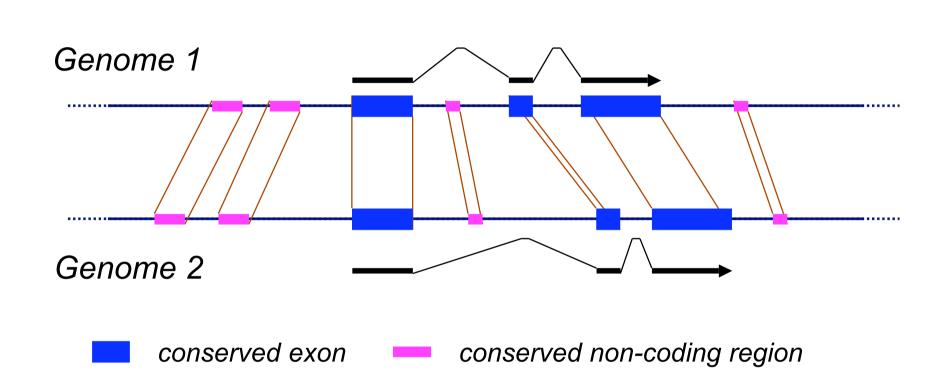
Regulatory Sequence Analysis

Applications of comparative genomics to the analysis of regulatory sequences

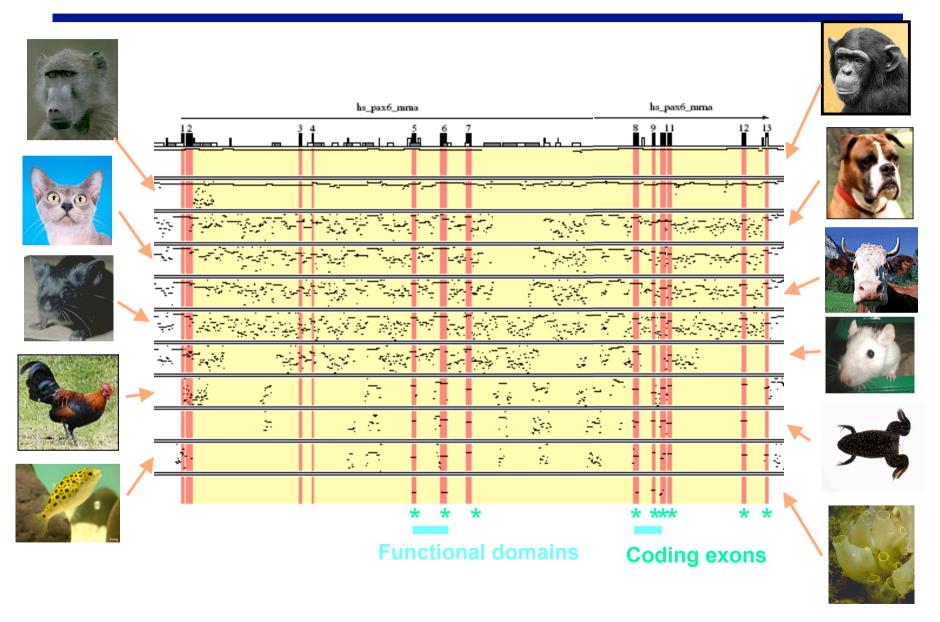
Jacques van Helden Jacques.van.Helden@ulb.ac.be

