

# Appendix 1.

## Frequencies of putative strong promoters observed over 32 prokaryotic genomes

<i>n</i>	<i>a</i>	<i>s</i> (Mbp)	<i>at</i>	<i>g</i>	<i>CI</i>		<i>CII</i>	
					<i>sp<sub>CI</sub></i>	<i>p<sub>1CI</sub></i> (%)	<i>sp<sub>CII</sub></i>	<i>p<sub>1CII</sub></i> (%)
<i>Aquifex aeolicus</i> v5	Others_AA	1.55	57.62	<b>1522</b>	<b>214</b>	14.1	<b>333</b>	21.9
<i>Bacillus subtilis</i> 168	Firm_BS	4.21	59.27	<b>3979</b>	<b>796</b>	20.0	<b>1225</b>	30.8
<i>Borrelia burgdorferi</i> b31	Spiro_BB	0.91	73.07	<b>850</b>	<b>43</b>	5.1	<b>49</b>	5.8
<i>Brucella melitensis</i> 16m chr1	Proteo_BM	3.28	45.56	<b>2059</b>	<b>25</b>	1.2	<b>93</b>	4.5
<i>Chlamydomonas reinhardtii</i> ar 39	Chla_CPn	1.22	62.55	<b>1069</b>	<b>23</b>	2.2	<b>30</b>	2.8
<i>Clostridium perfringens</i> str13	Firm_CPe	3.03	74.93	<b>2532</b>	<b>797</b>	31.5	<b>835</b>	33.0
<i>Deinococcus radiodurans</i> r1 chr1	Others_DR	3.05	34.01	<b>2521</b>	<b>22</b>	0.9	<b>204</b>	8.1
<i>Escherichia coli</i> k12	Proteo_EC	4.63	52.34	<b>4173</b>	<b>96</b>	2.3	<b>254</b>	6.1
<i>Haemophilus influenzae</i> rd kw20	Proteo_HI	1.83	64.22	<b>1673</b>	<b>31</b>	1.9	<b>37</b>	2.2
<i>Helicobacter pylori</i> j99	Proteo_HP	1.64	62.95	<b>1478</b>	<b>31</b>	2.1	<b>34</b>	2.3
<i>Listeria innocua</i>	Firm_LI	3.01	64.25	<b>2962</b>	<b>713</b>	24.1	<b>946</b>	31.9
<i>Listeria monocytogenes</i> strain EGD	Firm_LM	2.94	63.97	<b>2837</b>	<b>707</b>	24.9	<b>926</b>	32.6
<i>Mycobacterium leprae</i> tn	Atb_ML	3.26	49.55	<b>2670</b>	<b>31</b>	1.2	<b>122</b>	4.6
<i>Mycobacterium tuberculosis</i> h37rv	Atb_MT	4.41	35.47	<b>3909</b>	<b>32</b>	0.8	<b>290</b>	7.4
<i>Mycoplasma genitalium</i> G37	Molli_MGe	0.58	69.54	<b>441</b>	<b>4</b>	0.9	<b>4</b>	0.9
<i>Mycoplasma pneumoniae</i> M129	Molli_MPn	0.81	61.74	<b>644</b>	<b>18</b>	2.8	<b>26</b>	4.0
<i>Neisseria meningitidis</i> mc58	Proteo_NM	2.27	51.90	<b>1954</b>	<b>50</b>	2.6	<b>92</b>	4.7
<i>Oceanobacillus iheyensis</i> hte831	Firm_OI	3.63	66.41	<b>3398</b>	<b>882</b>	26.0	<b>1049</b>	30.9
<i>Pseudomonas aeruginosa</i> pa01	Proteo_PAe	6.26	35.28	<b>5565</b>	<b>45</b>	0.8	<b>353</b>	6.3
<i>Rickettsia prowazekii</i> madrid e	Proteo_RPM	1.11	73.16	<b>796</b>	<b>6</b>	0.8	<b>6</b>	0.8
<i>Salmonella typhimurium</i> lt2	Proteo_ST	4.85	51.22	<b>4334</b>	<b>117</b>	2.7	<b>342</b>	7.9
<i>Shewanella oneidensis</i> mr1	Proteo_SO	4.96	56.69	<b>4501</b>	<b>118</b>	2.6	<b>167</b>	3.7
<i>Sinorhizobium meliloti</i> 1021	Proteo_SM	3.65	39.50	<b>3272</b>	<b>55</b>	1.7	<b>385</b>	11.8
<i>Staphylococcus aureus</i> mw2	Firm_SA	2.82	69.77	<b>2610</b>	<b>560</b>	21.5	<b>622</b>	23.8
<i>Streptococcus pneumoniae</i> r6	Firm_SPn	2.03	62.39	<b>1861</b>	<b>213</b>	11.5	<b>285</b>	15.3
<i>Streptomyces coelicolor</i> a3 (2)	Atb_SC	8.66	29.15	<b>4665</b>	<b>22</b>	0.5	<b>347</b>	7.4
<i>Thermoanaerobacter tengcongensis</i>	Firm_TT	2.68	63.83	<b>2588</b>	<b>581</b>	22.5	<b>715</b>	27.6
<i>Thermotoga maritima</i>	Others_TM	1.86	54.59	<b>1790</b>	<b>305</b>	17.0	<b>707</b>	39.5
<i>Treponema pallidum</i> nichols	Spiro_TPN	1.13	47.00	<b>980</b>	<b>41</b>	4.2	<b>125</b>	12.8
<i>Vibrio cholerae</i> n16961 chr1	Proteo_VC	4.03	54.80	<b>2618</b>	<b>30</b>	1.2	<b>72</b>	2.8
<i>Xanthomonas campestris</i> atcc 33913	Proteo_XC	5.07	35.55	<b>4120</b>	<b>6</b>	0.2	<b>97</b>	2.4
<i>Yersinia pestis</i>	Proteo_YP	4.6	55.36	<b>4090</b>	<b>52</b>	1.3	<b>136</b>	3.3

