriptional network

Pattern / motif characterization by Disconnectivity

Pattern	ID	<i>E. coli</i>			<i>S. cerevisiae</i>			Mammals		
		Freq	Z-Score	\overline{Dis}	Freq	Z-Score	\overline{Dis}	Freq	Z-Score	\overline{Dis}
	6	4777	11.23	0.0039	11892	14.54	0.0018	1916	-0.39	0.0023
	12	160	-11.21	0.0189	295	-14.25	0.0135	1068	-1.67	0.011
	14	-	-	-	18	-1.30	0.0063	73	-10.47	0.0079
										

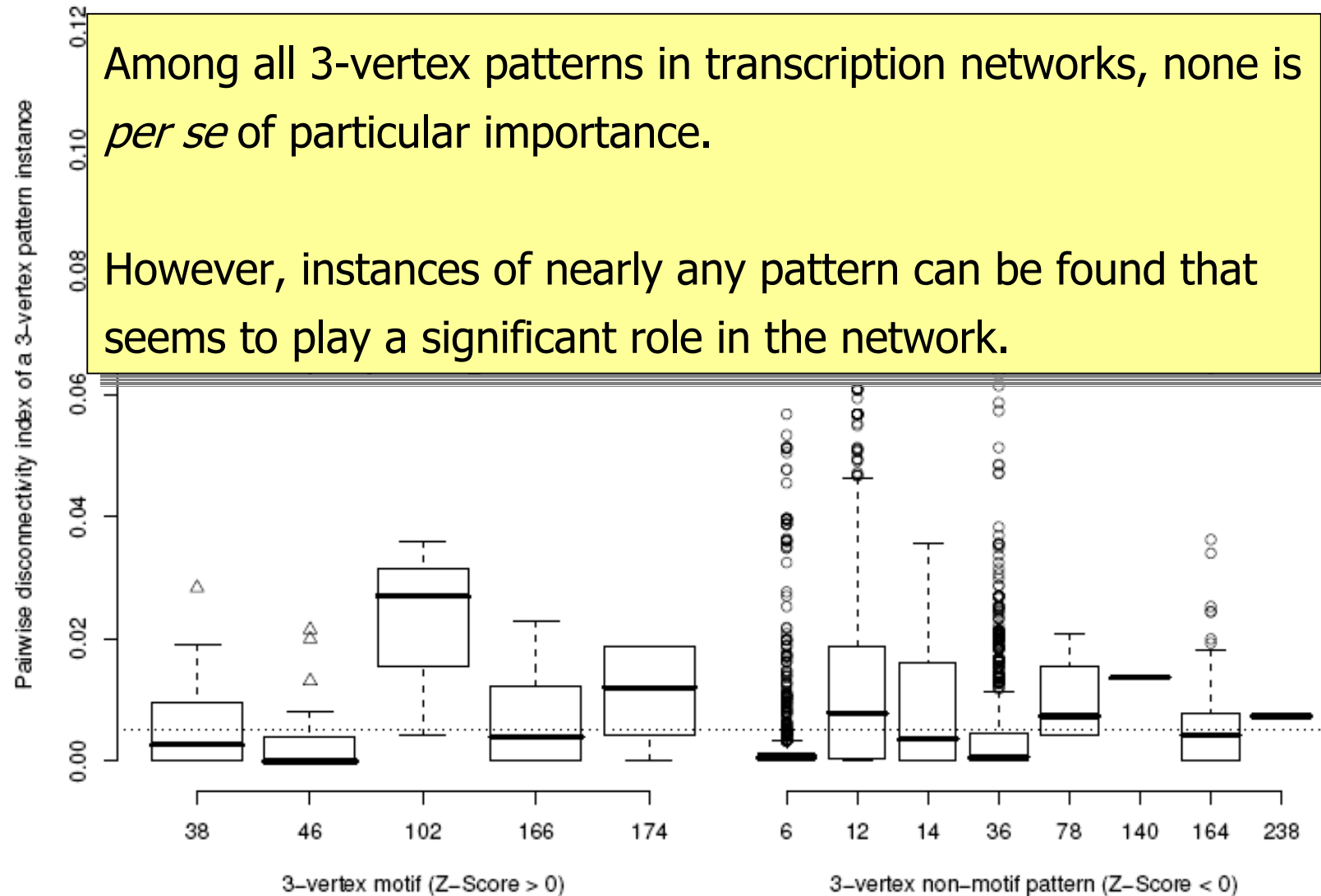
In transcription networks, the only 3-vertex motif appearing in all networks analyzed is the feed-forward loop (confirmed).

