

## *Regulatory Sequence Analysis*

# ***Regulatory regions and regulatory elements***

Jacques.van.Helden@ulb.ac.be

Université Libre de Bruxelles, Belgique

Laboratoire de Bioinformatique des Génomes et des Réseaux (BiGRe)

<http://www.bigre.ulb.ac.be/>

# Genomic sequences

- A genome  $G$  contains a set of  $n$  chromosomes.
  - $G = \{S_1, S_2, \dots, S_i, \dots, S_n\}$
- Each chromosome is a molecule of deoxyribonucleic acid (DNA), a polymer of 4 nucleotides
  - A Adenosine
  - C Cytidine
  - G Guanosine
  - T Thymidine
- Each chromosome is represented as a sequence ( $S_i$ ) of a text written in a 4-letter alphabet ( $A$ )
  - $A = \{A, C, G, T\}$
  - $S_i = (s_{i1}, s_{i2}, \dots, s_{ij}, \dots, s_{iL_i})$
  - $L_i$  is the length of the  $i^{th}$  chromosome

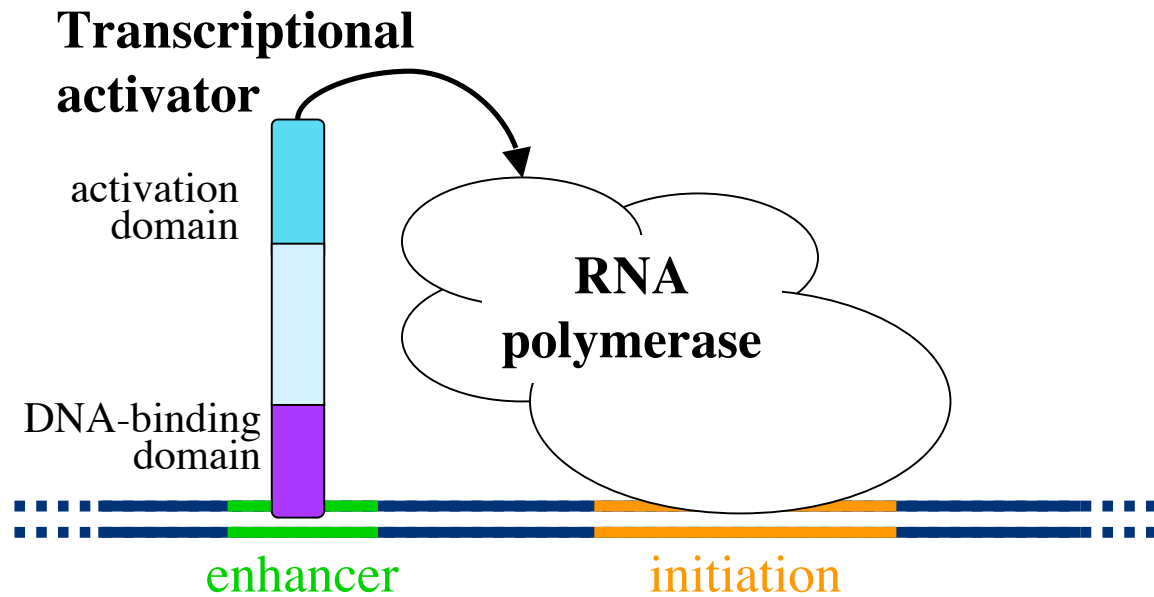
# The non-coding genome

Organism	Year	Size Mb	Genes	genes/Mb	coding %	non-coding %	repetitive %	Transcribed %
<i>Mycoplasma genitalium</i>	1995	0.6	481	802	90	10		
<i>Haemophilus influenzae</i>	1995	1.8	1 717	954	86	14		
<i>Escherichia coli</i>	1997	4.6	4 289	932	87	13		
<i>Saccharomyces cerevisiae</i>	1996	12	6 286	524	72	28		
<i>Arabidopsis thaliana</i>	2001	120	27 000	225	30	70		
<i>Caenorhabditis elegans</i>	1998	97	19 000	196	27	73		
<i>Drosophila melanogaster</i>	2000	165	16 000	97	15	85		
<i>Homo sapiens</i>	2001	3 200	31 000	10	3	97	46	28

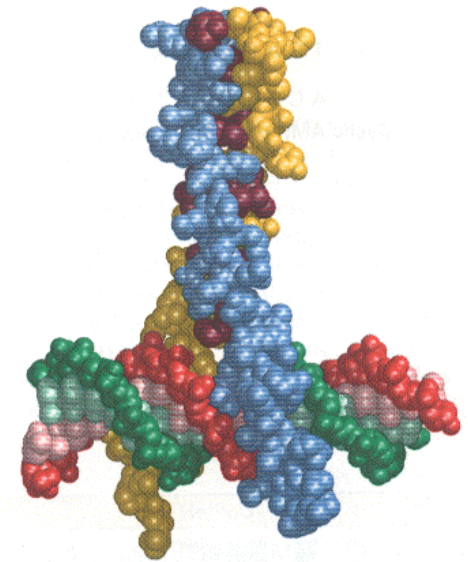
# Genome sizes - some examples

Nom d'espèce	Nom commun	Année de publication	Taille du génome Mb	Nombre de gènes	Distance moyenne entre gènes Kb	Fraction couverte par des gènes codants %	Fraction non-codante %	Fraction répétitive %	Fraction transcrite %	Remarques
Bactérie										
<i>Mycoplasma genitalium</i>	<i>Mycoplasma</i>	1995	0.6	481	1.2	90	10			Petit génome (intracellulaire)
<i>Haemophilus influenzae</i>		1995	1.8	1 717	1.0	86	14			Premier génome bactérien séquencé
<i>Escherichia coli</i>	Entérobactérie	1997	4.6	4 289	1.1	87	13			
Levures										
<i>Saccharomyces cerevisiae</i>	Levure du boulanger	1996	12	6 286	1.9	72	28			Premier génome eucaryote
Animaux										
<i>Caenorhabditis elegans</i>	Ver nématode	1998	97	19 000	5	27	73			Premier génome de métazoaire
<i>Drosophila melanogaster</i>	Mouche à vinaigre	2000	165	16 000	10	15	85			
<i>Ciona intestinalis</i>			174	14 180	12					
<i>Danio rerio</i>	Poisson zèbre		1 527	18 957	81					
<i>Xenopus laevis</i>	Xénope (amphibien)		1 511	18 023	84					
<i>Gallus gallus</i>	Poule		2 961	16 736	177					
<i>Ornithorynchus anatinus</i>	Ornithorynque		1 918	17 951	107					
<i>Mus musculus</i>	Souris	2002	3 421	23 493	146					
<i>Pan troglodytes</i>	Chimpanzé		2 929	20 829	141					
<i>Homo sapiens</i>	Humain	2001	3 200	21 528	149	2	98	46	28	Version "brouillon"
1000 génomes humains		> 2008								Projet annoncé en janvier 2008
Plantes										
<i>Arabidopsis thaliana</i>	Arabette	2001	120	27 000	4	30	70			Premier génome de plante
<i>Oryza sativa</i>	Riz		390	37 544	10					
<i>Zea mais</i>	Maïs		2 500	50 000	50			50		Nb de gènes approximatif
<i>Triticum aestivum</i>	Blé		16 000							Génome hexaploïde
<i>Lilium</i>	Lys		120 000							
<i>Psilotum nudum</i>			250 000							

