

## ORIGINAL PAPER

Sorel Fitz-Gibbon · Anthony J. Choi  
Jeffrey H. Miller · Karl O. Stetter · Melvin I. Simon  
Ronald Swanson · Ung-Jin Kim

## A fosmid-based genomic map and identification of 474 genes of the hyperthermophilic archaeon *Pyrobaculum aerophilum*

Received: October 18, 1996 / Accepted: October 30, 1996

**Abstract** We have constructed a physical map of the approximately 1.7-Mb genome of the hyperthermophilic archaeon *Pyrobaculum aerophilum*. Derived from a 12× coverage genomic fosmid library with an average insert size of 36 Kb, the map consists of a single circular contig of 96 overlapping fosmid clones with 211 markers ordered along them. One hundred of the sequence markers have strong similarities to known genes. Many overlaps were also checked using restriction fingerprint analysis. This map is an important step in the elucidation of the sequence of the entire genome of *Pyrobaculum aerophilum*. To this end we have determined more than 95% of the genome with 15000 random sequences. Each sequence has been screened against the public sequence databases to identify similarities to known genes. We report here a list of the 474 putative genes we have identified.

**Key words** *Pyrobaculum aerophilum* · Genome · Hyperthermophile · Archaea · Eocyte

Communicated by J. Wiegel

S. Fitz-Gibbon · A.J. Choi · J.H. Miller (✉)  
Department of Microbiology and Molecular Genetics and The  
Molecular Biology Institute, University of California at Los Angeles,  
1602 Molecular Sciences Building, 405 Hilgard Avenue, Los Angeles,  
CA 90095-1489, USA  
Tel. +1-310-825-8460; Fax +1-310-206-3088  
e-mail: jhmiller@ewald.mbi.ucla.edu

K.O. Stetter  
Archaeenzentrum, Regensburg University, 93053 Regensburg,  
Germany

M.I. Simon · R. Swanson<sup>1</sup> · U.-J. Kim  
Division of Biology, 147-75, California Institute of Technology,  
Pasadena, CA 91125, USA

Present address:

<sup>1</sup> Recombinant Biocatalysis Inc., Philadelphia, PA, USA

### Introduction

*Pyrobaculum aerophilum* is a rod-shaped hyperthermophilic archaeon isolated from a boiling marine water hole (Völkl et al. 1993). Due to its distinctive characteristics, *P. aerophilum* is both an interesting and a suitable organism for evolutionary and comparative studies. It has an extremely high optimal growth temperature of 98°C, and yet, in contrast to other hyperthermophiles, it is facultatively aerobic. Furthermore, *P. aerophilum* can be plated to form colonies within four days with up to 100% efficiency (Völkl et al. 1993). Since strictly anaerobic metabolism predominates at the upper temperature limits of life (Robb and Place 1995), cultivation of hyperthermophiles has been difficult and consequently they are only poorly described. Based on 16S rRNA analysis, hyperthermophiles are the living organisms with the greatest similarity to ancient organisms (Stetter 1992; Woese 1987). Thus hyperthermophiles are of great interest and, due to its ease of cultivation, *P. aerophilum* is an ideal candidate for development as a model hyperthermophilic organism. We plan the construction of a genetic system which will allow for mutagenesis and reverse genetic studies, and for the analysis of gene function. Determining the entire genome sequence will simplify the development of a genetic system for *P. aerophilum* and will provide valuable data not only for research on hyperthermophily but also for phylogeny research and the elucidation of the least studied domain of life, Archaea.

Our strategy for sequencing the entire genome of *P. aerophilum* has been primarily a random “shotgun” approach. We have sequenced random fragments from the genome, used computer algorithms to compare the sequences, and then assembled the overlapping sequences into contigs. While the “shotgun” sequencing approach (Fleischmann et al. 1995) is an efficient way to obtain most of the sequence, it becomes impractical to complete the sequence by this method. To close the final gaps in the sequence, a directed approach is desirable. We have developed a genome map in order to facilitate a directed approach to sequence completion, to confirm our sequence

contigs, and to aid in resolving repetitive and otherwise difficult regions.

## Results

We have constructed a genomic map consisting of 211 markers on 96 overlapping large clones. The large inserts (~36 Kb) are carried in the fosmid vector pFOS1. This vector uses the *E. coli* F-factor origin of replication and thus is kept at a low copy number. This vector has been shown to maintain large insert DNA fragments with relatively high stability (Kim et al. 1992). The entire genomic fosmid library consists of 768 clones, which corresponds to a 12× genomic coverage based on an estimated genome size of 1.7 Mb. This level of redundancy was sufficient to cover the whole genome without gaps. By hybridizing the library with various probes we were able to place markers on the map and determine the overlaps of the fosmids.

### Construction and screening of fosmid library

Genomic DNA was partially digested with *Sau3AI* before ligation to the *Bam*HI cloning site of the fosmid arms. The ligation reaction contained a large molar excess of fosmid arms over genomic DNA in order to minimize the chance of chimera formation. Packaging was done using the Stratagene GigaPack XL kit which preferentially packages large inserts. The average insert size for this library was 36 Kb as determined by *Not*I digestion and pulsed field gel electrophoresis (not shown). The entire library was gridded onto single small (8 cm × 12 cm) nylon membranes using the Biomek 1000 laboratory workstation. The library was screened by a variety of probes that are described in detail later.

### Fosmid contig assembly and verification

Three methods were used to identify overlapping fosmids. First, small inserts from the pUC18 clones were amplified by PCR and hybridized to high-density colony membranes of the large insert fosmid library. As part of our large-scale random sequencing project, both ends of these pUC18 inserts had previously been sequenced and compared to the public sequence databases. Of the 146 markers on the map, 112 were putative genes. These include ribosomal genes, DNA and RNA polymerase genes, tRNA genes, various repair genes, *recA* homologs, and numerous metabolic genes (Fig. 1). Analysis of these data revealed several ambiguous regions, possibly due to repetitive sequences in the markers. These regions were clarified by a fingerprinting method consisting of computer assisted analysis of the banding patterns of fosmids digested with *Ban*I and *Msp*I and labeled at the *Ban*I ends (Coulson et al. 1986; Sulston et al. 1988). Twenty-four overlaps were confirmed by this method (Fig. 1), yielding 10 clear contigs. The final 10 gaps were closed using RNA probes transcribed from the T7 and

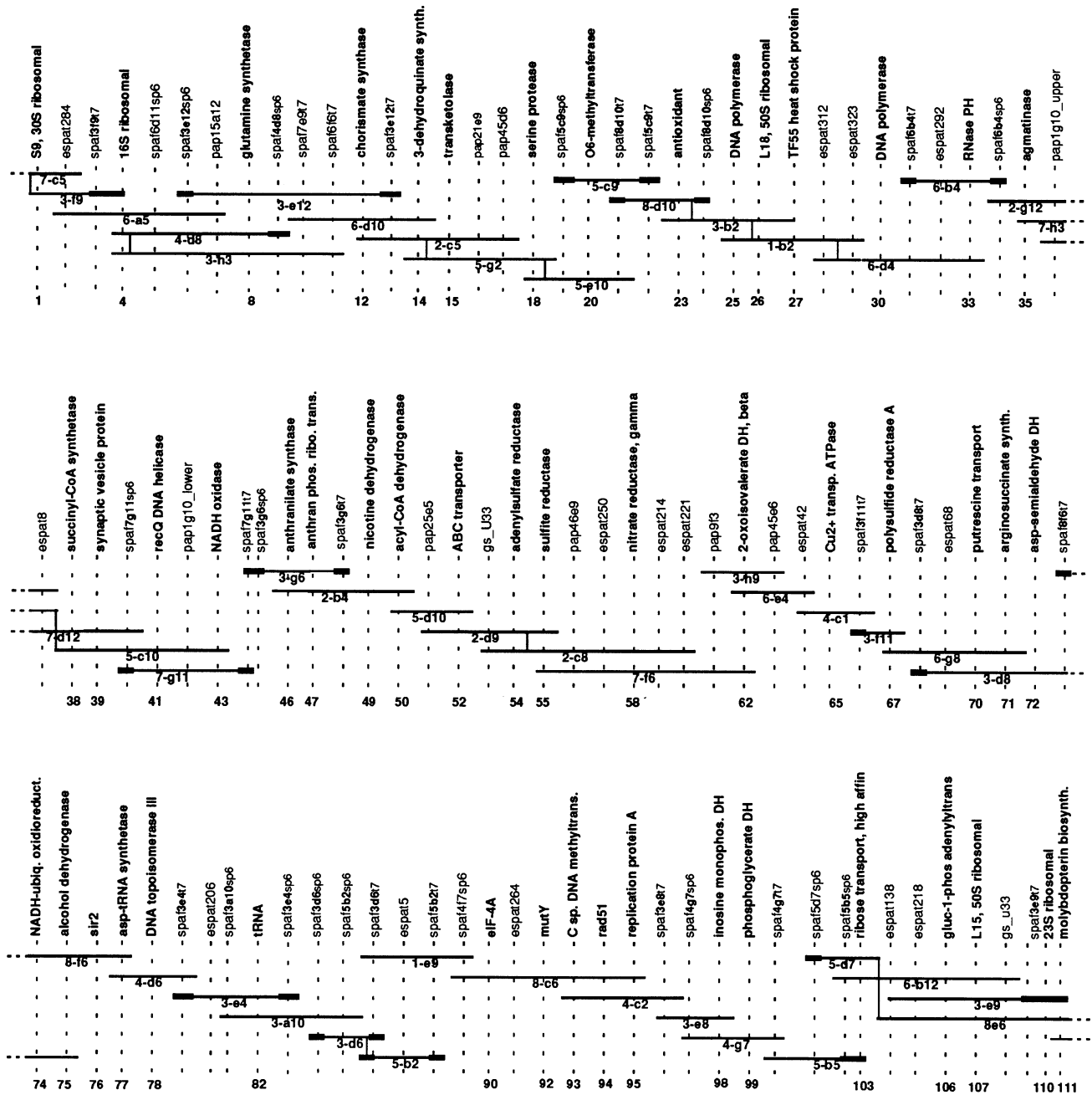
Sp6 promoters flanking the fosmid vector cloning site. These RNA probes were radiochemically labeled and hybridized back to the whole fosmid library. Prior to generating the probe, the fosmid DNA was digested with *Hinc*II, which yields DNA fragments averaging approximately 1000 bp. *Hinc*II was chosen over a less frequent cutter, to keep the length of the labeled transcript short. This was important for minimizing the inclusion of repetitive sequences in the riboprobes. We found larger probes often gave ambiguous results. Using this method we successfully identified fosmids spanning the 10 remaining gaps and confirmed 38 other overlaps (Fig. 1). In three cases, a second round of probing was necessary to cross the gaps, using probes from the ends of the fosmids identified in the first round. An overlapping set of 96 fosmids was identified, giving a redundancy of approximately 2×. From the average size of the 96 clones in the map and the overall redundancy that they represent, we estimate the genome size of this microbe to be 1.7 Mb.

### Genome sequence

The map reported here will be used as an aid to determining the entire genome sequence of *P. aerophilum*. The sequencing effort has so far yielded more than 15 000 random single-pass sequences which were determined using fluorescent multiplex automated sequencers (ABI 373, Applied Biotechnology). Base calls were made from the ABI trace data using the Phred software package (Brent Ewing). This software not only assigns bases from the trace data but also assigns a quality value for each base call. This factor becomes extremely important when the sequences are assembled, obviating the need for extensive manual editing of contigs. The Phrap software package (Phil Green) makes use of the quality information provided by the Phred software, while other assembly programs allow regions of low-quality data to disrupt the consensus sequence. The Phrap software then uses a novel method of assembly to create contigs which are a mosaic of the highest quality parts of reads (Table 1).

