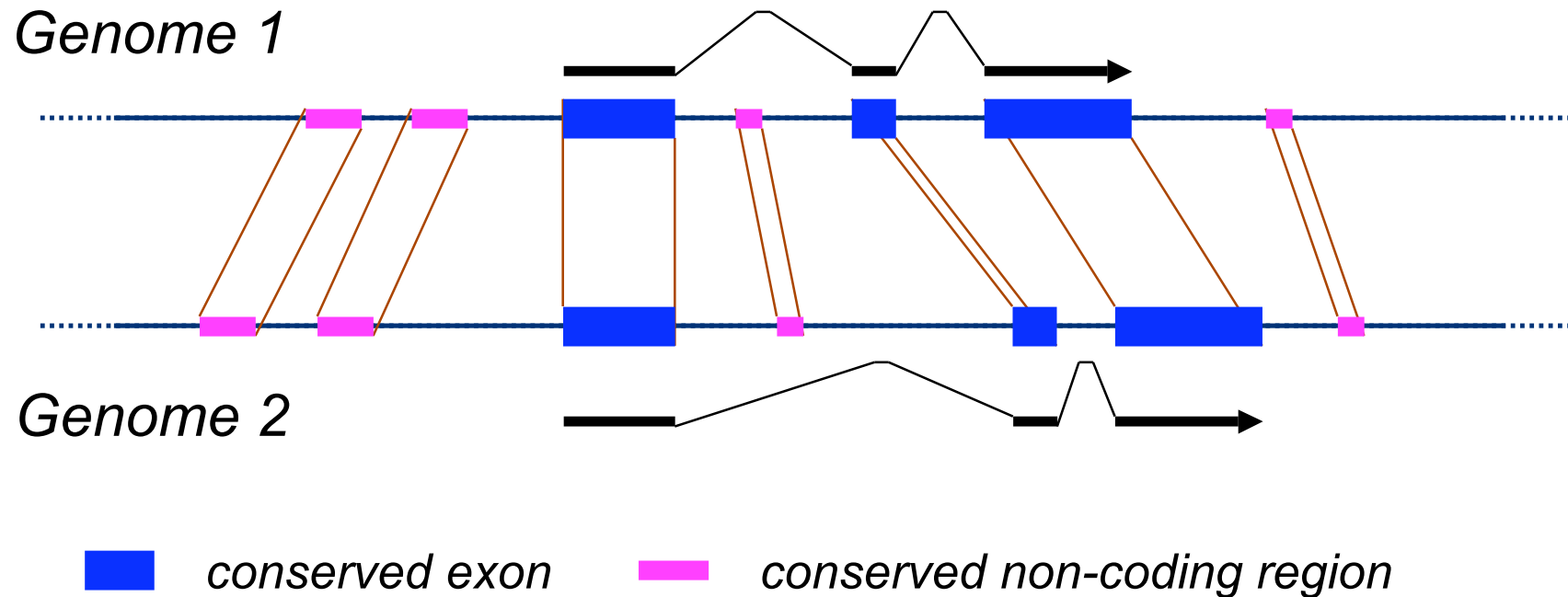


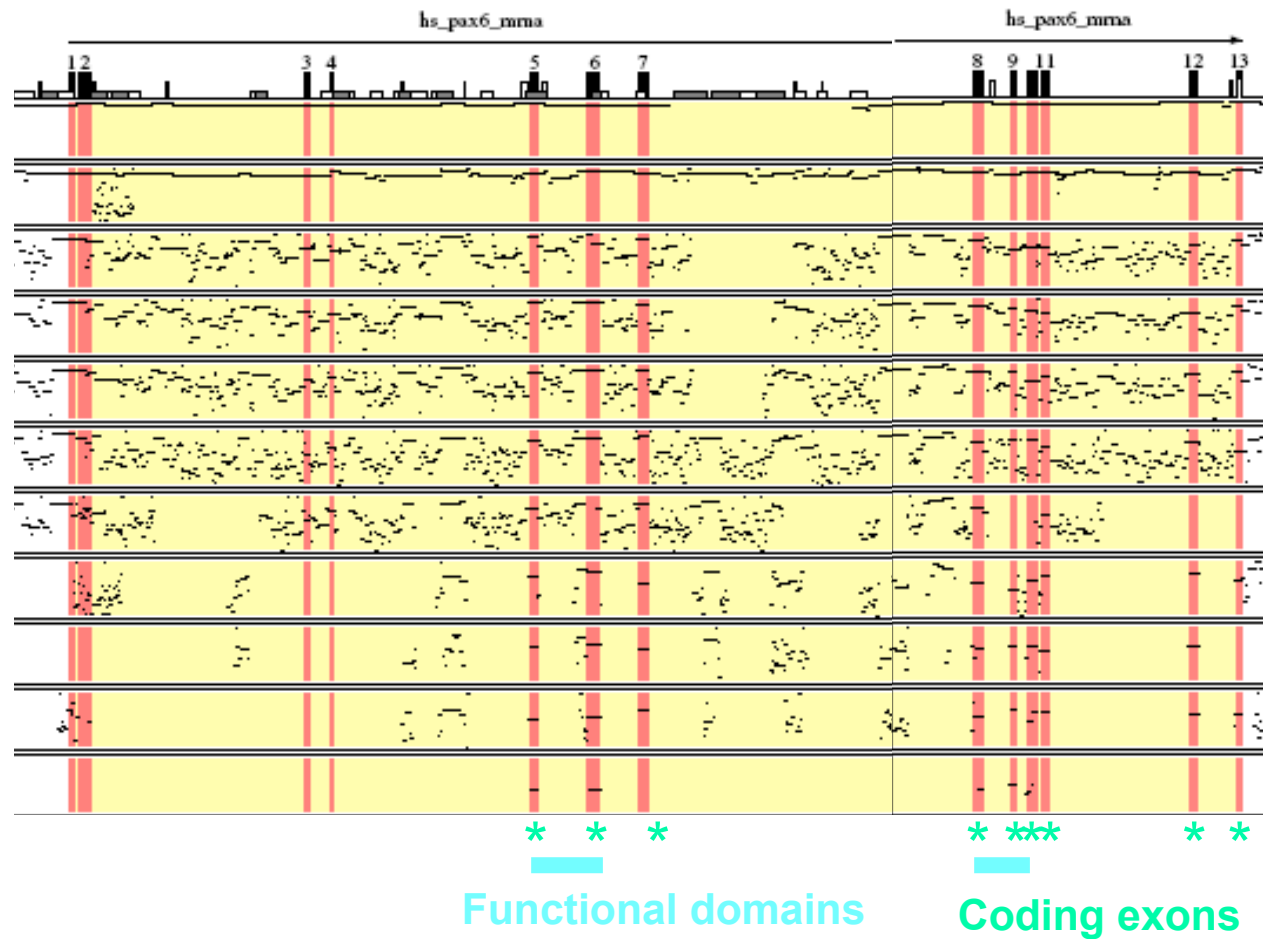
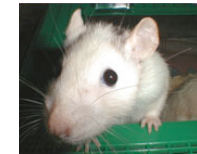
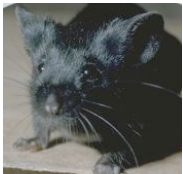
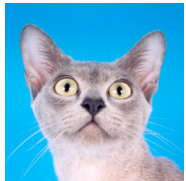
***Applying comparative genomics
to detect cis-acting elements***