# Transcriptional activation

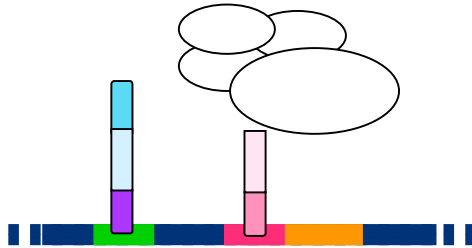


GCN4  
(leucine-zipper)  
binding to DNA

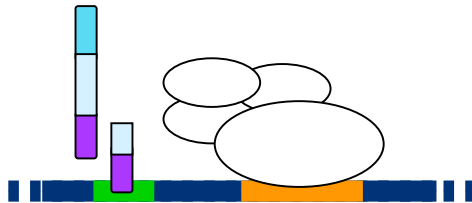


Source: L.Stryer, (1995).  
Biochemistry. p1003

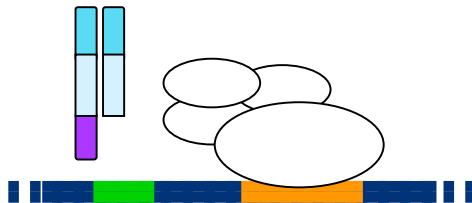
# Transcriptional repression



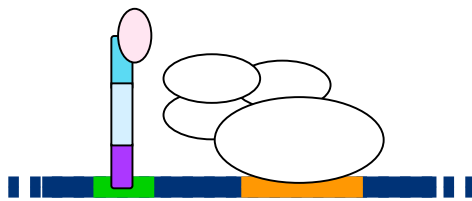
Prevent RNA polymerase from accessing DNA



