

Regulatory Sequence Analysis

***Regulatory regions
and regulatory elements***

Jacques van Helden
Jacques.van.Helden@ulb.ac.be

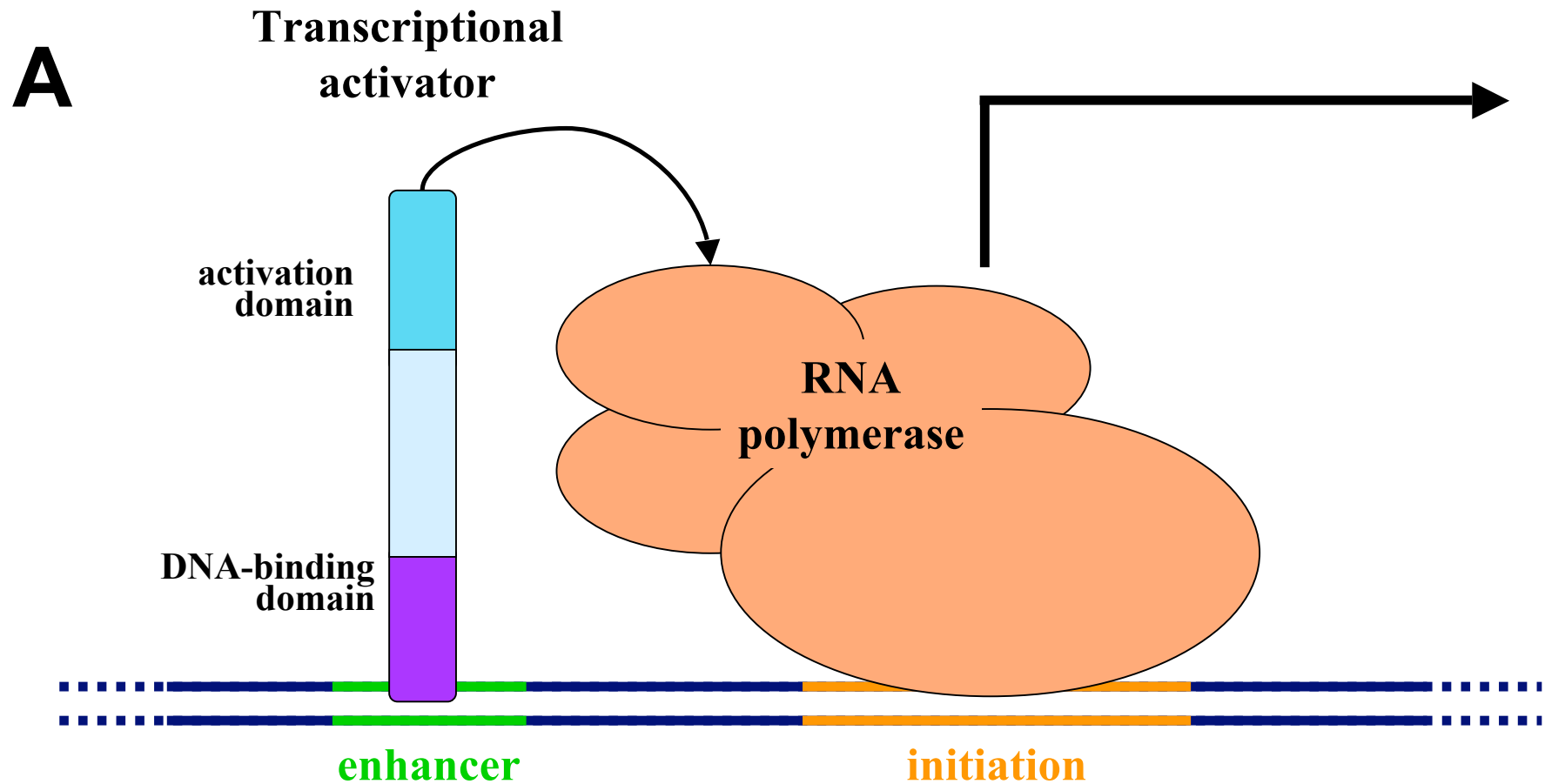
The non-coding genome

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

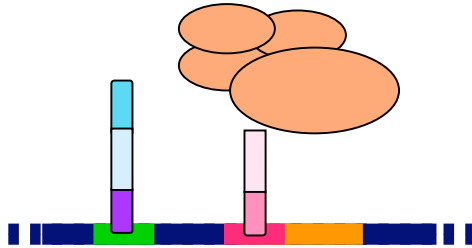
The genome challenge



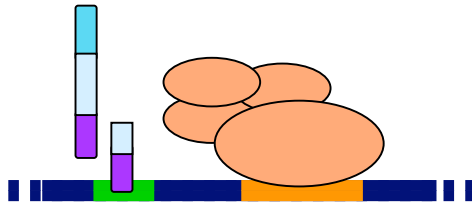
Transcriptional activation



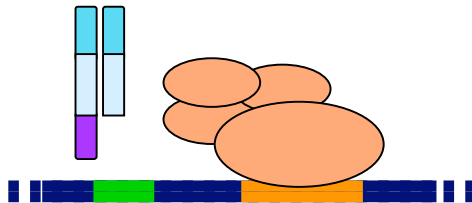
Transcriptional repression



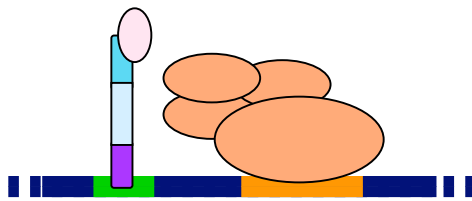
Prevent RNA polymerase from accessing DNA



