

UDC 577.218.577.577.121.9

Computer Analysis of Regulatory Signals in Bacterial Genomes. Fnr Binding Sites

A. V. Gerasimova, D. A. Rodionov, A. A. Mironov, and M. S. Gelfand

State Research Center for Biotechnology GosNII Genetika, Moscow, 113545 Russia;

E-mail: misha@imb.imb.ac.ru

Received March 16, 2001

Abstract—Comparative approach to computer analysis of regulatory signals allows one to predict new signals in bacterial genomes with high accuracy. A prediction is reliable whenever candidate signals are consistently observed in several related genomes. We applied comparative approach to the analysis of the Fnr regulon of gamma-proteobacteria. Responding to changes in the aerobic/anaerobic state of the medium, the transcriptional factor Fnr regulates expression of many genes. We predicted Fnr binding sites in 12 genes regulated by Fnr, and identified 17 new operons as potential members of the Fnr regulon of *Escherichia coli*. In addition, we described the Fnr regulon of other gamma-proteobacteria.

Key words: Fnr, computer analysis, aerobic/anaerobic regulation, *Escherichia coli*

INTRODUCTION

The transcriptional factor Fnr regulates aerobic/anaerobic-dependent gene expression in gamma-proteobacteria. In *Escherichia coli*, expression of more than 120 Fnr-regulated genes depends on alternation of aerobic and anaerobic conditions [1, 2]. Fnr activates expression of several genes of anaerobic enzymes, such as nitrate and nitrite reductase (anaerobic respiration) and pyruvate-formate lyase (formate-acetyltransferase) (anaerobic fermentation). Moreover, Fnr can repress some genes of respiratory enzymes, such as cytochrome *d* oxidase and NADH dehydrogenase. The aim of this study was to analyze the Fnr regulon and find new potentially Fnr-regulated genes in the *Escherichia coli* genome and in the less studied genomes of other gamma-proteobacteria.

METHODS

We considered the genomes of seven related gamma-proteobacteria. Full nucleotide sequences of the *Escherichia coli* [3], *Haemophilus influenzae* [4], *Vibrio cholerae* [5], and *Pseudomonas aeruginosa* [6] genomes were extracted from GenBank (<http://www.ncbi.nlm.nih.gov/GenBank/>) [7]. Preliminary nucleotide sequences of *Salmonella typhi*, *Klebsiella pneumoniae*, and *Yersinia pestis* were extracted from <http://www.sanger.ac.uk/Projects/Microbes/> and <http://genome.wustl.edu/gsc/Projects/bacteria.shtml>. Bacterial genomes were analyzed using the software package Genome Explorer [8]. The training set was compiled using information on the Fnr sites obtained from the database of regulatory sites DPInteract (<http://arep.med.harvard.edu/dpinteract/>) [9]. The rec-

ognition rule was constructed using the program SignalX [8]. Functional annotation of genomes was done using the program BLASTA (<http://www.ncbi.nlm.nih.gov/BLAST/>) [10] and the data bank of protein amino acid sequences SWISSPROT (<http://expasy.hcuge.ch/sprot/>) [11]. Alignments of intergene regions in related organisms were constructed using the program Menteric (<http://globin.cse.psu.edu/enterix/menteric/menteric.html>) [12].

RESULTS AND DISCUSSION

Basing on of the training set, we constructed a weight matrix using the program SignalX:

<i>a</i>	<i>c</i>	<i>g</i>	<i>t</i>
−0.09	−0.09	−0.32	0.49
−0.29	0.05	−0.29	0.52
−0.26	−0.26	0.55	−0.03
0.40	−0.15	−0.37	0.12
−0.39	0.07	−0.06	0.38
−0.08	−0.01	0.02	0.07
0.07	−0.26	0.00	0.19
0.19	0.00	−0.26	0.07
0.07	0.02	−0.01	−0.08
0.38	−0.06	0.07	−0.39
0.12	−0.37	−0.15	0.40
−0.03	0.55	−0.26	−0.26
0.52	−0.29	0.05	−0.29
0.49	−0.32	−0.09	−0.09

Table 1. Known and potential Fnr-regulated operons in the genomes of *Escherichia coli* (*E. c.*), *Salmonella typhi* (*S. t.*), *Klebsiella pneumoniae* (*K. p.*), *Yersinia pestis* (*Y. p.*), *Haemophilus influenzae* (*H. i.*), *Vibrio cholerae* (*V. c.*), and *Pseudomonas aeruginosa* (*P. a.*)