Phylogenetic footprinting to define regulatory regions

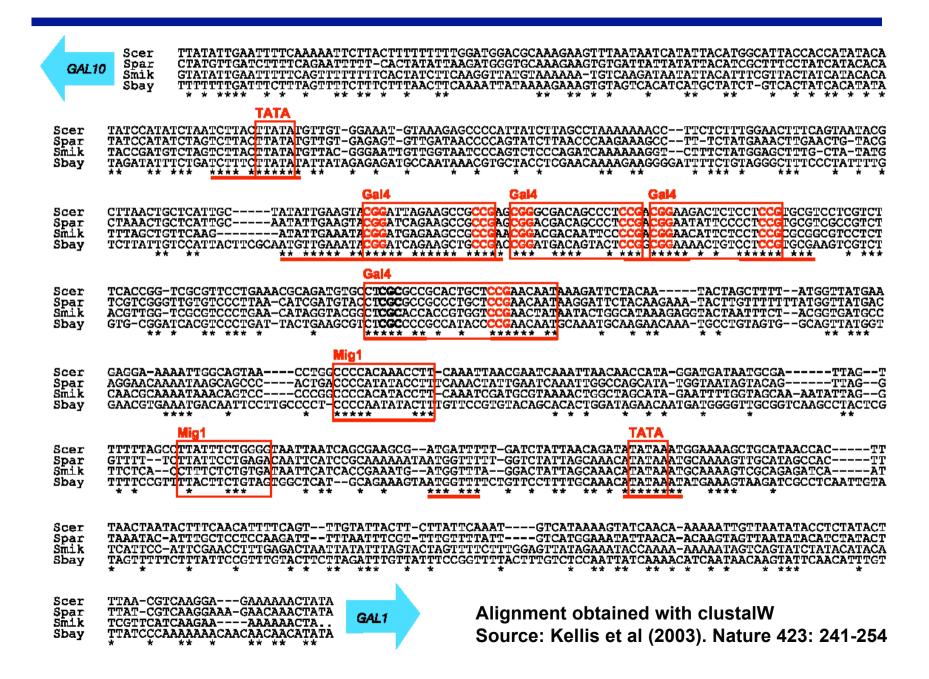


- Within non-coding sequences, regulatory elements evolve slower than their surrounding.
- Conserved non-coding sequences contain a high concentration in regulatory elements.