Competition for factor binding site

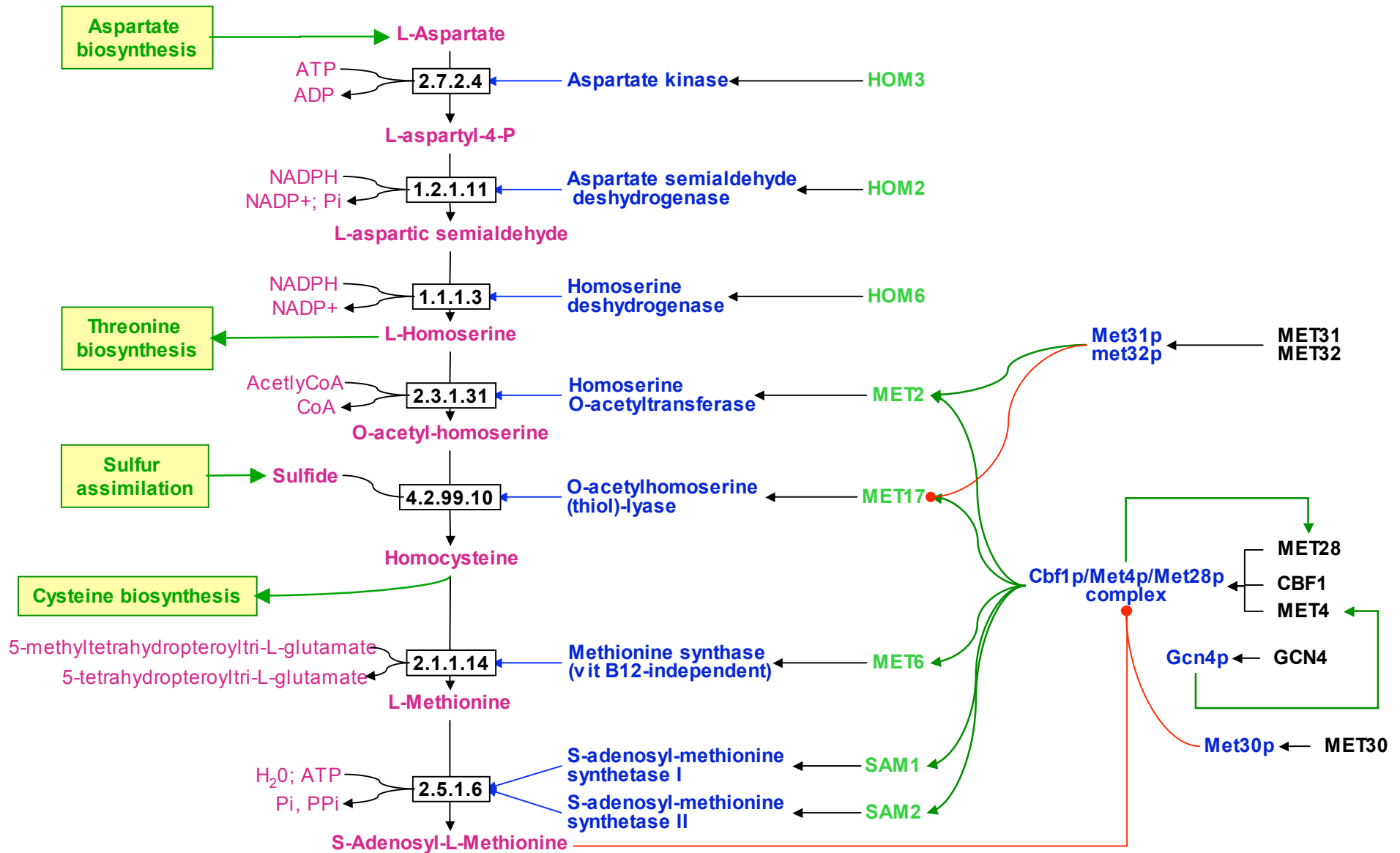


Factor titration

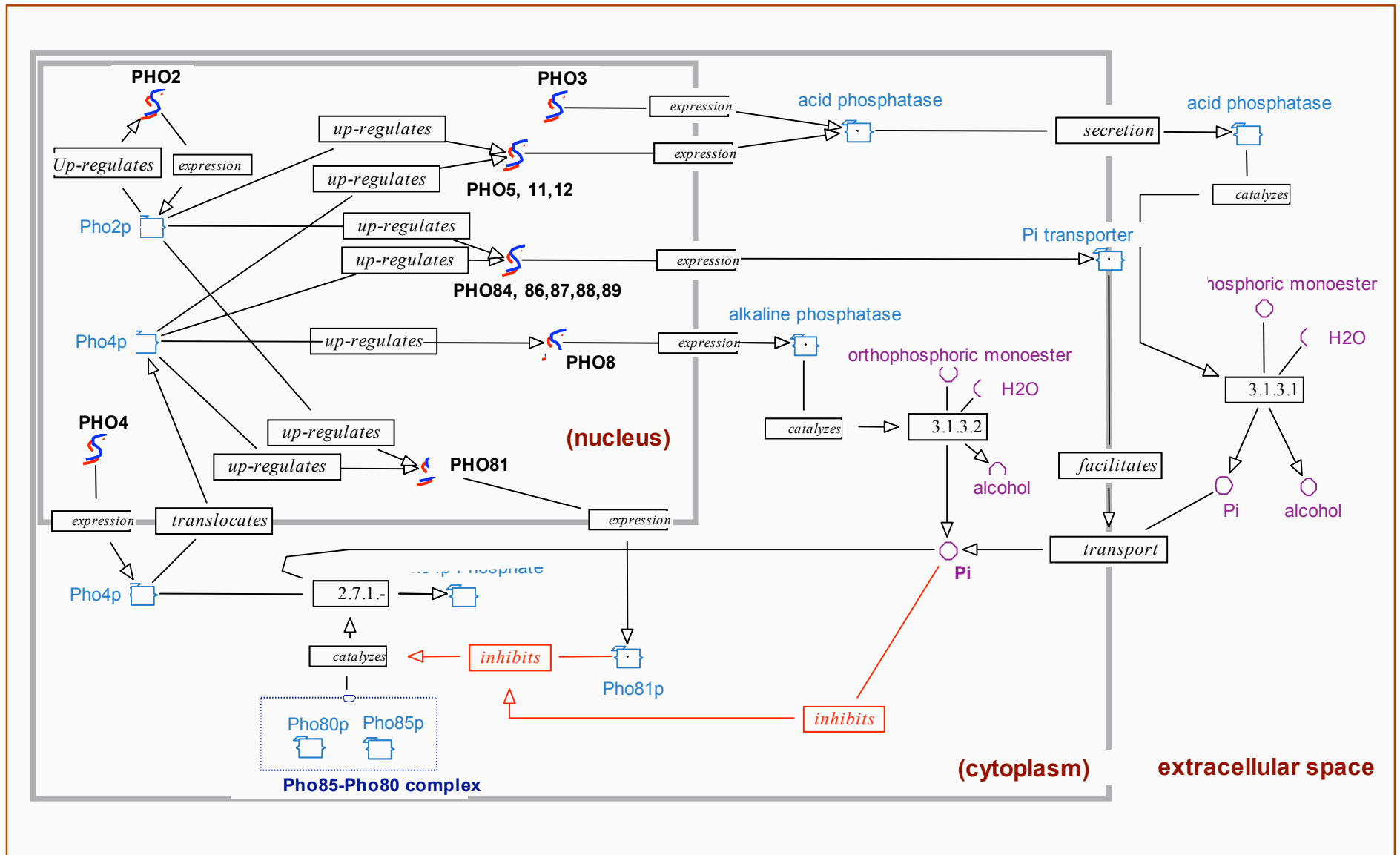


Prevent transcription factor from interacting with RNA-polymerase (bind with activation domain)

# Methionine Biosynthesis in *S.cerevisiae*



# Phosphate utilization in yeast





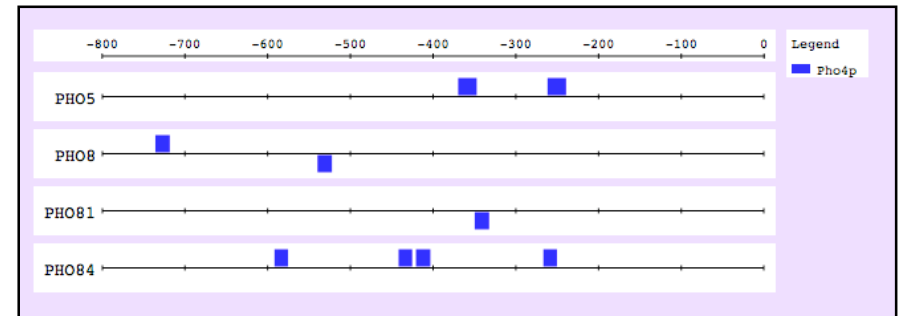
*Interface between the yeast Pho4p protein and one of its binding sites*



# Transcription factor binding site (TFBS)

Gene	Ft_type	Factor	Strand	left	right	Sequence
PHO5	site	Pho4p	D	-370	-347	TAAATTAG <b>CACGTTTT</b> CGCATAGA
PHO5	site	Pho4p	D	-262	-239	TGGCACTCAC <b>CACGTGGG</b> ACTAGCA
PHO8	site	Pho4p	R	-540	-522	ATCGCTGC <b>CACGTGG</b> CCCGA
PHO8	site	Pho4p	D	-736	-718	ATATTAAGCGTGCGGGTAA
PHO81	site	Pho4p	R	-350	-332	TTATTTCG <b>CACGTGCC</b> ATAA
PHO84	site	Pho4p	D	-592	-575	TTACGC <b>CACGTTGG</b> TGCTG
PHO84	site	Pho4p	D	-421	-403	TTTCCAG <b>CACGTGGG</b> CGG
PHO84	site	Pho4p	D	-442	-425	TAGTTCC <b>CACGTGG</b> ACGTG
PHO84	site	Pho4p	DR	-879	-874	aaaagtgt <b>CACGTG</b> ataaaaat
PHO84	site	Pho4p	D	-267	-250	TAATACGC <b>CACGTTTT</b> TAA

- A *transcription factor binding site (TFBS)* is a **location** within a sequence, where a transcription factor binds specifically.
- A site can be
  - characterized experimentally (known site)
  - inferred by some algorithm (predicted site)
- Example
  - binding sites for the yeast transcription factor Pho4p. Coordinates are relative to the start codon.



# Alignment of transcription factor binding sites

## Binding sites for the yeast Pho4p transcription factor

(Source : Oshima et al. Gene 179, 1996; 171-177)

Gene	Site Name	Sequence	Affinity
PHO5	UASp2	---aCtCaCA <b>CACGTGG</b> ACTAGC-	high
PHO84	Site D	---TTTCCA <b>GCACGTGG</b> GCGGA--	high
PHO81	UAS	----TTATG <b>GCACGTGC</b> GAATAA--	high
PHO8	Proximal	GTGATCGCT <b>GCACGTGG</b> CCCGA---	high
group 1	consensus	----- <b>gCACGTGg</b> g-----	high
PHO5	UASp1	--TAAATT <b>GCACGTTT</b> TCGC----	medium
PHO84	Site E	----AATAC <b>GCACGTTT</b> TTAATCTA	medium
group 2	consensus	----- <b>cgCACGTTt</b> t-----	medium
Degenerate consensus		----- <b>GCACGTKKk</b> -----	high-med

### IUPAC ambiguous nucleotide code

A	A	Adenine
C	C	Cytosine
G	G	Guanine
T	T	Thymine
R	A or G	puRine
Y	C or T	pYrimidine
W	A or T	Weak hydrogen bonding
S	G or C	Strong hydrogen bonding
M	A or C	aMino group at common position
K	G or T	Keto group at common position
H	A, C or T	not G
B	G, C or T	not A
V	G, A, C	not T
D	G, A or T	not C
N	G, A, C or T	aNy

### Non-binding sites

PHO5	UASp3	--TAATTTG <b>GCAT</b> <u>I</u> GTGCGATCTC--	No binding
PHO84	Site C	-----ACGTCC <b>CACGTG</b> <u>A</u> ACTAT--	No binding
PHO84	Site A	-----TTTAT <b>CACGTG</b> <u>A</u> CACTTTTT	No binding
PHO84	Site B	-----TTAC <b>GCACGTT</b> <u>G</u> GTGCTG--	No binding
PHO8	Distal	---TTACCC <b>GCACG</b> <u>C</u> TTAATAT---	No binding

## *From binding sites to count matrix*

- The TRANSFAC database contains 8 binding sites for the yeast transcription factor Pho4p
  - 5/8 contain the core of high-affinity binding sites (CACGTG)
  - 3/8 contain the core of medium-affinity binding sites (CACGTT)

R06098	\TCAC <b>CACGT</b> GGGA\
R06099	\GGC <b>CACGT</b> GCAG\
R06100	\TGAC <b>CACGT</b> GGGT\
R06102	\CAG <b>CACGT</b> GGGG\
R06103	\TTC <b>CACGT</b> GCGA\
R06104	\ACG <b>CACGTT</b> GGT\
R06097	\CAG <b>CACGTTT</b> TC\
R06101	\TACC <b>CACGTTT</b> TC\

# Count matrix

## Alignment of Pho4p binding sites (TRANSFAC annotations)

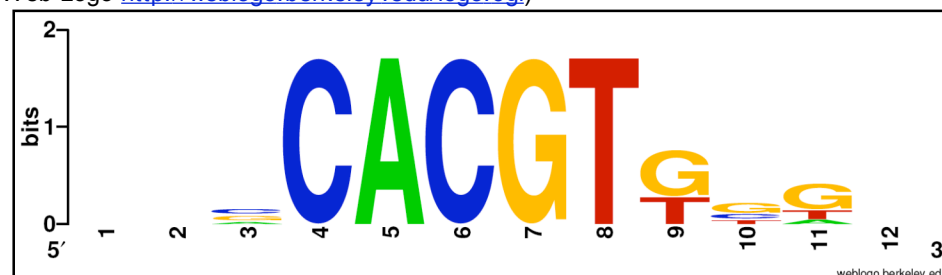
R06098	T	C	A	C	A	C	G	T	G	G	G	A
R06099	G	G	C	C	A	C	G	T	G	C	A	G
R06100	T	G	A	C	A	C	G	T	G	G	G	T
R06102	C	A	G	C	A	C	G	T	G	G	G	G
R06103	T	T	C	C	A	C	G	T	G	C	G	A
R06104	A	C	G	C	A	C	G	T	T	G	G	T
R06097	C	A	G	C	A	C	G	T	T	T	T	C
R06101	T	A	C	C	A	C	G	T	T	T	T	C