Competition for factor binding site

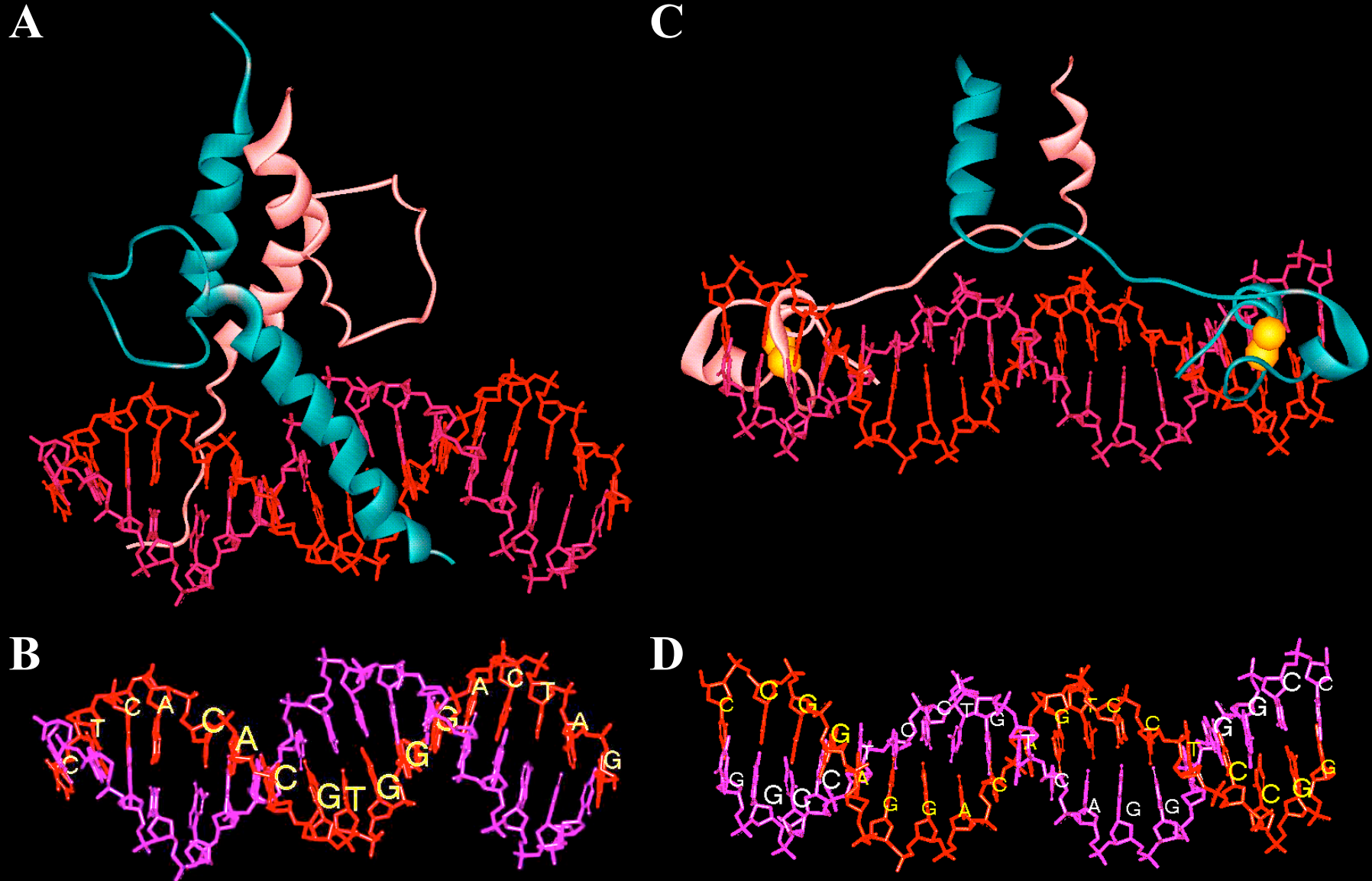


Factor titration

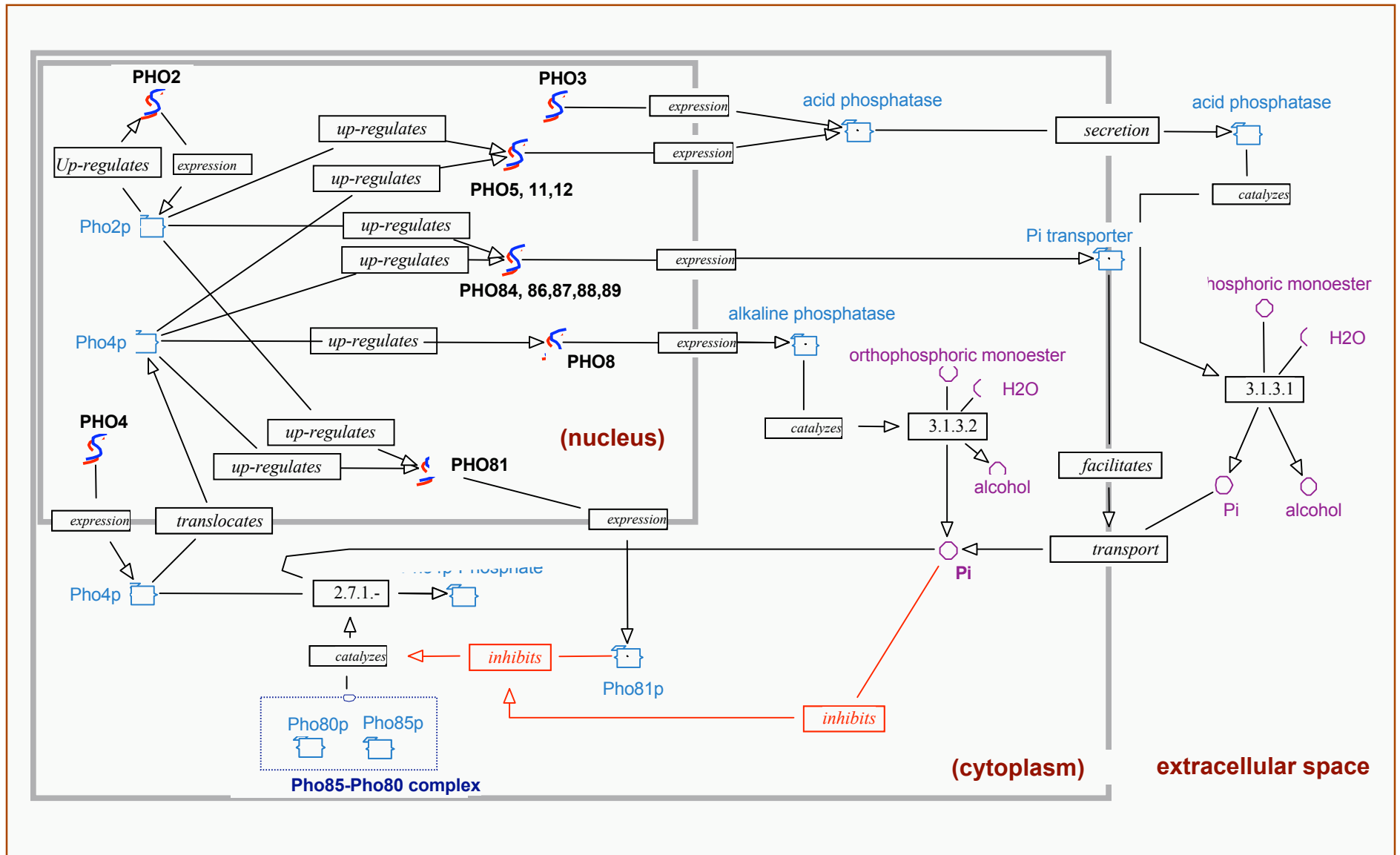


Prevent transcription factor from interacting with RNA-polymerase (bind with activation domain)

