

Complete Genome Sequence of an Aerobic Thermoacidophilic Crenarchaeon, *Sulfolobus tokodaii* strain7

Yutaka KAWARABAYASI,^{1,2,*} Yumi HINO,¹ Hiroshi HORIKAWA,¹ Koji JIN-NO,¹ Mikio TAKAHASHI,¹ Mitsuo SEKINE,¹ Sin-ichi BABA,¹ Akiho ANKAI,¹ Hiroki KOSUGI,¹ Akira HOSOYAMA,¹ Shigehiro FUKUI,¹ Yoshimi NAGAI,¹ Keiko NISHIJIMA,¹ Rie OTSUKA,¹ Hidekazu NAKAZAWA,¹ Minako TAKAMIYA,¹ Yumiko KATO,¹ Takio YOSHIKAWA,¹ Toshihiro TANAKA,¹ Yutaka KUDOH,¹ Jun YAMAZAKI,¹ Norihiro KUSHIDA,¹ Akio OGUCHI,¹ Ken-ichi AOKI,¹ Sayaka MASUDA,¹ Masao YANAGII,¹ Masami NISHIMURA,¹ Akihiko YAMAGISHI,³ Tairo OSHIMA,³ and Hisasi KIKUCHI¹

National Institute of Technology and Evaluation, 2-49-10 Nishihara, Shibuya, Tokyo 151-0066, Japan,¹ National Institute of Advanced Industrial Science and Technology, 1-1 Higashi, Tsukuba, Ibaraki 305-0046, Japan,² and Tokyo University of Pharmacy and Life Science, Horinouchi, Hachioji, Tokyo 192-0392, Japan³

(Received 23 July 2001; revised 2 August 2001)

Abstract

The complete genomic sequence of an aerobic thermoacidophilic crenarchaeon, *Sulfolobus tokodaii* strain7 which optimally grows at 80°C, at low pH, and under aerobic conditions, has been determined by the whole genome shotgun method with slight modifications. The genomic size was 2,694,756 bp long and the G+C content was 32.8%. The following RNA-coding genes were identified: a single 16S–23S rRNA cluster, one 5S rRNA gene and 46 tRNA genes (including 24 intron-containing tRNA genes). The repetitive sequences identified were SR-type repetitive sequences, long dispersed-type repetitive sequences and Tn-like repetitive elements. The genome contained 2826 potential protein-coding regions (open reading frames, ORFs). By similarity search against public databases, 911 (32.2%) ORFs were related to functional assigned genes, 921 (32.6%) were related to conserved ORFs of unknown function, 145 (5.1%) contained some motifs, and remaining 849 (30.0%) did not show any significant similarity to the registered sequences. The ORFs with functional assignments included the candidate genes involved in sulfide metabolism, the TCA cycle and the respiratory chain. Sequence comparison provided evidence suggesting the integration of plasmid, rearrangement of genomic structure, and duplication of genomic regions that may be responsible for the larger genomic size of the *S. tokodaii* strain7 genome. The genome contained eukaryote-type genes which were not identified in other archaea and lacked the CCA sequence in the tRNA genes. The result suggests that this strain is closer to eukaryotes among the archaea strains so far sequenced.

The data presented in this paper are also available on the internet homepage (http://www.bio.nite.go.jp/E-home/genome_list-e.html/).

Key words: aerobic thermoacidophilic crenarchaeon; genome sequencing; whole genome shotgun method; comparative analysis; plasmid

1. Introduction

Among the complete genome sequences of thermophilic archaea reported,^{1–5} two genomes were sequenced by our group: *Pyrococcus horikoshii* OT3¹ and *Aeropyrum pernix* K1.² These two species are hyper-thermophilic and

optimally grow at over 95°C. While *P. horikoshii* OT3 is an anaerobic euryarchaeon, *A. pernix* K1 is an aerobic crenarchaeon. The entire genomic sequence of *A. pernix* K1 was the first and the only complete genomic data of crenarchaeota. It is also the only known genomic sequence of a strictly aerobic hyperthermophile. To obtain more information about the genomic sequence data of *A. pernix* K1, it is useful to compare the data with data from closely related species.

We therefore selected *Sulfolobus tokodaii* strain7, a species of genus *Sulfolobus* in crenarchaeon for sequencing. This strain was isolated from Beppu hot springs

Communicated by Michio Oishi

* To whom correspondence should be addressed. Tel. +81-3-3481-8951, Fax. +81-3-3481-1962, E-mail: kyutaka@nite.go.jp

† The entire genome sequence has been deposited in DDBJ/GenBank/EMBL databases under accession numbers AP000981–AP000990.

in Kyushu, Japan in 1983.⁶ This species grows under aerobic conditions, as does *A. pernix* K1. The optimal growth temperature of *S. tokodaii* strain7 is 80°C and the optimal pH is between 2 and 3. The genomic data of this strain is expected to provide the information on not only the thermostability of proteins but also the characteristics of cells living in an acidic environment. The genomic size of this strain is approximately 2.7 Mbp,⁷ 1 Mbp larger than the other two species previously determined by our group: *P. horikoshii* OT3 and *A. pernix* K1. In addition to above genomic features, *S. tokodaii* strain7 has no extra-chromosomal genetic unit⁷ and is able to convert hydrogen sulfide to sulfate (personal communication from A. Yamagishi and T. Oshima).

To determine the entire sequence of this genome, the shotgun libraries with short and long inserts were constructed both from the entire genomic DNA and from the restriction fragments of genomic DNA digested with *Bss*H II. The raw sequencing data from shotgun clones were assembled using the software PhredPhrap and Consed.⁸⁻¹⁰ The remaining sequencing gaps were filled by walking of long insert shotgun clones.

2. Materials and Methods

2.1. Bacterial strains and genomic DNA

The strain, *S. tokodaii* strain7 deposited in the Japan Collection of Microorganisms (JCM number 10545) was used for genome sequencing in this study. These cells were inoculated into 100 ml of the *Sulfolobus* culture medium¹¹ which was prepared in 500-ml Erlenmeyer flasks, and cultured at 80°C with vigorous shaking. The genomic DNA was isolated principally based on the method of Yamagishi and Oshima.¹² *Escherichia coli* DH10B was used for the preparation of plasmid clones.

2.2. Construction of shotgun clones

The *S. tokodaii* strain7 genomic DNA was sonicated for 5, 10, and 20 sec at the L position of a sonicator Biorupter (Cosmo Bio, Tokyo, Japan), followed by size-fractionation by agarose gel electrophoresis. The fractions from 0.8 to 1.2 kb and from 2.0 to 2.5 kb were independently cloned into the *Hinc* II site of pUC118 (Takara Shuzo, Kyoto, Japan). They are referred to as short-fragment and long-fragment shotgun libraries, respectively.

For the construction of the shotgun libraries from restriction fragments, the genomic DNA was prepared in the agarose-plug as described by Kondo et al.,⁷ and after complete digestion of DNA in the plugs with *Bss*H II, the digests were resolved in 0.4% low-melting-point agarose gel by electrophoresis. The bands were cut out and the gel was dissolved by incubation at 50°C with Agarase (FMC, Rockland, ME, USA) overnight. The resul-

tant solution containing DNA fragments was sonicated, and the short-fragment shotgun library was prepared as above, except for *Bss*H II fragment A (945 kb) from which both short- and long-fragment shotgun libraries were constructed.

2.3. Sequencing

Plasmid DNA was prepared by an Autogen 740 automatic DNA preparation system (Autogen, Framington, MA, USA). The sequencing reaction was performed using two kinds of cycle sequencing kits, a dye-primer cycle sequencing kit and a dye-terminator cycle sequencing kit. The sequence data were detected by ABI-DNA sequencers (377XL; Perkin-Elmer ABI, Foster City, CA, USA). Also the dye-terminator cycle sequencing kit was used for filling the sequencing gaps. The long-fragment shotgun clones with inserts covering gap regions were used for filling of gaps by walking with synthesized primers.