<i>E. c.</i>	<i>S. t.</i>	<i>K. p.</i>	<i>Y. p.</i>	<i>V. c.</i>	<i>H. i.</i>	<i>P. a.</i>	Function	Metabolic pathway	Fnr influence on gene expression
First group									
<i>narXL/narK GJI</i>	<i>narXL/ /narKGJI</i>	<i>narXL/ /narKGJI</i>	##	##	##	<i>narX/narK narG</i>	Nitrate reductase, nitrate transporter	Anaerobic respiration	Activates
<i>nirBD; nirC</i>	<i>NirBD</i>	<i>nirBD</i>	<i>nirBD</i>	<i>nirBD</i>	#	#	Nitrite reductase	Denitrification	Activates
<i>cydAB</i>	<i>cydAB</i>	<i>cydAB</i>	<i>cydAB</i>	<i>cydAB</i>	<i>cydAB</i>	#	Cytochrome <i>d</i> oxidase	Aerobic respiration	Represses
<i>ndh</i>	<i>Ndh</i>	<i>ndh</i>	<i>ndh</i>	<i>ndh</i>	<i>ndh</i>	<i>ndh</i>	NADH dehydrogenase	Aerobic respiration	Represses
<i>nrfABC-DEFG</i>	<i>nrfABC-DEFG</i>	#	#	#	<i>nrfABCD</i>	#	Formate-dependent nitrite reductase	Anaerobic respiration	Activates
<i>fdnGHI</i>	<i>FdnG</i>	<i>fdnG?</i>	<i>#, fdnGHI</i>	##	<i>fdnGHI</i>	<i>#, fdnGHI</i>	Formate dehydrogenase	Anaerobic respiration	Activates
<i>focA-pflB</i>	<i>FocA-pflB</i>	<i>focA-pflB</i>	<i>focA-pflB</i>	<i>focA-pflB</i>	<i>focA-pflB</i>	#	Formate transporter, pyruvate-formate lyase	Anaerobic fermentation	Activates
<i>ansB</i>	<i>AnsB</i>	<i>ansB</i>	<i>ansB</i>	#	<i>ansB</i>	<i>ansB</i>	L-Asparaginase	Catabolism of asparagine	Activates
Second group									
<i>pdhR-aceEF-lpdA</i>	<i>pdhR-aceEF-lpdA</i>	<i>pdhR-aceEF-lpdA</i>	<i>pdhR-aceEF-lpdA</i>	<i>pdhR-aceEF-lpdA</i>	# <i>aceEF, lpdA</i>	<i>pdhR-aceEF-lpdA</i>	Pyruvate dehydrogenase	Pyruvate metabolism	Represses
<i>feoAB</i>	<i>FeoAB</i>	<i>feoAB</i>	<i>feoAB</i>	<i>feoB</i>	#	<i>feoAB</i>	Ferrous iron transporter		
<i>nrdDG</i>	<i>NrdDG</i>	<i>nrdD, nrdG</i>	<i>nrdDG</i>	<i>nrdDG</i>	<i>nrdD, nrdG</i>	#	Anaerobic ribonucleoside-triphosphate reductase	Nucleoside metabolism	Activates
<i>dmsABC</i>	<i>dmsABC</i>	<i>dmsA, dmsBC</i>	<i>dmsABC</i>	#	<i>dmsABC</i>	##	Anaerobic dimethyl sulfoxide reductase	Anaerobic growth on various sulfoxides	Activates
<i>dcuA</i>	<i>dcuA</i>	<i>dcuA</i>	<i>dcuA</i>	<i>dcuA</i>	#	#	C4-dicarboxylate anaerobic transporter	The Krebs cycle (transport)	Activates
<i>fnr</i>	<i>fnr</i>	<i>fnr</i>	<i>fnr</i>	<i>fnr</i>	<i>fnr</i>	<i>fnr</i>	Fnr transcription regulator		Represses
<i>arcA</i>	<i>arcA</i>	<i>arcA</i>	<i>arcA</i>	<i>arcA</i>	<i>arcA</i>	#	ArcA transcription regulator		Activates
<i>yfiD</i>	<i>yfiD</i>	<i>yfiD</i>	<i>yfiD</i>	<i>yfiD</i>	<i>yfiD</i>	#	?		Represses
<i>tdcABC</i>	<i>tdcABC</i>	<i>tdcABC</i>	#	#	#	#	tdcABC operon transcriptional activator	Catabolism of threonine	Activates
<i>dcuC</i>	<i>dcuC</i>	<i>dcuC</i>	#	<i>dcuC</i>	#	#	C4-dicarboxylate anaerobic transporter	The Krebs cycle (transport)	Activates

Table 1. (Contd.)