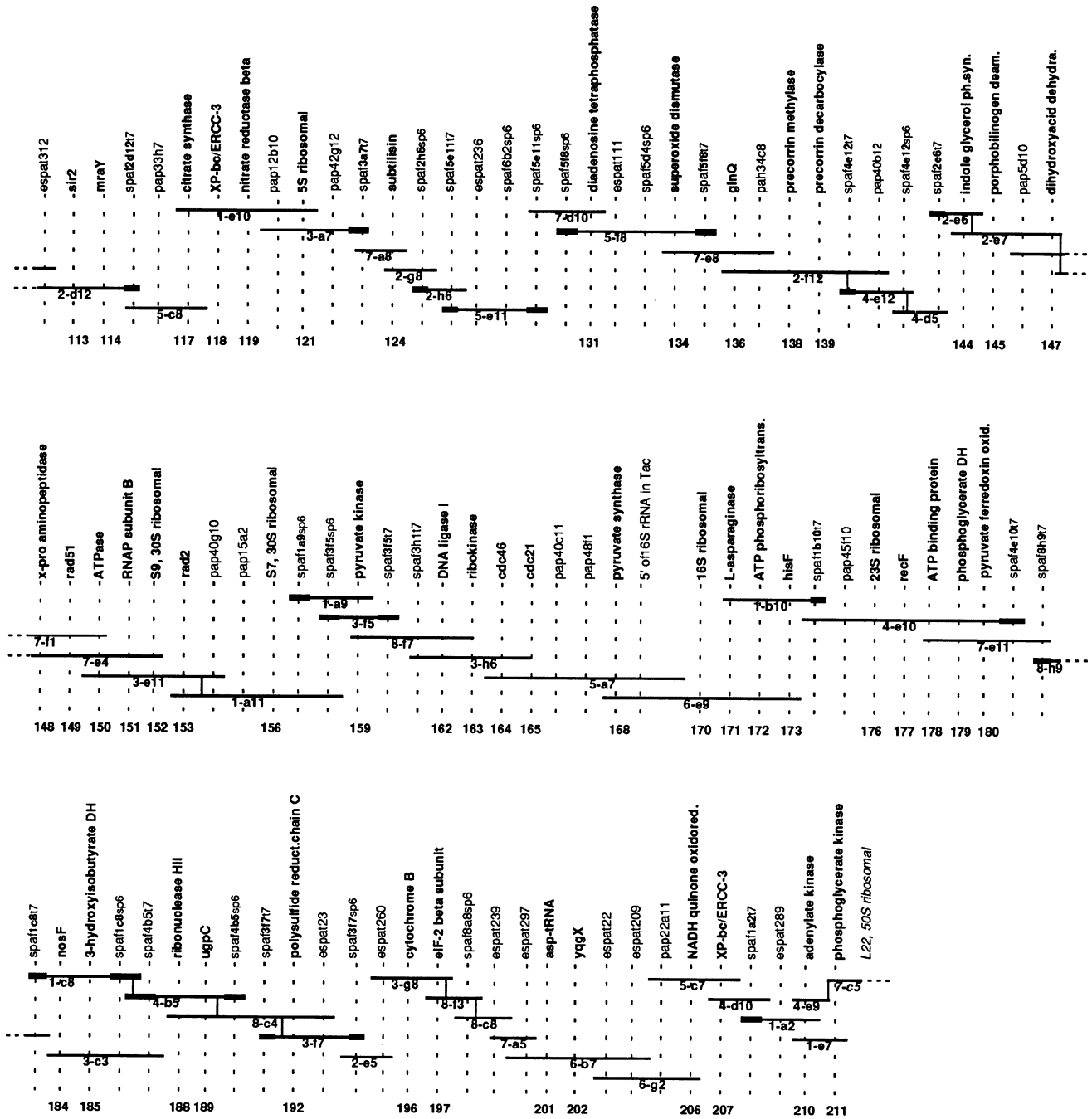
After having completed the random sequencing phase of the project, we have slightly over 800 contigs and thus approximately 800 gaps in our sequence. Poisson analysis suggests an average gap size of less than 200 bases long. Several methods are being used in parallel in order to move quickly to a finished sequence. Where possible, new sequences are being derived from the adjacent regions of partially sequenced pUC18 clones. In most cases an oligo primer has to be designed and synthesized for this purpose. On a small number of clones, a deletion subcloning method (unpublished) is being used to reach the unsequenced region at the center of the insert while still using the standard priming sites on the vector.

Of the 800 gaps in our sequence there are approximately 200 gaps which are not represented in the pUC18 library. The fosmid map and templates will be used to extend the sequences across these gaps. The fosmid map includes 146 sequence markers which have been assembled into the se-



**Fig. 1.** The horizontal lines represent the inserts of fosmid clones covering the entire circular genome of *Pyrobaculum aerophilum*. Fosmid designations are indicated immediately below the horizontal lines. Vertical dashed lines represent sequence markers placed along the

genome via hybridization experiments. Where the sequence tag has similarity to a known gene, the gene name is written in **bold**. All markers were tentatively identified as the listed genes by their similarity to RNA or amino acid sequences in the public databases.



**Fig. 1.** (continued) Details such as blast scores, accession numbers and full names are shown in Table 1a by the corresponding number. Riboprobes are labeled "spaf...t7" or "spaf...sp6". **Bold regions of**

*the horizontal lines* represent the regions from which the riboprobes were transcribed. *Vertical solid lines* represent the 24 overlaps confirmed by fingerprinting

**Table 1a.** Similarity data for 100 putative genes which were placed on the fosmid map (Fig. 1). All matches have a spurious match probability of less than  $10^{-5}$  except for eukaryotic IF-4A (90). The numbers (*Map #*) correspond to the markers' position on the fosmid map (Fig. 1)

Map #	Hit identification	Probability (-E)	Accession	Organism
1	S9, 30S ribosomal	30	sp P39468	<i>Sulfolobus acidocaldarius</i>
4	16S ribosomal	54	gb M35966	<i>Thermoproteus tenax</i>
8	glutamine synthetase	35	sp P36205	<i>Thermotoga maritima</i>
12	chorismate synthase	37	pir S13070	<i>Synechocystis sp.</i>
14	3-dehydroquinate synthase	62	sp P07639	<i>Escherichia coli</i>
15	transketolase	26	gi 1149710	<i>Clostridium perfringens</i>
18	serine protease	26	sp P29139	<i>Bacillus polymyxa</i>
20	O6-methyltransferase	9	sp P16455	<i>Homo sapiens</i>
23	antioxidant	55	gi 1045502	<i>Sulfolobus sp.</i>
25	DNA polymerase	150	gi 807830	<i>Pyrodictium occultum</i>
26	L18, 50S ribosomal	52	sp P14033	<i>Methanococcus vannielii</i>
27	TF55 heat shock protein	79	sp P28488	<i>Sulfolobus shibatae</i>
30	DNA polymerase	200	gi 807828	<i>Pyrodictium occultum</i>
33	RNase PH	15	sp P37939	<i>Mycobacterium leprae</i>
35	agmatinase	16	sp P37819	<i>Streptomyces clavuligerus</i>
38	succinyl-CoA synthetase alpha&beta	86	sp P25126	<i>Thermus flavus</i>
39	synaptic vesicle protein	37	pir S34961	<i>Rattus norvegicus</i>
41	DNA Helicase, <i>recQ</i>	22	sp P30015	<i>Escherichia coli</i>
43	NADH oxidase	20	sp P37061	<i>Enterococcus faecalis</i>
46	anthranilate synthase	91	sp Q06128	<i>Sulfolobus solfataricus</i>
47	anthranilate phosphoribosyltransferase	69	gi 149037	<i>Bacillus pumilus</i>
49	nicotine dehydrogenase	44	pir S37570	<i>Arthrobacter nicotinovorans</i>
50	acyl-CoA dehydrogenase	33	sp P46703	<i>Mycobacterium leprae</i>
52	ABC transporter	11	sp P39326	<i>Escherichia coli</i>
54	adenylyl-sulphate reductase alpha-subunit	44	gi 1183905	<i>Desulfovibrio vulgaris</i>
55	sulfite reductase	56	pir S27479	<i>Archaeoglobus fulgidus</i>
58	nitrate reductase gamma subunit	24	sp P42177	<i>Bacillus subtilis</i>
62	2-oxoisovalerate dehydrogenase, beta	82	sp P37941	<i>Bacillus stearothermophilus</i>
65	Cu <sup>2+</sup> transporting ATPase	31	sp P32113	<i>Enterococcus hirae</i>
67	polysulfide reductase chain <i>psrA</i>	29	sp P31075	<i>Wolinella succinogenes</i>
70	putrescine transport	10	sp P31136	<i>Escherichia coli</i>
71	argininosuccinate synthase	109	sp P13256	<i>Methanococcus vannielii</i>
72	aspartate-semialdehyde dehydrogenase	60	sp P41394	<i>Leptospira interrogans</i>
74	NADH-ubiquinone oxidoreductase	40	sp P42026	<i>Bos taurus</i>
75	alcohol dehydrogenase	28	sp P39462	<i>Sulfolobus solfataricus</i>
76	<i>sir2</i>	36	gi 845686	<i>Staphylococcus aureus</i>
77	asp-tRNA synthetase-prokaryote	62	gi 1146247	<i>Bacillus subtilis</i>
78	DNA topoisomerase III	41	gi 1292913	<i>Homo sapiens</i>
82	tRNA	12	emb X05070	<i>Thermoproteus tenax</i>
90	eukaryotic IF-4A	5	gi 1072122	<i>Homo sapiens</i>
92	<i>mutY</i> = adenine glycosylase	37	sp P29588	<i>Methanobacterium thermo.</i>
93	cytosine specific DNA methyltransferase	14	gi 1165245	<i>Neisseria gonorrhoeae</i>
94	<i>rad51</i>	24	pir S37673	<i>Homo sapiens</i>
95	replication protein A	7	sp P05682	<i>Agrobacterium rhizogenes</i>
98	inosine monophosphate dehydrogenase	7	pir B48868	<i>Methanosarcina thermophila</i>
99	phosphoglycerate dehydrogenase	55	gi 1146196	<i>Bacillus subtilis</i>
103	ribose transport, high affinity	23	gi 290599	<i>Escherichia coli</i>
106	glucose-1-phosphate adenylyltransferase	21	gi 699151	<i>Mycobacterium leprae</i>
107	L15, 50S ribosomal	8	sp P14032	<i>Methanococcus vannielii</i>
110	23S ribosomal RNA	228	emb X05480	<i>Desulfurococcus mobilis</i>
111	molybdopterin biosynthesis	44	gi 1002858	<i>Anabaena sp.</i>
113	<i>sir2</i>	21	gi 845686	<i>Staphylococcus aureus</i>
114	phospho-N-acetylmuramoyl-pentapeptide-transferase	7	sp Q03521	<i>Bacillus subtilis</i>
117	citrate synthase	64	sp P39120	<i>Bacillus subtilis</i>