The raw sequence data were first analyzed with the software Phred.⁷ After the elimination of contaminated sequences derived from *E. coli* or vector DNAs, the treated data were assembled into contigs using the software Phrap⁸ and Consed.⁹ The assembled sequences were split into 30-kbp segments and the sequence in each segment was re-assembled and edited by Sequencher (Gene Codes, Ann Arbor, MI, USA).

2.4. Computational analysis

The criteria used for the assignment of potential protein-coding regions were similar to those used in the previous paper.¹ However, the open reading frames (ORFs) which had neither similarity nor motif sequences completely included within longer ORFs were not assigned as potential protein-coding regions in this work. Similarity search of the assigned ORFs was mainly performed by the Smith-Waterman algorithm.¹³ The databases used for similarity search were GenBank release 109, EMBL release 56.0, Swiss-Prot release 38.0, PIR release 62.0, and Owl release 31.4.

3. Results and Discussion

3.1. Determination of entire genome sequence

The physical map of circular genomic DNA of *S. tokodaii* strain7 has been reported by Kondo et al.⁷ According to their observation that the genomic DNA generates six fragments upon digestion with *Bss*H II, these six fragments were used for the construction of fragment-specific shotgun libraries.

The entire genome sequence was determined by the whole genome shotgun method¹⁴ with slight modifications. The two kinds of shotgun libraries, one with approximately 1-kb inserts (short-fragment shotgun library) and the other with approximately 2.2-kb inserts

(long-fragment shotgun library), were prepared from the genomic DNA and from *Bss*H II fragment A. From the other *Bss*H II fragments, only short-fragment shotgun libraries were constructed. The clones from the short-fragment shotgun libraries were sequenced from one end whereas the clones from the long-fragment shotgun libraries were sequenced from both ends. A total of 61,000 readings of raw sequencing data were accumulated and assembled using PhredPhrap.^{8–10}

As indicated in the following section, a large number of long dispersed repetitive sequences were identified from the consensus sequences of primary contigs generated by PhredPhrap. To avoid the confusion of assembling by these repetitive sequences, the nucleotide sequence of these repetitive units were added in the cross-match database of PhredPhrap and were not used for contig assembling. The resulting contigs were separated into 30-kbp regions and the raw data contained in each region was re-assembled and edited by Sequencher.

To fill the remaining sequencing gaps, the clones carrying the end sequences of contigs were selected from the long-fragment shotgun libraries and sequenced by walking with synthesized primers. The nucleotide sequences of long dispersed repetitive units were also determined by walking the long-fragment shotgun clones which included each repetitive unit with synthesized primers. All the sequences were determined by co-incidence of sequences of at least two clones per base.

To confirm the authenticity of the genomic sequence constructed, the restriction pattern of each 15-kb fragment directly amplified from the genomic DNA by long PCR were compared with those deduced from the sequence data.

The total length of the genome finally confirmed was 2,694,756 bp. The nucleotide position was numbered from the one end of the *Bss*H II restriction site located on the *Not* I B and *Rsr* II F fragments in the physical map shown by Kondo et al.⁷

Distribution of ACGT along the strand of the entire genome was 33.4% A, 16.3% C, 16.5% G, and 33.8% T, resulting in a G + C content of 32.8%. Through the processes of genomic sequencing, no evidence was obtained for the presence of an extra-chromosomal unit.

3.2. RNA coding genes

The entire genomic sequence was subjected to similarity search against the rRNA sequences registered in the databases. It was found that the *S. tokodaii* strain7 genome contained a single 16S–23S rRNA cluster and one 5S rRNA. This organization is similar to those of *P. horikoshii* OT3¹ and *A. pernix* K1.²

Forty-five tRNA genes were identified by searching with tRNA scan,¹⁵ and one tRNA gene for Leu (GAG) was identified by similarity search. The species of tRNA genes identified are shown in Table 1. Forty-one of

these genes were discretely mapped, while the remaining six tRNA genes were mapped as clusters of two tRNA genes. As noted in other microorganisms,^{1,2,16} no tRNA genes containing A at the first position of the anticodon were identified. Similar to the archaea genomes already sequenced, the Met-tRNA gene occurred in triplicate. However, the sequences had no similarity to one another, suggesting that the different Met-tRNA species may be used for the translation initiation of different classes of genes.

Twenty-four tRNA genes were found to contain 11- to 57-bp long introns. These genes containing introns and the organization of the introns in each gene are summarized in Table 2. In all of the intron-containing tRNA genes except for the tRNA-Leu (GAG), tRNA-Glu (UUC) and tRNA-Glu (CUC) genes, the introns were identified 1 bp downstream from the 3' end of the anticodon triplet. The intron in the tRNA-Leu (GAG) gene had been inserted 4-bp upstream from the 5' end of anticodon region. The 17 bp-long intron portion was assigned within the D-loop of two tRNA-Glu genes. The insertion position of these introns was similar to that of the tRNA-Thr (UGU) gene identified in the *A. pernix* K1 genome.²

In contrast to tRNA genes in other archaea, most of the tRNA genes identified did not contain the 3' terminal CCA sequence. On the other hand, tRNA nucleotidyltransferase,¹⁷ that catalyzed the addition of the CCA sequence to tRNA transcripts, was found on this genome (ST0952). In bacteria, the 3' CCA is generally encoded on the tRNA genes. The role of *E. coli* tRNA nucleotidyltransferase is the repair of tRNA 3' ends.^{18–20} In eukaryotes, where CCA is rarely encoded by tRNA genes, tRNA nucleotidyltransferase is essential.²¹ Although the role of tRNA nucleotidyltransferase is unclear in *S. tokodaii* strain7, the structure of tRNA genes and the presence of a CCA-adding enzyme are similar to eukaryotic cells.

The genes for all tRNA synthetases have been identified in the assigned ORFs, except for those for tRNA^{Gln} and tRNA^{Asn}. Instead, similar to *Bacillus subtilis*²² and *Deinococcus radiodurans*,²³ the three subunits of Glu-tRNA^{Gln} amidotransferase (ST1283, ST0282, and STS140), which were necessary for transamidation of Glu-tRNA^{Gln} and Asp-tRNA^{Asn}, were identified.

3.3. Repetitive sequences identified

Three types of repeating units including Tn-like elements were found on the genome of *S. tokodaii* strain7.

The first type, named type A, is comprised of repetitive sequences composed of LR and SR segments, and is divided into two subtypes based on the sequences of SR segments. The subtype A-I repetitive unit contained 24 bp-long well-conserved SR segments (GMTAATCCWTAATGGAATTGAAAG) and a

Table 1. Summary of assigned tRNA genes.

UUU (Phe)		UCU (Ser)		UAU (Tyr)		UGU (Cys)	
UUC (Phe)	⊙	UCC (Ser)	○	UAC (Tyr)	⊙	UGC (Cys)	⊙
UUA (Leu)	○	UCA (Ser)	⊙	UAA (End)		UGA (End)	
UUG (Leu)	⊙	UCG (Ser)	⊙	UAG (End)		UGG (Trp)	⊙
CUU (Leu)		CCU (Pro)		CAU (His)		CGU (Arg)	
CUC (Leu)	⊙	CCC (Pro)	⊙	CAC (His)	○	CGC (Arg)	⊙
CUA (Leu)	⊙	CCA (Pro)	○	CAA (Gln)	○	CGA (Arg)	○
CUG (Leu)	⊙	CCG (Pro)	○	CAG (Gln)	○	CGG (Arg)	○
AUU (Ile)		ACU (Thr)		AAU (Asn)		AGU (Ser)	
AUC (Ile)	○	ACC (Thr)	○	AAC (Asn)	○	AGC (Ser)	○
AUA (Ile)		ACA (Thr)	⊙	AAA (Lys)	⊙	AGA (Arg)	⊙
AUG (Met)	⊙ *	ACG (Thr)	⊙	AAG (Lys)	⊙	AGG (Arg)	⊙
GUU (Val)		GCU (Ala)		GAU (Asp)		GGU (Gly)	
GUC (Val)	○	GCC (Ala)	○	GAC (Asp)	⊙	GGC (Gly)	○
GUA (Val)	○	GCA (Ala)	○	GAA (Glu)	⊙	GGA (Gly)	○
GUG (Val)	○	GCG (Ala)	○	GAG (Glu)	⊙	GGG (Gly)	○

tRNA genes assigned are indicated by ○ and those with introns by ⊙
 *; All three tRNA^{Met} possessed the introns