Phylogenetic footprinting to define regulatory regions

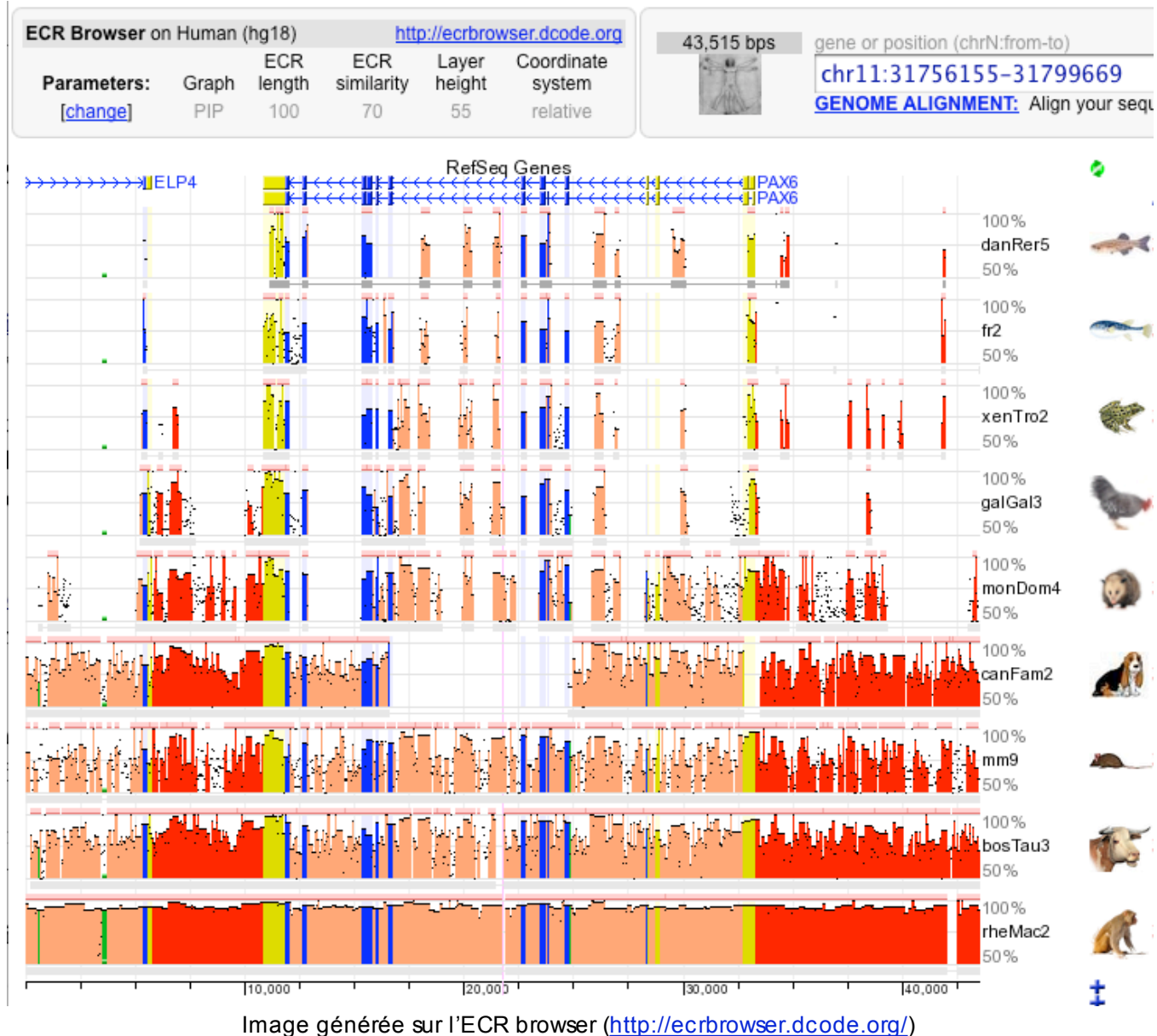


- Within non-coding sequences, regulatory elements evolve slower than their surrounding.
- Conserved non-coding sequences contain a high concentration in regulatory elements.

