

Appendix 5.

Comparing observations in bacterial genomes with expectations in randomly generated genomes

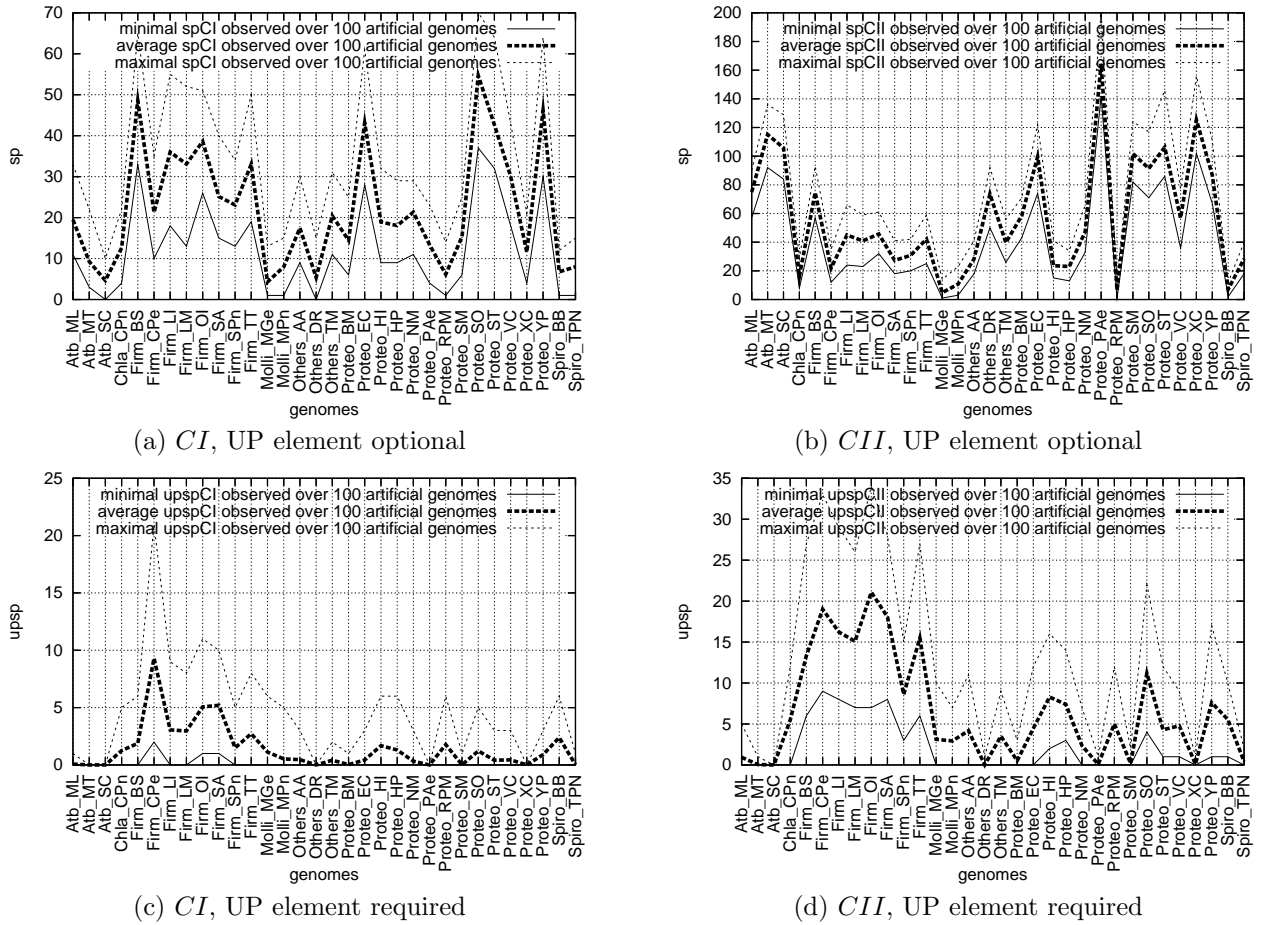


Figure 5.1. Minimal, average and maximal values observed over 100 randomly generated genomes, for sp and $upsp$ respectively; 32 bacterial genomes are considered (see Figure 1 for genome nomenclature). For each such bacterial genome, 100 artificial genomes are generated at random, which have same proportions of A, C, T, G nucleotides and same total number of genes encoding proteins as the bacterial genome. sp_{CI} denotes the number of genes with a putative **S**trong **P**romoter identified under constraint set *CI*. sp_{CII} is defined similarly for constraint set *CII* (see text, Subsection "Genome analysis upon request" for the definition of *CI* and *CII* constraints). $upsp_{CI}$ denotes the number of genes with an **U**P element in their putative **S**trong **P**romoter, and identified under constraint set *CI*. $upsp_{CII}$ is defined similarly.

genome name	abbreviation	<i>CI</i>		<i>CII</i>		<i>CI</i>		<i>CII</i>	
		UP element optional		UP element required		UP element optional		UP element required	
genomes for which Z-score value is above threshold									
		5	2	5	2	5	2	5	2
<i>Mycobacterium leprae tn</i>	Atb_ML		ML	ML	ML				
<i>Mycobacterium tuberculosis h37rv</i>	Atb_MT		MT	MT	MT	-	-	MT	MT
<i>Streptomyces coelicolor a3 (2)</i>	Atb_SC		SC	SC	SC	-	-	-	-
<i>Aquifex aeolicus vf5</i>	Others_AA	AA	AA	AA	AA	AA	AA	AA	AA
<i>Deinococcus radiodurans r1</i>	Others_DR		DR	DR	DR	-	-		
<i>Thermotoga maritima</i>	Others_TM	TM	TM	TM	TM	TM	TM	TM	TM
<i>Brucella melitensis 16m</i>	Proteo_BM		BM	BM	BM			BM	BM
<i>Escherichia coli k12</i>	Proteo_EC	EC	EC	EC	EC		EC	EC	EC