Table 2. The size and position of introns in the intron-containing tRNA genes.

tRNA-ID	tRNA species	anticodon	gene size (nt)	Intron position* start	Intron position* end	Intron size (nt)	mature tRNA (nt)
STrRNA02	tRNA-Glu	TTC	91	22**	38	17	74
STrRNA03	tRNA-Lys	TTT	99	39	63	25	74
STrRNA04	tRNA-Lys	CTT	101	39	65	27	74
STrRNA06	tRNA-Cys	GCA	88	38	51	14	74
STrRNA08	tRNA-Met	CAT	98	39	62	24	74
STrRNA11	tRNA-Thr	CGT	98	39	62	24	74
STrRNA15	tRNA-Met	CAT	93	39	57	19	74
STrRNA16	tRNA-Tyr	GTA	93	39	57	19	74
STrRNA18	tRNA-Ser	CGA	111	39	64	26	85
STrRNA19	tRNA-Leu	CAG	101	41	56	16	85
STrRNA20	tRNA-Glu	CTC	91	22**	38	17	74
STrRNA21	tRNA-Ser	TGA	95	39	49	11	84
STrRNA23	tRNA-Phe	GAA	91	39	55	17	74
STrRNA25	tRNA-Leu	TAG	101	41	56	16	85
STrRNA26	tRNA-Leu	GAG	104	34**	52	19	85
STrRNA27	tRNA-Pro	GGG	94	41	58	18	76
STrRNA29	tRNA-Ile	GAT	90	39	54	16	74
STrRNA31	tRNA-Thr	TGT	92	40	56	17	75
STrRNA32	tRNA-Arg	GCG	88	40	52	13	75
STrRNA33	tRNA-Leu	CAA	97	40	52	13	84
STrRNA34	tRNA-Met	CAT	87	40	51	12	75
STrRNA36	tRNA-Arg	CCT	99	40	63	24	75
STrRNA40	tRNA-Trp	CCA	131	39	95	57**	74
STrRNA46	tRNA-Arg	TCT	91	40	55	16	75

*: Nucleotide positions from the 5' end of tRNA genes.

** : These numerals indicated the unusual position or size of introns.

variable 36- to 169-bp sequence, followed by 223 bp-, 244 bp-, or 310 bp-long LR segments. This type of repetitive unit was present at three different positions with 48, 104, and 119 repeats of SR segments (coordinates: 2,667,399–2,664,097; 2,678,286–2,671,319; and 2,693,179–2,685,053).

The subtype A-II repetitive unit was composed of 73 and 113 repeats of 25-bp-long well-conserved SR seg-

ments (GATGAATCCCAAAGGAATTGAAAG) and a variable 33- to 45-bp spacer sequence, followed by 228-bp-long LR segments. These repetitive units were identified at two different positions (coordinates: 30,837–25,997 and 32,475–39,835).

The second type of repetitive sequences was Tn-like elements. This was extracted by analysis of ORFs that are present in repetitive sequences. Out of nine ORFs

Table 3. The lengths and coordinations of the Tn-like elements and dispersed repetitive sequences.

type of repetitive sequence	length (nt)	coordinates	direction
Tn-like element	1779	57900 - 59678	<
		660704 - 662475	<
		1034923 - 1036702	>
		1077993 - 1079772	>
		1157568 - 1159339	>
		1907731 - 1909502	>
		2562560 - 2564337	<
truncated Tn-like element	345	120380 - 120715	>
		715738 - 716074	>
		740655 - 740991	<
		760848 - 761184	<
		1115939 - 1116276	<
		1146090 - 1146426	>
		1623632 - 1623968	<
		1777801 - 1778137	>
		1784191 - 1784527	<
dispersed repetitive unit			
subtype I	1459	1309109 - 1310567	>
		1367448 - 1368906	<
subtype II	1322	823496 - 824817	<
		885027 - 886348	>
		1164946 - 1166267	<
subtype III	844	1922728 - 1924049	<
		163506 - 164343	<
		269897 - 270739	<
		833260 - 834103	<
subtype IV	518	1164031 - 1164874	<
		2269440 - 2270271	<
		2394814 - 2395331	>
		2396410 - 2396927	<

assigned as transposases, seven ORFs with 136 amino acid residues were located within the highly conserved 1779-bp-long sequences. Thus, these regions were assigned as a Tn-like element. However, the inverted repeats, which were generally identified at the border of Tn elements were not found. In addition to this complete Tn-like element, its truncated form was also identified. The end region of each truncated element was identical to the end sequence of the Tn-like element. The size of this element was 352 bp long and the transposase was not found in this truncated element. The position and direction of these Tn-like repetitive elements and truncated form of the Tn-like element are summarized in Table 3.

The third type of repetitive units comprised four subtypes of different repeating units from 518 bp to 1459 bp, and appeared dispersed along the entire genome. The length and coordinate of these four different repetitive sequences are summarized in Table 3. These long repetitive sequences often interfered with the assembly of raw sequencing data, but the problem could be minimized by registration of repetitive sequences into a cross-match database in PhredPhrap. The biological significance of these repeats, except for the Tn-like element, is not known.

3.4. Assignment and similarity search of potential protein-coding regions

The assignment of potential protein-coding regions was performed as the previous paper¹ with some modifications. The ORFs which consist of longer than 100 codons starting with ATG or GTG were designated by a two-letter code (ST) plus a four-digit number indicating the ORF position. The ORFs which were entirely included within longer ORF on either strand and showed neither similarity nor motif sequences were not assigned as potential protein-coding regions. There were a total of 2558 ST-class ORFs.

Subsequently, shorter ORFs consisting of 50 to 99 codons were extracted from the regions where no ST-class ORF was assigned plus a neighboring 150-bp region overlapping the next ST-class ORFs. Among these short ORFs, those which possessed significant similarity to either the registered sequences in databases or to protein motifs were taken as potential protein-coding regions and named with a three-letter code (STS) plus a three-digit number indicating the ORF position. The organization of all ORFs assigned among the entire genome is shown in Fig. 1.

Thus, a total of 2826 ORFs were assigned along the entire genomic sequence (ST-class ORFs = 2558, STS-class ORFs = 268). The average size of the ORFs was

Table 4. Summary of functional classification of ORFs.

functional categories	number of ORFs
Amino acid biosynthesis	77
Purines, pyrimidines, nucleosides, and nucleotides	55
Fatty acid and phospholipid metabolism	41
Biosynthesis of cofactors, prosthetic groups, and carriers	53
Central intermediary metabolism	34
Energy metabolism	184
Transport and binding proteins	92
DNA metabolism	45
Transcription	47
Protein synthesis	109
Protein fate	42
Regulatory functions	12
Cell envelope	21
Cellular processes	34
Other categories	23
Unknown function	42
Conserved hypothetical protein	921
	total 1832

267 amino acid residues, and the longest one consisted of 1933 residues (ST0620). The 2826 assigned ORFs occupy 83.9% of the entire genome. The G + C content in these coding regions was 33.52%, and that in non-coding regions was 28.51%.

For the assignment of gene function, the products name in public databases with a Zscore greater than or equal to 20, or with over 30% amino acid sequence identity along the entire coding region, were taken for ST-class ORFs. The name of gene products with a Zscore greater than or equal to 12 were also used for STS-class ORFs. Among ST- and STS-class ORFs, there were 911(32.2%) ORFs with assigned function, 921 (32.6%) showed significant similarity to registered sequences with unknown function, and 145 (5.1%) contained some motifs. The remaining 849 ORFs showed no significant similarity to the sequences in public databases and protein motif sequence. The predicted products of ORFs with known functions are summarized in Table 11, in which all products predicted are classified according to functional categories.