Table 1a. Continued

Map #	Hit identification	Probability (-E)	Accession	Organism
118	DNA excision repair helicase, XP-bc/ERCC-3	14	gi 902048	<i>Arabidopsis thaliana</i>
119	nitrate reductase beta	157	sp P42176	<i>Bacillus subtilis</i>
121	5S ribosomal RNA	40	gb M16530	<i>Sulfolobus</i> sp.
124	subtilisin	61	sp P29139	<i>Bacillus polymyxa</i>
131	diadenosine tetraphosphatase (mutT-like)	18	gi 1054947	<i>Sus scrofa</i>
134	superoxide dismutase (Mn)	19	sp Q08713	<i>Sulfolobus acidocaldarius</i>
136	glutamine transport protein <i>glnQ</i>	46	sp P10346	<i>Escherichia coli</i>
138	precorrin methylase	38	sp Q05630	<i>Salmonella typhimurium</i>
139	precorrin decarboxylase	18	sp Q05632	<i>Salmonella typhimurium</i>
144	indole-3-glycerol-phosphate synthase	19	sp Q06121	<i>Sulfolobus solfataricus</i>
145	porphobilinogen deaminase	52	sp P16616	<i>Bacillus subtilis</i>
147	dihydroxyacid dehydrase	43	sp P31959	<i>Clostridium pasteurianum</i>
148	X-PRO aminopeptidase	41	gi 1303901	<i>Bacillus subtilis</i>
149	<i>rad51</i>	91	gi 1054624	<i>Xenopus laevis</i>
150	ATPase, alpha, membrane associated	126	sp P09639	<i>Sulfolobus acidocaldarius</i>
151	RNA polymerase B	94	sp P11513	<i>Sulfolobus acidocaldarius</i>
152	S9, 30S ribosomal	30	sp P39468	<i>Sulfolobus acidocaldarius</i>
153	<i>rad2</i>	24	sp P26793	<i>Saccharomyces cerevisiae</i>
156	S7, 30S ribosomal	55	sp P41206	<i>Desulfurococcus mobilis</i>
159	pyruvate kinase	30	gi 1041097	<i>Bacillus psychrophilus</i>
162	DNA ligase	32	sp Q02093	<i>Desulfurolobus ambivalens</i>
163	DNA ligase	122	sp Q02093	<i>Caenorhabditis elegans</i>
164	DNA replication licensing factor (MCM2, <i>cdc46</i> )	98	gi 852053	<i>Drosophila melanogaster</i>
165	<i>cdc21</i>	5	sp P33991	<i>Homo sapiens</i>
168	pyruvate synthase	43	pir  S22396	<i>Halobacterium halobium</i>
170	16S ribosomal	47	gb M38637	<i>Thermoplasma acidophilum</i>
171	L-asparaginase	34	gi 496102	<i>Lupinus albus</i>
172	ATP phosphoribosyltransferase	21	gi 1117923	<i>Yarrowia lipolytica</i>
173	<i>hisF</i> (cyclase) (histidine biosynthesis)	70	sp Q02133	<i>Lactococcus lactis</i>
176	23S ribosomal RNA	26	gb M86626	<i>Pyrobaculum occultum</i>
177	<i>recF</i>	8	sp P24900	<i>Salmonella typhimurium</i>
178	ATP-binding protein	23	gi 1184189	<i>Escherichia coli</i>
179	phosphoglycerate dehydrogenase	37	sp P37666	<i>Escherichia coli</i>
180	pyruvate ferredoxin oxidoreductase	55	gi 1197358	<i>Pyrococcus furiosus</i>
184	Cu transporting ATP-binding <i>nosF</i>	11	sp P19844	<i>Pseudomonas stutzeri</i>
185	3-hydroxyisobutyrate DH	47	gi 1001605	<i>Synechocystis</i> sp.
188	RNase HII	9	pir  S53S08	<i>Saccharomyces cerevisiae</i>
189	nucleotide-binding protein <i>ugpC</i>	80	sp P10907	<i>Escherichia coli</i>
192	polysulfide reductase chain C, <i>psrC</i>	28	sp P31076	<i>Wolinella succinogenes</i>
196	cytochrome <i>b</i>	19	gi 927525	<i>Sulfolobus acidocaldarius</i>
197	eIF-2 beta subunit	6	sp P41375	<i>Drosophila melanogaster</i>
201	asp-tRNA	9	gb L07321	<i>Thermococcus celer</i>
202	<i>yqgX</i>	11	gi 1303871	<i>Bacillus subtilis</i>
206	NADH-quinone oxidoreductase	25	sp P29925	<i>Paracoccus denitrificans</i>
207	DNA excision repair helicase, XP-bc/ERCC-3	16	gi 902048	<i>Arabidopsis thaliana</i>
210	adenylate kinase	29	gi 1086550	<i>Methanococcus jannaschii</i>
211	phosphoglycerate kinase	69	sp P20971	<i>Methanothermobacter ferredoxin</i>

Approximately 15 000 sequences were screened individually against the nonredundant set of public databases using BLAST (Altschul et al. 1990; States et al. 1991) with default settings. The spurious match probabilities (as determined by the BLAST algorithm) are reported as orders of magnitude, e.g.,  $3.6E-24$  becomes 24. All of the matches reported have spurious match probabilities of less than  $10^{-5}$ ; most of the matches reported have spurious match probabilities of less than  $10^{-10}$ . Detailed analyses of matches have not been performed. Duplicate copies of genes are not known or reported. The identification and

categorization of the matches should be considered tentative. Assignments are based upon Monica Riley's categories (Riley 1993; see also Bult et al. 1996). The cited names of the organisms are as they appear in the data entry of the gene bank and thus do not necessarily reflect the present correct nomenclature of the listed organisms. Upon completion of the sequence, extensive analysis will be performed and reported. We are developing a web site on which our data will be publicly available. E-mail to sorel@ewald.mbi.ucla.edu for web site information or for help with specific data.

**Table 1b.** Similarity data for 474 *Pyrobaculum aerophilum* random sequence tags

ID #	Hit identification	Probability (-E)	Accession	Organism
<b>Amino acid biosynthesis</b>				
1	glutamate-1-semialdehyde 2,1 aminomutase	24	gp X53695	<i>Synechococcus</i> sp.
2	NMHase	27	gp A17961	<i>Arthrobacter</i> sp.
3	ornithine aminotransferase	12	gp X81802	<i>Bacillus subtilis</i>
4	S-adenosyl-L-homocysteine hydrolase	81	gp Z50174	<i>Sulfolobus solfataricus</i>
5	2-isopropylmalate synthase	42	sp P05342	<i>Lactococcus lactis</i>
6	2-oxoisovalerate dehydrogenase, beta	43	sp P37941	<i>Bacillus subtilis</i>
7	3-phosphoshikimate 1-carboxyvinyltrans.	13	sp P24497	<i>Klebsiella pneumoniae</i>
8	anthranilate phosphoribosyltransferase	45	gp M65060	<i>Methanobacterium thermo.</i>
9	anthranilate synthase I	68	sp Q08653	<i>Thermotoga maritima</i>
10	anthranilate synthase II	52	gp M98048	<i>Sulfolobus solfataricus</i>
11	chorismate mutase	9	sp P27603	<i>Pseudomonas stutzeri</i>
12	chorismate synthase	34	gp X67516	<i>Synechocystis</i> sp.
13	dihydroxy-acid dehydratase	31	sp P39522	<i>Saccharomyces cerevisiae</i>
14	<i>hisF</i>	45	sp P26721	<i>Azospirillum brasilense</i>
15	<i>hisH</i>	11	sp Q02132	<i>Lactococcus lactis</i>
16	histidinol-phosphate transaminase	13	pir A30270	<i>Escherichia coli</i>
17	imidazoleglycerol-phosphate dehydratase	20	sp P06633	<i>Saccharomyces cerevisiae</i>
18	indole-3-glycerol-phosphate synthase	22	sp Q06121	<i>Sulfolobus solfataricus</i>
19	ATP phosphoribosyltransferase	15	sp Q02129	<i>Lactococcus lactis</i>
20	ketol-acid reductoisomerase	51	pir A47037	<i>Corynebacterium glutam.</i>
21	kynurenine/alpha-aminoacidipate aminotrans.	22	gp Z50144	<i>Rattus norvegicus</i>
22	N-acetylornithine aminotransferase	47	pir S44189	<i>Anabaena</i> sp.
23	phosphoribosylanthranilate isomerase ( <i>trpF</i> )	7	gp X17149	<i>Vibrio parahaemolyticus</i>
24	phosphoribosylpyrophosphate synthetase	15	gp M76553	<i>Leishmania donovani</i>
25	tryptophan synthase, alpha	5	gp M65060	<i>Methanobacterium thermo.</i>
26	tryptophan synthase, beta	6	sp P19868	<i>Bacillus stearothermophilus</i>
27	acetylglutamate kinase	8	sp P36840	<i>Bacillus subtilis</i>
28	argininosuccinate synthetase (Ass)	51	sp P13256	<i>Methanococcus vannielii</i>
29	aspartate semialdehyde dehydrogenase	35	pir A44846	<i>Leptospira interrogans</i>
30	aspartate aminotransferase	12	gp D50624	<i>Streptomyces virginiae</i>
31	aspartate beta-semialdehyde dehydrogenase	27	gp M77500	<i>Leptospira borgpetersenii</i>
32	aspartate kinase-homoserine dehydrogenase	15	gp L33912	<i>Zea mays</i>
33	ATP phosphoribosyltransferase	20	gp U07830	<i>Schizosaccharomyces pombe</i>
34	cystathionine gamma-synthase	17	sp P00935	<i>Escherichia coli</i>
35	L-isoaspartyl protein carboxyl methyltrans.	32	gp M63493	<i>Escherichia coli</i>
36	lysine-sensitive aspartokinase III	20	sp P08660	<i>Escherichia coli</i>
37	dihydrodipicolinate synthase ( <i>mosA</i> )	15	sp Q07607	<i>Rhizobium meliloti</i>
38	succinyl-diaminopimelate desuccinylase	8	sp P24176	<i>Escherichia coli</i>
39	homoserine kinase ( <i>thrB</i> )	16	gp Y00522	<i>Calothrix</i> sp.
40	3-isopropylmalate dehydrogenase	37	gp U07980	<i>Bos taurus</i>
41	acetolactate synthase, large subunit	41	sp P37251	<i>Bacillus subtilis</i>
42	branched-chain-amino-acid-transaminase	39	pir S30668	<i>Escherichia coli</i>
43	branched-chain amino acid transport BRAF	17	sp P21629	<i>Pseudomonas aeruginosa</i>
44	branched-chain amino acid transport BRAG	23	sp P21630	<i>Pseudomonas aeruginosa</i>
45	ketol-acid reductoisomerase	56	gp L03181	<i>Bacillus subtilis</i>
46	argininosuccinate lyase	14	sp P11447	<i>Escherichia coli</i>
47	argininosuccinate synthetase	16	gp M21315	<i>Methanococcus vannielii</i>
48	ferredoxin-dependent glutamate synthase	32	gp U03006	<i>Spinacia oleracea</i>
49	glutamate-1-semialdehyde 2,1 aminomutase	25	pir A35789	<i>Hordeum vulgare</i>
50	glutamate dehydrogenase	34	gp L19995	<i>Thermococcus litoralis</i>
51	glutamine synthetase	46	gp X60160	<i>Thermotoga maritima</i>
52	L-glu-D-fructose-6-phosphate amidotrans.	31	gp U21932	<i>Bacillus subtilis</i>
53	NADP-specific glutamate dehydrogenase	52	sp P39475	<i>Sulfolobus shibatae</i>
54	aspartate carbamoyltransferase	14	sp P19936	<i>Serratia marcescens</i>
55	hydroxymethyltrans.,3-methyl-2-oxobut.	46	sp P31057	<i>Escherichia coli</i>
56	serine hydroxymethyltransferase	24	gp U02131	<i>Mycoplasma genitalium</i>
57	serine O-acetyltransferase	12	sp P05796	<i>Escherichia coli</i>
58	thioredoxin reductase	19	pir S38988	<i>Eubacterium acidaminophilum</i>
<b>Biosynthesis of cofactors, prosthetic groups, and carriers</b>				
59	3-oxoacyl [acyl-carrier protein] reductase	21	sp P27582	
60	<i>adgA</i>	18	gp X59399	<i>Rhodobacter capsulatus</i>
61	cobyric acid a,c-diamide synthase	19	sp P21632	<i>Pseudomonas denitrificans</i>
62	dihydroflavonol 4-reductase	24	pir S38474	<i>Lycopersicon esculentum</i>
63	mRNA for dihydroflavonol-4-reductase	12	gp X15536	<i>Antirrhinum majus</i>
64	<i>pqqF/A/B/C</i>	9	gp X87299	<i>Pseudomonas fluorescens</i>
65	precorrin 3 methylase	43	sp Q05590	<i>Salmonella typhimurium</i>
66	precorrin-8W decarboxylase	10	sp Q05632	<i>Salmonella typhimurium</i>
67	GTP cyclohydrolase I	27	gp X85954	<i>Campylobacter jejuni</i>