**Table 1.1** Frequencies of genes harbouring a potentially strong  $\sigma_{70}$  promoter, under two constraint sets, in 32 prokaryotic genomes. *n*: micro-organism name; *a*: abbreviation for micro-organism name (Atb: *Actinobacteria*, Chla: *Chlamydia*, Firm: *Firmicutes* (among which Molli.: *Mollicutes*), "Others" group, Proteo: *Proteobacteria*, Spiro: *Spirochaetales*); *s*: genome size; *at*: average percentage of nucleotides T and A over the 350 bp-long region upstream of start codon, computed over the total number *g* of genes encoding proteins in the genome; *sp<sub>CI</sub>*: number of genes harbouring a putative **Strong Promoter** identified under constraint set *CI*, *sp<sub>CII</sub>*: *idem*, under the more relaxed constraint set *CII* (see text, Subsection "Genome analysis upon request" for the definition of *CI* and *CII* constraints);  $p_{1CI} = 100 \times sp_{CI}/g$ ,  $p_{1CII} = 100 \times sp_{CII}/g$ .

<i>n</i>	<i>a</i>	<i>at</i>	<i>g</i>	<i>CI</i>		<i>CII</i>		$p1_{CI} \times p2_{CI}$ (%)	$p1_{CII} \times p2_{CII}$ (%)
				<i>upsp<sub>CI</sub></i>	<i>p2<sub>CI</sub></i> (%)	<i>upsp<sub>CII</sub></i>	<i>p2<sub>CII</sub></i> (%)		
Aquifex aeolicus vf5	Others_AA	57.62	<b>1522</b>	<b>19</b>	8.9	<b>73</b>	21.9	1.3	4.8
Bacillus subtilis 168	Firm_BS	59.27	<b>3979</b>	<b>115</b>	14.5	<b>416</b>	34.0	2.9	10.5
Borrelia burgdorferi b31	Spiro_BB	73.07	<b>850</b>	<b>25</b>	58.1	<b>40</b>	81.6	2.9	4.7
Brucella melitensis 16m	Proteo_BM	45.56	<b>2059</b>	<b>0</b>	0	<b>4</b>	4.3	0	0.2
Chlamydomonada pneumoniae ar 39	Chla_CPn	62.55	<b>1069</b>	<b>6</b>	26.1	<b>14</b>	46.7	0.6	1.3
Clostridium perfringens str13	Firm_CPe	74.93	<b>2532</b>	<b>511</b>	64.1	<b>730</b>	87.4	20.2	28.8
Deinococcus radiodurans r1	Others_DR	34.01	<b>2521</b>	<b>0</b>	0	<b>1</b>	0.5	0	0.0
Escherichia coli k12	Proteo_EC	52.34	<b>4173</b>	<b>3</b>	3.1	<b>22</b>	8.7	0.1	0.5
Haemophilus influenzae rd kw20	Proteo_HI	64.22	<b>1673</b>	<b>8</b>	25.8	<b>20</b>	54.1	0.5	1.2
Helicobacter pylori j99	Proteo_HP	62.95	<b>1478</b>	<b>7</b>	22.6	<b>17</b>	50.0	0.5	1.2
Listeria innocua	Firm_LI	64.25	<b>2962</b>	<b>145</b>	20.3	<b>501</b>	53.0	4.9	16.9
Listeria monocytogenes strain EGD	Firm_LM	63.97	<b>2837</b>	<b>147</b>	20.8	<b>488</b>	52.7	5.2	17.2
Mycobacterium leprae tn	Atb_ML	49.55	<b>2670</b>	<b>0</b>	0	<b>1</b>	0.8	0	0.0
Mycobacterium tuberculosis h37rv	Atb_MT	35.47	<b>3909</b>	<b>0</b>	0	<b>2</b>	0.7	0	0.1
Mycoplasma genitalium G37	Molli_MGe	69.54	<b>441</b>	<b>1</b>	25.0	<b>1</b>	25.0	0.2	0.2
Mycoplasma pneumoniae M129	Molli_MPn	61.74	<b>644</b>	<b>2</b>	11.1	<b>10</b>	38.5	0.3	1.6
Neisseria meningitidis mc58	Proteo_NM	51.90	<b>1954</b>	<b>3</b>	6.0	<b>13</b>	14.1	0.2	0.7