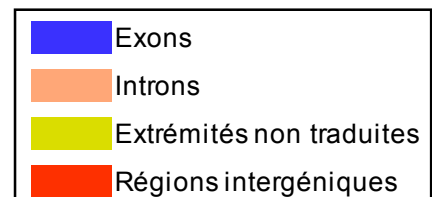
Phylogenetic footprints for the *pax6* gene



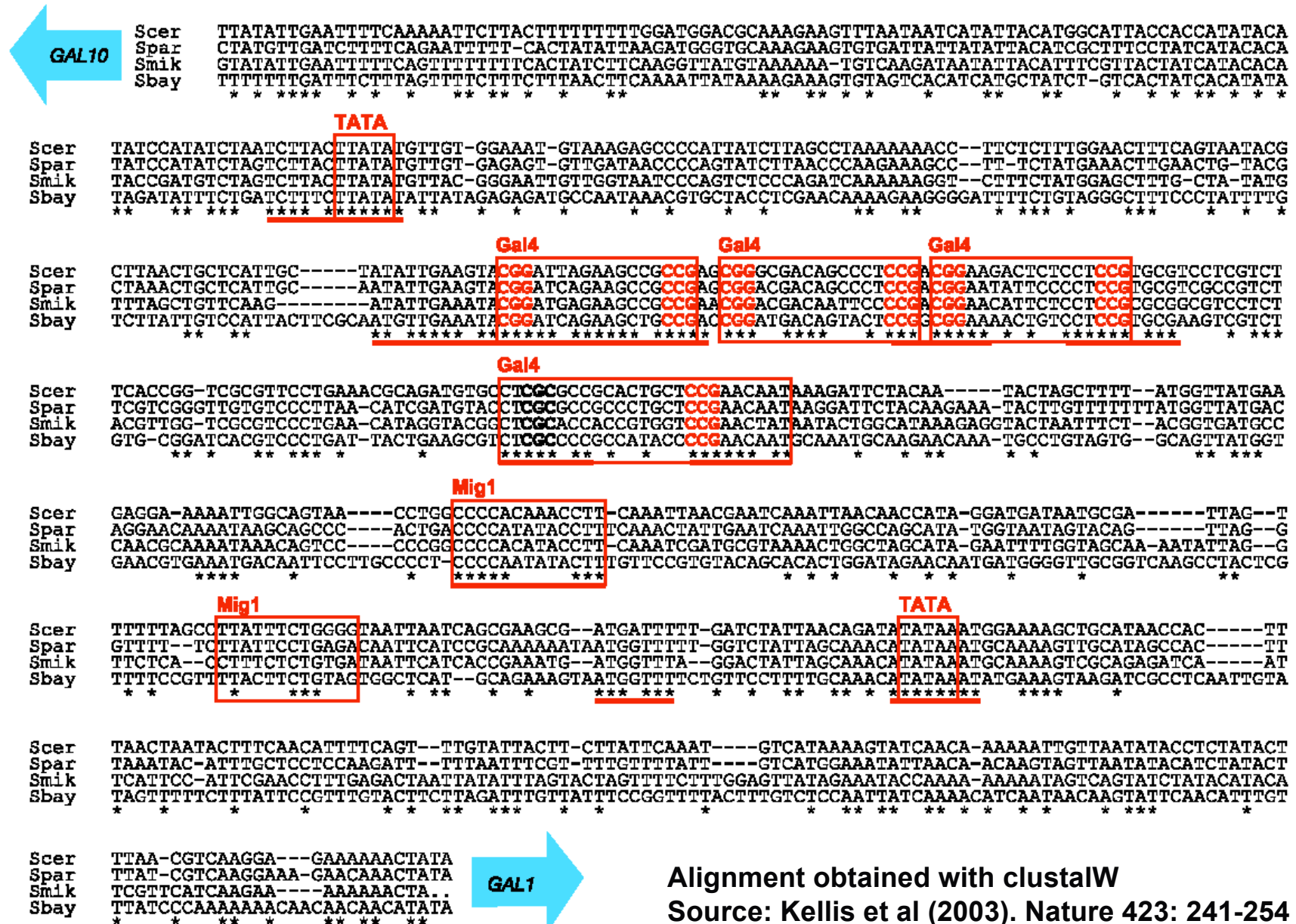
Pourcentages de positions identiques (PIP) dans la région chromosomique de Pax6



- La génomique comparative permet d'améliorer la localisation des gènes.
- Alignement de la région génomique contenant le gène Pax6, entre le génome humain et une série d'organismes de plus en plus distants évolutivement (de bas en haut).
- Les blocs de séquences conservées reflètent souvent la présence de fragments codants.
- Cependant, il existe également des segments conservés dans les régions non-codantes.