## Count matrix (TRANSFAC matrix F\$PHO4\_01)

Residue\position	1	2	3	4	5	6	7	8	9	10	11	12
<b>A</b>	1	3	2	0	8	0	0	0	0	0	1	2
<b>C</b>	2	2	3	8	0	8	0	0	0	2	0	2
<b>G</b>	1	2	3	0	0	0	8	0	5	4	5	2
<b>T</b>	4	1	0	0	0	0	0	8	3	2	2	2
<b>Sum</b>	8	8	8	8	8	8	8	8	8	8	8	8

## Tom Schneider's sequence logo

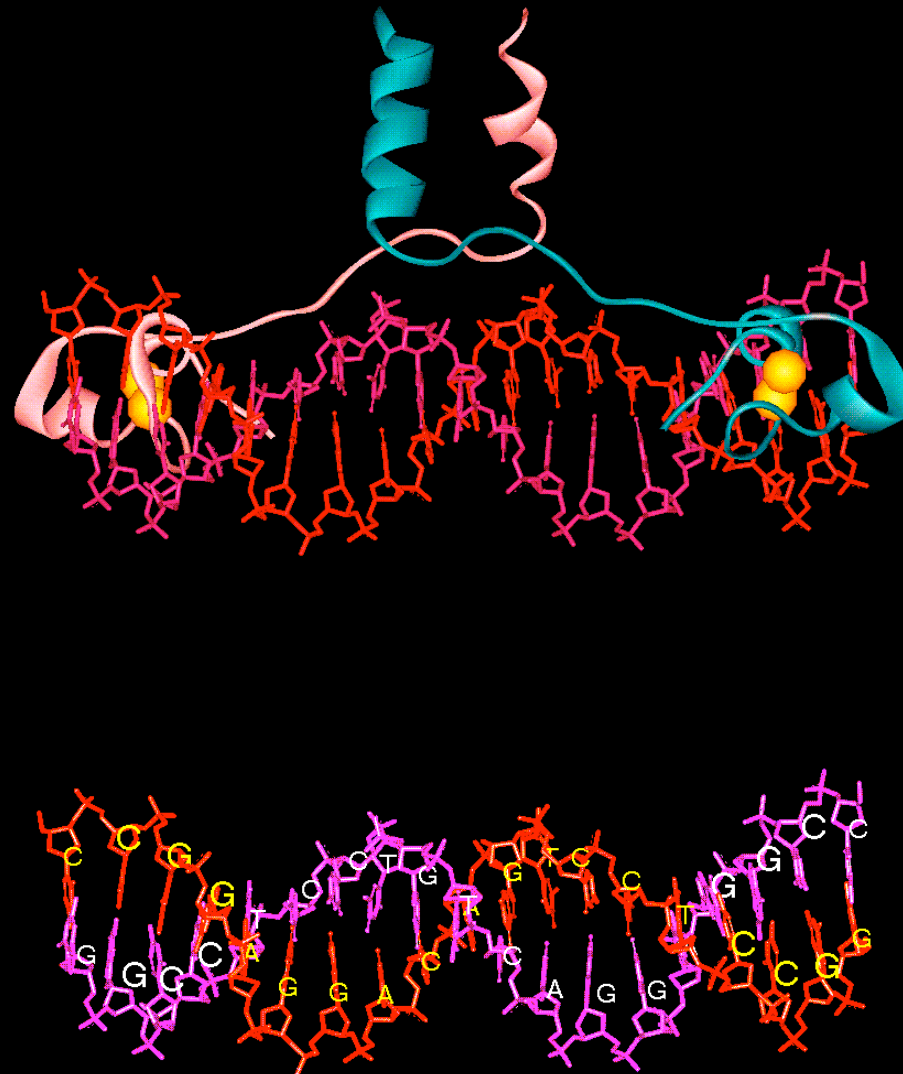
(generated with Web Logo <http://weblogo.berkeley.edu/logo.cgi>)



## Motif / pattern

- We use the term ***motif*** (or ***pattern***) in the sense of a model used to represent the specificity of binding for a transcription factor.
- A motif can be described using different formalisms.
  - Consensus string
    - nucleotide alphabet    **CACGTGGG**
    - IUPAC alphabet    **CACGTGKK**
    - regular expressions.    **CACGTG[GT][GT]**
  - Position-specific scoring matrix (PSSM)
  - Logo representation (Schneider, 1986)
  - Hidden Markov Models (HMM)

*Interface between the yeast Gal4p protein and one of its binding sites*

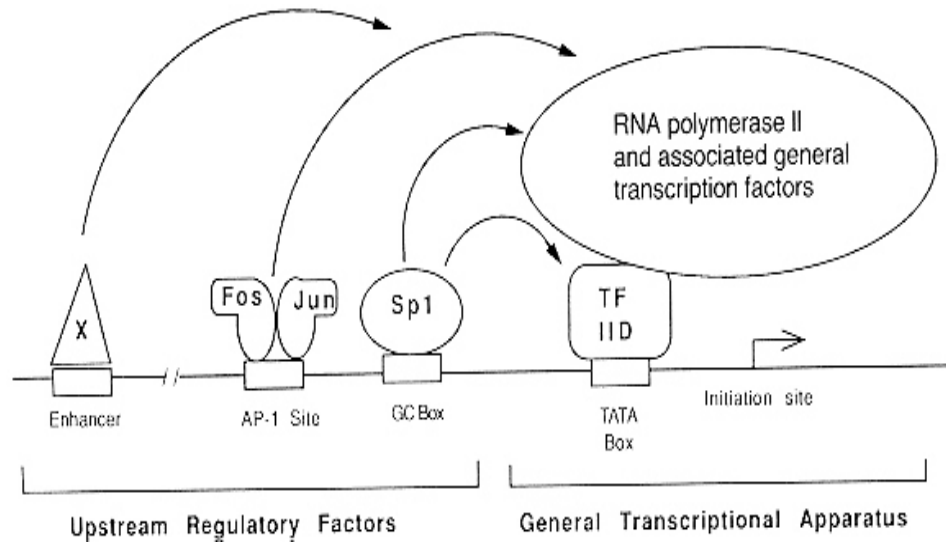


# *Characteristics of yeast regulatory elements*

- Short DNA sequences (5-30 bp)
  - Highly conserved core (5-8 bp), with partly conserved flanking nucleotides
  - Pair of very short oligonucleotides (3 nt) separated by a non-conserved segment (0-20 bp)
  
- In the yeast *Saccharomyces cerevisiae*
  - Located upstream the regulated gene
  - Strand-insensitive
    - Activity does not depend on the strand
  - Within 800 bp from the start codon
    - Activity does not depend on precise position



# Cis-regulatory modules (CRM)



- In higher organisms, some non-coding regions (typically 100-200 bp) contain closely packed binding sites for distinct transcription factors.
- These regions are called ***cis-regulatory modules (CRMs)***
- CRMs play the role of integrating devices.
- Depending on the combination of transcription factors present in the cell, they will activate or repress the expression of a target gene.