The ORFs showing similarities to the registered genes with known functions were classified according to the functional categories and are provided in Table 4; the product names are listed in Table 11. As indicated in Table 4, this genome contained a large number of genes related to energy metabolism. This may be due to the heterotrophic feature of this microorganism, which use a large number of compounds as the energy source.

The codon usages of the proteins encoded by the 2826 ORFs are summarized in Table 5. The codons with A or U at the third position appeared to be more frequently used, probably reflecting the relatively low G + C content of this strain.

3.5. Other features

3.5.1. Genes related to respiration

Since *S. tokodaii* strain7 is aerobic, the genes involved in the TCA cycle were searched, and all of the genes in this cycle were identified on the genome. Two copies of the genes coding for citrate synthase and two subunits of 2-oxoacid:ferredoxin oxidoreductase,²⁴ which plays the same role in archaea as alpha-ketoglutarate dehydrogenase, were identified in the genome. This result suggests that *S. tokodaii* strain7 has a complete TCA cycle system similar to that found in mitochondria of eukaryotes.

On the other hand, some genes in the respiratory chain which participate in the production of ATP were not identified by similarity searching. In particular, cytochrome *c* was not identified in this microorganism. This feature is similar to that in *A. pernix* K1. As cytochrome *c* has an important role in electron transfer to oxygen in eukaryotes, it is likely that this microorganism uses a different molecule for the same function or possesses a different pathway. Participation of Rieske iron-sulfur protein and/or sulfocyanin in place of cytochrome *c* has been suggested in *Sulfolobus* species.

3.5.2. ORFs related to sulfide metabolism

S. tokodaii strain7 is known to oxidize hydrogen sulfide to sulfate intracellularly, and this feature has been applied for the treatment of industrial waste water (personal communication from A. Yamagishi and T. Oshima). By similarity search of the genes relating to this pathway, a total of eight ORFs were detected as genes related to sulfide metabolism (from hydrogen sulfide to sulfate). These ORFs and their products are summarized in Table 6. These enzymes seem to be enough to oxidize sulfur from hydrogen sulfide to sulfate. Although

Table 5. Codon usages of the predicted proteins coded for by the 2826 ORFs.

UUU (Phe)	24003	31.53	UCU (Ser)	14563	19.13	UAU (Tyr)	27402	36.00	UGU (Cys)	3569	4.69
UUC (Phe)	10453	13.73	UCC (Ser)	3715	4.88	UAC (Tyr)	9701	12.74	UGC (Cys)	1395	1.83
UUA (Leu)	36642	48.14	UCA (Ser)	12912	16.96	UAA (Stop)	1532	2.01	UGA (Stop)	812	1.07
UUG (Leu)	7685	10.10	UCG (Ser)	2509	3.30	UAG (Stop)	482	0.63	UGG (Trp)	7716	10.14
CUU (Leu)	14179	18.63	CCU (Pro)	10889	14.31	CAU (His)	7283	9.57	CGU (Arg)	1103	1.45
CUC (Leu)	4243	5.57	CCC (Pro)	3211	4.22	CAC (His)	2636	3.46	CGC (Arg)	355	0.47
CUA (Leu)	12512	16.44	CCA (Pro)	13048	17.14	CAA (Gln)	11842	15.57	CGA (Arg)	700	0.92
CUG (Leu)	2916	3.83	CCG (Pro)	2622	3.44	CAG (Gln)	3964	5.21	CGG (Arg)	271	0.36
AUU (Ile)	30581	40.17	ACU (Thr)	16723	22.00	AAU (Asn)	26487	34.80	AGU (Ser)	12679	16.66
AUC (Ile)	6932	9.11	ACC (Thr)	3368	4.42	AAC (Asn)	10515	13.81	AGC (Ser)	4328	5.69
AUA (Ile)	37755	49.60	ACA (Thr)	12969	17.04	AAA (Lys)	39298	51.63	AGA (Arg)	19962	26.22
AUG (Met)	15580	20.47	ACG (Thr)	3096	4.07	AAG (Lys)	21379	28.09	AGG (Arg)	9091	11.94
GUU (Val)	22315	29.32	GCU (Ala)	19352	25.42	GAU (Asp)	27892	36.64	GGU (Gly)	18010	23.66
GUC (Val)	4210	5.53	GCC (Ala)	3927	5.16	GAC (Asp)	7224	9.49	GGC (Gly)	4105	5.39
GUA (Val)	22471	29.52	GCA (Ala)	15919	20.91	GAA (Glu)	35726	46.93	GGA (Gly)	19959	26.22
GUG (Val)	6357	8.35	GCG (Ala)	2969	3.90	GAG (Glu)	17583	23.10	GGG (Gly)	5585	7.34

Numerals in the second column are the sum of codons present, and those in the third column are the frequency of occurrence per thousand in ORFs of *S. tokodaii* strain7

Table 6. List of ORFs related to H₂S oxidizing reactions.

ORF-ID	Length of ORF (a. a.)	predicted product	predicted function
ST0615	384	sulfide dehydrogenase	H ₂ S = Sulfar + H ₂
ST0971	390	sulfide dehydrogenase (flavocytochrome c) flavoprotein chain	H ₂ S = Sulfar + H ₂
ST1010	208	sulfite oxidase	SO ₃ ²⁻ + O ₂ + H ₂ O = SO ₄ ²⁻ + H ₂ O ₂
ST1839	270	thiosulfate reductase electron transport protein	H ₂ S = S ₂ O ₃ ²⁻
ST2564	293	thiosulfate sulfurtransferase	S ₂ O ₃ ²⁻ + cyanide = SO ₃ ²⁻ + thiocyanate
ST2566	628	sulfite reductase	H ₂ S + 3 Oxidized Ferredoxin + 3 H ₂ O = SO ₃ ²⁻ + 3 Reduced Ferredoxin
ST2567	239	phosphoadenosine phosphosulfate reductase	SO ₃ ²⁻ + H ₂ = SO ₃ ²⁻ + H ₂ O
ST2568	412	sulfate adenyltransferase	ATP + SO ₄ ²⁻ = PPi + AMP-SO ₄ ²⁻

the products of these ORFs have not yet been isolated, confirmation of their activities would provide the valuable information for sulfide metabolism in this microorganism and facilitate the improvement of this system for industrial applications.

3.5.3. Possible integration of plasmid into chromosome

It has previously been reported that *Sulfolobus* sp. NOB8H2 contains 41,229 bp long plasmid, pNOB8.^{25,26} However, no extra-chromosomal genetic unit was detected in the analysis in the present study, whereas these two strains belong to the same genus. To clarify the possibility that a plasmid had been integrated into this genome, we searched ORF homologues of pNOB8. As indicated in Table 7, a total of 37 gene families, consisting of 66 ORFs, were identified in the *S. tokodaii* strain7 genome, although their organization and the number of genes included in each gene family are variable in the genome. These results imply that a plasmid from an ancestral species was integrated into the *S. tokodaii* strain7 genome, and following rearrangement and duplication of the genome resulted in the variation of the order and

copy number of genes originally encoded by the plasmid. Peng et al. previously showed that pHEN7, a plasmid isolated from *S. islandicus*, was integrated into the host chromosome by integrase encoded by the plasmid itself.²⁷ The length of the plasmid is very short, 7.83 kb, and the integrase is highly similar to that found in SVV1 virus.²⁸ Two kinds of integrases were identified in the genome of *S. tokodaii* strain7, one was partially similar to that in SSV1 virus and the other was highly similar to the integrase/recombinase possessed in bacteria. The mechanism of integration in *S. tokodaii* strain7 may be different from that in *S. islandicus*.

Out of the 66 ORFs identified, the functions of only three genes were assigned, and the remaining 63 ORFs were assigned as genes with unknown function. These genes with unknown function may be indispensable for *S. tokodaii* strain7, as these genes have been maintained in the genome until now.