<i>E. c.</i>	<i>S. t.</i>	<i>K. p.</i>	<i>Y. p.</i>	<i>V. c.</i>	<i>H. i.</i>	<i>P. a.</i>	Function	Metabolic pathway	Fnr influence on gene expression
Third group									
<i>b2503</i>	<i>b2503</i>	<i>b2503</i>	#	#	#	#	?		
<i>ompW (yciD)</i>	<i>yciD</i>	<i>yciD</i>	<i>yciD</i>	<i>yciD</i>	#	#	Outer membrane protein W (porin)		
<i>pyrG</i>	<i>pyrG</i>	<i>pyrG</i>	<i>pyrG</i>	<i>pyrG</i>	<i>pyrG</i>	<i>pyrG</i>	CTP synthetase	Nucleoside metabolism	
<i>upp</i>	<i>upp</i>	<i>upp</i>	<i>upp</i>	<i>upp</i>	<i>upp</i>	<i>upp</i>	Uracil phosphoribosyl-transferase	Nucleoside metabolism	
<i>b0780/moaA</i>	<i>b0780/moaA</i>	<i>b0780/moaA</i>	<i>b0780/moaA</i>	<i>b0780/moaA</i>	#/ <i>moaA</i>	# #	Molybdenum cofactor of protein A biosynthesis		
<i>ccpR (yhjA)</i>	<i>yhjA</i>	<i>yhjA</i>	#	<i>yhjA</i>	#	<i>yhjA</i>	Cytochrome <i>c551</i> peroxidase	Antiperoxide protection	
<i>ppsA</i>	<i>ppsA</i>	<i>ppsA</i>	<i>ppsA</i>	<i>ppsA</i>	#	<i>ppsA</i>	Phosphoenolpyruvate synthase	Gluconeogenesis	
<i>b0873-b0872</i>	<i>b0873</i>	<i>b0873</i>	<i>b0873</i>	#	#	#	?, NADH oxidoreductase		
<i>fadL</i>	<i>fadL</i>	<i>fadL</i>	<i>fadL</i>	<i>fadL</i>	<i>fadL</i>	#	Long-chain fatty acid transporter	Fatty acid metabolism	
<i>wrbA</i>	<i>wrbA</i>	<i>wrbA</i>	<i>wrbA</i>	<i>wrbA</i>	#	<i>wrbA</i>	Trp repressor binding protein		
<i>sfhB</i>	<i>sfhB</i>	<i>sfhB</i>	<i>sfhB</i>	<i>sfhB</i>	<i>sfhB</i>	<i>sfhB</i>	?		
<i>gltX</i>	<i>gltX</i>	<i>gltX</i>	<i>gltX</i>	<i>gltX</i>	<i>gltX</i>	<i>gltX</i>	Glutamyl-tRNA synthetase		
<i>mltA</i>	<i>mltA</i>	<i>mltA</i>	<i>mltA</i>	<i>mltA</i>	<i>mltA</i>	<i>mltA</i>	Murein transglycosylase A		
<i>mtlA</i>	<i>mtlA</i>	<i>mtlA</i>	<i>mtlA</i>	<i>mtlA</i>	#	#	Mannitol-specific transporter	Sugar metabolism	
<i>b1973</i>	<i>b1973</i>	<i>b1973</i>	#	#	#	#	?		
<i>yjiO</i>	<i>yjiO</i>	<i>yjiO</i>	#	#	#	#	Antibiotic transporter	Drug resistance	
<i>ycfC</i>	<i>ycfC</i>	<i>ycfC</i>	<i>ycfC</i>	#	<i>ycfC</i>	<i>ycfC</i>	?		
<i>last</i>	<i>last</i>	<i>last</i>	#	#	#	#	?		

Note: Boldface means that the found Fnr signal is conserved in this genome; # means that there is no orthologous gene for this organism.

Table 2. Known and potential Fnr boxes in the genomes of *Escherichia coli* (*E. c.*), *Salmonella typhi* (*S. t.*), *Klebsiella pneumoniae* (*K. p.*), *Yersinia pestis* (*Y. p.*), *Haemophilus influenzae* (*H. i.*), *Vibrio cholerae* (*V. c.*), and *Pseudomonas aeruginosa* (*P. a.*)