## Regulatory regions

<b>organism</b>	<b>coli</b>	<b>yeast</b>	<b>metazoan</b>
<b>location</b>	upstream overlap. Initiation	upstream	upstream downstream within introns
<b>distance range</b>	-400 to +50 bp	-800 to -1 bp	from several Kbs to several Mb !
<b>position effect</b>	often essential	often irrelevant	often irrelevant
<b>strand</b>	sensitive or symmetric	insensitive	insensitive
<b>most common core</b>	spaced pair of 3nt	~5-8 conserved bp	~5-8 conserved bp
<b>repeated sites</b>	rare	occasional	frequent
<b>cis-regulatory modules (CRMs)</b>			frequent

# Questions and approaches

- **Pattern matching**

- If we know the consensus for a given transcription factor, can we predict its binding sites in a DNA sequence ?

- **Matching a library of patterns**

- Can we scan a sequence for matches with the consensus of all the currently known transcription factors ?

- **Pattern discovery**

- Starting from a set of co-regulated genes, can we predict cis-acting elements involved in their transcriptional regulation ?

- **Phylogenetic footprinting**

- Can we detect regulatory signals by searching conserved elements in non-coding sequences of orthologous genes ?

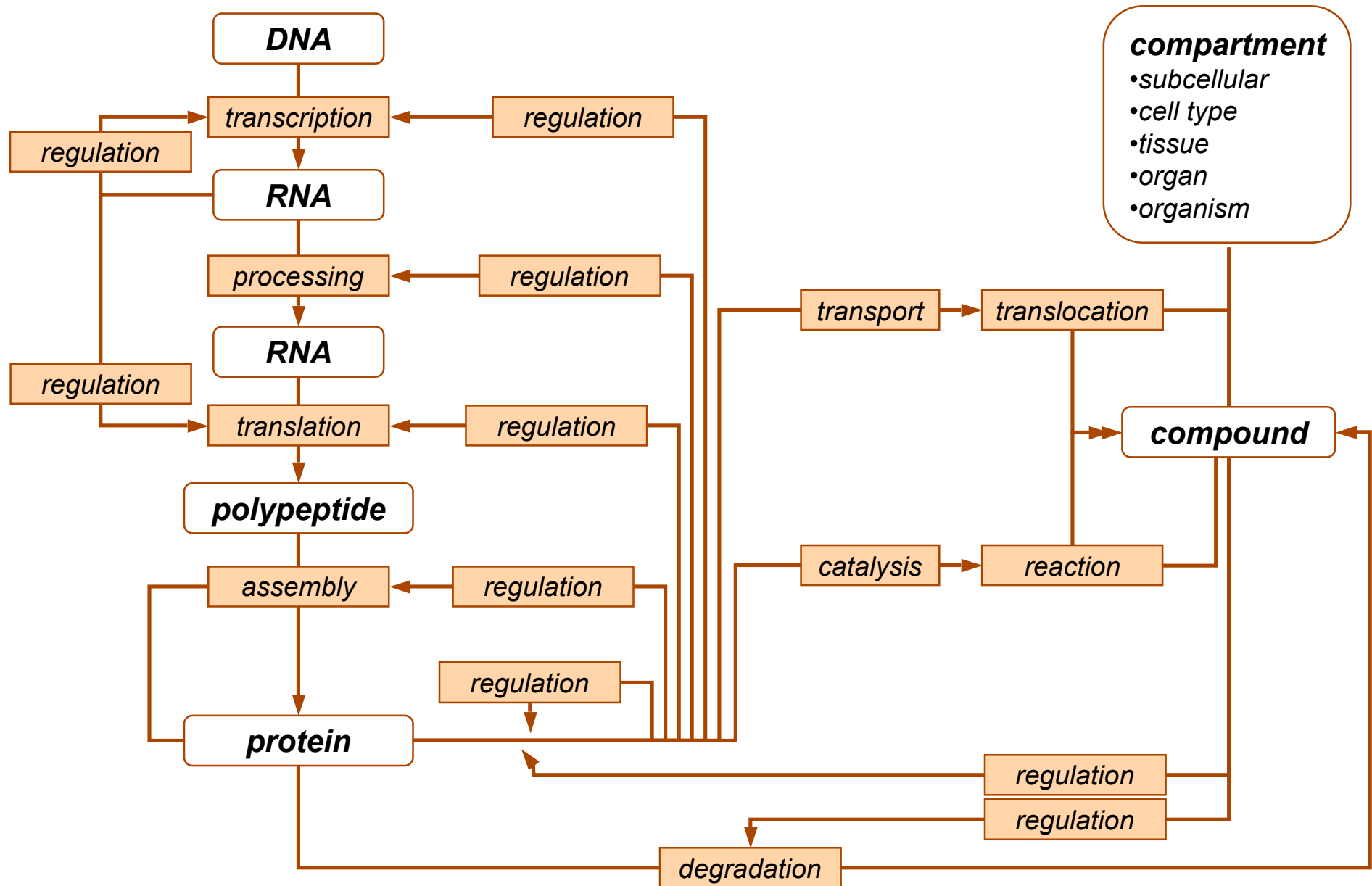
- **Network inference**

- Can we infer groups networks of regulation from cis-regulatory elements ?

- **Gene classification** on the basis of pattern scores

- Can we classify genes on the basis of the presence of regulatory motifs in their regulatory regions ?
- **Unsupervised classification (clustering):** regroup elements (genes) in clusters without a priori knowledge about these clusters. The clusters are “discovered” during the clustering process.
- **Supervised classification:** use pre-defined groups of genes (training sets) to train a program, and then use this program to assign new elements (genes) to one of the pre-defined groups.