Global alignment of intergenic regions



Another alignment in the same genomes

```
GAL80 (YML051W) upstream regions
Scer      ATGGCGCAAGTTTTCCGCTTTGTAATATATATTTATACCCCTTTCTTCTCTCCCCTGCAA
Spar      AGGGGCCAAAGCTCCCGCTCTGTAAAATATATTTATATCCCTTCCTTCTCTCCCCTGCAA
Smik      TAGGGACAAAGCCCGCCTTTTGTAAATATACTTATACCTCTCCTTCTCTCCCCTGCAA
Sbay      .....
          **   ***           *  *  *****  *****  *****  **   *  *****
          .....

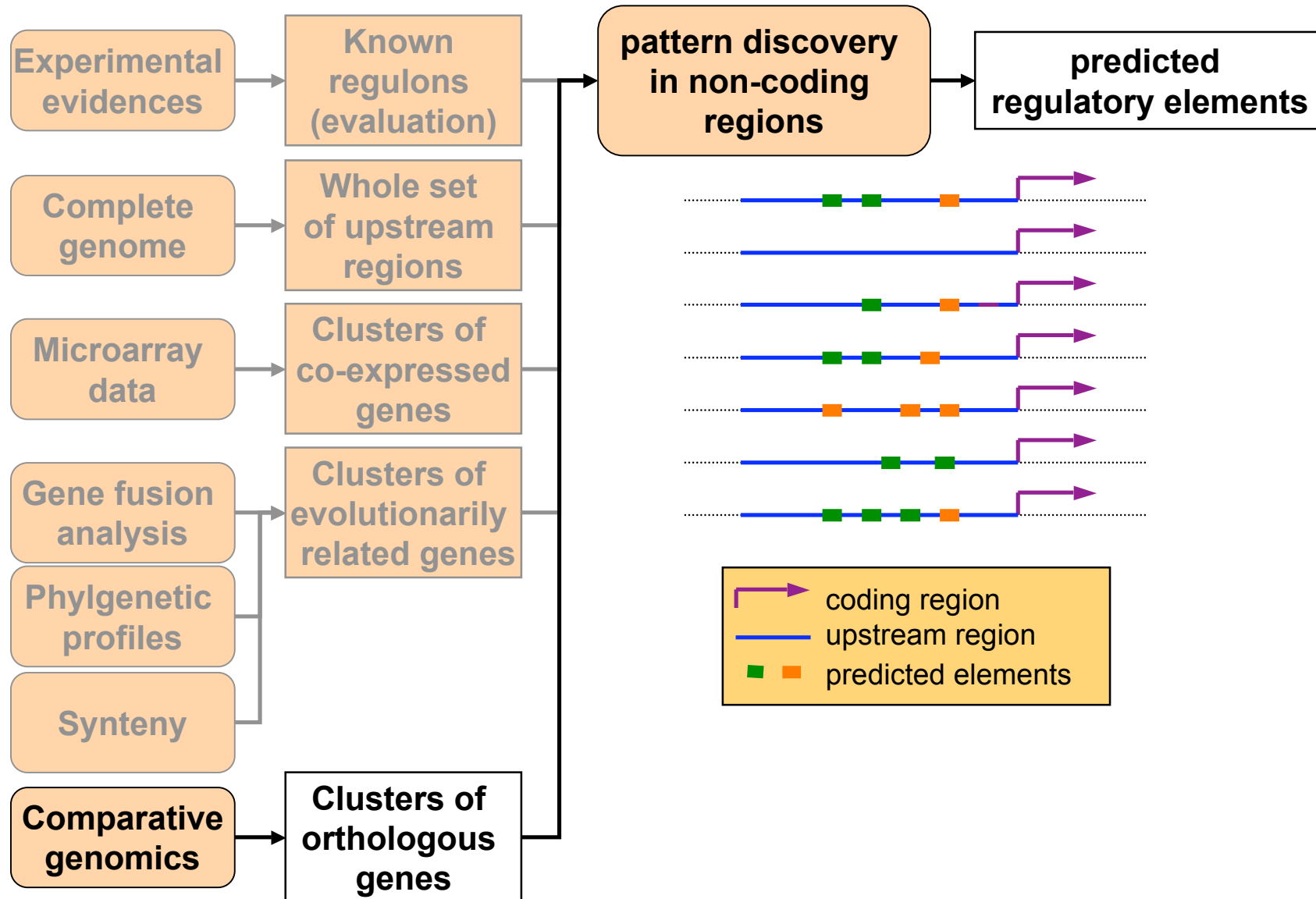
Scer      TATAATAGTTTAAATTCTAATATTAATAATA---TCCTATATTTTCTTCATTTACGGGCGC
Spar      TATAATAGTTTAAATTCTAATATTAATAATA---TCCTATATTTTCTTACC-ACGGGCGC
Smik      CATAATAGTTAACTCCTAATATTAATAATAATATCCTACAATTTCTTAGC-ACGGGGGC
Sbay      .....
          ***** *  *  ***** ***** ***** *  ***** *  *  ***** **
          .....

Scer      ACTCTCGCCCGAACGACCTCAAAATGTCTGCTACATTCATAATAACCAAAGCTCATAAC
Spar      ACTCTCGCCCGAACGACCTCAAAATGCTTGCTACATTCATAATAATCAAAGCTTATAAC
Smik      ACTCTCGCCCGAACGACCTCAAAACGCTTGCTACATCCATAATATTCAGAACTACATCAC
Sbay      .....
          ***** ***** *  ***** ***** ***** ** **      ** **
          .....

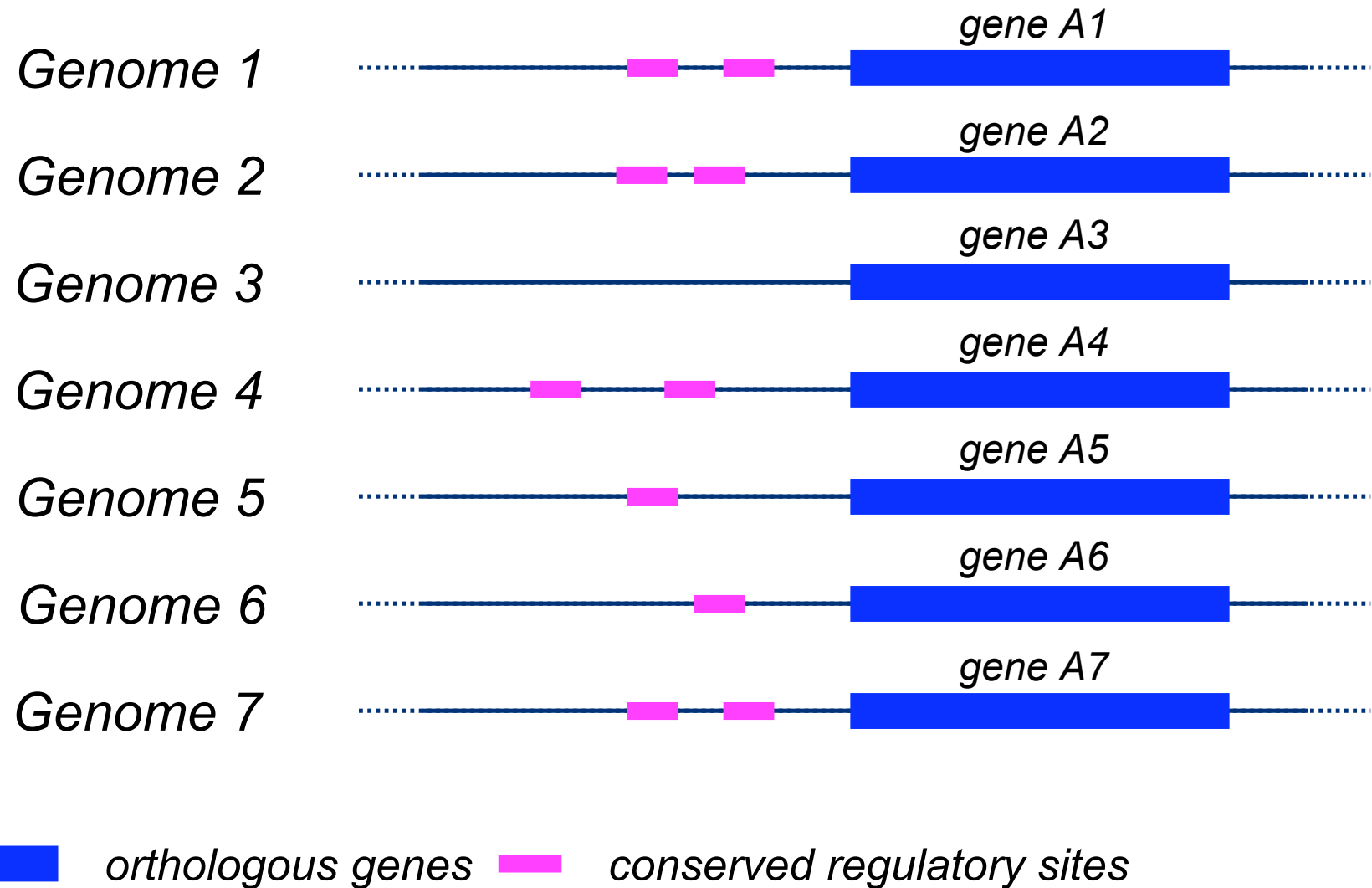
Scer      TTTTTTTTTT----TGAACCTGAATATATATACATCACATATCACTGCTGGTCCTTGCCGA
Spar      TTTTTTTTTTCTTTGTACCTGAATATATATACATCTCATGTCACTGCTGGTCCTTGCCGG
Smik      TTTTTTTTTT-----GTACATAAAAATATATAC--CACATGTCACTGCTGATCCTTGCTGA
Sbay      .....
          *****          *  *  *  *  *****  *  *  ***** ***** *
          .....

Scer      CCAGCGTATACAATCTCGATAGTTGGTTT-C-CCGTTCTTTCCACTCCCGTCATGGACTA
Spar      CCAGCGTATACAACCTCGATAGCTGGTTTTT-C-CCGTTCTTCCACTCCTGTCATGGACTA
Smik      CGAGCGTATACAAGCTCGATAGCTGGTCTTTACCGTGCCATTCCCTGCCGTCATGGACTA
Sbay      .....
          *  ***** ***** ***** *  ***** *  *  *  *****
```