Table 1b. *Continued*

ID #	Hit identification	Probability (-E)	Accession	Organism
68	GTP cyclohydrolase II ( <i>ribA</i> )	14	sp P17620	<i>Bacillus subtilis</i>
69	porphobilinogen deaminase	23	gp M57676	<i>Bacillus subtilis</i>
70	uroporphyrin-III C-methyltransferase	14	sp P29928	<i>Bacillus megaterium</i>
71	lipoate biosynthesis A	19	gp U32688	<i>Haemophilus influenzae</i>
72	thiamine biosynthetic	7	gp U17350	<i>Zea mays</i>
73	<i>chLE/N</i>	22	gp M21151	<i>Escherichia coli</i>
74	molybdopterin biosynth. ( <i>moeA</i> )	15	sp P12281	<i>Escherichia coli</i>
75	molybdenum cofactor biosynthesis B	32	pir S31880	<i>Escherichia coli</i>
76	riboflavin synthase, beta	14	gp U32810	<i>Haemophilus influenzae</i>
<b>Cell envelope</b>				
77	A10	26	gp L21027	<i>Mus musculus</i>
78	erythromycin biosynth. sensory transductn	19	sp P39623	<i>Bacillus subtilis</i>
79	gamma-aminobutyrate permease	9	gp U31756	<i>Bacillus subtilis</i>
80	laminin receptor homolog	25	gp S35960	<i>Homo sapiens</i>
81	transmembrane	24	gp U00039	<i>Escherichia coli</i>
82	ankyrin	17	gp L35601	<i>Drosophila melanogaster</i>
83	integral membrane phosphoprotein band 7.2b	41	gp U17297	<i>Mus musculus</i>
84	polysulfide reductase A	7	sp P31075	<i>Wolinella succinogenes</i>
85	proline/betaine transporter	37	sp P30848	<i>Escherichia coli</i>
86	oligopeptide permease	7	gp X05491	<i>Salmonella typhimurium</i>
87	membrane spanning	26	gp X77636	<i>Bacillus subtilis</i>
88	<i>n</i> -acetyl-gamma-glutamyl-phosphate reduct	12	sp Q07906	<i>Bacillus stearothermophilus</i>
<b>Cellular processes</b>				
89	carbon dioxide binding	25	gi 600730	<i>Brucella melitensis</i>
90	carbonic anhydrase	11	sp P17582	<i>Escherichia coli</i>
91	ERV operon: <i>frvX</i>	15	sp P32153	<i>Escherichia coli</i>
92	<i>sua5</i>	43	gp Z38002	<i>Saccharomyces cerevisiae</i>
93	<i>ftsYEX</i> (cell division control)	30	gp X04398	<i>Escherichia coli</i>
94	( <i>fc3</i> ) <i>cpn60</i>	6	emb X75420	<i>Plasmodium falciparum</i>
95	heat-shock	54	gp D29672	<i>Pyrococcus sp.</i>
96	signal recognition particle	42	sp P27414	<i>Sulfolobus acidocaldarius</i>
<b>Central intermediary metabolism</b>				
97	2-hydroxyhepta-2,4-diene-1,7-dioate isom.	19	sp P37352	<i>Escherichia coli</i>
98	alpha-amylase ( <i>amyA</i> )	12	gp L13279	<i>Escherichia coli</i>
99	enoyl-CoA hydratase	26	gp X79899	<i>Clostridium difficile</i>
100	glyoxysomal citrate synthase	16	gp D38132	<i>Cucurbita sp.</i>
101	inorganic pyrophosphatase	28	gp X83728	<i>Nicotiana tabacum</i>
102	malate oxidoreductase	39	sp P16468	<i>Bacillus stearothermophil.</i>
103	nodulation ATP-binding I	8	pir S27496	<i>Bradyrhizobium japonicum</i>
104	pyruvate-flavodoxin oxidoreductase ( <i>nifj</i> )	10	sp P03833	<i>Klebsiella pneumoniae</i>
105	ribokinase	14	sp P05054	<i>Escherichia coli</i>
106	desulfoviridin, gamma	18	gp L05610	<i>Desulfovibrio vulgaris</i>
<b>Energy metabolism</b>				
107	3-hydroxybutyryl-CoA dehydrogenase	11	pir A43723	<i>Clostridium acetobutylicum</i>
108	3-hydroxyisobutyrate dehydrogenase	25	sp P23523	<i>Escherichia coli</i>
109	3-oxoacyl-[acyl-carrier-protein] reductase	14	sp P28643	<i>Cuphea lanceolata</i>
110	4-coumarate-CoA ligase	20	gp U12012	<i>Pinus taeda</i>
111	5-oxopent-3-ene-1,2,5-tricarbox. decarbox.	20	gp X75028	<i>Escherichia coli</i>
112	acetyl-CoA synthetase	36	gp M63968	<i>Methanotherx soehngenii</i>
113	adenylsulfate reductase	28	pir S18928	<i>Archaeoglobus fulgidus</i>
114	subtilisin	61	sp P29139	<i>Bacillus polymyxa</i>
115	aminopeptidase P	25	gi 1046027	<i>Mycoplasma genitalium</i>
116	<i>N</i> -acetylglutamate-g-semialdehyde DH ( <i>argC</i> )	15	sp Q07906	<i>Bacillus stearothermophil.</i>
117	ATPase	87	gp J04836	<i>Methanosarcina barkeri</i>
118	<i>atrP</i>	13	gp X86160	<i>Escherichia coli</i>
119	beta-lactamase	9	sp P00811	<i>Escherichia coli</i>
120	carbamoylphosphate synthetase, large	31	gi 1304392	<i>Sulfolobus solfataricus</i>
121	cellulase	26	gp L32742	<i>Caldocellum saccharolyticum</i>
122	cephalosporin C acylase	9	gp A17015	<i>Pseudomonas diminuta</i>
123	chloroplast ATPase	26	gp X60752	chloroplast <i>Odontella sine.</i>
124	dihydrolipoamide dehydrogenase	27	sp P11959	<i>Bacillus stearothermophil.</i>
125	DR-nm23	16	gp U29656	<i>Homo sapiens</i>
126	enolase	17	sp P29201	<i>Haloarcula</i>
127	<i>fadBA</i> operon (fatty acid oxidizing)	8	gb M74164	<i>Escherichia coli</i>
128	formate dehydrogenase	9	gp X54057	<i>Wolinella succinogenes</i>
129	fumarate reductase	25	pir S34619	<i>Thermoplasma acidophilum</i>
130	glucosamine-fruc.-6-phos. aminotransferase	24	sp P39754	<i>Bacillus subtilis</i>



Table 1b. Continued

ID #	Hit identification	Probability (-E)	Accession	Organism
131	glyceraldehyde 3-phosphate dehydrogenase	49	sp P19315	<i>Methanobacterium formic.</i>
132	indolepyruvate decarboxylase	19	sp P2323	<i>Enterobacter cloacae</i>
133	L-isoaspartyl protein carboxyl methyltrans.	24	sp P24206	<i>Escherichia coli</i>
134	lactaldehyde dehydrogenase	11	gp M64541	<i>Escherichia coli</i>
135	lipoic acid metabolism ( <i>lipB</i> )	37	sp P30976	<i>Escherichia coli</i>
136	membrane-associated ATPase alpha	49	gp J03218	<i>Sulfolobus acidocaldarius</i>
137	membrane-associated ATPase beta	58	sp P13052	<i>Sulfolobus acidocaldarius</i>
138	<i>mmgC</i>	15	gp U29084	<i>Bacillus subtilis</i>
139	NADH oxidase	25	sp P37061	<i>Enterococcus faecalis</i>
140	NADH-plastoquinone oxidoreductase	26	sp P27724	<i>Synechocystis sp.</i>
141	<i>nap57</i>	35	gp Z34922	<i>Rattus norvegicus</i>
142	nitrate reductase, alpha ( <i>narZ</i> )	68	sp P19319	<i>Escherichia coli</i>
143	nitrate reductase, beta ( <i>narY</i> )	58	sp P42176	<i>Bacillus subtilis</i>
144	nitrate reductase, gamma	9	gp Z49884	<i>Bacillus subtilis</i>
145	nitrate transporter	9	sp P38044	<i>Synechococcus sp</i>
146	periplasmic divalent cation tolerance ( <i>cutA</i> )	17	sp P36654	<i>Escherichia coli</i>
147	phospho-2-dehydro-3-deoxyheptonate aldol.	29	pir S21418	<i>Bacillus subtilis</i>
148	porphobilinogen deaminase	25	gp M95623	<i>Homo sapiens</i>
149	porphobilinogen synthase ( <i>hem B</i> )	48	pir S42531	<i>Synechococcus sp.</i>
150	precorrin decarboxylase	9	sp Q05632	<i>Salmonella typhimurium</i>
151	precorrin methylase	22	sp Q05630	<i>Salmonella typhimurium</i>
152	precorrin-3 methylase	17	sp P21922	<i>Pseudomonas denitrificans</i>
153	glucose-1-phosphate cytidyltransferase ( <i>PsaIp</i> )	11	gp U19608	<i>Saccharomyces cerevisiae</i>
154	thiosulfate sulfurtransferase	21	sp P16385	<i>Saccharopolyspora erythr.</i>
155	pyruvate phosphate dikinase	41	gp U02529	<i>Entamoeba histolytica</i>
156	pyruvate synthase	57	pir S22397	<i>Halobacterium halobium</i>
157	pyruvate, orthophosphate dikinase	44	sp P11155	<i>Zea mays</i>
158	ribose 5-phosphate isomerase A	10	gp U10438	<i>Caenorhabditis elegans</i>
159	sucrase-isomaltase	13	gp L25926	<i>Rattus norvegicus</i>
160	succinate dehydrogenase (ubiq.) iron-sulfur	15	sp P21914	<i>Drosophila melanogaster</i>
161	sulfite reductase	18	pir S27479	<i>Archaeoglobus fulgidus</i>
162	enoyl-CoA hydratase	19	gp Z27079	<i>Caenorhabditis elegans</i>
163	tungsten formylmethanofuran dehydrogenase	7	gp X87970	<i>Methanobacterium thermo.</i>
164	tartrate dehydratase	9	sp P05847	<i>Escherichia coli</i>
165	<i>ttuD</i>	16	gp U32375	<i>Agrobacterium vitis</i>
166	UDP-glucose 4-epimerase	29	sp P45602	<i>Klebsiella pneumoniae</i>
167	uroporphyrinogen III methyltransferase	11	gp U05002	<i>Bacillus megaterium</i>
168	x-aconitate hydratase	41	sp P09339	<i>Bacillus subtilis</i>
169	anaerobic dimethyl sulfoxide reductase A	13	gp U32785	<i>Haemophilus influenzae</i>
170	NADH dehydrogenase	45	sp P42026	<i>Bovine</i>
171	Carbamate kinase	31	sp P13982	<i>Pseudomonas aeruginosa</i>
172	cystathionine gamma-lyase	44	gp S52028	<i>Homo sapiens</i>
173	L-asparaginase	20	sp P30362	<i>Lupinus arboreus</i>
174	4-hydroxybutyrate dehydrogenase	14	gp L36817	<i>Alcaligenes eutrophus</i>
175	formate DH, nitrate-inducible, alpha	12	pir S18213	<i>Wolinella succinogenes</i>
176	iron sulfur	26	gp M96826	<i>Methanothermobacter feravidus</i>
177	L-lactate dehydrogenase	33	gp M93720	<i>Plasmodium falciparum</i>
178	adenosine triphosphatase	15	gp L33259	<i>Helicobacter pylori</i>
179	ATP sulfurylase	7	gp U07353	<i>Penicillium chrysogenum</i>
180	aldehyde:ferredoxin oxidoreductase	8	gp X79777	<i>Pyrococcus furiosus</i>
181	cytochrome <i>b</i>	14	sp P39480	<i>Sulfolobus acidocaldarius</i>
182	cytochrome <i>c</i> biogenesis	24	gp Z22517	<i>Bradyrhizobium japonicum</i>
183	cytochrome <i>c</i> oxidase assembly factor ( <i>cyoE</i> )	30	gp U00013	<i>Mycobacterium leprae</i>
184	cytochrome <i>c</i> oxidase polypeptide I	35	sp P08681	<i>Mitochondria Chlamydomonas</i>
185	cytochrome oxidase I	57	gp U08900	<i>Mitochondria Trichoniscus pusil.</i>
186	cytochrome oxidase II	15	gi 155085	<i>Thermus thermophilus</i>
187	D-lactate dehydrogenase (cytochrome)	14	gp Z67750	<i>Saccharomyces cerevisiae</i>
188	electron transfer flavoprotein alpha	40	gp U17110	<i>Clostridium acetobutylicum</i>
189	electron transfer flavoprotein beta	20	sp P38975	<i>Paracoccus denitrificans</i>
190	ferredoxin oxidoreductase	20	gp X64521	<i>Halobacterium halobium</i>
191	ferredoxin-nitrite reductase	14	pir S30920	<i>Nicotiana tomentosiformis</i>
192	formaldehyde:ferredoxin oxidoreductase	57	gp X83963	<i>Thermococcus litoralis</i>
193	molybdenum- iron-sulfur flavoprotein ( <i>codH</i> )	12	gp X82447	<i>Oligotropha carboxidovora.</i>
194	polysulfide reductase B	43	sp P31076	<i>Wolinella succinogenes</i>
195	polysulfide reductase C	5	sp P31077	<i>Wolinella succinogenes</i>
196	quinol oxidase polypeptide I/II	22	sp P39481	<i>Sulfolobus acidocaldarius</i>
197	sulfide dehydrogenase	7	sp Q06530	<i>Chromatium vinosum</i>
198	thiosulfate reductase	40	gp L32188	<i>Salmonella typhimurium</i>
199	alcohol dehydrogenase	26	sp P12311	<i>Bacillus stearothermophil.</i>
200	alcohol dehydrogenase family member Ke 6	19	pir A48154	<i>Mus musculus</i>