Transcription factor-DNA interfaces



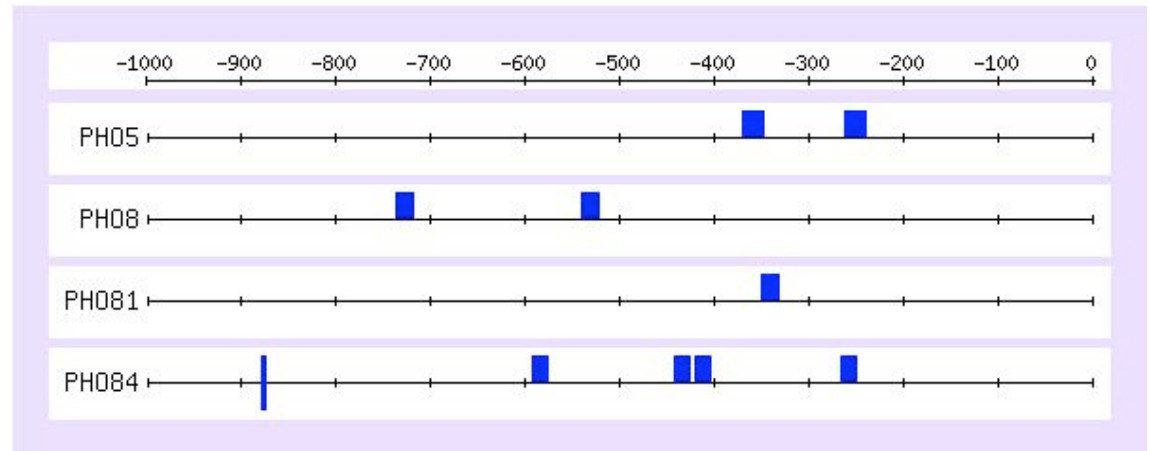
Phosphate utilisation in yeast



Transcription factor binding site (TFBS)

Gene	Ft_type	Factor	Strand	left	right	Sequence
PH05	site	Pho4p	D	-370	-347	TAAATTAG CACGTTTT CGCATAGA
PH05	site	Pho4p	D	-262	-239	TGGCACTCAC CACGTGGG ACTAGCA
PH08	site	Pho4p	D	-540	-522	TCGGGCCACGTG CAGCGAT
PH08	site	Pho4p	D	-736	-718	ATATTAAGCGTGCGGGTAA
PH081	site	Pho4p	D	-350	-332	TTAT GGCACGTG CGAATAA
PH084	site	Pho4p	D	-592	-575	TTACG CACGTTGGT GCTG
PH084	site	Pho4p	D	-421	-403	TTTCCAG CACGTGGGG C GG
PH084	site	Pho4p	D	-442	-425	TAGTT CCACGTGG ACGTG
PH084	site	Pho4p	DR	-879	-874	aaaagtgt CACGTG ataaaaaat
PH084	site	Pho4p	D	-267	-250	TT AAAAACGTG CGTATTA

- A **transcription factor binding site (TFBS)** is a location within a sequence.
- A site can be
 - experimentally characterized (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 transcription factor Pho4p

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

Gene	Site Name	Sequence	Affinity
PHO5	UASp2	---aCtCaCACACGTGGGACTAGC-	high
PHO84	Site D	---TTTCCAGCACGTGGGCGGA--	high
PHO81	UAS	----TTATGACACGTGCGAATAA--	high
PHO8	Proximal	GTGATCGCTGACGTGGCCCGA---	high
PHO5	UASp3	--TAATTTGCGATGTGCGATCTC--	low
PHO84	Site C	-----ACGTCCACGTGGAACATAT--	low
PHO84	Site A	-----TTTATCACGTGACACTTTTT	low
group 1	consensus	-----gCACGTGggac-----	high-low
PHO5	UASp1	--TAAATTAGCACGTTTTCGC----	medium
PHO84	Site E	-----AATACGCACGTTTTAATCTA	medium
PHO84	Site B	-----TTACGCACGTGGGTGCTG--	low
PHO8	Distal	---TTACCCGCACGCTTAATAT---	low
group 2	consensus	-----cgCACGTTt-----	med-low
Degenerate consensus		-----GCACGTTKKk-----	