# Molecular networks (shamefully simplified)



## *Regulatory Sequence Analysis*

# ***Supplementary material***

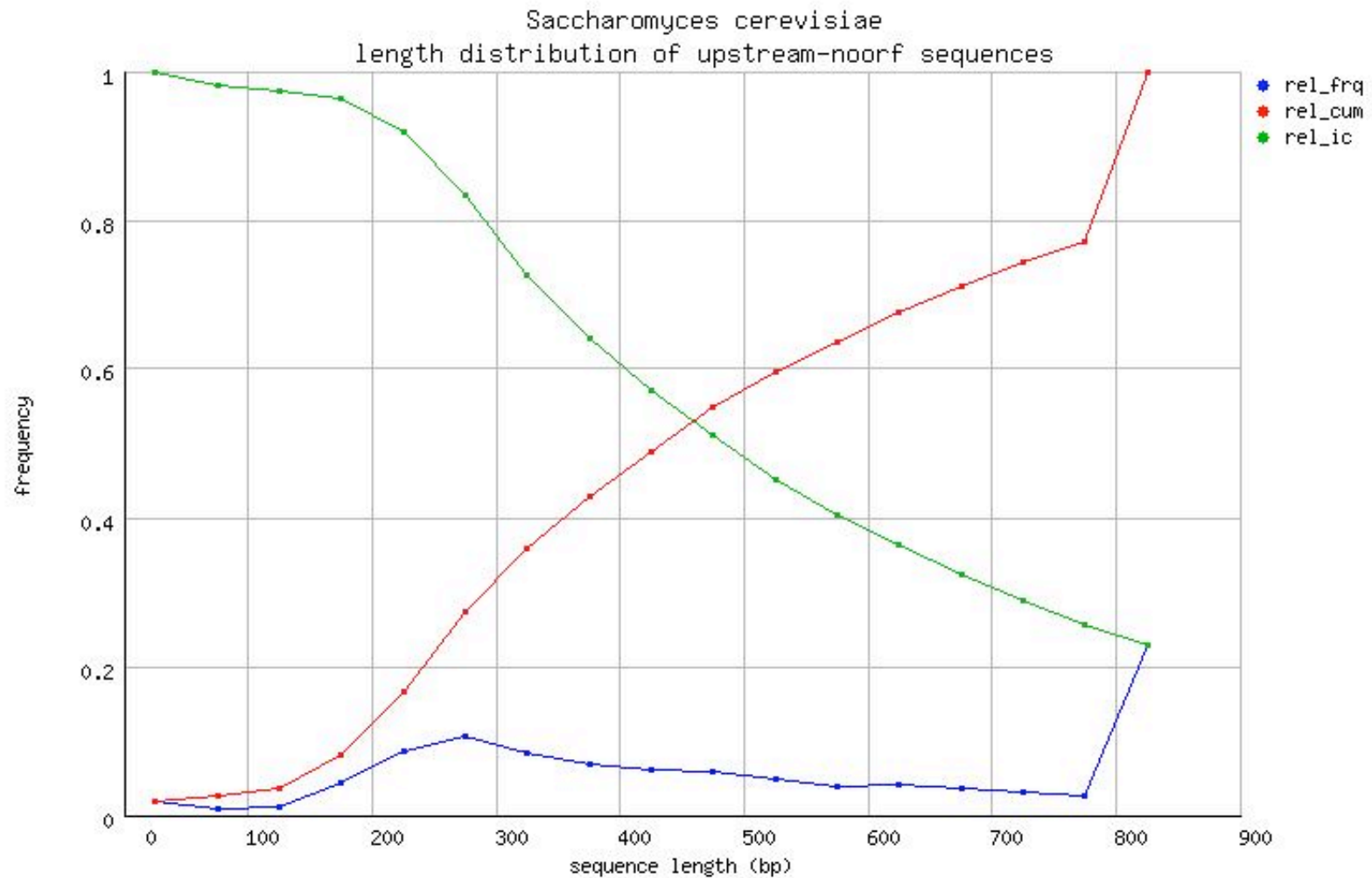
Jacques.van.Helden@ulb.ac.be

Université Libre de Bruxelles, Belgique

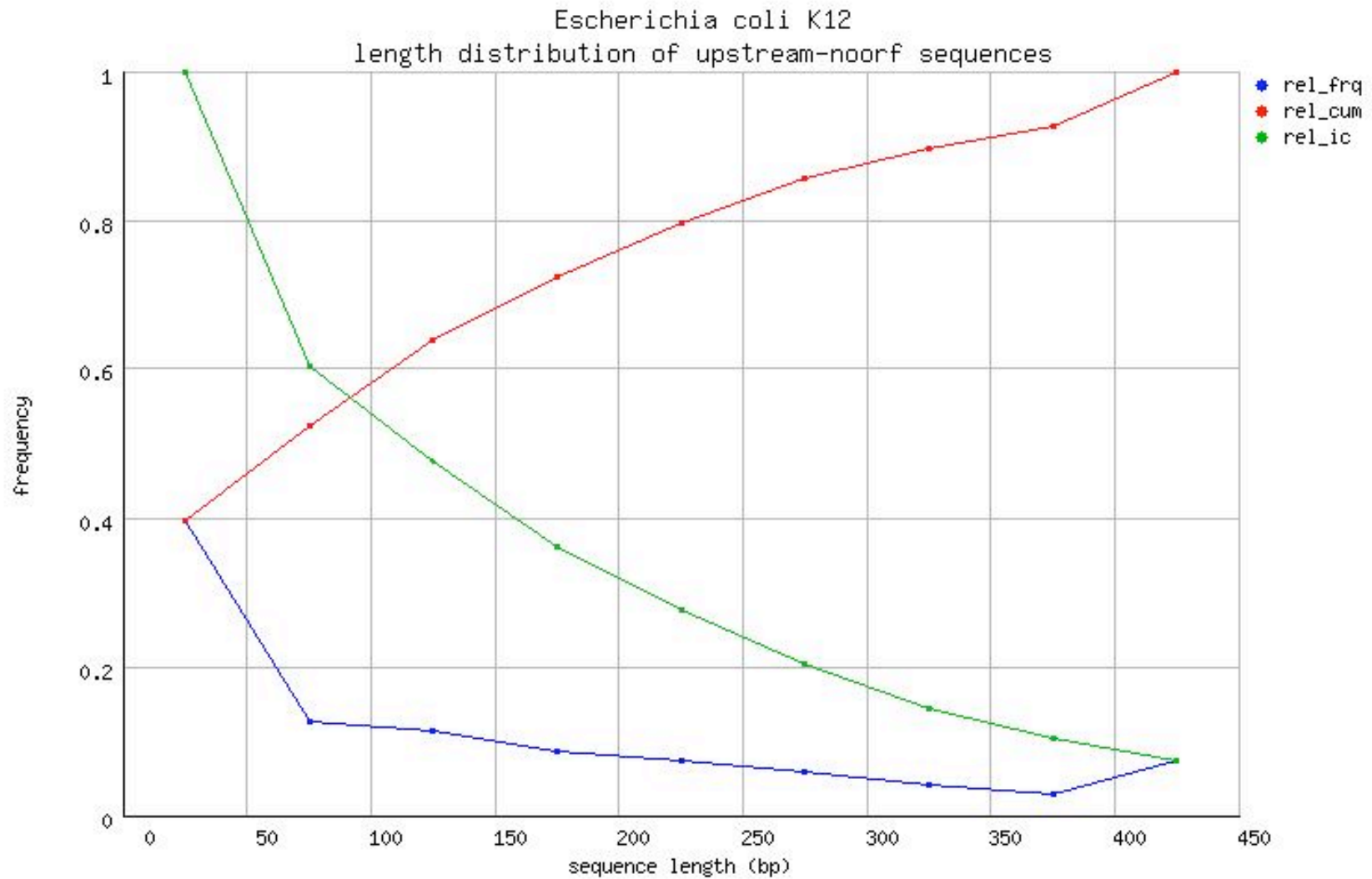
Laboratoire de Bioinformatique des Génomes et des Réseaux (BiGRe)

<http://www.bigre.ulb.ac.be/>

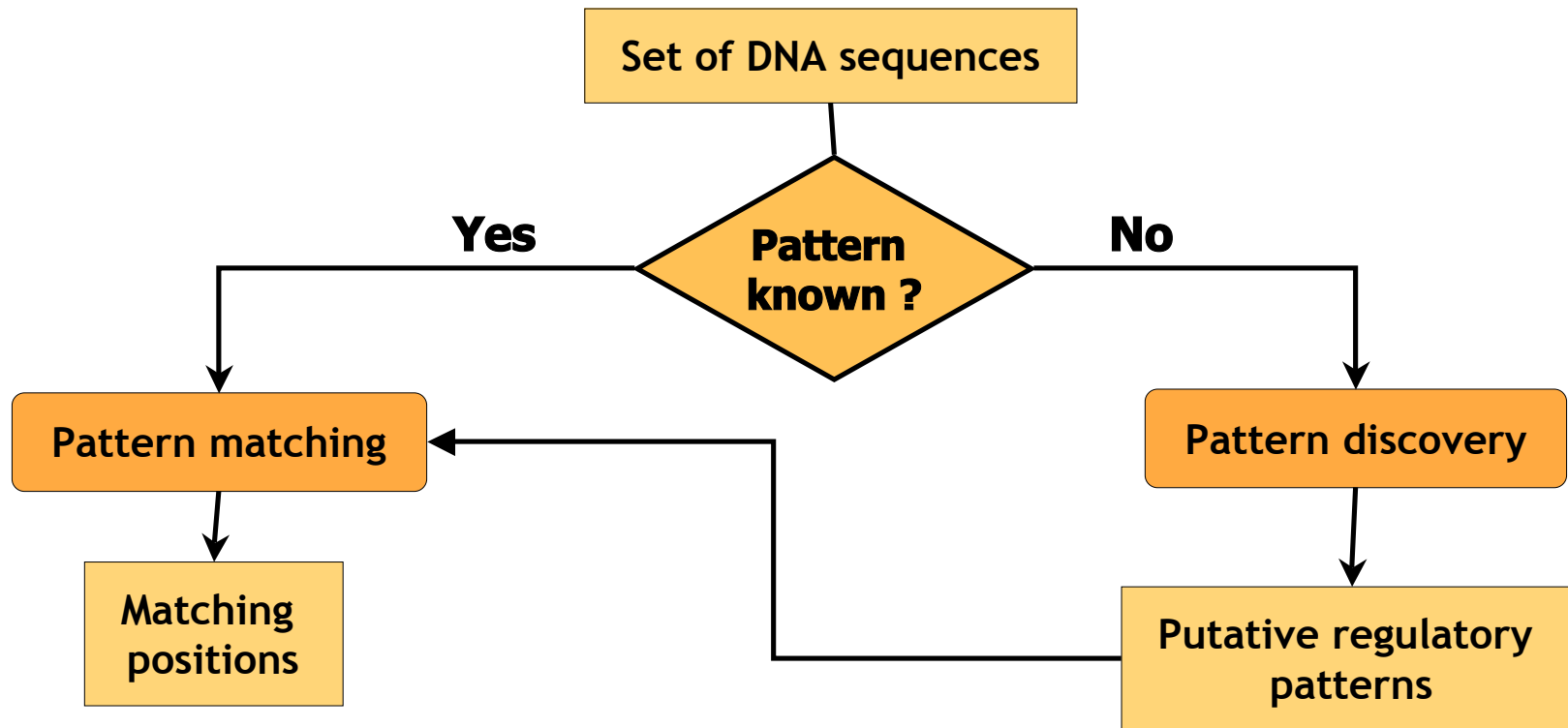
## Distribution of upstream sequence lengths *Saccharomyces cerevisiae*