Table 1b. *Continued*

ID #	Hit identification	Probability (-E)	Accession	Organism
201	aldehyde dehydrogenase, cytosolic	34	sp P13601	<i>Rattus norvegicus</i>
202	(salicyl)aldehyde dehydrogenase	31	gp U19817	<i>Arthrobacter globiformis</i>
203	3-phosphoglycerate kinase	21	gp M55529	<i>Methanothermus fervidus</i>
204	4-aminobutyrate aminotransferase	13	sp P40829	<i>Mycobacterium leprae</i>
205	fructose 1,6-biphosphate aldolase (class II)	21	gp X14436	<i>Escherichia coli</i>
206	<i>otsA</i>	27	gp U15187	<i>Mycobacterium leprae</i>
207	<i>pgk</i>	49	gp X80178	<i>Sulfolobus solfataricus</i>
208	phosphoenolpyruvate-utilizing	27	gp S74619	<i>Staphylothermus marinus</i>
209	phosphoglycerate kinase	17	sp P20971	<i>Methanothermus fervidus</i>
210	pyruvate kinase	12	gp U12980	<i>Saccharomyces cerevisiae</i>
211	pyruvate, water dikinase	59	gp U08376	<i>Pyrococcus furiosus</i>
212	fumarate dehydrogenase	46	gp X75402	<i>Sulfolobus solfataricus</i>
213	fumarate hydratase class II	41	sp P39461	<i>Sulfolobus solfataricus</i>
214	glucose dehydrogenase	7	gp D90044	<i>Bacillus megaterium</i>
215	transketolase	20	gp X67688	<i>Homo sapiens</i>
216	lipoamide dehydrogenase	15	sp P35484	<i>Acholeplasma laidlawii</i>
217	phosphonopyruvate decarboxylase	11	gp D37809	<i>Streptomyces hygroscopic.</i>
218	pyruvate dehydrobenase E1	24	sp P35488	<i>Acholeplasma laidlawii</i>
219	galactokinase	7	sp Q01415	<i>Homo sapiens</i>
220	aconitate hydratase	50	sp P37032	<i>Legionella pneumophila</i>
221	carboxyphosphonoenolpyruvate mutase	36	gp X67953	<i>Streptomyces hygroscopic.</i>
222	citrate synthase	34	gp X55282	<i>Thermoplasma acidophilum</i>
223	citrate synthase II	42	gp U05257	<i>Bacillus subtilis</i>
224	hydroxymethylglutaryl-CoA reductase	54	gp X54658	<i>Hevea brasiliensis</i>
225	malic acid	24	gb M19485	<i>Bacillus stearothermophil.</i>
226	phosphorylated isocitrate dehydrogenase	57	pdb 4ICD	<i>Escherichia coli</i>
227	succinyl-CoA synthetase	49	sp P25126	<i>Thermus aquaticus</i>
<b>Fatty acid and phospholipid metabolism</b>				
228	acetyl-CoA synthetase	71	sp P27095	<i>Methanotherx soehngenii</i>
229	acyl-CoA dehydrogenase	19	sp P15650	<i>Rattus norvegicus</i>
230	acyl-coA ligase(luciferase)	49	sp P29212	<i>Escherichia coli</i>
231	beta-ketothiolase	33	gp L37761	<i>Acinetobacter sp.</i>
232	enoyl-coA hydratase, mitochondrial	14	gp X15958	<i>Rattus norvegicus</i>
233	glyoxysomal beta-ketoacyl-thiolase	7	gp X93015	<i>Brassica napus</i>
234	geranylgeranyl pyrophosphate synthase	29	gp M87280	<i>Sulfolobus acidocaldarius</i>
235	geranyltransferase	12	sp Q08291	<i>Bacillus stearothermophil.</i>
236	hexaprenyl pyrophosphate synthetase	10	gp Z26494	<i>Saccharomyces cerevisiae</i>
237	HMG-CoA reductase	25	sp P34135	<i>Haloferax volcanii</i>
238	isoprenyl diphosphate synthase	12	gp S75695	<i>Methanobacterium thermo.</i>
239	long chain fatty acid coA ligase	12	gp U32686	<i>Haemophilus influenzae</i>
240	medium chain fatty acid coA ligase	31	sp Q00594	<i>Pseudomonas oleovorans</i>
241	thiolase	28	pir JC4032	<i>Clostridium acetobutylicum</i>
<b>Purines, pyrimidines, nucleosides, and nucleotides</b>				
242	diadenosine tetraphosphatase ( <i>apaH</i> )	12	gp X04711	<i>Escherichia coli</i>
243	mannosyltransferase B	6	gp D43637	<i>Escherichia coli</i>
244	deoxycytidine triphosphate deaminase	41	sp Q02103	<i>Desulfurolobus ambivalens</i>
245	thioredoxin reductase	36	gp L04500	<i>Eubacterium acidaminophi.</i>
246	CTP synthase ( <i>pyrG</i> )	62	gp U00021	<i>Mycobacterium leprae</i>
247	rRNA (adenosine- <i>N</i> 6, <i>N</i> 6-)-dimethyltransf.	11	sp P06992	<i>Escherichia coli</i>
248	adenylate kinase	18	gp U39882	<i>Methanococcus jannaschii</i>
249	amidophosphoribosyltransferase	11	gp D28868	<i>Arabidopsis thaliana</i>
250	glu phosphoribosylpyrophosphate amidotrans.	7	sp P00497	<i>Bacillus subtilis</i>
251	GMP synthase (glutamine-hydrolyzing)	44	sp P29727	<i>Bacillus subtilis</i>
252	<i>guaA</i>	33	gp U00015	<i>Mycobacterium leprae</i>
253	phosphoribosylamine-glycine ligase	17	sp P12039	<i>Bacillus subtilis</i>
254	phosphoribosylformylglycinamide synthase I	21	sp P12041	<i>Bacillus subtilis</i>
255	phosphoribosylformylglycinamide synthase II	29	gp M85265	<i>Lactobacillus casei</i>
256	phosphoribosylglycinamide formyltransf.	11	pir S37105	<i>Arabidopsis thaliana</i>
257	pur operon encoding purine biosynthesis	20	gp J02732	<i>Bacillus subtilis</i>
258	<i>purL</i>	12	gp U15182	<i>Mycobacterium leprae</i>
259	carbamoyl-phosphate synthetase	47	gp J05503	<i>Mesocricetus auratus</i>
260	dihydroorotate synthetase	19	gp U09990	<i>Methanobacterium thermo.</i>
261	dihydroorotate dehydrogenase	6	sp P32747	<i>Schizosaccharomyces pom.</i>
262	orotate phosphoribosyltransferase	19	pir A30492	<i>Bacillus subtilis</i>
263	phosphoribosyl-amp 1,6 cyclohydrolase	25	gp X82010	<i>Rhodobacter sphaeroides</i>
264	carbamyl phosphate synthetase	31	gp X87371	<i>Saccharomyces cerevisiae</i>
265	<i>thiA</i>	17	gp U26178	<i>Bacillus subtilis</i>
266	UDP galactose 4-epimerase	21	gp M94964	<i>Klebsiella pneumoniae</i>
267	UMP synthase	24	gp U22260	<i>Caenorhabditis elegans</i>