Gene	Genome	Nucleotide sequence of the Fnr box	Fnr		CRP*		
			position	weight	position	weight	
First group							
<i>narK</i>	<i>E. c.</i>	TTGATTTAcATCAA	-74	5.15		-	
	<i>E. c.</i>	aTGATaaAtATCAA	-112	4.20			
	<i>S. t.</i>	TTGATTTAtATCAA	-74	5.05	-3	3.62	
	<i>K. p.</i>	TTGATaTAAATCAA	-74	5.05		-	
	<i>P. a.</i>	TTGATTcctATCAA	-75	4.41	-79	3.67	
<i>narG</i>	<i>E. c.</i>	TTGATcggtATCAA	-106	4.66		-	
	<i>S. t.</i>	TTGATcggtATCAA	-106	4.66		-	
	<i>K. p.</i>	TTGATcgctATCAA	-107	4.59		-	
<i>nirB</i>	<i>E. c.</i>	TTGATTTAcATCAA	-73	5.15		-	
	<i>S. t.</i>	TTGATTTAcATCAA	-70	5.15		-	
	<i>K. p.</i>	TTGATTTAcATCAA	-70	5.15		-	
	<i>Y. p.</i>	TTGATTTAcATCAA	-69	5.15	28	3.61	
<i>ndh</i>	<i>E. c.</i>	TTGATTaAcATCAA	-151	5.03	-155	3.83	
	<i>S. t.</i>	TTGATgcAcATCAA	-151	4.65	-155	3.88	
	<i>S. t.</i>	TTGtTgTtAATtAA	-70	3.94			
	<i>K. p.</i>	TTGATgcAcATCAA	-149	4.65	-150	3.88	
	<i>Y. p.</i>	TTGATaTAtATCAA	-146	4.9		3.76	
	<i>V. c.</i>	TTGATaaAtATCAA	-200	4.78		3.83	
	<i>cydA</i>	<i>E. c.</i>	TTGATaTtATCAA	-346	4.78		-
		<i>E. c.</i>	TTGtTcTcgATCAA	-294	4.57		
<i>S. t.</i>		TTGATTTtAATCAA	-347	5.08		3.67	
<i>S. t.</i>		TTGtccgtgATCAA	-295	4.14	-555	4.09	
<i>K. p.</i>		TTGATTTAtATCAA	-346	5.05		-	
<i>K. p.</i>		TTGATcaccgTCgA	-246	3.98			
<i>Y. p.</i>		TTGtTcTAAATCAA	-315	4.84	-424	3.92	
<i>Y. p.</i>		TTGtgcTAgATCAA	-264	4.48	-387	3.92	
<i>V. c.</i>		TTGATTTAgATCAA	-221	5.12	-382	3.99	
<i>V. c.</i>		TTGATTTAgATCAt	-271	4.54			
<i>V. c.</i>		TTGATTgtttTCAA	-337	3.97			
<i>H. i.</i>		TTGATcTAAgTCAA	-293	4.81		3.68	
<i>nrfA</i>		<i>E. c.</i>	TTGATTaAAgaCAA	-142	4.49		-
		<i>S. t.</i>	TTGATTaAAgaCAA	-144	4.49	-781	3.5
	<i>H. i.</i>	TTtATTTAAaCAA	-117	4.44		4.12	
	<i>H. i.</i>	TTGATcaAgcTCAA	-67	4.48	87	3.91	
<i>ansB</i>	<i>E. c.</i>	TTGtTTaAcgTCAA	-70	4.44	-124	3.53	
	<i>K. p.</i>	TTGATTaAtgTcTcA	-45	3.81		-	
	<i>P. a.</i>	TTGcTgggcATCAA	-927	3.91		-	
<i>focA</i>	<i>E. c.</i>	aTGATcTAtATCAA	-73	4.39		3.59	
	<i>S. t.</i>	aTGATcTAtATCAA	-70	4.39		4.07	
	<i>K. p.</i>	cTGATgaAAgaCAA	-244	3.86		-	
	<i>Y. p.</i>	aTGATccAtATCAA	-78	3.94		3.81	
	<i>H. i.</i>	TTGtgaatAATCAA	-320	4.09	-437	3.82	
	<i>H. i.</i>				-217	3.53	
<i>fdnG</i>	<i>E. c.</i>	TTGAggTAggTCAA	-134	4.32		-	
	<i>K. p.</i>	cTGATcgAAAaCAA	-200	4.07		-	
	<i>H. i.</i>	aTGATcTAgATCAc	-210	3.65		4.96	
		TTtAacgAAATCAA	-82	3.58		3.45	

Table 2. (Contd.)