Phylogenetic footprints for the pax6 gene



Global alignment of intergenic regions

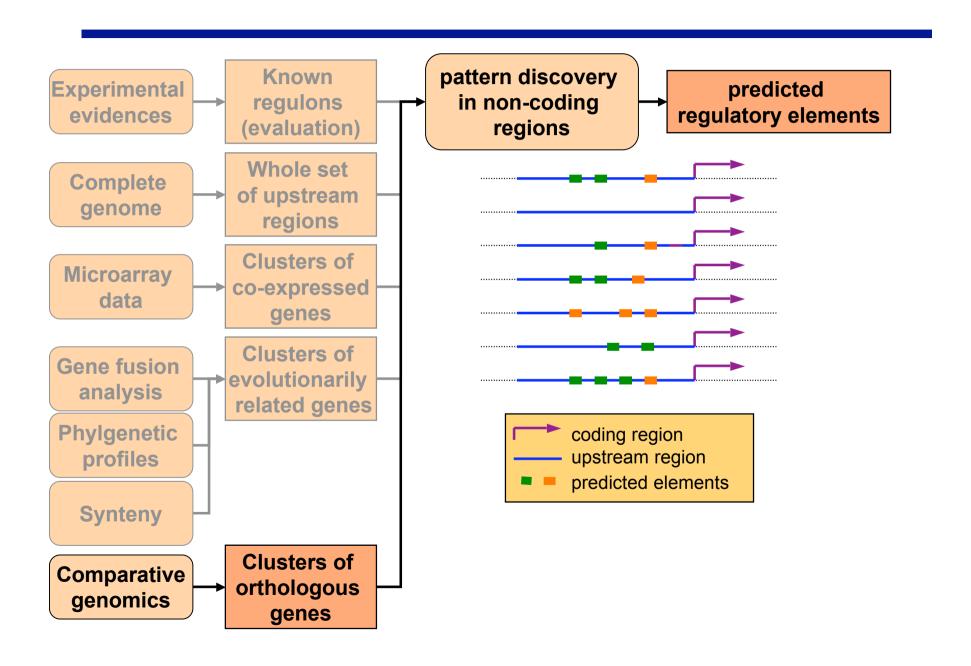


Another alignment in the same genomes

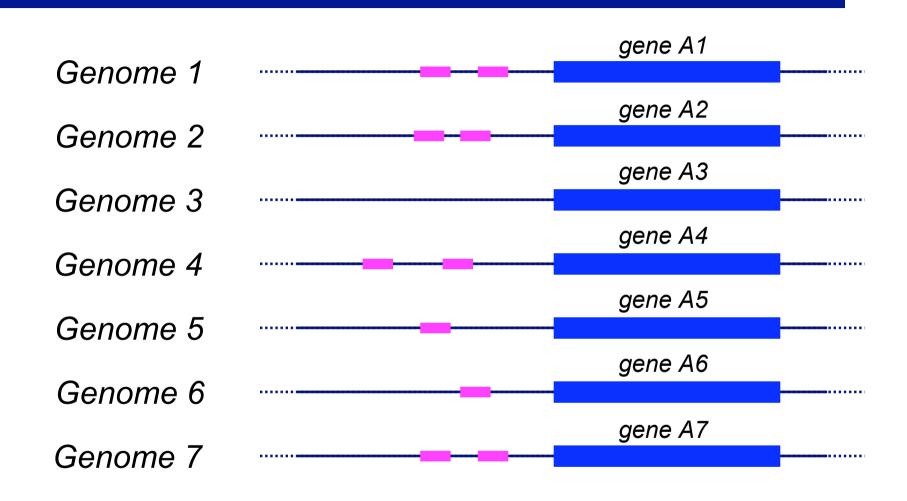
GAL80	(YML051W)	upstream regions
Scer		ATGGCGCAAGTTTTCCGCTTTGTAATATATATTTATACCCCTTTCTTCTCTCCCCTGCAA
Spar		AGGGGCCAAAGCTCCCGCTCTGTAAAATATATTTATATCCCTTCCTT
Smik		TAGGGACAAAGCCCGCCTTTTGTAATATATACTTATACCCTCTCCTTCTCTCCCCTGCAA
Sbay		** *** * * ***** **** ** * * *****
Scer		TATAATAGTTTAATTCTAATATTAATAATATCCTATATTTTCTTCATTTACCGGCGC
Spar		TATAATAGTTTAATTCTAATATTAATAATATCCTATATTTTCCTTACC-ACCGGCGC
Smik		CATAATAGTTAACTCCTAATATTAATAATATATCCTACAATTTCCTTAGC-ACCGGGGC
Sbay		******* * * ********* **** * **** * * ****
Scer		ACTCTCGCCCGAACGACCTCAAAATGTCTGCTACATTCATAATAACCAAAAGCTCATAAC
Spar		ACTCTCGCCCGAACGACCTCAAAATGCTTGCTACATTCATAATAATCAAAAGCTTATAAC
Smik		ACTCTCGCCCGAACGACCTCAAAACGCTTGCTACATCCATAATATTCAGAACTACATCAC
Sbay		
		******* ** ** ** ** ** ** ** ** ** ** *
Scer		TTTTTTTTTTGAACCTGAATATATATACATCACATATCACTGCTGGTCCTTGCCGA
Spar		TTTTTTTTTCCTTTGTACCTGAATATATATACATCTCATGTCACTGCTGGTCCTTGCCGG
Smik		TTTTTTTTTGTACATAAAATATATACCACATGTCACTGCTGATCCTTGCTGA
Sbay		

Scer		CCAGCGTATACAATCTCGATAGTTGGTTT-C-CCGTTCTTTCCACTCCCGTCATGGACTA
Spar		CCAGCGTATACAACCTCGATAGCTGGTTTTC-CCGTTCTTCCCACTCCTGTCATGGACTA
Smik		CGAGCGTATACAAGCTCGATAGCTGGTCTTTACCGTGCCATTCCCTGCCGTCATGGACTA
Sbay		
		* ****** *** **** * * * * * * * * * * *

Motifs in clusters of orthologous genes (COGs)