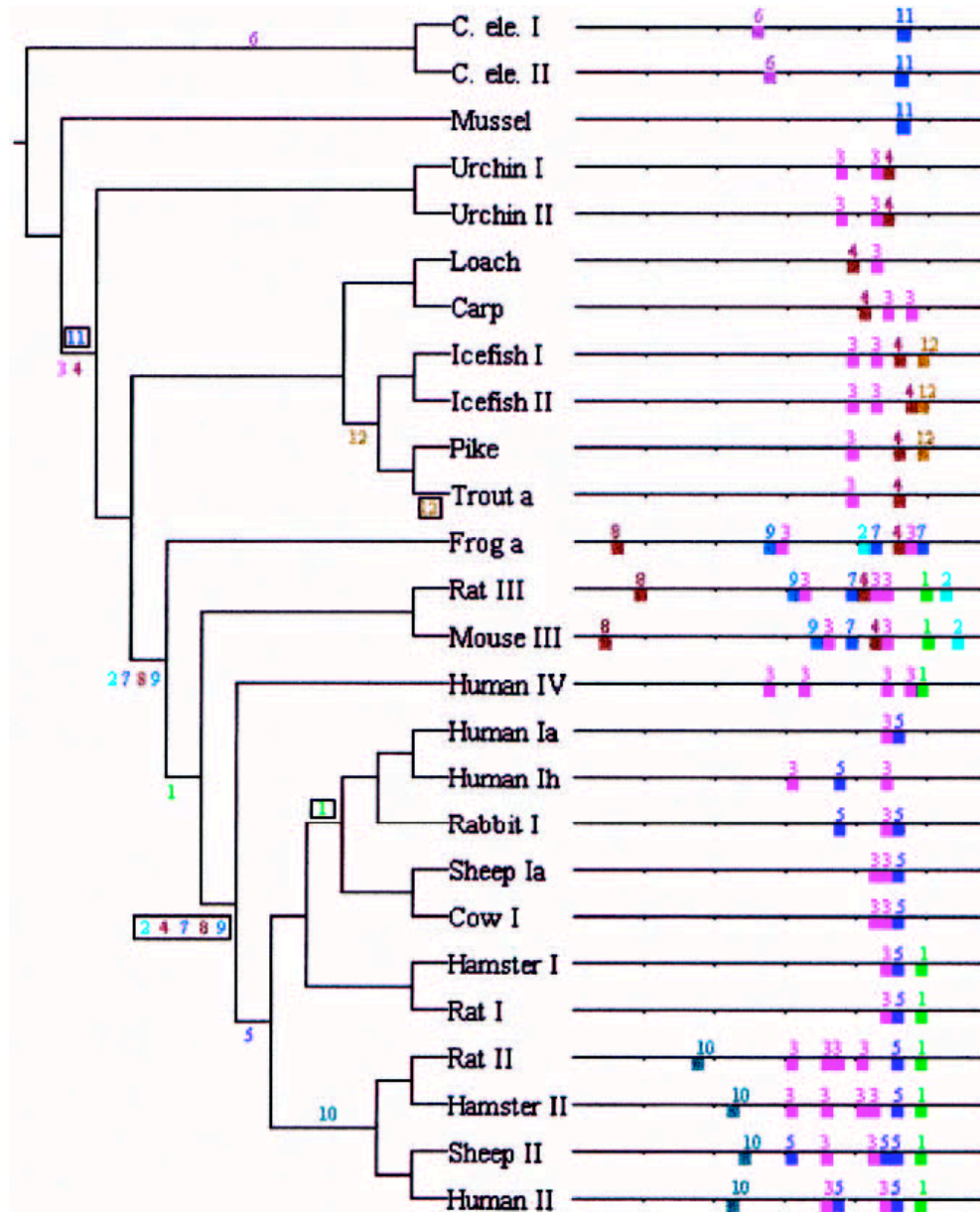
Motifs in clusters of orthologous genes (COGs)