Gene	Genome	Nucleotide sequence of the Fnr box	Fnr		CRP*	
			position	weight	position	weight
Second group						
<i>yfiD</i>	<i>E. c.</i>	TTGATTTAAATCAA	-120	5.20		3.82
	<i>E. c.</i>	TTGATgTAAAaCAA	-173	4.87		
	<i>S. t.</i>	TTGATTTAAATCAA	-112	5.20		3.76
	<i>S. t.</i>	TTGtTTTAcATCAA	-165	4.87		
	<i>K. p.</i>	TTGATaTAAATCAA	-107	5.05		-
	<i>K. p.</i>	TTGtTTTAcATCAA	-160	4.87		
	<i>Y. p.</i>	TTGATaTAAAaCAt	-128	4.72		3.55
	<i>Y. p.</i>	TTGATaTAAAaCAt	-181	4.19	-48	3.92
	<i>V. c.</i>	TTGATTTAggTCAA	-598	4.81		3.85
	<i>V. c.</i>	TTGATTTgtgTCAA	-438	4.29		3.5
	<i>H. i.</i>	TTaATTTAgATCAA	-143	4.31		4.38
	<i>H. i.</i>	TTcATTatAAaCAA	-230	3.87		
<i>narX</i>	<i>E. c.</i>	TTGATgTAAAaCAA	-278	5.15		-
	<i>E. c.</i>	TTGATaTttATCAt	-240	4.20		
	<i>S. t.</i>	TTGATaTAAATCAA	-278	5.05	-321	3.62
	<i>K. p.</i>	TTGATaTAAATCAA	-250	5.05		-
	<i>V. c.</i>	TTGtTTTggATCAA	-129	4.39		
	<i>P. a.</i>	TTGATaggAATCAA	-113	4.41		
<i>nrdD</i>	<i>E. c.</i>	TTGAgcTAcATCAA	-248	4.63		4.22
	<i>S. t.</i>	TTGtTcTAcATCAA	-247	4.79		-
	<i>K. p.</i>	TTGtTcTgggTCAA		4		-
	<i>Y. p.</i>	TTGtTcTAggTCAA	-193	4.45		-
	<i>V. c.</i>	TTGATcTAAATCAA	-126	5.12		-
	<i>H. i.</i>	TTGATaTtAATCAg	174	4.35	-207 -153	4.67 3.71
<i>arcA</i>	<i>E. c.</i>	TTGATaTAtgTCAA	-290	4.59		-
	<i>S. t.</i>	TTGATaTAtgTCAA	-289	4.59		-
	<i>K. p.</i>	TTGATaTAtgTCAA	-289	4.59	-558	3.54
	<i>Y. p.</i>	TTGATaTAtgTCAA	-296	4.59		3.6
	<i>V. c.</i>	TTGATgTAAATCAA	-254	5.15		-
	<i>H. i.</i>	TTGtTTTtATCgA	-264	4.18	-129	4.66
<i>b0621 (dcuC)</i>	<i>E. c.</i>	TTGATTTtAATCAg	-102	4.50		-
	<i>S. t.</i>	gTGATTTtAATCAg	-99	3.69	-361	3.63
	<i>K. p.</i>	TTGATTTgcATCAg	-105	4.21		-
	<i>V. c.</i>	TTGcTTTAgATCAt	-102	3.99		-
<i>fnr</i>	<i>E. c.</i>	TTGAcaaAtATCAA	-32	4.47		-
	<i>S. t.</i>	TTGAcaaAtATCAA	-31	4.47		-
	<i>K. p.</i>	TTGAccaAtATCAA	-29	4.54		-
	<i>K. p.</i>	cTGtTTTtAATCAA	-131	4.22		-
	<i>Y. p.</i>	TTGAcgcAtATCAA	-33	4.24		3.51
	<i>V. c.</i>	TTGAcgTAcATCAA	-33	4.79		-
	<i>H. i.</i>	TTGcgTTAgATCAA	-40	4.13		3.81
	<i>E. c.</i>	aTGATTTcAATCAA	-116	4.43	-132	3.88
<i>pdhR</i>	<i>S. t.</i>	cTGATTTcAATCAA	-117	4.43		-
	<i>K. p.</i>	cTGATTTcAATCAA	-133	4.43		-
	<i>Y. p.</i>	aTGATTTcggTCAA	-133	4.04	-50	3.77
	<i>V. c.</i>	aTGATTTAggTCAA	-121	4.23		3.76
	<i>E. c.</i>	aTGtTgTAAATCAA	-132	4.29		-
<i>aceE</i>	<i>V. c.</i>	aTGATTTAggTCAA	-941	4.23		3.76
	<i>E. c.</i>	TTGtTaaAAAaCAA	-24	4.37		-
<i>dcuA</i>	<i>S. t.</i>	TTGtTaaAcAaCAA	-23	4.32	-2	3.62
	<i>K. p.</i>	TTGtTaaAAAaCAA	-22	4.37		-
	<i>Y. p.</i>	TTGAaTgAAATCAA	-405	4.24		-

Table 2. (Contd.)