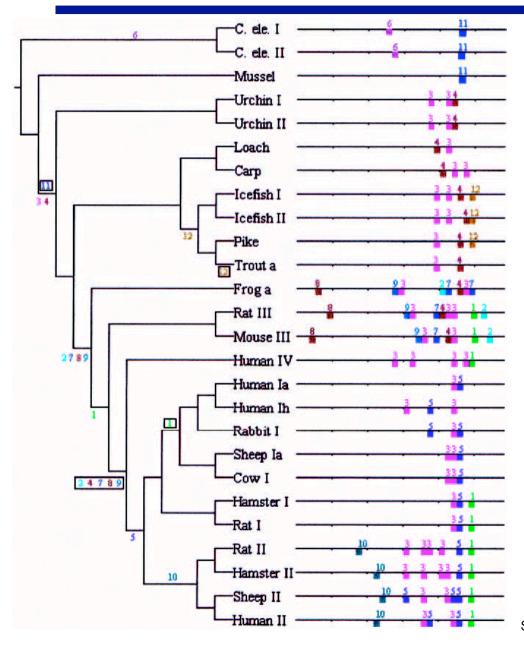


Phylogenetic footprinting to predict regulatory sites



orthologous genes — conserved regulatory sites

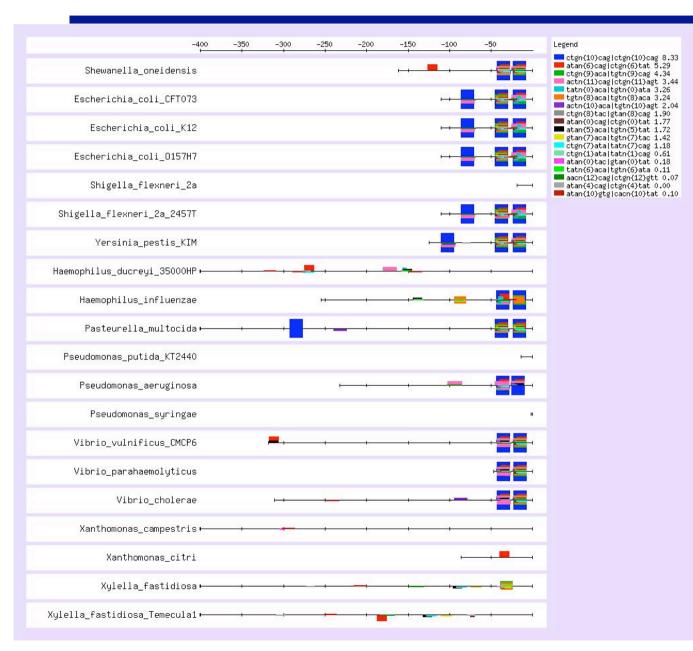
Footprinter example metallothionein



- 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.

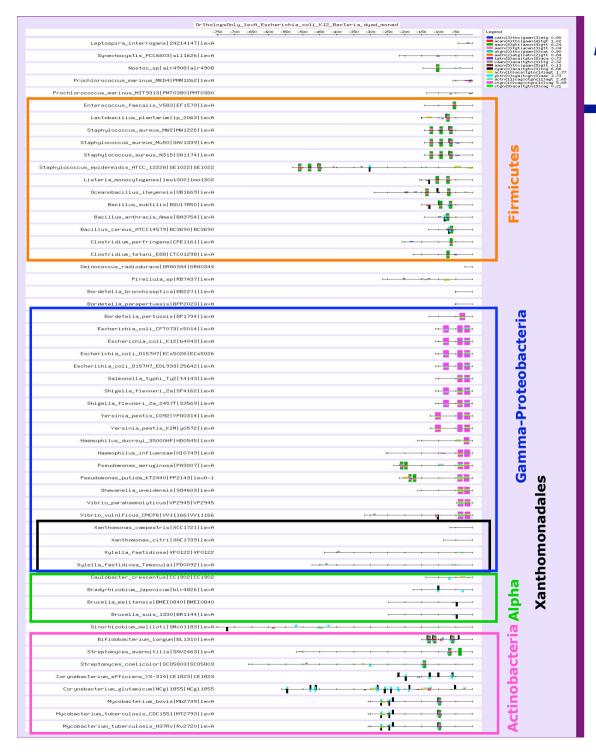
Source: Blanchette and Tompa (2002). Genome Research. 12, 739–748.

Pattern discovery in upstream regions of COGs



- Sequences
 - Upstream sequences of the cluster of orthologs for the gene lexA in all Gammaproteobacteria
- Pattern discovery
 - dyad-analysis
- A very highly significant signal is detected. It is found in two conserved positions in most genes of the group, and in 3 positions for certain bacteria.
- This motif corresponds to the known lexA binding site

Janky & van Helden, in prep.