**Table 1b.** *Continued*

ID #	Hit identification	Probability (-E)	Accession	Organism
268	uracil phosphoribosyl transferase	17	gp U10246	<i>Toxoplasma gondii</i>
269	cytidine deaminase	28	sp P19079	<i>Bacillus subtilis</i>
270	HPRT	8	gb M88110	<i>Plasmodium falciparum</i>
271	uridine phosphorylase	45	sp P12758	<i>Escherichia coli</i>
272	galactose-1-phosphate uridylyltransferase	21	sp P31764	<i>Haemophilus influenzae</i>
<b>Regulatory functions</b>				
273	adenosylhomocysteinase	36	sp P10760	<i>rat</i>
274	ATP synthase	11	gp M22402	<i>Sulfolobus acidocaldarius</i>
275	carbon starvation A	14	sp P15078	<i>Escherichia coli</i>
276	GTP binding ( <i>hflX</i> )	18	gp U00019	<i>Mycobacterium leprae</i>
277	hydrobenase expression	9	sp P31905	<i>Alcaligenes eutrophus</i>
278	indoleacetamide hydrolase	12	sp P06618	<i>Pseudomonas syringae</i>
279	GTP-binding	10	gp Z49068	<i>Caenorhabditis elegans</i>
280	lactose operon repressor	6	sp P03023	<i>Escherichia coli</i>
281	mg11	14	sp U15635	<i>Mus musculus</i>
282	pleiotropic regulatory	44	gp M29002	<i>Bacillus stearothermophilus</i>
283	repressor	7	gp M81646	<i>Agrobacterium tumefacie.</i>
284	signal recognition particle (SRP54)	11	sp P37106	<i>Arabidopsis thaliana</i>
<b>Replication</b>				
285	DMC1/ <i>rad51/recA</i>	43	sp P25453	<i>yeast</i>
286	<i>recF</i>	9	sp P24900	<i>S. typhimurium</i>
287	<i>recQ</i>	25	sp P30015	<i>E coli</i>
288	DNA repair helicase ( <i>rad25</i> )	7	sp P19447	<i>yeast</i>
289	DNA repair ( <i>rad2</i> )	8	sp P39750	<i>Schizosaccharomyces pom.</i>
290	DNA repair, ionizing radiation (XRCC1)	5	sp P18887	<i>human</i>
291	<i>mutT</i>	8	gp U00021	<i>Mycobacterium leprae</i>
292	O6-methyltransferase	9	sp P16455	<i>human</i>
293	methyltransferase ( <i>uvrC</i> )	5	gp L29642	<i>Pseudomonas fluorescens</i>
294	cytosine sp. DNA methyltransferase	20	gp S86113	<i>Neisseria gonorrhoeae</i>
295	possible G-T mismatches repair	16	sp P29588	<i>Methanobacterium thermo.</i>
296	<i>umuD/lexA</i> -type , UV SOS operon	9	sp P04153	<i>Escherichia coli</i>
297	SIR2 (silent information regulator)	7	sp P06700	<i>yeast</i>
298	superoxide dismutase (Mn)	19	sp Q08713	<i>Sulfolobus acidocaldarius</i>
299	reverse gyrase	41	gp L10651	<i>Sulfolobus acidocaldarius</i>
300	DNA replication licensing factor	19	gp L41762	<i>Drosophila melanogaster</i>
301	DNA topoisomerase I	11	gp L27797	<i>Bacillus subtilis</i>
302	DNA topoisomerase III	21	sp P13099	<i>Saccharomyces cerevisiae</i>
303	DNA topoisomerase (ATP-hydrolysing)	42	pir S47332	<i>Erwinia carotovora</i>
304	modification methylase FNUDI	14	sp P34906	<i>Fusobacterium nucleatum</i>
305	heat shock	25	gi 473965	<i>Pyrococcus sp.</i>
306	TF55 heat shock	80	sp P28488	<i>Sulfolobus shibatae</i>
307	TF56	57	gp L34691	<i>Sulfolobus shibatae</i>
308a	DNA polymerase	34	gp M74198	<i>Thermococcus litoralis</i>
308b	DNA polymerase	59	gi 807828	<i>Pyrodictium occultum</i>
308c	DNA polymerase	63	gp D38573	<i>Pyrodictium occultum</i>
308d	DNA polymerase	60	gp D38574	<i>Pyrodictium occultum</i>
308e	DNA polymerase I	22	sp P26811	<i>Sulfolobus solfataricus</i>
309	replication factor C, small subunit	23	sp P35249	<i>Homo sapiens</i>
310	replication factor C, large subunit	12	sp P35601	<i>Mus musculus</i>
311	adenylylsulfate reductase	36	gp X63435	<i>Archaeoglobus fulgidus</i>
312	alkyl hydroperoxide reductase	49	gp U36479	<i>Sulfolobus sp.</i>
313	diadenosine tetraphosphatase	19	gi 1054947	<i>Sus scrofa</i>
314	DNA-ligase	122	sp Q02093	<i>Desulfurolobus</i>
315	endonuclease III	19	gp U11289	<i>Bacillus subtilis</i>
316	Holliday junction DNA helicase	5	gp U32716	<i>Haemophilus influenzae</i>
317	minichromosome maintenance	6	sp P29469	<i>Saccharomyces cerevisiae</i>
318	protein-L-isoaspartate (D-asp) O-methyltrans.	13	gp U09669	<i>Caenorhabditis elegans</i>
319	structure specific endonuclease	10	sp P26793	<i>Saccharomyces cerevisiae</i>
320	XPBara	9	gp U29168	<i>Arabidopsis thaliana</i>
321	YSA1	7	sp Q01976	<i>Saccharomyces cerevisiae</i>
<b>Transcription</b>				
322	RNA polymerase sigma factor ( <i>ntxA</i> )r	11	gp X69959	<i>Azorhizobium caulinodans</i>
323	TATA-binding	57	gp U23419	<i>Sulfolobus shibatae</i>
324	RNase PH (tRNA nucleotidyltransferase)	15	gp L10328	<i>Escherichia coli</i>
325a	DNA-directed RNA polymerase	41	pir A28213	<i>Methanobacterium thermo.</i>
325b	DNA-directed RNA polymerase II, RpB10	16	sp P29199	<i>Haloarcula marismortui</i>
325c	DNA-directed RNA polymerase II, RpB3	17	pir D44126	<i>Halobacterium marismort.</i>
326	DNA-directed RNA polymerase A	81	sp P31813	<i>Thermococcus celer</i>

Table 1b. Continued

ID #	Hit identification	Probability (-E)	Accession	Organism
327	DNA-directed RNA polymerase B	65	sp P11513	<i>Sulfolobus acidocaldarius</i>
328	DNA-directed RNA polymerase C	38	gp X14818	<i>Sulfolobus acidocaldarius</i>
329	DNA-directed RNA polymerase E	29	sp P39466	<i>Sulfolobus acidocaldarius</i>
330	RNA helicase	20	gp L01622	<i>Escherichia coli</i>
331	transcription activator	11	gp M73546	<i>Bacillus subtilis</i>
332	transcription elongation factor S-II	26	sp Q07271	<i>Sulfolobus acidocaldarius</i>
333	transcription factor IIB	18	gp U20899	<i>Sulfolobus shibatae</i>
334	transcriptional activator ( <i>tenA</i> )	11	sp P25052	<i>Bacillus subtilis</i>
<b>Translation</b>				
335	5S rRNA	41	gb M16530	<i>Sulfolobus sp.</i>
336	16S rRNA	43	gb M38637	<i>Thermoplasma acidophilum</i>
337	23S rRNA	65	emb X05480	<i>Desulfurococcus mobilis</i>
338	tRNA/5S rRNA gene cluster	12	emb X00916	<i>Methanococcus vannielii</i>
339	ala-tRNA	16	emb X05069	<i>Thermoproteus tenax</i>
340	ala-tRNA synthetase	64	sp P35029	<i>Sulfolobus solfataricus</i>
341	asn-tRNA synthetase	29	sp P10723	<i>Brugia malayi (nematode)</i>
342	asp-tRNA	26	emb X07692	<i>Methanobacterium thermo.</i>
343	asp-tRNA synthetase	38	gp U32810	<i>Haemophilus influenzae</i>
344	cys-tRNA synthetase	59	gp X56234	<i>Escherichia coli</i>
345	glu-tRNA	7	gb L36898	<i>Saccharomyces cerevisiae</i>
346	glu-tRNA synthetase	17	gp X07466	<i>Homo sapiens</i>
347	gly-tRNA	96	emb X14835	<i>Thermofilum pendens</i>
348	gly-tRNA synthetase	27	gp X78993	<i>Saccharomyces cerevisiae</i>
349	his-tRNA synthetase	19	gp L36863	<i>Sulfolobus shibatae</i>
350	ile-tRNA synthetase	43	pir S21569	<i>Methanobacterium thermo.</i>
351	Iso-tRNA synthetase	27	sp P41368	<i>Staphylococcus aureus</i>
352	leu-tRNA	25	emb X05071	<i>Thermoproteus tenax</i>
353	leu-tRNA synthetase	22	sp P11325	mitochondrial – yeast
354	met-tRNA	9	gb M26978	<i>Methanothermus fervidus</i>
355	met-tRNA synthetase	13	pir A25424	<i>Saccharomyces cerevisiae</i>
356	pro-tRNA	8	gb J01365	<i>Saccharomyces cerevisiae</i>
357	pro glu-tRNA synthetase	28	gp M74104	<i>Drosophila melanogaster</i>
358	ser-tRNA	10	gb M97644	<i>Methanopyrus kandleri</i>
359	ser-tRNA synthetase	49	sp P26636	<i>Cricetulus griseus</i>
360	thr-tRNA synthetase	22	sp P18256	<i>Bacillus subtilis</i>
361	trp-tRNA	11	gb K02528	<i>Halobacterium volcanii</i>
362	tyr-tRNA synthetase	15	gp L12221	<i>Saccharomyces cerevisiae</i>
363	val-tRNA synthetase	18	sp P36420	<i>Lactobacillus casei</i>
364	elongation factor, EF 1 alpha	72	sp P35021	<i>Sulfolobus solfataricus</i>
365	elongation factor, EF 2	75	gp X69297	<i>Sulfolobus solfataricus</i>
366	Initiation factor, IF 4A	5	sp P10081	<i>Saccharomyces cerevisiae</i>
367	initiation factor, IF 5A	19	sp P28461	<i>Sulfolobus acidocaldarius</i>
368	L2, 50S ribosomal	16	sp P21479	<i>Methanococcus vannielii</i>
369	L3, 60/50S ribosomal	14	sp P35684	<i>Oryza sativa</i>
370	L4, 50S ribosomal	12	sp Q06845	<i>Halobacterium halobium</i>
371	L6,50S ribosomal	21	sp P14030	<i>Methanococcus vanelli</i>
372	L10, ribosomal, acidic	11	sp P35023	<i>Sulfolobus solfataricus</i>
373	L11, 50S ribosomal	21	sp P35025	<i>Sulfolobus solfataricus</i>
374	L12, ribosomal	24	gp X59038	<i>Sulfolobus solfataricus</i>
375	L13, 50S ribosomal	22	sp P29198	<i>Haloarcula marismortui</i>
376	L15, 50S ribosomal	28	gp U16148	<i>Thermoplasma acidophilum</i>
377	L18, 50S ribosomal	50	sp P14033	<i>Methanococcus vannielii</i>
378	L23, ribosomal	19	gp Y00772	<i>Methanococcus vannielii</i>
379	L24, 50S ribosomal	26	sp P14034	<i>Methanococcus vannielii</i>
380	L29, 50S ribosomal	10	sp P04457	<i>Bacillus stearothermophil.</i>
381	L30, 50S ribosomal	18	sp P14035	<i>Methanococcus vanelli</i>
382	L32e, 50S ribosomal	17	sp P14549	<i>Methanococcus vannielii</i>
383	L37a, 60S ribosomal	16	pir S24170	<i>Gallus gallus</i>
384	S3, 30S ribosomal	15	sp P20281	<i>Haloarcula marismortui</i>
385	S5, 30S ribosomal	40	sp P14036	<i>Methanococcus vannielii</i>
386	S6, ribosomal	10	sp P17116	<i>Escherichia coli</i>
387	S7, 30S ribosomal	43	sp P35026	<i>Sulfolobus solfataricus</i>
388	S8, 30S ribosomal	35	sp P14038	<i>Methanococcus vannielii</i>
389	S10, 30S ribosomal	12	sp P23357	<i>Haloarcula marismortui</i>
390	S11, 30S ribosomal	57	sp P39469	<i>Sulfolobus acidocaldarius</i>
391	S12, 30S ribosomal	62	sp P29161	<i>Thermococcus celer</i>
392	S13, 40S ribosomal	46	sp Q05761	<i>Zea mays</i>
393	S15, ribosomal	21	pir A35908	<i>Homo sapiens</i>
394	S17, ribosomal	14	pir A24028	<i>Rattus norvegicus</i>