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

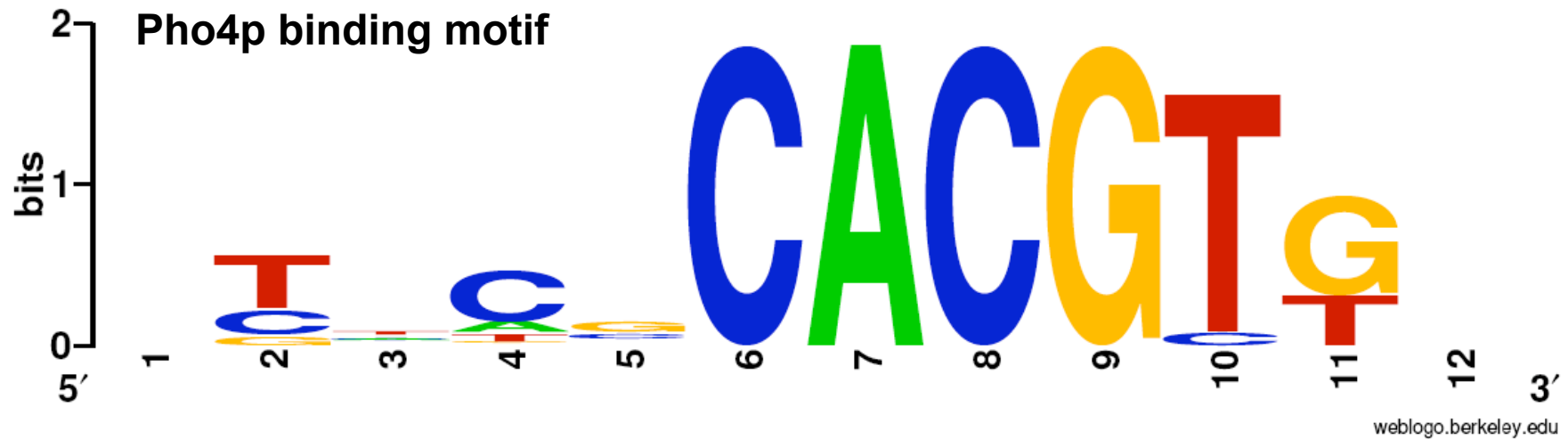
Regulatory sites : matrix description

Position-specific scoring matrix (PSSM)

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

Binding motif for the yeast Pho4p transcription factor
(Source : Transfac matrix F\$PHO4_01)

Sequence logo



Sequence logo

A A T G T A T G G

Rap1

G G T G G C A A A A

Rpn4

A A A T G A G T C A

Gcn4

G A A T T C A G A A

HSE

T G G G G G T A G G

Mig1

A A T T C A C G T G

Cbf1

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.

- String-based descriptions

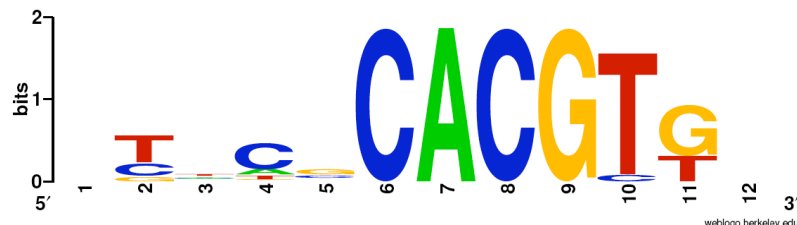
- nucleotide alphabet **CACGTGGG**
- IUPAC alphabet **CACGTGKK**
- regular expressions. **CACGTG[GT][GT]**

- Probabilistic descriptions

- Position-specific scoring matrix (PSSM)

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

- Hidden Markov Model (HMM) (not treated here)
- Logo representation (Schneider, 1986)