However, motifs *per se* do not play a more important role for the coherence of the network than non-motif patterns.

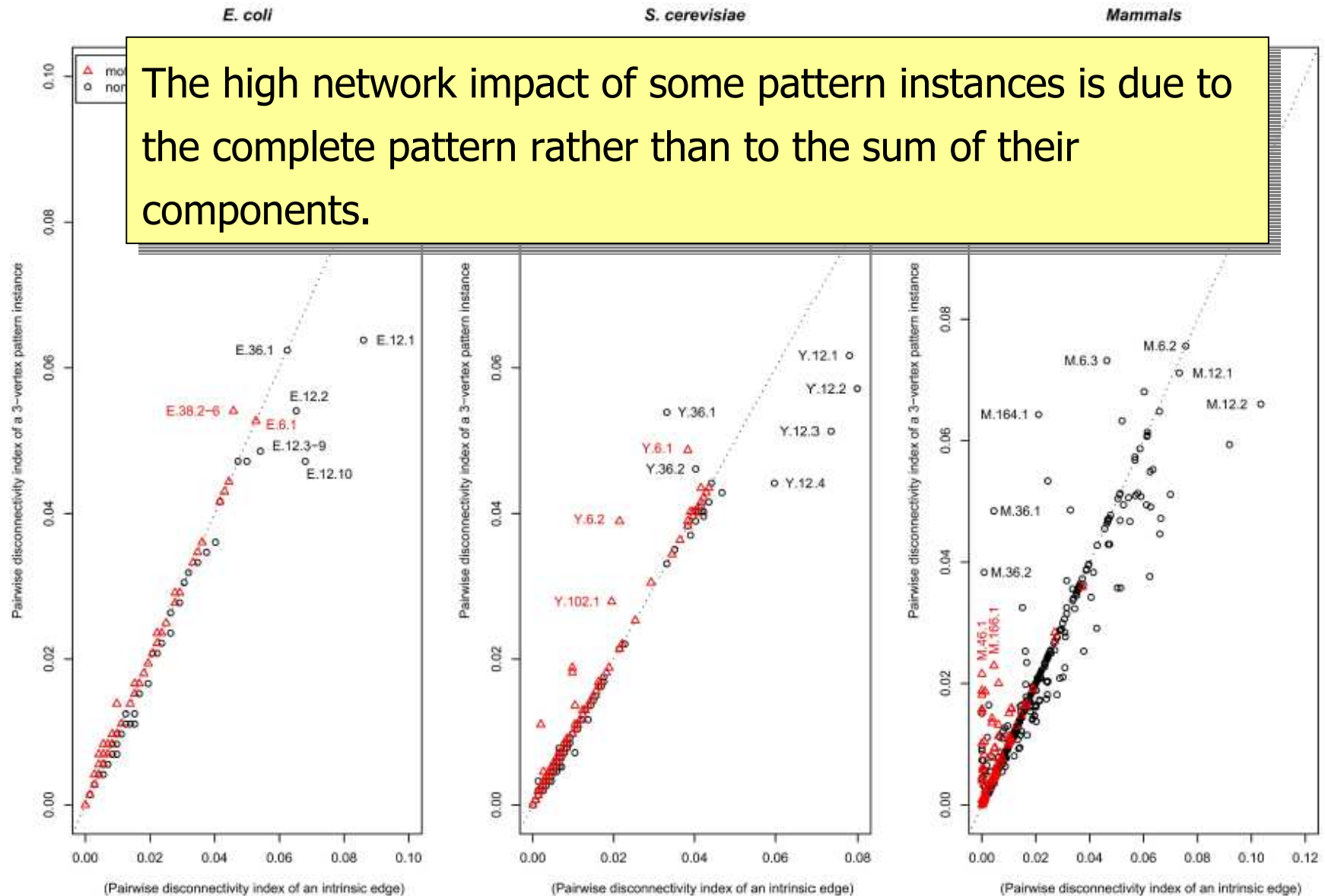
	164	-	-	-	-	-	-	197	-6.52	0.0051
	166	-	-	-	1	4.65	0.0052	20	7.12	0.0058
	174	-	-	-	-	-	-	6	7.31	0.0109
	238	-	-	-	-	-	-	1	-	0.0073

Among all 3-vertex patterns in transcription networks, none is *per se* of particular importance.