**Table 1b.** *Continued*

ID #	Hit identification	Probability (-E)	Accession	Organism
395	S18, 40S ribosomal	20	gp Z46260	<i>Saccharomyces cerevisiae</i>
396	S26e, 40S ribosomal	13	pir S47942	<i>Saccharomyces cerevisiae</i>
397	ribosomal alanine acetyltransferase	10	sp P09453	<i>Escherichia coli</i>
398	trmlp	26	gp Z48758	<i>Saccharomyces cerevisiae</i>
399	tRNA-splicing endonuclease $\beta$	6	sp P16658	<i>Saccharomyces cerevisiae</i>
400	ARD1 <i>N</i> -acetyl transferase	12	gp X77588	<i>Homo sapiens</i>
401	multifunctional aminoacyl-tRNA synthetase	26	sp P07814	<i>Homo sapiens</i>
402	ATP-dependent protease	10	gp L19301	<i>Myxococcus xanthus</i>
403	beta-type proteasome	25	gp U22157	<i>Methanosarcina thermophi.</i>
404	dipeptidase	17	gp D13142	<i>Sus scrofa</i>
405	<i>leuA</i>	11	gp U18657	<i>Haemophilus influenzae</i>
406	lysosomal alpha glucosidase	28	gp Y00839	<i>Homo sapiens</i>
407	prolidase	11	gp Z34896	<i>Lactobacillus delbrueckii</i>
408	proteasome, beta	18	sp P28061	<i>Thermoplasma acidophilum</i>
409	protein serine/threonine phosphatase 2A	14	gp X56261	<i>Saccharomyces cerevisiae</i>
410	TVG	29	gp D13178	<i>Thermoactinomyces vulga.</i>
411	sulfite reductase	10	pir S27478	<i>Archaeoglobus fulgidus</i>
<b>Transport and binding proteins</b>				
412	2'-5' oligoadenylate binding	25	pir S52166	<i>Homo sapiens</i>
413	<i>appB</i>	11	gp U20909	<i>Bacillus subtilis</i>
414	ATP-binding	25	gp Z25798	<i>Bacillus subtilis</i>
415	<i>cysA</i>	20	gp M32101	<i>Escherichia coli</i>
416	maltose-binding periplasmic	11	sp P19576	<i>Salmonella typhimurium</i>
417	<i>nataA/B</i>	31	gp U30873	<i>Bacillus subtilis</i>
418	p87	16	gp S47919	<i>Bos sp.</i>
419	alkylphosphonate uptake ( <i>phn</i> )A-Q	16	gp J05260	<i>Escherichia coli</i>
420	selenium-binding; ap56	36	sp P17563	<i>Mus musculus</i>
421	transport sec61, alpha	15	sp P38379	<i>Pyrenomonas salina</i>
422	arginine permease, substrate-binding	17	gp U08865	<i>Listeria monocytogenes</i>
423	D-lactate dehydrogenase	18	sp P32891	<i>Saccharomyces cerevisiae</i>
424	<i>dcxAD</i>	17	gp X56678	<i>Bacillus subtilis</i>
425	dipeptide transport system permease <i>dppB</i>	24	sp P37316	<i>Escherichia coli</i>
426	dipeptide transport system permease <i>dppC</i>	32	gp U17295	<i>Haemophilus influenzae</i>
427	glutamine permease	27	gp X14180	<i>Escherichia coli</i>
428	glutamine-binding periplasmic	19	sp P10344	<i>Escherichia coli</i>
429	branched-chain amino acid transport	46	sp P30294	<i>Salmonella typhimurium</i>
430	spoOk operon	26	gp M57689	<i>Bacillus subtilis</i>
431	oligopeptide perm. ( <i>oppC</i> , sporulation initiation)	20	sp P24139	<i>Bacillus subtilis</i>
432	oligopeptide transport permease <i>appB</i>	10	sp P42062	<i>Bacillus sbutilis</i>
433	oligopeptide transport permease <i>appC</i>	8	sp P42063	<i>Bacillus subtilis</i>
434	oligopeptide transport ATP binding <i>appF</i>	32	sp P42065	<i>Bacillus subtilis</i>
435	oligopeptide transport ATP-binding <i>oppF</i>	51	sp P08007	<i>Salmonella typhimurium</i>
436	phosphonates transport ATP-binding <i>phnC</i>	21	sp P16677	<i>Escherichia coli</i>
437	phosphonates transp. sys. permease ( <i>phnE</i> )	9	sp P16683	<i>Escherichia coli</i>
438	translocase SecY	23	gp X85020	<i>Sulfolobus acidocaldarius</i>
439	proline/betaine transporter	24	gp U14003	<i>Escherichia coli</i>
440	putrescine transport system permease ( <i>potI</i> )	10	sp P31136	<i>Escherichia coli</i>
441	spermidine/putrescine-binding periplasmic	10	gp U32731	<i>Haemophilus influenzae</i>
442	phosphate transport ATP-binding <i>pstB</i>	45	sp P07655	<i>Escherichia coli</i>
443	lactose transport ATP-binding <i>lack</i>	13	sp Q01937	<i>Agrobacterium radiobacter</i>
444	sterol carrier -2	24	gp M75884	<i>Homo sapiens</i>
445	sugar-binding	7	gp M77351	<i>Streptococcus mutans</i>
446	cation efflux system CZCD	11	sp P13512	<i>Alcaligenes eutrophus</i>
447	copper transport atp-binding <i>nosf</i>	17	sp P19844	<i>Pseudomonas stutzeri</i>
448	mercuric reductase	14	gp X73112	<i>Pseudomonas fluorescens</i>
449	<i>n</i> -ferritin repressor (FRP)	48	gb M95815	<i>Oryctolagus cuniculus</i>
450	Na <sup>(+)</sup> -ATPase B	17	sp Q08637	<i>Enterococcus hirae</i>
451	periplasmic divalent cation tolerance	14	gp Z36905	<i>Escherichia coli</i>
452	potassium/copper transporting ATPase	32	sp P32113	<i>Enterococcus hirae</i>
453	AP56	39	gp S56599	<i>Mus sp.</i>
454	rhodanese-like ( <i>cysA</i> )	41	gp M29612	<i>Saccharopholyspora erythr.</i>
455	<i>SN</i> -glycerol-3-phos. transport ATP-binding ( <i>ugpC</i> )	56	sp P10907	<i>Escherichia coli</i>
456	sulfate permease ( <i>cysA</i> )	19	sp P14788	<i>Anacystis nidulans</i>
<b>Other categories</b>				
457	7-alpha-hydroxysteroid hydrogenase	8	pir A42468	<i>Eubacterium sp.</i>
458	acylphosphatase	13	sp P07032	<i>Gallus gallus</i>
459	cheA/W/Y	22	gp U30501	<i>Thermotoga maritima</i>

**Table 1b.** *Continued*

ID #	Hit identification	Probability (-E)	Accession	Organism
460	FixA (nitrogen fixation)	8	sp P31573	<i>Escherichia coli</i>
461	FixB (nitrogen fixation)	41	gp M91817	<i>Clostridium acetobutylic.</i>
462	FixC (nitrogen fixation)	30	sp P31575	<i>Escherichia coli</i>
463	FixX (nitrogen fixation)	33	pir S49190	<i>Azotobacter vinelandii</i>
464	homocitrate synthase ( <i>nifV</i> , nitrogen fixation)	39	sp P05342	<i>Azotobacter vinelandii</i>
465	phosphomannomutase	15	gp M60873	<i>Pseudomonas aeruginosa</i>
466	ABC transporter ( <i>pstC-1</i> )	14	gp Z47982	<i>Mycobacterium tuberculosis.</i>
467	<i>sec61</i>	12	gp X77805	<i>Pyrenomonas salina</i>
468	threonine dehydratase biosynthetic	6	sp Q02145	<i>Lactococcus lactis</i>
469	vacuolar H <sup>+</sup> phosphatase	49	gp M81892	<i>Arabidopsis thaliana</i>
470	multidrug resistance	18	gp Z49126	<i>Caenorhabditis elegans</i>
471	phenylacrylic acid decarboxylase	38	gp L09263	<i>Saccharomyces cerevisiae</i>
472	daunorubicin-doxorubicin polyketide synth.	16	gp L35560	<i>Streptomyces peucetius</i>
473	pyrroline-5-carboxylate reductase	6	pir JC2078	<i>Thermus aquaticus</i>
474	amylomaltase	43	gp J01796	<i>Streptococcus pneumoniae</i>

quence contigs. The distance between these sequence markers is 11.5Kb on average, while the average distance between the 200 sequence gaps will be 8.5Kb, or in other words the average contig size will be 8.5Kb. Thus, a randomly placed contig is likely to fall across a sequence marker which consequently maps that contig to particular fosmid clones. The few contigs which do not cross a sequence marker will be mapped by hybridization to the 96 fosmid clones.

## Discussion

The mapping and sequencing of the *Pyrobaculum aerophilum* genome will provide valuable resources to researchers in many disciplines. We report here a list of 474 putative genes and a genome map. Access to these gene sequences now will provide researchers with a valuable tool for a multitude of studies.

A few whole genome sequencing projects have already been completed (Bult et al. 1996; Fleischmann et al. 1995; Fraser et al. 1995) and several others are in progress (Burland et al. 1995; Charlebois et al. 1996; Coulson 1996; Levy 1994). These efforts are concentrated on key organisms dispersed throughout the domains of life and will be important in elucidating the evolutionary path of modern organisms.

*Pyrobaculum aerophilum* is a member of the crenarchaea, which were thought to be a branch of the archaea but whose position in the phylogenetic tree of life has been a subject of controversy. It is uncertain whether the archaea are truly monophyletic (sharing a common ancestor distinct from other groups) or whether the crenarchaea in fact only share a distinct common ancestor with the eukaryotes and not with the rest of the archaea (Rivera and Lake 1992). Tracing the evolution of a few single genes is not sufficient to resolve these questions since

many factors complicate the analyses including gene duplications, exon shuffling, and horizontal transfers (i.e., transfer of a gene from one organism to another). Accumulating large amounts of sequence data from representative organisms is the most efficient way to unscramble these evolutionary puzzles.

Complete genome sequences from representative organisms also provide researchers with the information needed to quickly isolate their gene of interest from virtually any organism. Furthermore, researchers will frequently be able to trace the evolutionary history of their favorite gene directly from the sequence databases. Quickly determining the absence of a particular gene in an organism is in itself a powerful tool. Currently, this takes weeks of hybridization experiments and still leaves the researcher uncertain of the gene's absence.

Although a detailed analysis comparing this sequence to other completed sequences will wait for the final edited sequence, there are already some surprises. One example is in the area of repair systems. We have identified tags to members of many of the major repair systems found in *E.coli* and the human, including those involved in excision repair, repair of oxidative damage, and methylated bases. *Pyrobaculum aerophilum* is a particularly interesting organism for the study of repair systems because of its remarkable similarity, at the amino acid sequence level, to eukaryotes. For instance, we have found homologs to genes involved in human repair disorders, such as xeroderma pigmentosum. Interestingly, despite the fact that we have sequenced more than 95% of the *Pyrobaculum* genome, we have not detected any genes involved in either the mismatch repair systems or deamination of cytosines. The latter is expected to be enhanced at high temperatures. The reported complete genome sequence of *Methanococcus jannaschii* (Bult et al. 1996) also fails to list genes involved in either of these repair pathways. Therefore, either the Archaea have a different way of repairing replication mismatches, or we must entertain the fascinating possibility that they have high mutation rates in vivo.

## Materials and methods

Laboratory procedures were performed as described in the laboratory manual (Sambrook et al. 1989) unless otherwise mentioned.