3.5.4. ORFs related to eukaryote-type gene families

Many genes assigned on archaeal genomes are often identified in eukaryotes or eubacteria.^{29–31} During sim-

Table 7. *S. tokodaii* ORFs orthologous to genes coded in plasmid pNOB8.

genes in pNOB8	Length (a. a.)	number of identified ORFs	<i>S. tokodaii</i> ORF-ID orthologous to pNOB8 genes						
ORF01	116	3	ST1310	ST0251	ST1317				
ORF02	188	1	ST0159						
ORF03	81	5	ST0758	ST0236	STS156	STS229	ST1340		
ORF04	422	2	ST1308	ST1318					
ORF05	537	2	ST1307	ST2280					
ORF06	72	1	STS202						
ORF07	50	3	STS147	STS180	STS202				
ORF08	406	3	ST1056	ST2325	ST2008				
ORF09	87	1	STS029						
ORF10	1025	1	ST1326						
ORF17	253	2	ST2505	ST1338					
ORF18	94	7	STS156	STS035	STS149	STS148	ST1340	STS157	ST0236
ORF19	97	4	STS156	STS149	ST0173	STS035			
ORF20	108	1	ST1342						
ORF22	164	1	ST1344						
ORF23	69	1	STS162						
ORF24	72	2	STS163	STS072					
ORF25	92	1	STS164						
ORF26	101	3	ST1345	STS045	ST2503				
ORF27	439	4	ST1346	ST0313	ST2498	ST0255			
ORF28	80	4	STS146	STS224	ST0837	ST1890			
ORF30	248	3	ST1850	ST1311	ST0254				
ORF31	630	1	ST1312						
ORF34	86	1	ST1479						
ORF37	52	1	STS159						
ORF38	604	1	ST1315						
ORF39	165	1	ST1316						
ORF40	65	2	STS155	STS151					
ORF41	110	3	ST0251	ST1310	ST1317				
ORF42	205	1	ST2077						
ORF43	74	1	ST1343						
ORF44	93	7	STS033	ST1486	STS229	STS148	ST1060	ST1339	ST0766
ORF45	470	2	ST1308	ST1318					
ORF46	315	1	ST1320						
ORF48	142	1	ST1321						
ORF50	152	1	ST1316						
ORF52	81	2	STS154	STS152					
total		81							

Table 8. *S. tokodaii* ORFs with similarity to the eukaryote-type genes.

ORF-ID	Length of ORFs (aa)	predicted product	distribution
ST0059	462	selenium-binding protein 2	human, mouse, plant, <i>C. elegans</i>
ST0155	353	flug protein.	<i>A. nidulans</i>
ST0233	288	hypothetical protein.	<i>S. cerevisiae</i>
ST0467	548	DNA replication licensing factor	human, mouse, frog, fruit fly, <i>S. pombe</i>
ST0779	583	acylamino-acid-releasing enzyme.	human, rat, pig, <i>C. elegans</i>
ST0940	767	oligosaccharyl transferase sit3 subunit.	human, mouse
ST0945	254	hypothetical protein.	<i>S. cerevisiae</i> , <i>C. elegans</i>
ST1110	386	nonspecific lipid-transfer protein.	mouse, rat
ST1257	108	s100 calcium-binding protein a3.	human, mouse
ST1350	363	nonspecific lipid-transfer protein.	human
ST1603	125	hypothetical protein.	<i>S. cerevisiae</i>
ST1745	565	acylamino-acid-releasing enzyme.	human, rat, pig, <i>C. elegans</i>
ST2082	360	sterol-regulatory element-binding protein.	human, <i>C. griseus</i>
ST2367	297	serine/threonine protein phosphatase pp2a catalytic subunit.	plant, human, mouse, fruit fly, <i>S. cerevisiae</i>

ilarity search of the *S. tokodaii* strain7 ORFs against public databases, we found that 14 ORFs show similarity to gene families that are only identified in eukaryotes. The homologues of these ORFs are not detected either in other archaea or bacteria. These ORFs and predicted gene products are given in Table 8. In addition to these eukaryotic-type genes, this strain con-

tained 37 archaea-specific ORFs and 53 ORFs which are present both in archaea and eukaryotes. Identification of eukaryote-specific genes together with the absence of the CCA sequence in tRNA genes suggests that this strain is closer to eukaryote among archaea that have had their genomes completely sequenced.

Table 9. Homologous ORFs identified by sequence comparison among the 2826 ORFs of *S. tokodaii* strain7.

number of homologous ORFs* in one group	number of groups	total ORFs
12	1	12
11	1	11
10	1	10
9	4	36
8	6	48
7	11	77
6	14	84
5	25	125
4	56	224
3	94	282
2	281	562
total	494	1471

*: ORFs with Zscore higher than 8 in SW search and amino acid identity higher than 30% in 70% of the entire region were taken.

Table 10. List of the same ORFs identified on the *S. tokodaii* strain7 genome.

predicted function	length of ORFs (a. a.)	ORFs included in one same gene families
transposase	136	ST0043 ST2553 ST1048 ST1096 ST1904 ST0680 ST1158
insertion element protein	242	ST0142 ST1907 ST1139 ST1082 ST0671
not assigned	235	ST0143 ST0857 ST1165 ST2259 ST1981
not assigned	330	ST0847 ST0907 ST1167 ST1916 ST2230
not assigned	240	ST0162 ST1719 ST0798 ST1987
not assigned	347	ST0041 ST1098 ST1050
not assigned	108	ST0042 ST1097 ST1049
not assigned	350	ST0055 ST1957 ST1769
not assigned	381	ST0152 ST2431 ST2008
not assigned	168	ST0254 ST1311 ST1850
not assigned	434	ST0678 ST1160 ST1906
not assigned	144	ST0679 ST1159 ST1905
IS hypothetical protein	202	ST1092 ST1939 ST2495
not assigned	328	ST0046 ST0099
ferredoxin	104	ST0163 ST1175
not assigned	123	ST0856 ST1718
IS element DNA-binding protein	314	ST1122 ST1908
insertion element protein	100	ST1156 ST2015
not assigned	470	ST1304 ST1360
transposase	340	ST1779 ST2430
not assigned	353	ST1952 ST2181
sulfocyanin	220	ST2393 ST2394

3.5.5. Duplication of ORFs

By sequence comparison among the ORFs, 1471 ORFs were grouped into 494 families (Table 9). The result can be interpreted that ORFs in each group were generated by duplication of an ancestral sequence. Such homologous ORFs were present in the genome as either tandem repeats of a single ORF or the repeats of a single or a cluster of ORFs at different locations.

Within these families, we detected 22 families that are composed of ORFs which have over 95% length identity and over 95% amino acid identity. The gene families detected are summarized in Table 10. In addition to transposase and IS-element related proteins, sixteen kinds of proteins which possessed the same sequence were detected at different positions in the genome. This result

suggests that rearrangement or duplication of the genome may continue to occur.

3.5.6. Other notable genes

The structure of the operon encoding H⁺-ATPase (ST1435, ST1436, ST1437, ST1438, ST1439, STS172) of this strain has already been reported by Denda et al.³²⁻³⁵ Data comparison indicated that the length of the delta subunit is 80 amino acids shorter than that of ST1435. Because some subunits which were present in *E. coli*³⁶ or *Syneccoccus*³⁷ were not identified in the operon cloned from *S. tokodaii* strain7, Denda et al. suggested two possibilities: One is that the operon encodes all of the genes necessary for the whole ATPase complex and the other is that the operon encoding the subunits of

ATPase were split into two (or more) independent operons. Our data suggest that one more ORF (ST1434) upstream of ST1435 may be included in this operon as the subunit of the H⁺-ATPase. It seems reasonable to conclude that this operon in *S. tokodaii* strain7 contains all of the genes of the whole ATPase complex.

The genes containing inteins, defined as the self-splicing portions of a polypeptide sequence,^{38,39} have been identified in the genome of *P. horikoshii* OT3¹ and *A. pernix* K1.² However, no genes with the intein portion were identified in the *S. tokodaii* strain7 genome. Like all other aerobic organisms, the gene for superoxide dismutase^{2,40-42} was also present in this organism (ST2283).