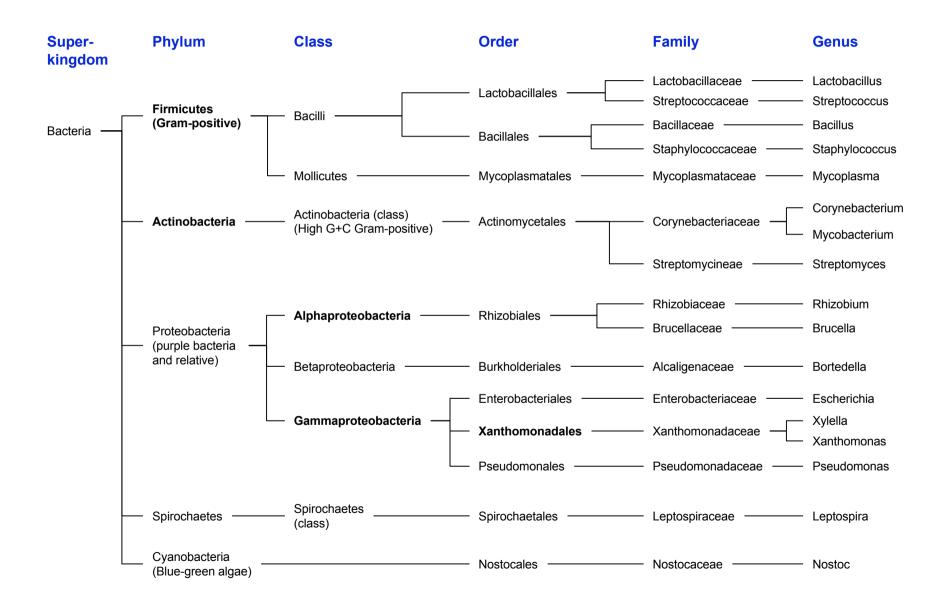


Evolution of cis-acting elements

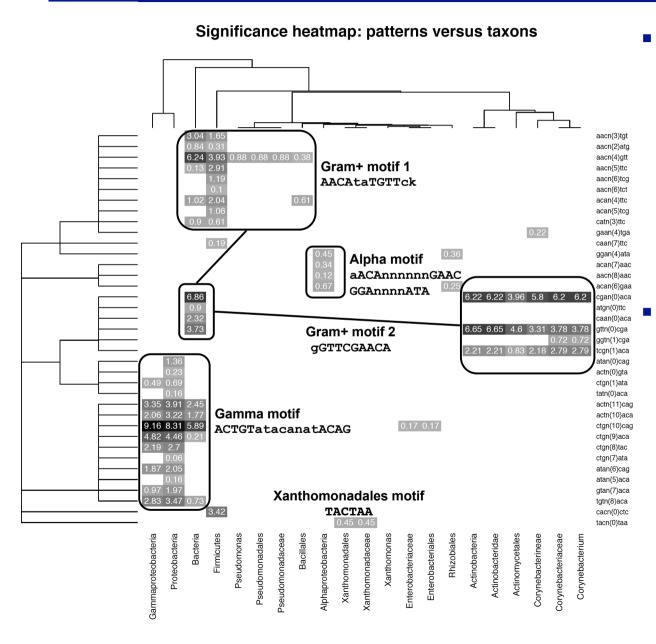
- Pattern discovery was applied
 - to upstream sequences the of 56 bacterial orthologs of lexA
 - at all the taxonomical levels.
- Note: the figure only shows the motifs disovered at one level (all bacterai analyzed together).
- Several motifs are discovered.
- Several taxon-specific motifs are discovered
 - Gamma-prtoeobacteria
 - Xanthomonadales
 - Gram-positive bacteria
 - Firmicutes
 - Actinobacteria
 - Alpha-proteobacteria
- These motifs correspond to those documented in the biological literature.

Janky & van Helden, in prep.