## Distribution of upstream sequence lengths *Escherichia coli* K12



# *Pattern matching vs pattern discovery*





## Other examples of sequence logos

ATGTAI GG

Rap1

GGTGGC AAAA

Rpn4

AAI GA TCA

Gcn4

GAA T C GAA

HSE

IG GGGG A

Mig1

AAI ICACGTG

Met4/Cbf1

## *Typical situations : pattern discovery*

- Selected sequence set
  - e.g. family of 20 co-regulated genes, obtained from DNA chip experiment  
→ identify putative regulatory sites
- Genome-scale pattern discovery
  - e.g. all upstream sequences  
→ identify transcription initiation signals
  - e.g. all downstream sequences  
→ identify 3' maturation signals

## *Typical situations : pattern matching*

- Selected genes, selected patterns
  - e.g. 10 genes known to be regulated by a factor  
→ search matching positions
- Selected genes, library of patterns
  - → infer putative action of any previously known transcription factor
- All genes, selected patterns
  - → classify all the genes of a genome according to putative regulatory properties

# Met4p binding sites

gene	start	end	sequence
MET3	-367	-349	GAAAAG <b>TCACGTG</b> TAATTT
MET3	-384	-366	AAAAGG <b>TCACGTG</b> ACCAGA
MET14	-235	-217	CTAATTT <b>TCACGTG</b> ATCAAT
MET16	-185	-167	ATCATT <b>TCACGTG</b> GCTAGT
ECM17	-311	-293	ATTTCAT <b>TCACGTG</b> CGTATT
ECM17	-339	-321	.TTTGTC <b>TCACGTG</b> ATATTTC
MET10	-255	-237	.CCACAC <b>TCACGTG</b> AGCTTAT
MET10	-237	-219	.TAGAAG <b>TCACGTG</b> ACCACAA
MET2	-360	-342	GTATTTT <b>TCACGTG</b> ATGCGC
MET2	-554	-536	TAATAAT <b>TCACGTG</b> ATATTT
MET17	-306	-288	.AAATGG <b>TCACGTG</b> AAGCTGT
MET17	-332	-314	TTGAGG <b>TCACATG</b> ATCGCA
MET6	-540	-522	GCCACAT <b>TCACGTG</b> CACATT
MET6	-502	-484	AATATTT <b>TCACGTG</b> ACTTAC
SAM2	-329	-311	.TCTACC <b>TCACGTG</b> ACTATAA
SAM2	-381	-363	.TCTTCAT <b>TCATGTG</b> ATTCATC

A	13	11	3	3	2	0	16	0	1	0	0	12
C	1	0	0	3	0	16	0	15	0	0	0	0
G	1	1	4	4	4	0	0	0	15	0	16	4
T	1	4	9	6	10	0	0	1	0	16	0	0

## Met31p binding sites

gene	start	end	sequence
MET14	-202	-182	CCTC <b>AAAAA</b> ATGTGGCAATGG
MET2	-313	-293	TGC <b>AAAAA</b> ATGTGGATGCAC
MET17	-227	-207	TCATG <b>AAA</b> ACTGTGTAAACATA
MET6	-313	-293	GTCGC <b>AAA</b> ACTGTGGTAGTCA
SAM2	-306	-286	GCTTG <b>AAA</b> ACTGTGGCGTTTT
SAM1	-283	-263	ACAGG <b>AAA</b> ACTGTGGTGGCGC
MET19	-173	-153	ATAAGC <b>AA</b> ACTGTGGTTCAT
MUP3	-188	-168	CGG <b>AAAAA</b> ACTGTGGCGTCGC
MET8	-184	-164	GG <b>AAAAA</b> AAATGTGAAAATCG
MET1	-232	-212	CATAAT <b>AA</b> ACTGTGAACGGAC
MET3	-259	-239	ACAAAG <b>CCACAGTTTT</b> ACAAC
MET28	-159	-139	CTAAC <b>CCACAGTTTT</b> GGGCG
MET8	-434	-414	TCTTGT <b>CCGCAGTTTT</b> ATCTG
MET30	-168	-148	GGGAAG <b>CCACAGTTT</b> GCGCGG
MET6	-405	-385	CTATCGA <b>ACTCGTTT</b> AGTCGC

A	5	11	14	14	14	2	0	0	0	0	2	5
C	2	2	0	0	0	11	0	0	1	0	0	5
G	5	0	0	0	0	0	0	14	0	14	11	1
T	2	1	0	0	0	1	14	0	13	0	1	3

## Pho4p binding sites

gene	start	end	sequence
PHO5	-260	-242	..GCACTCA <b>CACGTGGG</b> ACTA
PHO5	-260	-245	..GCACTCA <b>CACGTGGGA</b>
PHO5	-262	-239	TGGCACTCA <b>CACGTGGG</b> ACTAGCA
PHO8	-540	-522	...TCGGGC <b>CACGTGC</b> AGCGAT
PHO8	-736	-718	..ttaccgc <b>CACGCTT</b> aatat
PHO81	-350	-332	...TTATGG <b>CACGTGCG</b> AATAA
PHO84	-421	-403	..TTCCAG <b>CACGTGGG</b> GCGG
PHO84	-442	-425	...TAGTTC <b>CACGTGG</b> ACGTG
PHO84	-879	-874	.aaaagtgt <b>CACGTG</b> ataaaaat
PHO84	-267	-250	..taatacg <b>CACGTTTTT</b> aa
PHO84	-592	-575	....TTACG <b>CACGTT</b> GGTGCTG
PHO5	-368	-349	...AATTAG <b>CACGTTTT</b> CGCATA
PHO5	-369	-354	..AAATTAG <b>CACGTTT</b> CTC
PHO5	-370	-347	.TAAATTAG <b>CACGTTTT</b> CGCATAGA