The entire genomic sequence of *S. tokodaii* strain7 is the third sequence from aerobic thermophilic crenarchaeota. The comparison of this genome with other archaea or thermophilic microorganisms may provide important information about the difference between euryarchaeota and crenarchaeota, as well as differences of the thermostability of proteins and the origin or evolution of eukaryotes. More detailed analysis of gene organization and gene structure in comparison with other archaeal genomes is under investigation.

Acknowledgements: We thank Dr. P. Green for donation of the assembling software PhredPhrap and Mr. S. Suharnan for helpful discussions. This work was supported by the Ministry of Economy, Trade and Industry.

References

- Kawarabayasi, Y., Sawada, M., Horikawa, H. et al. 1998, Complete sequence and gene organization of the genome of a hyper-thermophilic archaeobacterium, *Pyrococcus horikoshii* OT3, *DNA Res.*, **5**, 55-76 & 147-155.
- Kawarabayasi, Y., Hino, Y., Horikawa, H. et al. 1999, Complete genome sequence of an aerobic hyper-thermophilic crenarchaeota, *Aeropyrum pernix* K1, *DNA Res.*, **6**, 84-76 & 147-155.
- Bult, C. J., White, O., Olsen, G. J. et al. 1996, Complete genome sequence of the methanogenic archaeon, *Methanococcus jannaschii*, *Science*, **273**, 1058-1073.
- Smith, D. R., Doucette-Stamm, L. A., Deloughery, C. et al. 1997, Complete genome sequence of *Methanobacterium thermoautotrophicum* ΔH: Functional analysis and comparative genomics, *J. Bacteriol.*, **179**, 7135-7155.
- Klenk, H. P., Clayton, R. A., Tomb, J. F. et al. 1997, The complete genome sequence of the hyperthermophilic, sulphate-reducing archaeon *Archaeoglobus fulgidus*, *Nature*, **390**, 364-370.
- Suzuki, T., Iwasaki, T., Uzawa, T. et al. *Sulfolobus tokodaii* sp. nov. (f. *Sulfolobus* sp. strain7), a new member of the genus *Sulfolobus* isolated from Beppu hot springs, Japan, *Extremophiles*, in press.
- Kondo, S., Yamagishi, A., and Oshima, T. 1993, A physical map of the sulfur-dependent archaeobacterium *Sulfolobus acidocaldarius* 7 chromosome, *J. Bac.*, **175**, 1532-1536.
- Ewing, B., Hillier, L., Wendl, M. C., and Green, P. 1998, Base-calling of automated sequencer traces using Phred. I. Accuracy assessment, *Genome Res.*, **8**, 175-185.
- Ewing, B. and Green, P. 1998, Base-calling of automated sequencer traces using Phred. II. Error probabilities, *Genome Res.*, **8**, 186-194.
- Gordon, D., Abajian, C., and Green, P. 1998, Consed: A graphical tool for sequence finishing, *Genome Res.*, **8**, 195-202.
- Brock, T. D., Brock, K. M., Belly, R. T., and Weiss, R. L. 1972, *Sulfolobus*: A new genus of sulfur-oxidizing bacteria living at low pH and high temperature, *Arch. Mikrobiol.*, **84**, 54-68.
- Yamagishi, A. and Oshima, T. 1990, Circular chromosomal DNA in the sulfur-dependent archaeobacterium *Sulfolobus acidocaldarius*, *Nucleic Acids Res.*, **18**, 1133-1136.
- Smith, T. F. and Waterman, M. S. 1981, Identification of common molecular subsequences, *J. Mol. Biol.*, **147**, 195-197.
- Fleishmann, R. D., Adams, M. D., White, O. et al. 1995, Whole-genome random sequencing and assembly of *Haemophilus influenzae* Rd, *Science*, **269**, 496-512.
- Lowe, T. M. and Eddy, S. R. 1997, tRNAscan-SE: A program for improved detection of transfer RNA genes in genomic sequence, *Nucleic Acids Res.*, **25**, 955-964.
- Nakamura, Y. and Tabata, S. 1997, Codon-anticodon assignment and detection of codon usage trends in seven microbial genomes, *Microbial & Comparative Genomics*, **2**, 299-312.
- Li, F., Wang, J., and Steitz, T. A. 2000, *Sulfolobus shibatae* CCA-adding enzyme forms a tetramer upon binding two tRNA molecules: A scrunching-shuttling model of CCA specificity, *J. Mol. Biol.*, **304**, 483-492.
- Cudny, H., Lupski, J. R., Godson, G. N., and Deutscher, M. P. 1986, Cloning, sequencing, and species relatedness of the *Escherichia coli* cca gene encoding the enzyme tRNA nucleotidyltransferase, *J. Biol. Chem.*, **261**, 6444-6449.
- Deutscher, M. P., Foulds, J., and McClain, W. H. 1974, Transfer ribonucleic acid nucleotidyl-transferase plays an essential role in the normal growth of *Escherichia coli* and in the biosynthesis of some bacteriophage T4 transfer ribonucleic acids, *J. Biol. Chem.*, **249**, 6696-6699.
- Zhu, L. and Deutscher, P. M. 1987, tRNA nucleotidyltransferase is not essential for *Escherichia coli* viability, *EMBO J.*, **6**, 2473-2477.
- Aebi, M., Kirchner, G., Chen, J.-Y., Vijayraghavan, U., Jacobson, A., Martin, N. C., and Abelson, J. 1990, Isolation of a temperature-sensitive mutant with altered tRNA nucleotidyltransferase and cloning of the gene encoding tRNA nucleotidyltransferase in the yeast *Saccharomyces cerevisiae*, *J. Biol. Chem.*, **265**, 16216-16220.
- Curnow, A. W., Hong, K. W., Yuan, R., Kim, S. I., Martins, O., Winkler, W., Henkin, T. M., and Söll, D. 1997, Glu-tRNA^{Gln} amidotransferase: A novel heterotropic enzyme required for correct decoding of glutamine codons during translation, *Proc. Natl. Acad. Sci.*