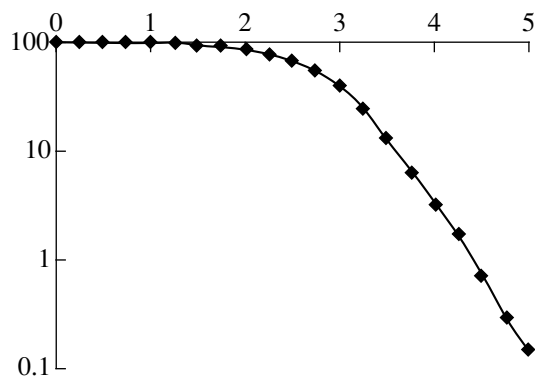
Gene	Genome	Nucleotide sequence of the Fnr box	Fnr		CRP*	
			position	weight	position	weight
<i>feoA</i>	<i>E. c.</i>	TTGAgaccgAaCAA	-204	4.18		3.59
	<i>S. t.</i>	TTGAgccAtATCAA	-214	4.08	-413	3.85
	<i>K. p.</i>	TTGcagcAtATCAA	-204	3.98		
	<i>Y. p.</i>	TTGATaTccAatAA	-488	3.72		-
	<i>P. a.</i>		-134	3.46		-
<i>dmsA</i>	<i>E. c.</i>	TTGATaccgAaCAA	-272	4.05		-
	<i>S. t.</i>	TcGATaTAtATCAg	-154	3.85		3.63
	<i>K. p.</i>	aTGATaatcATCgA	-589	3.71		-
	<i>Y. p.</i>	TTGATTccAgaCAA	-308	3.97		-
	<i>H. i.</i>	TTGATTTggcTCAA	-87	4.23		3.68
<i>tdcA</i>	<i>E. c.</i>	TTGAcaaAAATCAg	-182	4.04	-81	4.2
	<i>S. t.</i>	TTGATTgAAATCAg	-182	4.43	-81	4.2
	<i>K. p.</i>	TTGATTTtAATCAA	-521	5.08	-80	3.93
	<i>K. p.</i>				-52	3.61
Third group						
<i>b2503</i>	<i>E. c.</i>	TTGATaTAtATCAA	-140	4.90		-
	<i>S. t.</i>	TTGAcTTAAATCAA	-141	4.89	-91	3.75
	<i>K. p.</i>	TTGATTactATCAA	-134	4.74		-
<i>upp</i>	<i>E. c.</i>	TTGAcTaAAgTCAA	-67	4.46		-
	<i>S. t.</i>	TTGATccAggTCAA	-74	4.28		3.51
	<i>K. p.</i>	TTGATaTAcgTCAA	-73	4.69	-45	3.65
<i>b0780</i>	<i>Y. p.</i>	TTGATcTgAATCAg	-69	4.09		-
	<i>E. c.</i>	TTGATaTAcATCAt	-146	4.42		-
	<i>S. t.</i>	TTGATaTAtATCAt	-144	4.32		-
<i>moaA</i>	<i>K. p.</i>	TTGATgTcggTCAg	-92	3.99		-
	<i>Y. p.</i>	TTGtTTaAAATCAc	-188	3.99	-195	3.51
	<i>E. c.</i>	aTGATgTAtATCAA	-264	4.42		-
	<i>S. t.</i>	aTGATaTAtATCAA	-265	4.32		-
	<i>K. p.</i>	cTGAccgAcATCAA	-255	3.99		-
<i>yciD</i>	<i>Y. p.</i>	gTGATTTtAAaCAA	-356	3.99	-550	3.65
	<i>H. i.</i>	aTGATTTAAATCAA	-343	4.62		3.98
	<i>H. i.</i>	cTGATTTcATCAA	-195	4.45		-
	<i>E. c.</i>	TTGATTTAAATCAc	-163	4.39		4.56
	<i>S. t.</i>	TTaATcTggATCAA	-118	3.78	-83	3.57
<i>pyrG</i>	<i>K. p.</i>	TTGATTTcAcTCAt	-284	3.99	-138	3.64
	<i>k. p.</i>	TTaATcTggATCAA	-117	3.78	-82	3.56
	<i>Y. p.</i>	aTGATccAgATCAA	-124	4.01		3.92
	<i>V. c.</i>	TTGATTTccATCAA	-91	4.96	-144	3.75
	<i>E. c.</i>	TTGATTTgcgTCAA	-175	4.39		-
<i>fadL</i>	<i>S. t.</i>	TTGATTTAcgTCAA	-174	4.84		-
	<i>K. p.</i>	cTGATTTAcgTCAA	-230	4.26		-
	<i>V. c.</i>	TTGATTTgAAgCAA	-134	4.2		3.6
	<i>H. i.</i>	TTGAcTTAgATCAA	-372	4.81	-440	4.79
	<i>H. i.</i>				-376	3.68
<i>fadL</i>	<i>E. c.</i>	aTGATcTAAAaCAA	-238	4.26		3.76
					-264	4.42
	<i>S. t.</i>	aTGATcTAAAaCAA	-237	4.26		3.69
					-263	4.36
	<i>K. p.</i>	TTGATTTAggaaAA	-78	3.97		-
	<i>Y. p.</i>	TTtTTgAgATCAA	-177	4.07		-
<i>fadL</i>	<i>V. c.</i>	TTGATcTgATgAA	225	4.11		-
	<i>H. i.</i>	TTtATTTAtAaCAA	-30	4.19	-185	3.64

Table 2. (Contd.)

Gene	Genome	Nucleotide sequence of the Fnr box	Fnr		CRP*	
			position	weight	position	weight
<i>ppsA</i>	<i>E. c.</i>	TcGATgTccAaCAA	-2	4.16		-
	<i>S. t.</i>	TcGATgTccAaCAA	-1	4.16		-
	<i>K. p.</i>	TcGATgTccAaCAA	-1	4.16		-
	<i>Y. p.</i>	TcGATgTccAaCAA	47	4.16		-
<i>b0873</i>	<i>E. c.</i>	TTGcgcTAAATCAA	-103	4.13		-
	<i>S. t.</i>	TTGcgcTAAATCAA	-102	4.13		-
	<i>K. p.</i>	TTGcgcTAAATCAA	-100	4.13		-
	<i>Y. p.</i>	TTGcgTcAAATCAA	-113	3.76		-
<i>wrbA</i>	<i>E. c.</i>	TTGtTaTAAATCAA	-114	4.77		3.65
	<i>S. t.</i>	TTGtTaTAAATCAA	-110	4.77		-
	<i>K. p.</i>	TTGtTaTAAATCAA	-112	4.77		-
<i>sfhB</i>	<i>E. c.</i>	TTGAcTTccTCAA	-127	4.28		-
	<i>S. t.</i>	TTGAcTTccTCAA	-123	4.28		-
	<i>Y. p.</i>	TTGATTatccTCgA	-127	4.00		-
	<i>V. c.</i>	TTGAaTTAAcTCAA	15	3.99		-
	<i>H. i.</i>	TTGtTcTtgATaAA	-313	4.06		-
<i>gltX</i>	<i>E. c.</i>	TTcATgaAAATCAA	-2	4.22		-
	<i>S. t.</i>	TTcATgaAAATCAA	-1	4.22		-
	<i>K. p.</i>	TTcATgaAAATCAA	-4	4.22		-
<i>b1973</i>	<i>E. c.</i>	TcGtTTgtcATCAA	-138	4.09		-
	<i>S. t.</i>	TTGATaTcAAaaAA	-131	4	-26	3.68
	<i>K. p.</i>	TTGATTgAtATCgt	-166	3.81		-
<i>yjiO</i>	<i>E. c.</i>	TTGATTaAccgCAA	-281	4.04		-
	<i>S. t.</i>	TTGATTaAcATCAA	-289	5.03		3.85
	<i>K. p.</i>	TTGATTTtAaCAt	554	4.07		-
<i>yefC</i>	<i>E. c.</i>	TTtAcTTAAaCAA	-26	4.03		-
	<i>Y. p.</i>	TTGATggAAATaAA	-589	4.38		-
	<i>Y. p.</i>	cTGATccAAgTCAA	-67	3.78		-
	<i>H. i.</i>	TTtATTTggATCAA	-357	4.09		-
<i>last</i>	<i>E. c.</i>	TTGAcaTAtATCAA	-359	4.59		-
	<i>S. t.</i>	TTGAcaTAtATCAA	-359	4.59		-
	<i>K. p.</i>	TTGAcaTAtATCAA	-356	4.59	-95	3.54
<i>mtlA</i>	<i>E. c.</i>	TTGATaTcAcaCAA	-155	4.14		4.74
	<i>S. t.</i>				-165	4.67
	<i>K. p.</i>	gTGATcTtAATCAA	-328	3.72		-
	<i>Y. p.</i>	gTGATaaAtATCAA	-220	4.19		-
	<i>V. c.</i>	aTGATTTtgAaCAA	-287	3.97		-
	<i>E. c.</i>	TcGcTaTtAATCAA	-136	4.14	-163	3.69
<i>mtlA</i>	<i>V. c.</i>	TTGATggAttTCAA	-33	3.91		-
	<i>H. i.</i>	TTGATggAttTCAA	87	4.04		-