## *IUPAC ambiguous nucleotide code*

<b>A</b>	<b>A</b>	<b>Adenine</b>
<b>C</b>	<b>C</b>	<b>Cytosine</b>
<b>G</b>	<b>G</b>	<b>Guanine</b>
<b>T</b>	<b>T</b>	<b>Thymine</b>
<b>R</b>	<b>A or G</b>	<b>puRine</b>
<b>Y</b>	<b>C or T</b>	<b>pYrimidine</b>
<b>W</b>	<b>A or T</b>	<b>Weak hydrogen bonding</b>
<b>S</b>	<b>G or C</b>	<b>Strong hydrogen bonding</b>
<b>M</b>	<b>A or C</b>	<b>aMino group at common position</b>
<b>K</b>	<b>G or T</b>	<b>Keto group at common position</b>
<b>H</b>	<b>A, C or T</b>	<b>not G</b>
<b>B</b>	<b>G, C or T</b>	<b>not A</b>
<b>V</b>	<b>G, A, C</b>	<b>not T</b>
<b>D</b>	<b>G, A or T</b>	<b>not C</b>
<b>N</b>	<b>G, A, C or T</b>	<b>aNy</b>

## Pho4p binding specificity - matrix descriptions

**C**

C	Pho4p											
A	14	0	5	7	6	0	26	0	0	0	0	3
C	2	8	5	16	6	26	0	26	0	1	0	4
G	4	2	1	1	12	0	0	0	26	0	16	12
T	6	16	15	2	2	0	0	0	0	25	10	7

D

Pho4p.cacgtg												
A	2	17	0	0	0	0	2	1	8	5	5	13
C	16	0	18	0	0	0	6	3	4	5	0	1
G	0	1	0	18	0	18	9	12	2	5	2	1
T	0	0	0	0	18	0	1	2	4	3	11	3

# E

E
A
C
G
T



## Regulatory sites : matrix description

### Position-specific scoring matrix (PSSM)

Pos	1	2	3	4	5	6	7	8	9	10
A	3	2	0	12	0	0	0	0	1	3
T	1	1	0	0	0	0	11	5	4	4
G	3	7	0	0	0	12	0	7	5	4
C	5	2	12	0	12	0	1	0	2	1

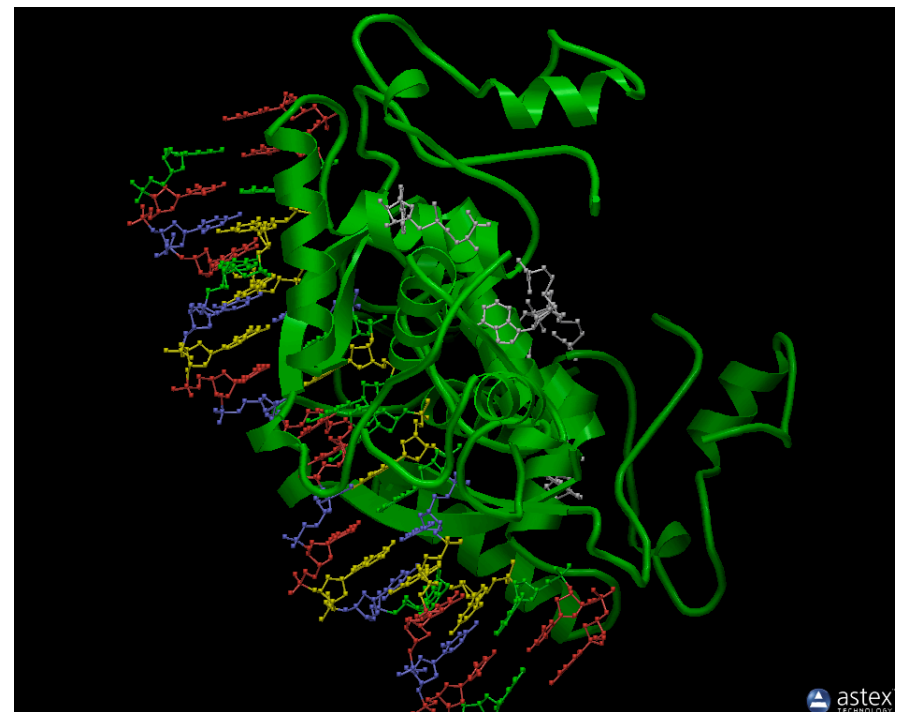
Binding motif for the yeast Pho4p transcription factor

Source : SCPD

<http://rulai.cshl.edu/cgi-bin/SCPD/getfactor?PHO4>

# Methionine repressor

- Crystal structure of the methionine repressor from *Escherichia coli*.
- In green: the MetJ protein forms a homodimer which is able to bind DNA.
- Nucleotide structure is coloured by type of nucleotide (A,C,G,T).
- In grey: the repressor is activated by binding of methionine molecules

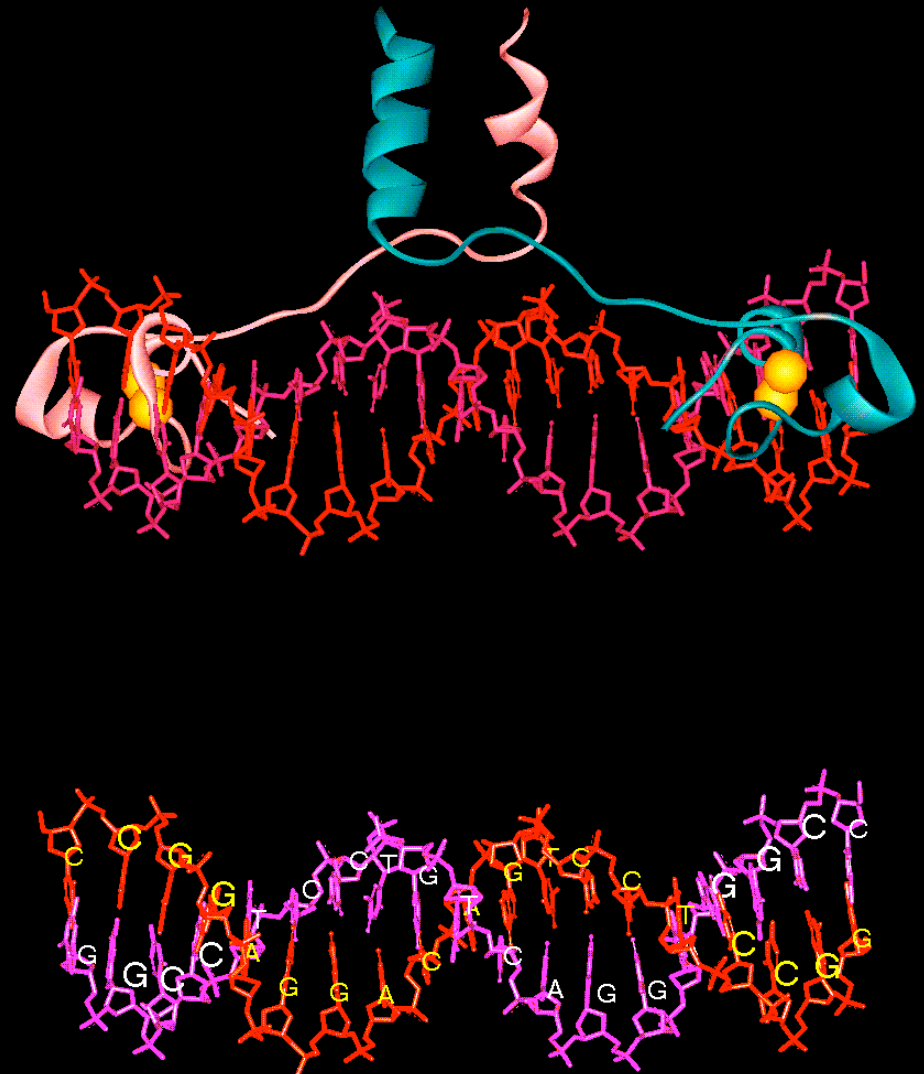


## Transcription factor-DNA interfaces

# Pho4p (yeast)



# Gal4p (yeast)



## The genome challenge



# *RNA polymerase*