<i>Haemophilus influenza rd kw20</i>	Proteo_HI			HI			HI		HI
<i>Helicobacter pylori j99</i>	Proteo_HP			HP					HP
<i>Neisseria meningitidis mc58</i>	Proteo_NM		NM	NM	NM		NM	NM	NM
<i>Pseudomonas aeruginosa pa01</i>	Proteo_PAe	PAe	PAe	PAe	PAe	-	-	PAe	PAe
<i>Sinorhizobium meliloti 1021</i>	Proteo_SM	SM	SM	SM	SM			SM	SM
<i>Shewanella oneidensis mr1</i>	Proteo_SO		SO	SO	SO		SO	SO	SO
<i>Salmonella typhimurium lt2</i>	Proteo_ST	ST	ST	ST	ST	ST	ST	ST	ST
<i>Vibrio cholerae n16961</i>	Proteo_VC		VC						
<i>Xanthomonas campestris atcc 33913</i>	Proteo_XC		XC	XC		-	-		
<i>Yersinia pestis</i>	Proteo_YP		YP	YP					YP

Table 5.1. Evaluation of the significance of $\sigma 70$ promoter frequencies for 18 bacterial species selected as non *Firmicutes* species with large genomes. For definition of *CI* and *CII* constraints, see Subsection "Genome analysis upon request". The significance is evaluated through the Z-score value (see text, Subsection "Empirical approach"). Two thresholds are considered (5 and 2). – means that the Z-score is not calculable. In a given column, the mention of a species points out statistical significance.

		Z-scores				
	UP element presence	<i>E. coli</i>	min	max	average	standard deviation
<i>CI</i>	optional	7.33	0.51	165.32	38.42	50.16
<i>CII</i>	"	15.38	2.02	172.62	51.12	56.87
<i>CI</i>	required	3.95	0.11	128.4	33.66	38.33
<i>CII</i>	"	8.43	0.11	239.71	59.94	82.06

Table 5.2. Evaluation of the significance of $\sigma 70$ promoter frequencies: comparison of *E. coli* with respect to 26 species with large genomes. For details about Z-scores, see Subsection "Comparing observations in bacterial genomes with expectations in genomes randomly generated".