Phylogenetic footprinting to predict regulatory sites



Footprinter example metallothionein

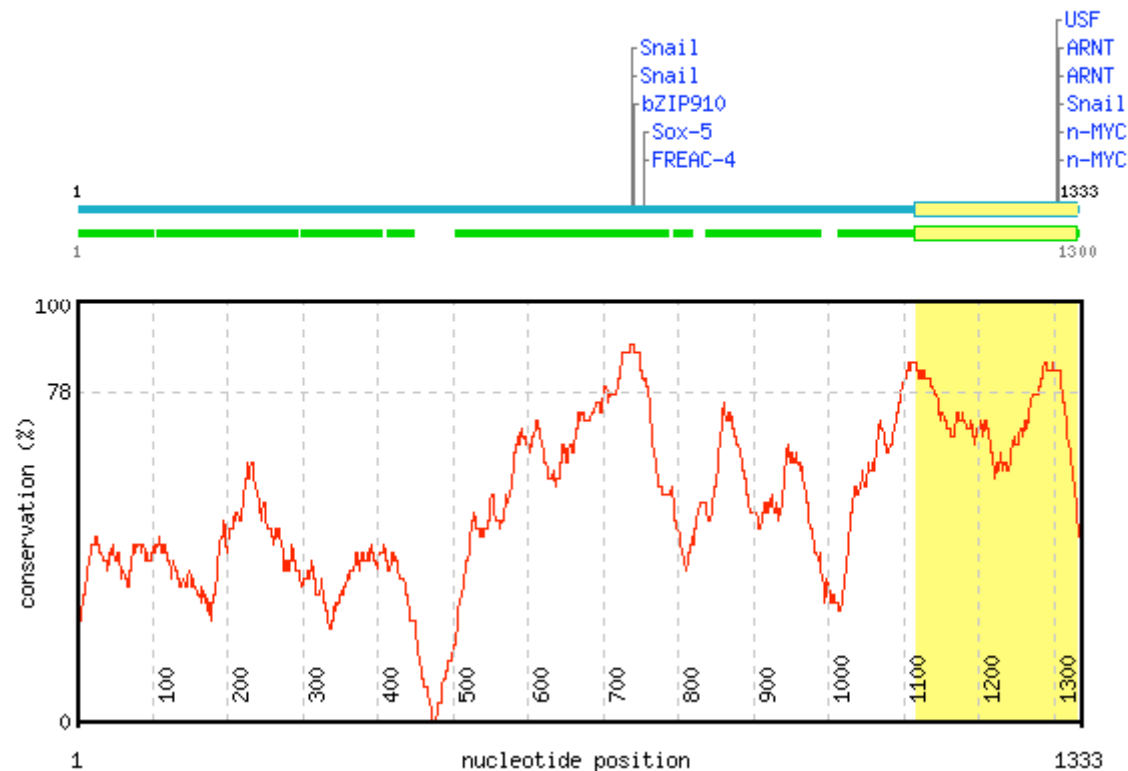


- 590 bp upstream of the same gene (methallothionein) in different species.
- 12 highly conserved motifs are detected.
- Each motif can be associated to a given internal node of the phylogenetic tree.

Cross-matches in promoters of orthologous genes

- Lenhard et al. (2003). J.Biology 2:13.
- 100 PSSM for known mammal transcription factors
- Searching for conserved matches in Human and mouse increases the selectivity by 85%.
- **Consite:** <http://mordor.cgb.ki.se/cgi-bin/CONSITE/consite/>