- USA, **94**, 11819–11826.
23. Curnow, A. W., Tumbula, D. L., Pelaschier, J. T., Kin, B., and Söll, D. 1998, Glutamyl-tRNA^{Gln} amidotransferase in *Deinococcus radiodurans* may be confined to asparagine biosynthesis, *Proc. Natl. Acad. Sci. USA*, **95**, 12838–12843.
 24. Zhang, Q., Iwasaki, T., Wakagi, T., and Oshima, T. 1996, 2-Oxoacid:Ferredoxin oxidoreductase from the thermoacidophilic archaeon, *Sulfolobus* sp. Strain7, *J. Biochem.*, **120**, 587–599.
 25. Schleper, C., Horz, I., Janekovic, D., Murphy, J., and Zilling, W. 1995, A multicopy plasmid of the extremely thermophilic archaeon *Sulfolobus* effects its transfer to recipients by mating, *J. Bacteriol.*, **177**, 4417–4426.
 26. She, Q., Phan, H., Garrett, R. A., Albers, S. V., Stedman, K. M., and Zilling, W. 1998, Genetic profile of pNOB8 from *Sulfolobus*: the first conjugative plasmid from an archaeon, *Extremophiles*, **2**, 417–425.
 27. Peng, X., Holz, I., Zilling, W., Garrett, R. A., and She, Q. 2000, Evolution of the family of pRN plasmids and their integrase-mediated insertion into chromosome of the crenarchaeon *Sulfolobus solfataricus*, *J. Mol. Biol.*, **303**, 449–454.
 28. Schleper, C., Kubo, K., and Zilling, W. 1992, The particle SVV1 from the extremely thermophilic archaeon *Sulfolobus* is a virus: demonstration of infectivity and of transfection with viral DNA, *Proc. Natl. Acad. Sci. USA*, **89**, 7645–7649.
 29. Makarova, K. S., Aravind, L., and Koonin, E. V. 1999, A superfamily of archaeal, bacterial, and eukaryotic proteins homologous to animal transglutaminases, *Protein Sci.*, **8**, 1714–1719.
 30. DiRuggiero, J., Brown, J. R., Bogert, A. P., and Robb, F. T. 1999, DNA repair systems in archaea: mementos from the last universal common ancestor?, *J. Mol. Evol.*, **49**, 474–484.
 31. Jain, R., Rivera, M. C., and Lake, J. A. 1999, Horizontal gene transfer among genomes: the complexity hypothesis, *Proc. Natl. Acad. Sci. USA*, **96**, 3801–3806.
 32. Denda, K., Konishi, J., Oshima, T., Date, T., and Yoshida, M. 1988, The membrane-associated ATPase from *Sulfolobus acidocaldarius* is distantly related to F1-ATPase as assessed from the primary structure of its a-subunit, *J. Biol. Chem.*, **263**, 6012–6015.
 33. Denda, K., Konishi, J., Oshima, T., Date, T., and Yoshida, M. 1988, Molecular cloning of the b-subunit of a possible non-F0F1 type ATP synthase from the acidothermophilic archaeobacterium, *Sulfolobus acidocaldarius*, *J. Biol. Chem.*, **263**, 17251–17254.
 34. Denda, K., Konishi, J., Oshima, T., Date, T., and Yoshida, M. 1989, A gene encoding the proteolipid subunit of *Sulfolobus acidocaldarius* ATPase complex, *J. Biol. Chem.*, **264**, 7119–7121.
 35. Denda, K., Konishi, J., Hajiro, K., Oshima, T., Date, T., and Yoshida, M. 1990, Structure of an ATPase operon of an acidothermophilic archaeobacterium, *Sulfolobus acidocaldarius*, *J. Biol. Chem.*, **265**, 21509–21513.
 36. Walker, J. E., Saraste, M., and Gay, N. J. 1984, The unc operon. Nucleotide sequence, regulation and structure of ATP-synthase, *Biochem. Biophys. Acta*, **768**, 164–200.
 37. Cozens, A. L. and Walker, J. E. 1987, The organization and sequence of the genes for ATP synthase subunits in the cyanobacterium *Synechococcus* 6301. Support for an endosymbiotic origin of chloroplasts, *J. Mol. Biol.*, **194**, 359–383.
 38. Hirata, R., Ohsumi, Y., Nakano, A., Kawasaki, H., Suzuki, K., and Anraku, Y. 1990, Molecular structure of a gene, VMA1, encoding the catalytic subunit of H⁺-translocating adenosine triphosphatase from vacuolar membranes of *Saccharomyces cerevisiae*, *J. Biol. Chem.*, **265**, 6726–6733.
 39. Kane, P. M., Yamashiro, C. T., Wolczyk, D. F., Neff, N., Goebel, M., and Stevens, T. H. 1990, Protein splicing converts the yeast TFP1 gene product to the 69-kD subunit of the vacuolar H⁺-adenosine triphosphatase, *Science*, **250**, 651–657.
 40. Yamano, S., Sako, Y., Nomura, N., and Maruyama, T. 1999, A cambialistic SOD in a strictly aerobic hyperthermophilic archaeon, *Aeropyrum pernix*, *J. Biochem. (Tokyo)*, **126**, 218–225.
 41. Knapp, S., Kardinahl, S., Hellgren, N., Tibbelin, G., Schafer, G., and Ladenstein, R. 1999, Refined crystal structure of a superoxide dismutase from the hyperthermophilic archaeon *Sulfolobus acidocaldarius* at 2.2 Å resolution, *J. Mol. Biol.*, **285**, 689–702.
 42. Kardinahl, S., Anemuller, S., and Schafer, G. 2000, The hyper-thermostable Fe-superoxide dismutase from the Archaeon *Acidianus ambivalens*: characterization, recombinant expression, crystallization and effects of metal exchange, *Biol. Chem.*, **381**, 1089–1101.

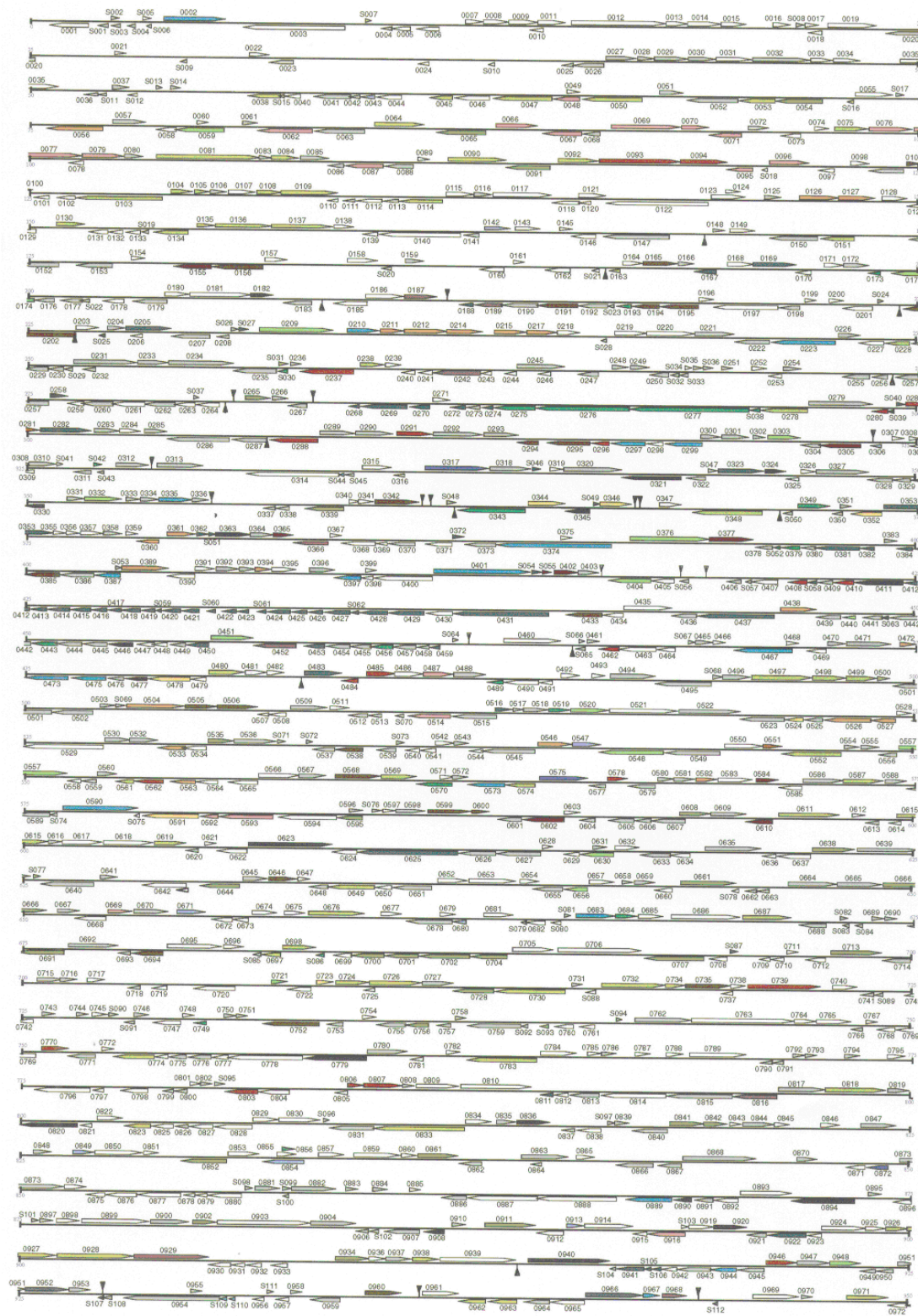


Figure 1. Gene map of the *Sulfolobus tokodaii* strain7 genome. A total of 2826 putatively identified potential protein-coding regions are shown, and the direction of transcription is indicated by arrows. Each line in the figure represents 100,000 bp of sequence in the *S. tokodaii* strain7 genome. Positions are given by numbers above or below the tic marks in each row. The ORFs are color coded by role category as described in the key. ORF ID numbers placed above or below the ORFs correspond to those in Table 11. The rRNA operon and tRNA genes, are labeled.