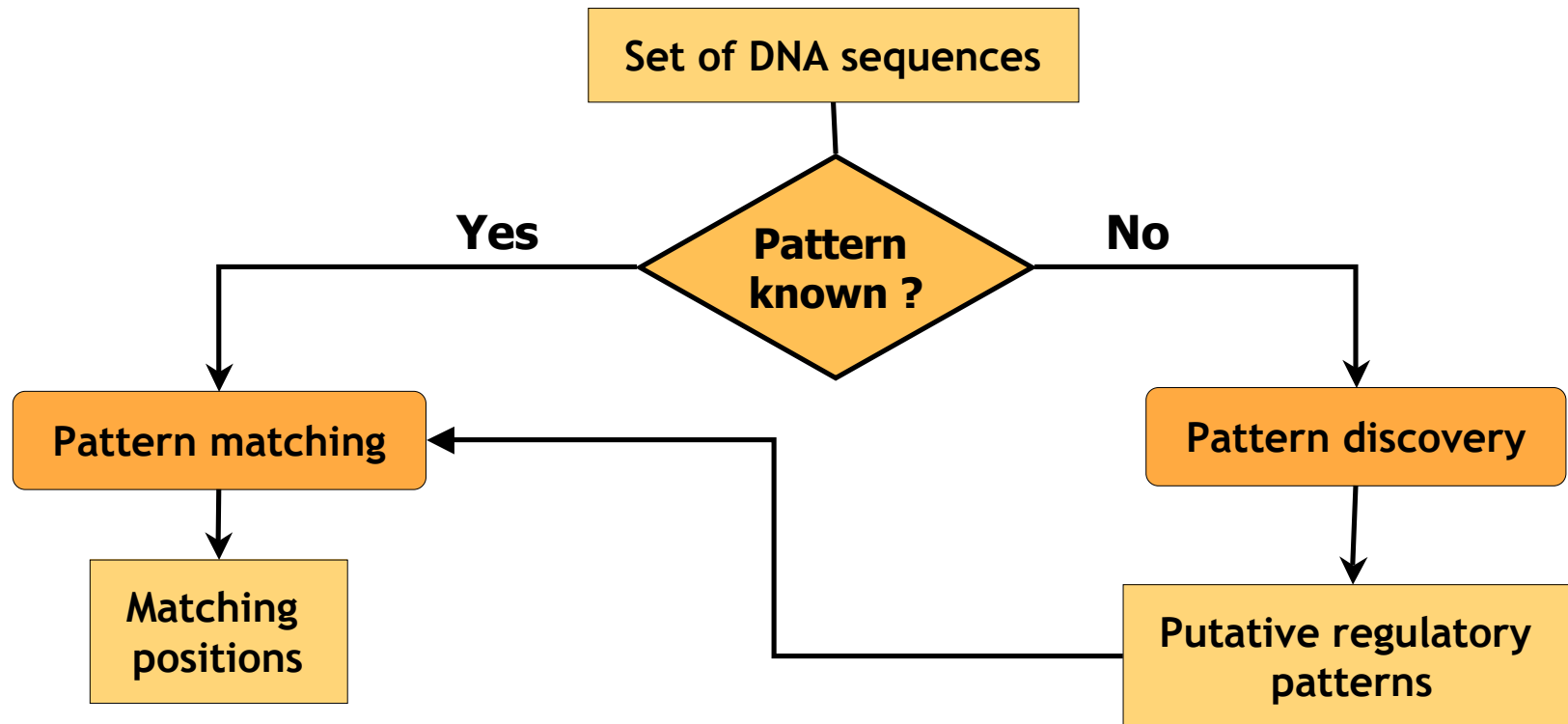
Characteristics of yeast regulatory sites

- Located upstream the regulated gene
- 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)
- Strand-insensitive
- Within 800 bp from the start codon
- Efficiency does not depend on
 - strand
 - position

Differences between species

organism	coli	yeast	higher organisms
location	upstream overlap. Initiation	upstream	upstream downstream within introns
distance range	-400 to +50 bp	-800 to -1 bp	over 100s of Kb
position effect	often essential	often irrelevant	often irrelevant
strand	sensitive or symmetric	insensitive	insensitive
most common core	spaced pair of 3nt	~5-8 conserved bp	~5-8 conserved bp
repeated sites	rare	occasional	frequent
composite elements			frequent

Pattern matching vs pattern discovery



Questions and approaches

1. If we know the consensus for a given transcription factor, can we predict its binding sites in a DNA sequence ?
 - **Pattern matching** against a sequence
2. Can we scan a sequence for matches with the consensus of all the currently known transcription factor ?
 - **Matching a library** of patterns against a sequence
3. Starting from a set of co-expressed genes, can we predict cis-acting elements involved in their transcriptional regulation ?
 - **Pattern discovery** within a sequence set
4. Can we detect regulatory signals by searching conserved elements in non-coding sequences of orthologous genes ?
 - **Phylogenetic footprinting**
5. Can we classify genes on the basis of the presence of regulatory motifs in their regulatory regions ?
 - **Gene classification** on the basis of pattern scores
 - **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.