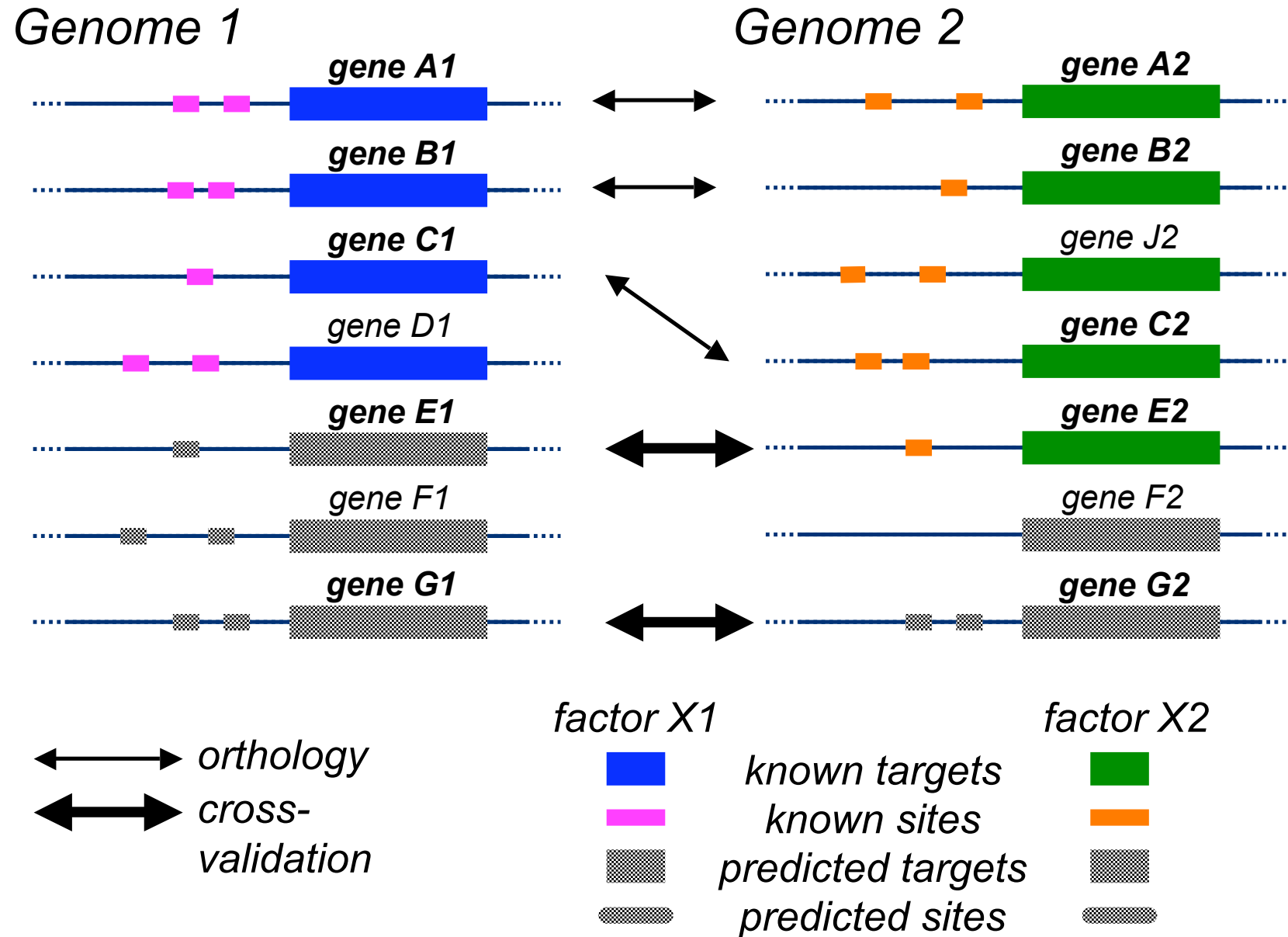
Conservation profile of *Human_IR*



Cross-validation of genome-scale pattern matching

- Genome-scale pattern matching raises many false positive
- Cross-validation :
 - ▣ gene *A* from genome *X* has a good match in its upstream sequence
 - ▣ ortholog *A'* from genome *Y* has a good match in its upstream sequence