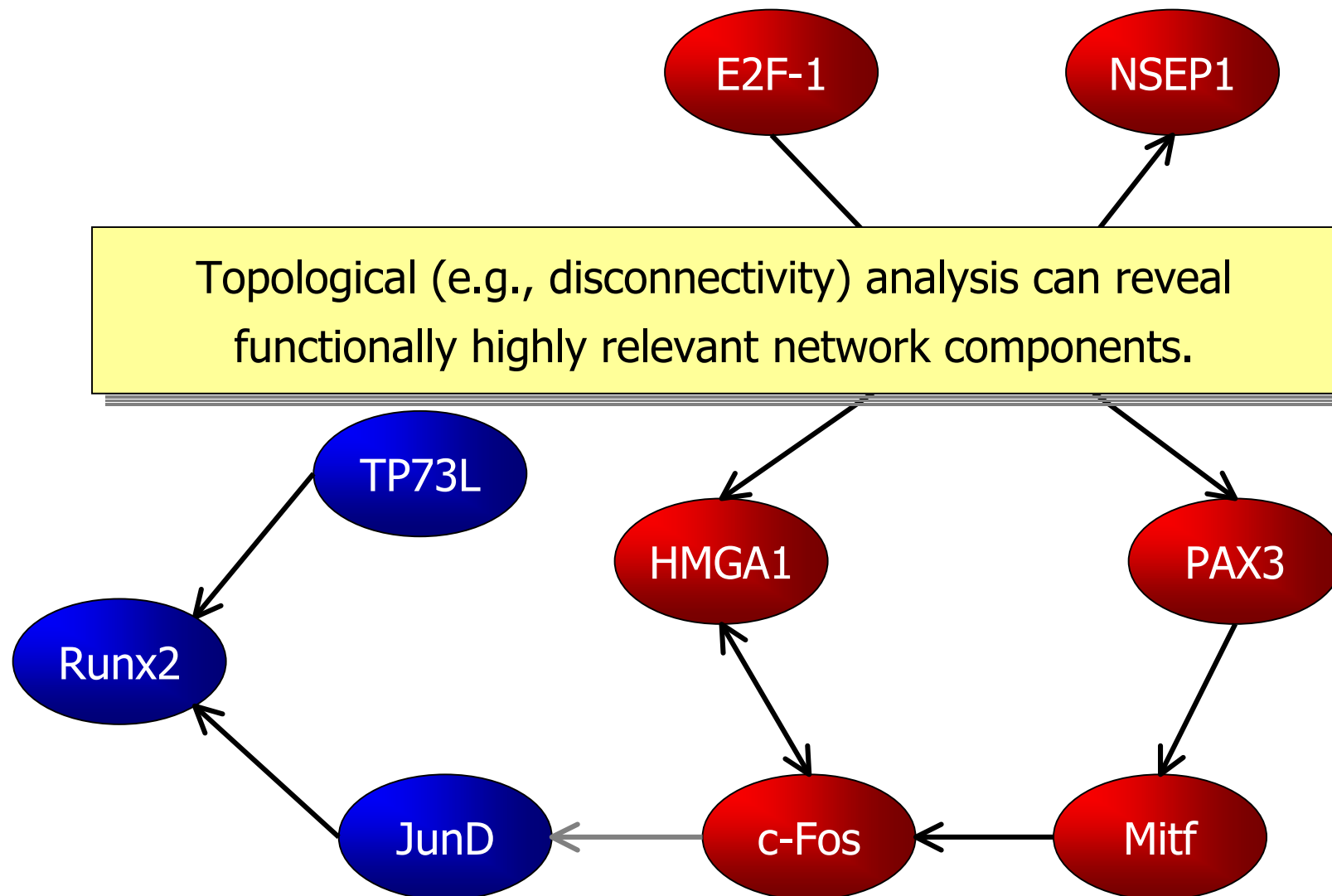
However, instances of nearly any pattern can be found that seems to play a significant role in the network.



Transcriptional network



The largest mammalian subnetwork according to motif analysis is functionally linked with cell-cycle control



- Through the past 20 years, our understanding of the mode how transcription factors (TFs) operate has made remarkable progress.
- New HTP technologies have increased the chances to come up with a genome-wide map of TF binding sites (TFBS), which, however, will always be a kind of snapshot of a selection of cellular states.
- Context-dependent approaches to TFBS prediction that also make use of comparative genomics may provide a major breakthrough.

- While TRANSFAC provides the knowledge base for this, state-of-the-art algorithms implemented in Explain™ reveal the syntax of promoters.
- The knowledge compiled in the TRANSFAC database may enable us to even predict the DNA-binding specificity of newly discovered TFs.
- Analyzing the architecture of the transcriptional network of cells will reveal key regulators that may be ideal candidates for diagnostic, therapeutic or biotechnological purposes.