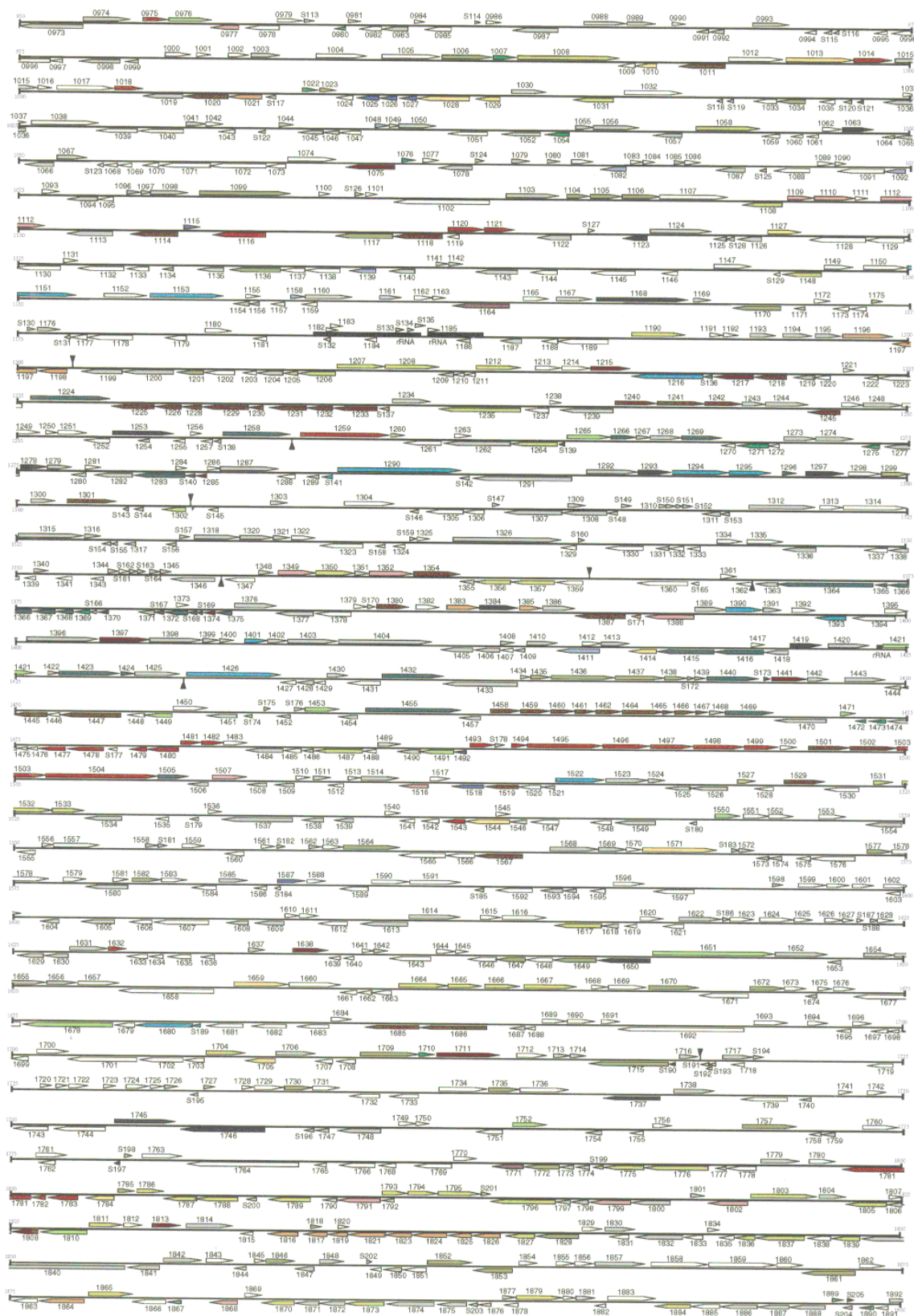


Figure 1. Continued.

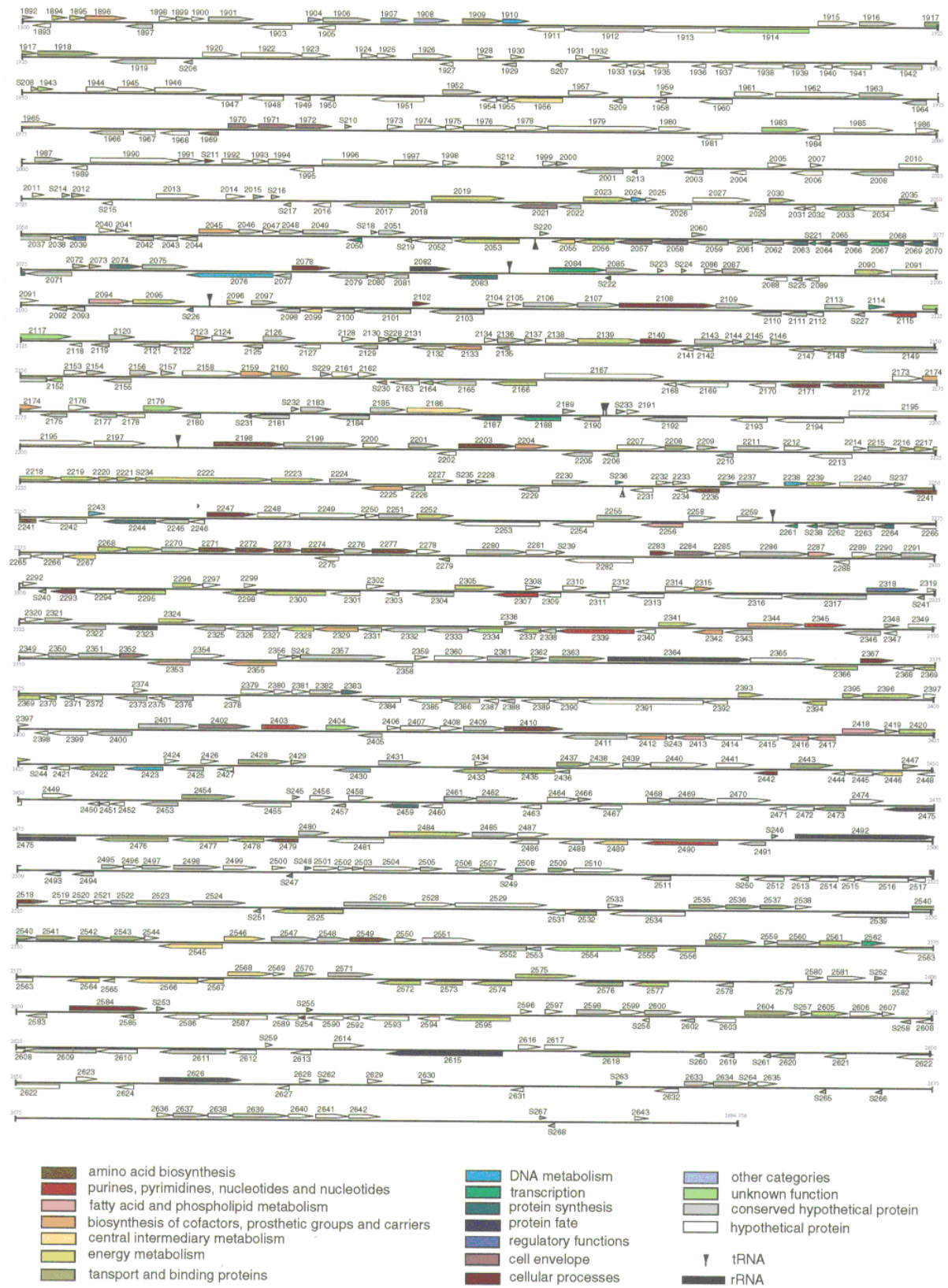


Figure 1. Continued.