Regulatory Sequence Analysis

Supplementary material

Jacques van Helden
Jacques.van.Helden@ulb.ac.be

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 TCACGTG TAATTT
MET3	-384	-366	AAAAGG TCACGTG ACCAGA
MET14	-235	-217	CTAATTT TCACGTG ATCAAT
MET16	-185	-167	ATCATT TCACGTG GCTAGT
ECM17	-311	-293	ATTTCA TCACGTG CGTATT
ECM17	-339	-321	.TTTGTC TCACGTG ATATTT
MET10	-255	-237	.CCACAC TCACGTG AGCTTAT
MET10	-237	-219	.TAGAAG TCACGTG ACCACAA
MET2	-360	-342	GTATTT TCACGTG ATGCGC
MET2	-554	-536	TAATAA TCACGTG ATATTT
MET17	-306	-288	.AAATGG TCACGTG AAGCTGT
MET17	-332	-314	TTGAGG TCACATG ATCGCA
MET6	-540	-522	GCCACAT TCACGTG CACATT
MET6	-502	-484	AATATTT TCACGTG ACTTAC
SAM2	-329	-311	.TCTACC TCACGTG ACTATAA
SAM2	-381	-363	.TCTTCA TCATGTG 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 AAAAA ATGTGGCAATGG
MET2	-313	-293	TGC AAAAA ATGTGGATGCAC
MET17	-227	-207	TCATG AA AACTGTGTAAACATA
MET6	-313	-293	GTCGC AA AACTGTGGTAGTCA
SAM2	-306	-286	GCTTG AA AACTGTGGCGTTTT
SAM1	-283	-263	ACAGG AA AACTGTGGTGGCGC
MET19	-173	-153	ATAAGC AA ACTGTGGTTCAT
MUP3	-188	-168	CGG AAAAA ACTGTGGCGTCGC
MET8	-184	-164	GG AAAAA AAATGTGAAAATCG
MET1	-232	-212	CATAAT AA ACTGTGAACGGAC
MET3	-259	-239	ACAAAG CCACAGTTTT ACAAC
MET28	-159	-139	CTAAC CCACAGTTTT GGGCG
MET8	-434	-414	TCTTGT CCGCAGTTTT ATCTG
MET30	-168	-148	GGGAAG CCACAGTTT GCGCGG
MET6	-405	-385	CTATCGA ACTCGTTTT 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 CACGTGGG ACTA
PHO5	-260	-245	..GCACTCA CACGTGGG A
PHO5	-262	-239	TGGCACTCA CACGTGGG ACTAGCA
PHO8	-540	-522	...TCGGGC CACGTGC AGCGAT
PHO8	-736	-718	..ttacccg CACG <u>C</u> TT aatat
PHO81	-350	-332	...TTATGG CACGTGCG AATAA
PHO84	-421	-403	..TTTCCAG CACGTGGG GCGG
PHO84	-442	-425	...TAGTTC CACGTGG ACGTG
PHO84	-879	-874	.aaaagtgt CACGTG ataaaaat
PHO84	-267	-250	..taatacg CACGTTTT aa
PHO84	-592	-575TTACG CACGTT GGTGCTG
PHO5	-368	-349	...AATTAG CACGTTTT CGCATA
PHO5	(?)	(?)	..AAATTAG CACGTTTT CGC
PHO5	-370	-347	.TAAATTAG CACGTTTT CGCATAGA

IUPAC ambiguous nucleotide code

A	A	A denine
C	C	C ytosine
G	G	G uanine
T	T	T hymine
R	A or G	pu R ine
Y	C or T	p Y rimidine
W	A or T	W weak hydrogen bonding
S	G or C	S trong hydrogen bonding
M	A or C	a M ino group at common position
K	G or T	K eto 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	a N y

Pho4p binding specificity - matrix descriptions

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

[illegible]