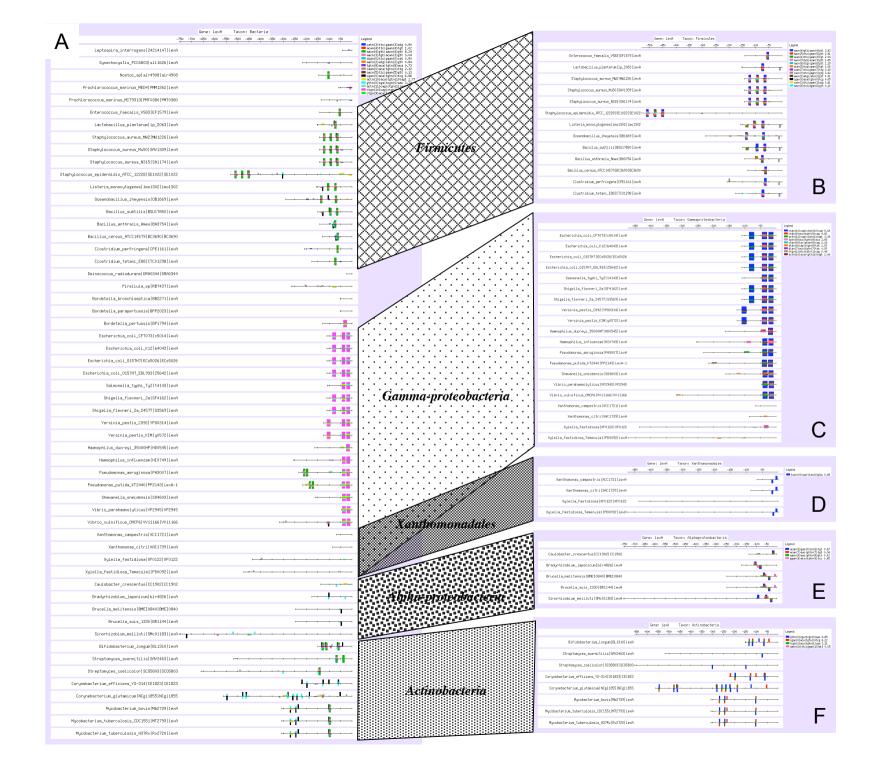
Bacteria taxonomy



Heat map - pattern sinificance per taxon



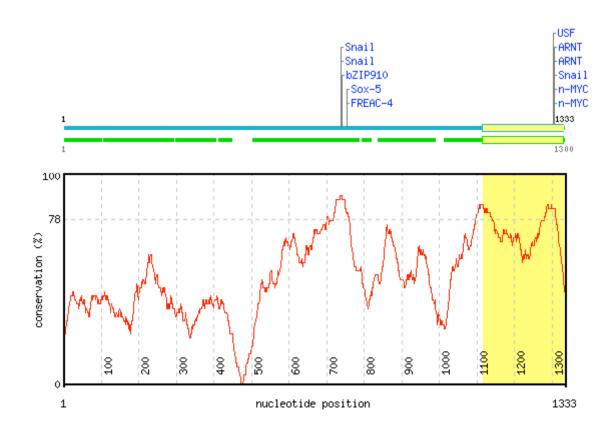
- The color in each cell indicates the significance of one pattern (row) in one taxon (column)
 - Rows are clustered by similarity between patterns.
 - Taxons are clustered by similarity between profiles of pattern significance.
- Several coherent groups of patterns are detected