### Genomic DNA extraction

*Pyrobaculum aerophilum* (type strain IM2; DSM 7523) was grown anaerobically and pelleted at Karl Stetter's laboratory, University of Regensburg, Germany, under the previously published conditions (Völkl et al. 1993). Approximately 0.5 g (wet weight) of the cell pellet was chilled with liquid nitrogen, ground with a mortar and pestle, and resuspended in 500 ml TE buffer. Lysozyme (5 ml, to give 50 mg/ml) was added to the cell suspension, followed by the addition of 30 ml 20% w/v SDS, and the lysis mixture was incubated at 37°C for 1 h. Lysed cells were phenol extracted once, ethanol precipitated, then resuspended in TE at a final concentration of 0.5 mg/ml.

### Fosmid library construction

For the construction of the genomic fosmid library, 10 µg of genomic DNA was resuspended in 500 µl *Sau3AI* buffer by gently swirling on a rotating platform for 3 h at 4°C followed by 10 min at 37°C. The DNA was restriction digested by adding 0.1 unit *Sau3AI* enzyme (New England Biolabs, Beverly, MA, USA) and incubating at 37°C for 15 minutes, followed by addition of another 0.1 U *Sau3AI* enzyme with additional incubation at 37°C for 15 minutes. Then, the DNA was extracted by chloroform, ethanol precipitated, and resuspended in 20 µl TE. Analysis on an agarose gel showed that most fragments were within 20–40 Kb. From this DNA stock a 3-µl aliquot was taken and dephosphorylated in a 25-µl reaction with 2 units HK-Phosphatase (Epicentre Technologies, Madison, WI, USA) at 30°C for 1 h, followed by heat inactivation of the phosphatase at 65°C for 15 min. Fosmid arms were prepared as described previously (Kim et al. 1992). Ligations were done with a large molar excess of vector DNA (50–100×) to minimize the probability of chimera formation. *Bam*HI digested fosmid arms (1 µg) were ligated with approximately 15 ng of the genomic DNA in a 15-µl reaction with 1 unit of T4 DNA ligase (Epicentre Technologies). The ligation product thus generated (3 µl) was packaged into phage particles using GigaPack XL in-vitro packaging extract (Stratagene, La Jolla, CA, USA) as instructed by the vendor. The packaged particles were suspended in 500 µl SM buffer. *E. coli* DH10B cells were grown in LB + 0.2% maltose + 10 mM MgSO<sub>4</sub> to an optical density at 600 nm (OD<sub>600</sub>) of 0.5–0.9, centrifuged at 4°C, and resuspended in 10 mM MgSO<sub>4</sub> to a final concentration of OD<sub>600</sub> = 1.0. Undiluted packaging lysate (25 µl) was mixed with 25 µl resuspended cells and kept at room temperature for 30 min. LB (200 µl at 37°C) was then added and kept at 37°C for 1 h with gentle shaking every 15 min. Cells were centrifuged and resuspended in 50 µl LB, which

was subsequently plated on LB plates containing 10 µg/ml chloramphenicol. Emerging colonies from the plates were inoculated into eight 96-well microtiter plates containing LB + 10% glycerol, which were then grown at 37°C overnight and kept frozen at –70°C until use.

### Construction of pUC18 library

Genomic DNA (10 µg) was diluted in 300 µl TE and sheared by sonicating with a Branson cell disrupter 200 sonicator at setting 3.5 for 2 s. This sheared DNA (180 µl) was mixed with 180 µl 2× *Bal31* buffer and kept at 30°C for 15 min. Then, 0.015 units of *Bal31* enzyme (New England Biolabs) was added and the mixture incubated at 30°C for 10 min. The mixture was then diluted to 660 µl with TE, extracted with phenol: chloroform and chloroform, and precipitated with ethanol. One fifth of the remaining DNA was run on a preparative 1.2% SeaKem GTG (FMC Bioproducts, Rockland, ME, USA) agarose gel. After brief staining with ethidium bromide, fragments of 1.5–2 Kb were excised and recovered by QiaexII (Qiagen, Chatsworth, CA, USA) binding following manufacturers instructions. Elution was in 20 µl 10 mM Tris pH 8.5. This entire volume was used for subcloning into the *Sma*I site of pUC18 (Ready-To-Go subcloning kit, Pharmacia Biotech, Piscataway, NJ, USA). These subclones were randomly sequenced from both ends (unpublished) using fluorescent multiplex automated sequencers (ABI 373, Applied Biotechnology, Foster City, CA, USA), and the public databases were searched with these sequences by BLAST (Altschul et al. 1990; States et al. 1991) to find similarities to previously known sequences.

### Generation of sequence marker probes by PCR amplification

Typically, 0.005 pmoles of template DNA, 20 pmoles of each primer, and 1 unit of TaqI polymerase (Gibco BRL, Gaithersburg, MD, USA) were combined in a 25-µl reaction. Polymerase chain reaction (PCR) products were isolated by running the PCR reactions on 1% SeaPlaque low melting point agarose gel (FMC Bioproducts, Rockland, ME, USA) in 1× TAE and excising DNA bands. The DNA fragments, in low melting point agarose, were radiochemically labeled using the DecaPrime kit (Ambion, Austin, TX, USA) as instructed by the vendor except that each reaction was scaled down to half the instructed values.

### Generation of riboprobes

Fosmid clone templates were prepared from 3 ml LB cultures using an AutoGen 740 automated miniprep machine (Integrated Separation Systems, Natick, MA, USA). Half of the DNA isolated from a 3-ml culture was used in a 100-µl digestion reaction with 5 units *Hinc*II restriction enzyme. Digested products were ethanol precipitated and resuspended in 5 µl TE or water. This concentrated template (1.5 µl) was used in the standard transcription reaction de-

scribed in the RiboScribe kit protocol (Epicentre Technologies), except that all volumes were scaled down to half, giving a final reaction size of 10  $\mu$ l. After the 2-h labeling period the entire contents of the reactions were used for hybridization.

### Hybridization

Fosmid library colony membranes were prepared by gridding the entire library, represented by eight 96-well microtiter plates, at high density onto single 8 cm  $\times$  12 cm nylon membranes using a Biomek 1000 laboratory workstation (Beckman Instruments, Fullerton, CA, USA). Labeled probes (PCR products or riboprobes, as described in the previous section) were transferred directly into hybridization chambers without purification. Colony hybridizations were incubated overnight in 1 M NaCl, 50 mM Tris-HCl pH 8.0, 5 mM EDTA pH 8.0, 10% PEG 8000, at 65°C. The membranes were washed in 2 $\times$  standard saline citrate (SSC), 0.1% sodium dodecyl sulfate (SDS) at 65°C for 30 min, followed by another wash in 0.2 $\times$  SSC, 0.1% SDS at 65°C for 30 min. The membranes were exposed to X-ray film from 3 h to overnight.

### Fingerprinting

Restriction fingerprint analysis of fosmid clones was done essentially according to the previously described modified protocol that was adapted from the original cosmid fingerprinting procedure (Coulson et al. 1986; Sulston et al. 1988). The fingerprint gels were digitized using a Phosphorimager (Molecular Dynamics, Foster City, CA, USA) and analyzed using the contigc program (<http://www.sanger.ac.uk>).

### Assembly of fosmid contig

Ordering and alignment of contigs was performed manually and based primarily on the results of hybridization of the markers to the fosmid clones. Two clones are considered overlapping if they are hit by one or more common probes. Ambiguities in the overlapping relations were resolved by fingerprint analysis or additional hybridization with riboprobes that are specific to clone ends.

**Acknowledgments** We thank Dr. Phil Green and Dr. David Gordon for use of their Phred and Phrap programs, Dr. David Mathog for helping us use the GCG DNA sequence analysis programs and Jack Gerson for critical reading of the manuscript. This work was supported by funding from the United States Department of Energy to M.I.S. and from the Office of Naval Research to J.H.M. (N00014-95-1-0938).

## References

- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) Basic local alignment search tool. *J Mol Biol* 215:403–410
- Bult CJ, White O, Olsen GJ, Zhou L, Fleischmann RD, Sutton GG, Blake JA, FitzGerald LM, Clayton RA, Gocayne JD, Kerlavage AR, Dougherty BA, Tomb JF, Adams MD, Reich CI, Overbeek R, Kirkness EF, Weinstock KG, Merrick JM, Glodek A, Scott JL, Geoghagen NSM, Weidman JF, Fuhrmann JL, Venter JC et al. (1996) Complete genome sequence of the methanogenic archaeon, *Methanococcus jannaschii* *Science* 273:1058–1073
- Burland V, Plunkett GR, Sofia HJ, Daniels DL, Blattner FR (1995) Analysis of the *Escherichia coli* genome VI: DNA sequence of the region from 92.8 through 100 minutes. *Nucleic Acids Res* 23:2105–2119
- Charlebois RL, Gaasterland T, Ragan MA, Doolittle WF, Sensen CW (1996) The *Sulfolobus solfataricus* P2 genome project. *FEBS Lett* 389:88–91
- Coulson A (1996) The *Caenorhabditis elegans* genome project. *C. elegans Genome Consortium. Biochem Soc Trans* 24:289–291
- Coulson A, Sulston J, Brenner S, Karn J (1986) Toward a physical map of the genome of the nematode *Caenorhabditis elegans*. *Proc Natl Acad Sci USA* 83:7821–7825
- Fleischmann RD, Adams MD, White O, Clayton RA, Kirkness EF, Kerlavage AR, Bult CJ, Tomb JF, Dougherty BA, Merrick JM et al. (1995) Whole-genome random sequencing and assembly of *Haemophilus influenzae* Rd *Science* 269:496–512
- Fraser CM, Gocayne JD, White O, Adams MD, Clayton RA, Fleischmann RD, Bult CJ, Kerlavage AR, Sutton G, Kelley JM et al. (1995) The minimal gene complement of *Mycoplasma genitalium* *Science* 270:397–403
- Kim UJ, Shizuya H, de Jong PJ, Birren B, Simon MI (1992) Stable propagation of cosmid sized human DNA inserts in an F factor based vector. *Nucleic Acids Res* 20:1083–1085
- Levy J (1994) Sequencing the yeast genome: an international achievement. *Yeast* 10:1689–1706
- Riley M (1993) *Microbiol Rev* 57:862
- Rivera MC, Lake JA (1992) Evidence that eukaryotes and eocyte prokaryotes are immediate relatives. *Science* 257:74–76
- Robb FT, Place AR (eds) (1995) *Archaea : A laboratory manual. Thermophiles*. Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, 3 vol
- Sambrook J, Fritsch EF, Maniatis T (1989) *Molecular cloning : a laboratory manual*, 2nd edn. Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, 3 vol
- States DJ, Gish W, Altschul SF (1991) Improved sensitivity of nucleic acid databases searches using application-specific scoring matrices. *Methods* 3:66–70
- Stetter KO (1992) Life at the upper temperature border. In: Tran Than Van JK, Mounolou JC, Schneider J, McKay C (eds) *Colloque Interdisciplinaire du Comite National de la Recherche Scientifique, Frontiers of Life*
- Sulston J, Mallett F, Staden R, Durbin R, Horsnell T, Coulson A (1988) Software for genome mapping by fingerprinting techniques. *Comput Appl Biosci* 4:125–32
- Völkl P, Huber R, Drobner E, Rachel R, Burggraf S, Trincone A, Stetter KO (1993) *Pyrobaculum aerophilum* sp. nov., a novel nitrate-reducing hyperthermophilic archaeum. *Appl Environ Microbiol* 59:2918–2926
- Woese CR (1987) Bacterial evolution. *Microbiol Rev* 51:221–271