Note: The cases of coincidence between the CRP and Fnr boxes are marked in bold.

* The search threshold for CRP boxes is equal to 3.5.



Weight distribution function for Fnr sites. Horizontal axis, the weight threshold value; vertical axis, percentage of genes selected at the given threshold, on a logarithmic scale.

The weight of each of the four nucleotides is given at each position of the matrix. The weights were defined by the following formula

$$W(b, k)$$

$$= 0.25 \sum_{i=A, C, G, T} \log[(N(b, k) + 0.5)/(N(i, k) + 0.5)],$$

where $N(b, k)$ is the count of nucleotide b at position k . The average weight W on the Bernoulli random sequence equals 0. The base of the logarithm is chosen such that the variance equals 1. In other words, W is the Z-statistics, and the probability of random occurrence of a signal can be assessed using the Gaussian distribution.

The score of a candidate site is the sum of positional weights of the constituent nucleotides:

$$S(b_1, \dots, b_n) = \sum_{k=1}^n W(b_k, k).$$

In these terms, the signal (Fnr box) is palindromic. Hence, the matrix obtained improves the already known consensus sequence TTGATnnnnATCAA [13]. Then, we searched upstream gene regions for potential Fnr boxes. As a result, candidate Fnr boxes were found upstream of 121 genes when the threshold was equal to 4.0. This choice of a threshold leads to the loss of some Fnr-regulated genes. We also searched for Fnr boxes using other values of the threshold, but the first choice appeared to be optimal. When the value of the threshold exceeds 4.0, many known sites are lost (underprediction), whereas when it is less than 4.0, almost all genes have potential Fnr boxes (overprediction). Indeed, consider the probability distribution of scores in the *E. coli* genome (figure). One can see that when the threshold equals 4.0, about 5% of all genes in the genome are selected. At the same time, it appears to be impossible to formally

assess the type 1 and type 2 errors, because we do not know exactly what genes compose the regulon. When constructing the weight matrix, we used sites from 9 out of 121 selected genes (obtained from the database of regulatory sites DPIPinteract). The fact that 12 genes are Fnr-regulated was experimentally confirmed. We applied the standard procedure of comparison between related genomes [14] and found out that the genomes of the considered bacteria (*S. typhi*, *K. pneumoniae*, *Y. pestis*, *H. influenzae*, *V. cholerae*, and *P. aeruginosa*) contain genes orthologous to *fnr*. This suggests that the Fnr regulon of these gamma-proteobacteria is conserved. Orthologs of 121 genes of *E. coli* with potential Fnr boxes upstream were identified in the genomes of other bacteria. We analyzed the upstream regions of the orthologous genes using the weight matrix. The 39 genes having upstream regions with potential Fnr boxes retained in at least three considered genomes (one of them is *E. coli*) are listed in Tables 1 and 2. These genes were divided into three groups. The first group includes 9 genes. Their sites were used to construct the recognition matrix. The second group consists of 12 genes that were experimentally demonstrated to be under regulation by Fnr but the corresponding Fnr boxes were not identified. The third group consists of 18 genes. They have upstream regions with potential Fnr boxes retained in at least two considered genomes in addition to *E. coli*.