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/

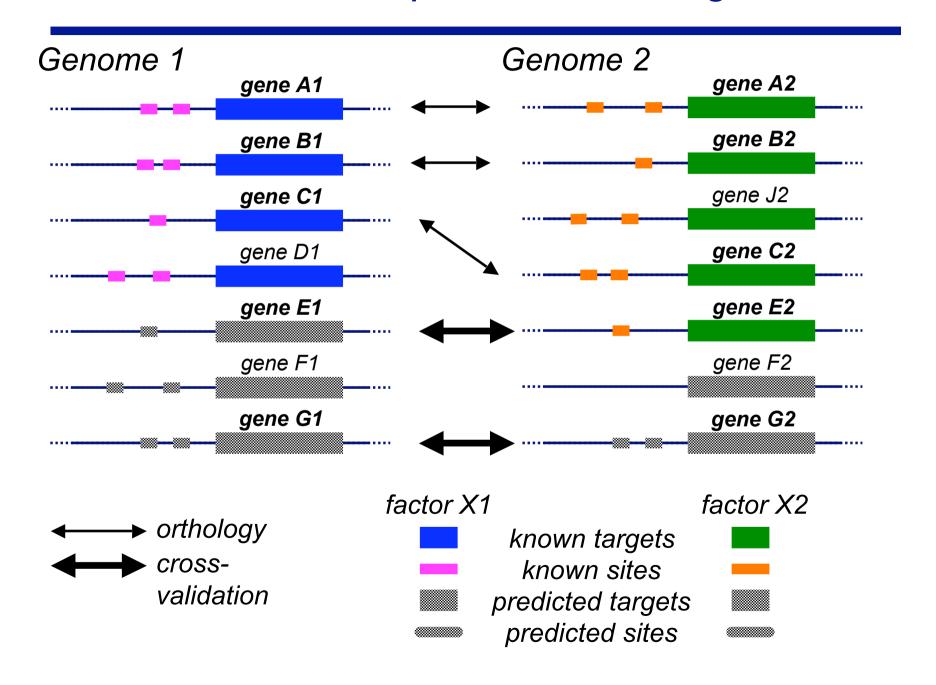
Conservation profile of Human IR



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

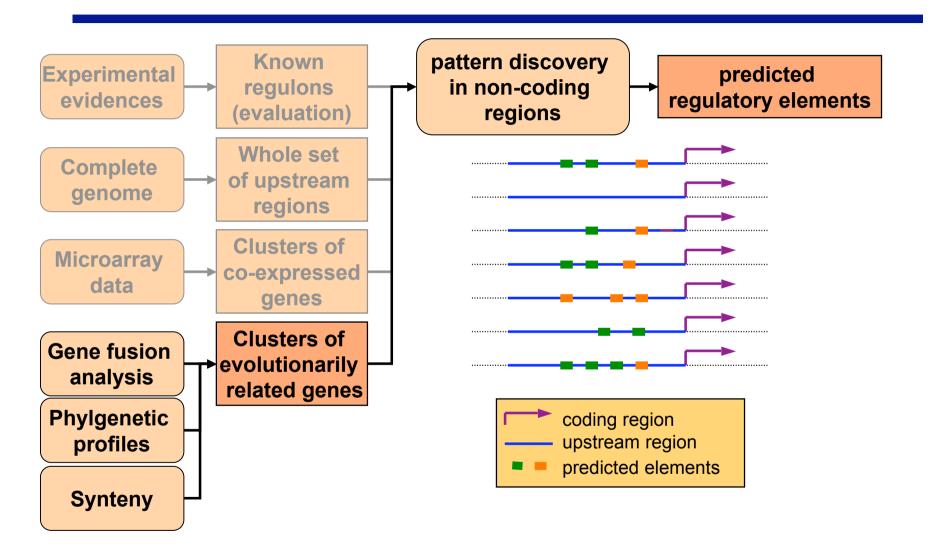
Cross-validation of pattern matching



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
Escherichia_coli_K12	EC21	cccctcaccctctt	aagagggtgaggggg	13.01
Escherichia_coli_K12	EC21	cccctcacccctt	aaggggtgaggggg	13.01
Escherichia_coli_K12	EC21	gccctcacccctc	gaggggtgagggc	13.01
Escherichia_coli_K12	EC21	ggggagagggtgagggga	tcccctcaccctctcccc	13.01
Escherichia_coli_K12	EC21	cctcaccctcaccctctcccctc	gaggggagagggtgagggtgagg	13.01
Escherichia_coli_K12	EC3	cccctcgccctt	aaggggcgaggggg	12.73
Escherichia_coli_K12	EC3	aagggcgaggggg	cccctcgccctt	12.73
Escherichia_coli_K12	EC3	gccctcgcccctc	gaggggcgagggc	12.73
Escherichia_coli_K12	EC3	cccctcaccctt	aaggggtgaggggg	12.73
Escherichia_coli_K12	EC3	cccctctccctt	aaggggagaggggg	12.73
Mycoplasma_pneumoniae	MP1	tataatact	agtattata	11.75
Mycoplasma_pneumoniae	MP1	cttaatactaat	attagtattaag	11.75
Escherichia_coli_K12	EC17	ccctctccctt	aagggagagggg	10.63
Escherichia_coli_K12	EC17	ccctctccctt	aaggggagagggg	10.63
Escherichia_coli_K12	EC17	ccctcgccctt	aagggcgagggg	10.63
Mycoplasma_pneumoniae	MP1	aataataag	cttattatt	10.4
Mycoplasma_pneumoniae	MP1	aataatattatt	aataatattatt	10.4
Mycoplasma_pneumoniae	MP1	taataataagnnnnnaataa	ttattnnnnncttattatta	10.4
Mycoplasma_pneumoniae	MP1	cttagtattatt	aataatactaag	10.4
Mycoplasma_pneumoniae	MP1	taataataagnnnnnaataa	ttattnnnnncttattatta	10.4
Mycoplasma_pneumoniae	MP1	aataatattaaga	tcttaatattatt	10.4
Mycoplasma_pneumoniae	MP1	cttagtatatataatatactaag	cttagtatatatatatatactaag	10.4
Mycoplasma_pneumoniae	MP1	taataataagnnnnnaataa	ttattnnnnncttattatta	10.4
Mycoplasma_pneumoniae	MP1	ctaatattatt	aataatattag	10.4
Mycoplasma_pneumoniae	MP1	taataataagnnnnnaataa	ttattnnnnncttattatta	10.4
Mycoplasma_pneumoniae	MP1	aataatattattnnngtactattataataag	cttattataatagtacnnnaataatattatt	10.4
Mycoplasma_pneumoniae	MP1	aataatattatc	gataatattatt	10.4

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.