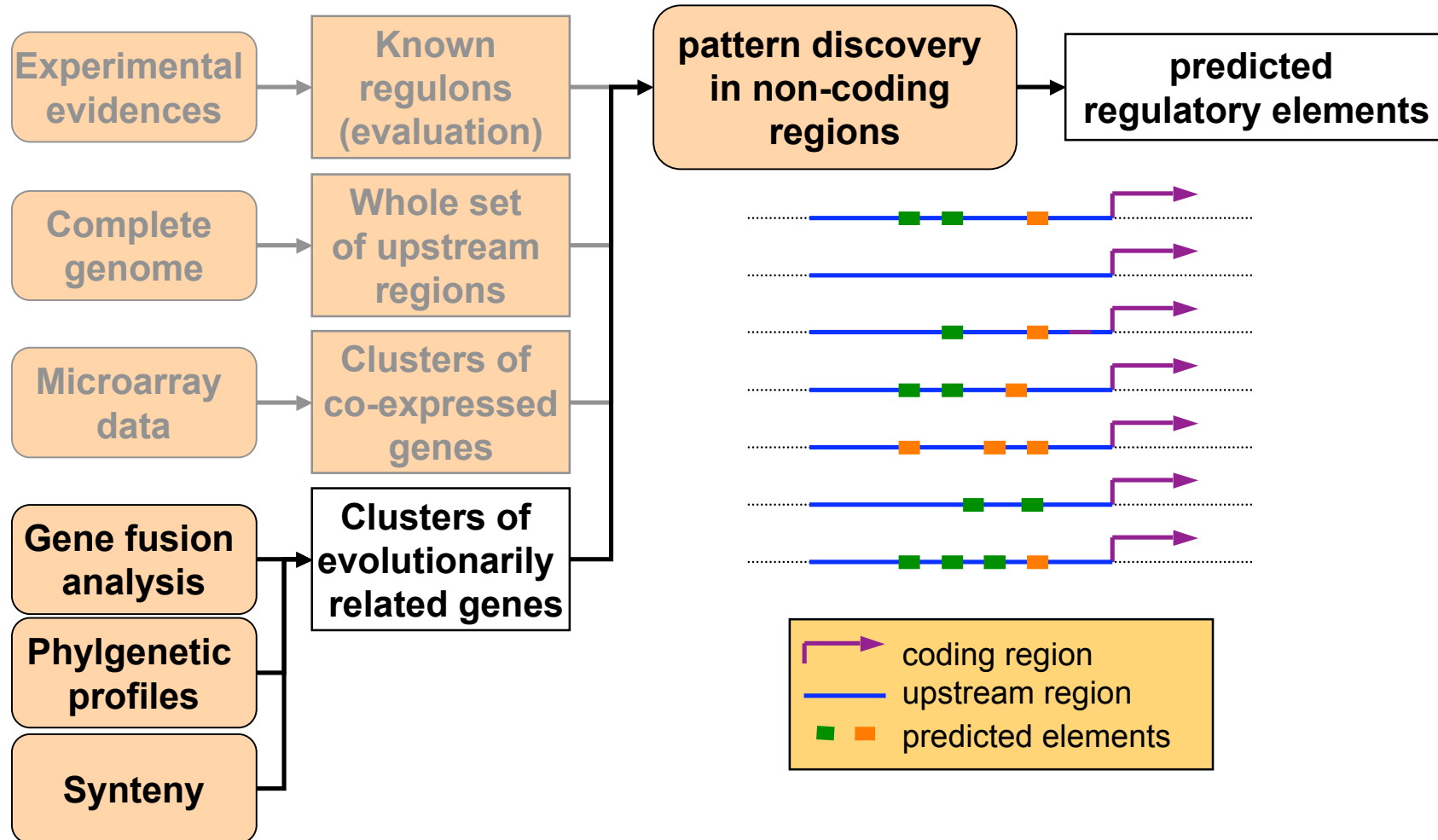
Cross-validation of pattern matching



Detection of functional clusters of genes

- Various methods of comparative genomics allows to detect clusters of functionally related genes
 - Operon conservation
 - Gene fusion analysis
 - Phylogenetic profiles (synteny)
- These functional clusters can be used to discover regulatory motifs in their upstream regions.

Clusters predicted from comparative genomics



Pattern discovery in predicted regulons

Organism	Cluster	pattern	reverse_complement	score
<i>Escherichia_coli_K12</i>	EC21	ccccctcaccctctt	aagagggtgaggggg	13.01
<i>Escherichia_coli_K12</i>	EC21	ccccctcaccctt	aaggggtgaggggg	13.01
<i>Escherichia_coli_K12</i>	EC21	gccctcaccctc	gaggggtgagggc	13.01
<i>Escherichia_coli_K12</i>	EC21	ggggagaggggtgagggga	tcccctcaccctctcccc	13.01
<i>Escherichia_coli_K12</i>	EC21	cctcaccctcaccctctcccctc	gaggggagaggggtgaggggtgagg	13.01
<i>Escherichia_coli_K12</i>	EC3	ccccctcgcccctt	aaggggcgaggggg	12.73
<i>Escherichia_coli_K12</i>	EC3	aagggcgaggggg	ccccctcgcccctt	12.73
<i>Escherichia_coli_K12</i>	EC3	gccctcgcccctc	gaggggcgagggc	12.73
<i>Escherichia_coli_K12</i>	EC3	ccccctcaccctt	aaggggtgaggggg	12.73
<i>Escherichia_coli_K12</i>	EC3	ccccctctcccctt	aaggggagaggggg	12.73
<i>Mycoplasma_pneumoniae</i>	MP1	tataatact	agtattata	11.75
<i>Mycoplasma_pneumoniae</i>	MP1	cttaataactaat	attagtattaag	11.75
<i>Escherichia_coli_K12</i>	EC17	cccctctccctt	aaggagagagggg	10.63
<i>Escherichia_coli_K12</i>	EC17	cccctctcccctt	aaggggagaggggg	10.63
<i>Escherichia_coli_K12</i>	EC17	cccctcgcccctt	aagggcgaggggg	10.63
<i>Mycoplasma_pneumoniae</i>	MP1	aataataag	cttattatt	10.4
<i>Mycoplasma_pneumoniae</i>	MP1	aataatattatt	aataatattatt	10.4
<i>Mycoplasma_pneumoniae</i>	MP1	taataataagnnnnnaataa	ttatnnnnncttattatta	10.4
<i>Mycoplasma_pneumoniae</i>	MP1	cttagtattatt	aataataactaag	10.4
<i>Mycoplasma_pneumoniae</i>	MP1	taataataagnnnnnaataa	ttatnnnnncttattatta	10.4
<i>Mycoplasma_pneumoniae</i>	MP1	aataatattaaga	tcttaatattatt	10.4
<i>Mycoplasma_pneumoniae</i>	MP1	cttagtatatataatataactaag	cttagtatatattatataactaag	10.4
<i>Mycoplasma_pneumoniae</i>	MP1	taataataagnnnnnaataa	ttatnnnnncttattatta	10.4
<i>Mycoplasma_pneumoniae</i>	MP1	ctaataattatt	aataatattag	10.4
<i>Mycoplasma_pneumoniae</i>	MP1	taataataagnnnnnaataa	ttatnnnnncttattatta	10.4
<i>Mycoplasma_pneumoniae</i>	MP1	aataatattatnnngtactattataataag	cttattataatagtacnnnaataatattatt	10.4
<i>Mycoplasma_pneumoniae</i>	MP1	aataatattatc	gataatattatt	10.4

gene clusters from McGuire et al. (2000). *Genome Res* 10(6), 744-57.

Phylogenetic footprinting resources

- CORG: a database for COmparative Regulatory Genomics
 - Dieterich et al. (2003), Nucleic Acids Res. 31:55-57.
 - <http://corg.molgen.mpg.de>
 - Systematic alignment of 15Kb upstream regions for each pair of mouse-human homologous genes (18.674 pairs).
 - 10.793 significant alignments ($P < 0.001$), containing 293.503 conserved non-coding blocks (CNB), covering 8% of the upstream sequences (<http://corg.molgen.mpg.de/stats.html>).

Summary - phylogenetic approaches

- Matching conserved sites for known transcription factors
 - Consite: <http://mordor.cgb.ki.se/cgi-bin/CONSITE/consite/>
 - Lenhard et al. (2003). J.Biology 2:13.
- Global alignment of promoters of orthologous genes
 - clustalW
 - e.g.: Kellis et al (2003). Nature 423: 241-254.
- Pattern discovery in promoters of orthologous genes
 - Footprinter: <http://bio.cs.washington.edu/software.html>
 - Blanchette and Tompa (2002). Genome Research. 12, 739–748.