Oceanobacillus ihayensis hte831	Firm_OI	66.41	<b>3398</b>	<b>217</b>	24.6	<b>576</b>	54.9	6.4	17.0
Pseudomonas aeruginosa pa01	Proteo_PAe	35.28	<b>5565</b>	<b>0</b>	0	<b>6</b>	1.7	0	0.1
Rickettsia prowazekii madrid e	Proteo_RPM	73.16	<b>796</b>	<b>2</b>	33.3	<b>6</b>	100	0.3	0.8
Salmonella typhimurium lt2	Proteo_ST	51.22	<b>4334</b>	<b>10</b>	8.6	<b>40</b>	11.7	0.2	0.9
Shewanella oneidensis mr1	Proteo_SO	56.69	<b>4501</b>	<b>10</b>	8.5	<b>38</b>	22.8	0.2	0.8
Sinorhizobium meliloti 1021	Proteo_SM	39.50	<b>3272</b>	<b>0</b>	0	<b>6</b>	1.6	0	0.2
Staphylococcus aureus mw2	Firm_SA	69.77	<b>2610</b>	<b>225</b>	40.2	<b>431</b>	69.3	8.6	16.5
Streptococcus pneumoniae r6	Firm_SPn	62.39	<b>1861</b>	<b>59</b>	27.7	<b>120</b>	42.1	3.2	6.5
Streptomyces coelicolor a3 (2)	Atb_SC	29.15	<b>4665</b>	<b>0</b>	0	<b>0</b>	0	0	0
Thermoanaerobacter tengcongensis	Firm_TT	63.83	<b>2588</b>	<b>150</b>	25.8	<b>407</b>	56.9	5.8	15.7
Thermotoga maritima	Others_TM	54.59	<b>1790</b>	<b>31</b>	10.2	<b>120</b>	17.0	1.7	6.7
Treponema pallidum nichols	Spiro_TPN	47.00	<b>980</b>	<b>1</b>	2.4	<b>8</b>	6.4	0.1	0.8
Vibrio cholerae n16961	Proteo_VC	54.80	<b>2618</b>	<b>2</b>	6.7	<b>5</b>	6.9	0.1	0.2
Xanthomonas campestris atcc 33913	Proteo_XC	35.55	<b>4120</b>	<b>0</b>	0	<b>0</b>	0	0	0
Yersinia pestis	Proteo_YP	55.36	<b>4090</b>	<b>3</b>	5.8	<b>20</b>	14.7	0.1	0.5

**Table 1.2** Frequencies of genes with a putative  $\sigma 70$  promoter harbouring an UP element, under two constraint sets, in 32 prokaryotic genomes. *n* micro-organism name; *a*: abbreviation for micro-organism name (see Table 1.1.); *at*: average percentage of nucleotides T and A over the 350 bp-long region upstream of start codon, computed over the total number *g* of genes in the genome; *upsp<sub>CI</sub>*: number of genes with an UP element harboured in the Strong Promoter region, identified under constraint set *CI*; *upsp<sub>CII</sub>*: *idem*, under the more relaxed constraint set *CII* (see text, Section "Genome analysis upon request" for the definition of *CI* and *CII* constraints);  $p2_{CI} = 100 \times upsp_{CI}/sp_{CI}$ ,  $p2_{CII} = 100 \times upsp_{CII}/sp_{CII}$ , with *sp<sub>CI</sub>* (resp. *sp<sub>CII</sub>*) the number of genes harbouring a Strong Promoter under constraints *CI* (resp. *CII*). The percentages indicated in the two rightmost columns are respectively  $p1_{CI} \times p2_{CI} = upsp_{CI}/g$  and  $p1_{CII} \times p2_{CII} = upsp_{CII}/g$ , that is the ratio of the number of genes identified with an UP element in their putative strong promoter, to the number *g* of genes encoding proteins in the genome considered.