In order to estimate the number of false positives, we used a statistical model. The estimates were as follows: $P = 5\%$ is the fraction of genes with candidate sites (figure), $g = 2/3$ is the fraction of orthologs in any two genomes, $n = 4000$ is the average number of genes in a genome. The expected number of false positives can be expressed as follows:

$$C_n^k n g^2 P^3 = (5 \times 6)/(1 \times 2) \times 4000 \times 4/9 \times 125 \times 10^{-6} = 3.3.$$

It is known that the Fnr protein is homologous to the regulator CRP, and their signals resemble one another [2]. Thus, it is possible that a number of the predicted sites are CRP binding sites. Using the weight matrix for CRP boxes [15], we obtained scores of the candidate CRP binding sites of genes from the Fnr regulon. Many known genes from the Fnr regulon (the first and the second groups of Table 2) have potential CRP boxes. In addition, it is known that the genes *ansB* and *tdcA* are under regulation by both CRP and Fnr [13, 16]. CRP also regulates the *mtlA* gene that belongs to the third group [17]. We observed numerous candidate CRP boxes upstream of this gene, and their scores were higher than those of candidate Fnr boxes. The hypothetical *b2503* gene was predicted to be the Fnr-regulated, because the alignment of the upstream regions of the orthologous genes revealed some conservatism in the area of the potential Fnr box.

Apparently, the *aldA* gene also belongs to the Fnr regulon. This gene does not have orthologs in the considered genomes and thus, using the formal approach, cannot be assigned to that regulon. However, there is a potential Fnr box upstream of *aldA* in the *E. coli* genome. In addition, this gene is known to be under regulation of ArcA and CRP [18] (the CRP binding site does not coincide with the predicted Fnr box). The ArcA enzyme controls aerobic respiration and regulates expression of several genes of the Fnr regulon, the *arcA* gene itself is regulated by Fnr [2] (Table 2). Detailed analysis of the ArcA regulon is the subject of another paper.

Thus, we have identified Fnr binding sites in upstream regions of 12 *E. coli* genes regulated by Fnr, and found 17 additional genes that may belong to the Fnr regulon. We have described the Fnr regulons of *S. typhi*, *K. pneumoniae*, *Y. pestis*, *H. influenzae*, *V. cholerae*, and *P. aeruginosa*. Currently we analyze the regulatory system Anr (the ortholog of Fnr) in the genomes of pseudomonads.

ACKNOWLEDGMENTS

This work was partially supported by the Russian Foundation for Basic Research (projects nos. 99-04-48247 and 00-15-99362), INTAS (99-1476), and the Howard Hughes Medical Institute (55000309).

REFERENCES

1. Bauer, C.E., Elsen, S., and Bird, T.H., *Annu. Rev. Microbiol.*, 1999, vol. 53, pp. 495–523.
2. Lynch, A.S. and Lin, E.C.C., *Escherichia coli and Salmonella. Cellular and Molecular Biology*, Neindhard, F.C., Ed., Washington DC: ASM Press, 1996, pp. 1526–1538.
3. Blattner, F.R., Plunkett, G., Bloch, C.A., *et al.*, *Science*, 1997, vol. 77, pp. 1453–1474.
4. Fleischmann, R.D., Adams, M.D., White, O., *et al.*, *Science*, 1995, vol. 269, pp. 496–512.
5. Heidelberg, J.F., Eisen, J.A., Nelson, W.C., *et al.*, *Nature*, 2000, vol. 406, pp. 477–483.
6. Stolver, C.K., Pham, X.Q., Erwin, A.L., *et al.*, *Nature*, 2000, vol. 406, pp. 959–964.
7. Benson, D.A., Karsch-Mizrachi, I., Lipman, D.J., *et al.*, *Nucleic Acids Res.*, 2000, vol. 28, pp. 15–18.
8. Mironov, A.A., Vinokurova, N.P., and Gelfand, M.S., *Mol. Biol.*, 2000, vol. 34, pp. 222–231.
9. Robison, K. and Church, G., *J. Mol. Biol.*, 1998, vol. 284, pp. 241–254.
10. Altschul, S.F., Madden, T.L., Schaffer, A.A., *et al.*, *Nucleic Acids Res.*, 1997, vol. 25, pp. 3389–3402.
11. Bairoch, A. and Apweiler, R., *Nucleic Acids Res.*, 2000, vol. 28, pp. 45–48.
12. Florea, L., Riemer, C., Schwartz, S., *et al.*, *Nucleic Acids Res.*, 2000, vol. 28, pp. 3486–3496.
13. Green, J., Irvine, A.S., Meng, W., and Guest, J.R., *Mol. Microbiol.*, 1996, vol. 19, pp. 125–137.
14. Mironov, A.A. and Gelfand, M.S., *Mol. Biol.*, 1999, vol. 33, pp. 127–132.
15. Gelfand, M.S., Novichkov, P.S., Novichkova, E.S., and Mironov, A.A., *Briefings in Bioinformatics*, 2000, vol. 1, pp. 357–371.
16. Chattopadhyay, S., Wu, Y., and Datta, P., *J. Bacteriol.*, 1997, vol. 179, pp. 4868–4873.
17. Ramseier, T.M. and Saier, M.H., *Microbiology*, 1995, vol. 141, pp. 1901–1907.
18. Pellicer, M.T., Lynch, A.S., De Wulf, P., *et al.*, *Mol. Gen. Genet.*, 1999, vol. 261, pp. 170–176.