



Title	Ecophysiology of <i>Thioploca ingrca</i> as revealed by the complete genome sequence supplemented with proteomic evidence
Author(s)	Kojima, Hisaya; Ogura, Yoshitoshi; Yamamoto, Nozomi; Togashi, Tomoaki; Mori, Hiroshi; Watanabe, Tomohiro; Nemoto, Fumiko; Kurokawa, Ken; Hayashi, Tetsuya; Fukui, Manabu
Citation	ISME journal, 9(5), 1166-1176 https://doi.org/10.1038/ismej.2014.209
Issue Date	2015-05
Doc URL	http://hdl.handle.net/2115/60178
Type	article (author version)
File Information	Main140919.pdf



[Instructions for use](#)

Ecophysiology of *Thioploca ingrica* as revealed by the complete genome sequence supplemented with proteomic evidence

5

Hisaya Kojima^{1*}, Yoshitoshi Ogura^{2,3}, Nozomi Yamamoto⁴, Tomoyuki Togashi⁵, Hiroshi Mori⁵, Tomohiro Watanabe¹, Fumiko Nemoto¹, Ken Kurokawa^{4,5}, Tetsuya Hayashi^{2,3}, and Manabu Fukui^{1*}

10

1. The Institute of Low Temperature Science, Hokkaido University Kita-19, Nishi-8, Kita-ku, Sapporo 060-0819, Japan

2. Division of Microbial Genomics, Department of Genomics and Bioenvironmental Science, Frontier Science Research Center, University of Miyazaki, Miyazaki 889-1692, Japan

15

3. Division of Microbiology, Department of Infectious Diseases, Faculty of Medicine, University of Miyazaki, Miyazaki 889-1692, Japan

4. Earth-Life Science Institute, Tokyo Institute of Technology, 2-12-1 E3-10, Ookayama, Meguro-ku, Tokyo 152-8550, Japan

5. Department of Biological Information, Graduate School of Bioscience and Biotechnology, Tokyo Institute of Technology, 2-12-1 M6-3, Ookayama, Meguro-ku, Tokyo 152-8550, Japan

20

Running title: Genomics and proteomics of *Thioploca ingrica*

Subject category: Integrated genomics and post-genomics approaches in microbial ecology

25

*Corresponding authors.

Postal address: The Institute of Low Temperature Science, Kita-19, Nishi-8, Kita-ku, Sapporo 060-0819, Japan

Tel: +81-11-706-5460 (H. Kojima)

Fax: +81-11-706-5460 (H. Kojima)

30

E-mail: kojimah@pop.lowtem.hokudai.ac.jp (H. Kojima)

my-fukui@pop.lowtem.hokudai.ac.jp (M. Fukui)

Abstract

Large sulfur-oxidizing bacteria, which accumulate a high concentration of nitrate, are
35 important constituents of aquatic sediment ecosystems. No representative of this group has been
isolated in pure culture, and only fragmented draft genome sequences are available for these
microorganisms. In this study, we successfully reconstituted the genome of *Thioploca ingrlica*
from metagenomic sequences, thereby generating the first complete genome sequence from this
group. The *Thioploca* samples for the metagenomic analysis were obtained from a freshwater
40 lake in Japan. A PCR-free paired-end library was constructed from the DNA extracted from the
samples and was sequenced on the Illumina MiSeq platform. By closing gaps within and
between the scaffolds, we obtained a circular chromosome and a plasmid-like element. The
reconstituted chromosome was 4.8 Mbp in length with a 41.2% GC content. A sulfur oxidation
pathway identical to that suggested for the closest relatives of *Thioploca* was deduced from the
45 reconstituted genome. A full set of genes required for respiratory nitrate reduction to dinitrogen
gas was also identified. We further performed a proteomic analysis of the *Thioploca* sample and
detected many enzymes/proteins involved in sulfur oxidation, nitrate respiration, and inorganic
carbon fixation as major components of the protein extracts from the sample, suggesting that
these metabolic activities are strongly associated with the physiology of *T. ingrlica* in lake
50 sediment.

Key words: *Thioploca*/ *Beggiatoaceae*/ complete genome sequence/ sulfur oxidation/
denitrification

55

Introduction

Although key microbial players in various biogeochemical processes have been identified, they are often not available in pure culture. One major group of such organisms is large sulfur-oxidizing bacteria, which have the capacity to accumulate a large amount of intracellular nitrate. In habitats where a sufficient supply of sulfide is available, they can form dense and widespread bacterial mats. Because of their large biomass and ecophysiological characteristics, they have been regarded as important constituents of aquatic ecosystems. In fact, they have been shown to have a considerable impact on nitrogen and phosphorus dynamics in marine sediments (Zopfi et al., 2001; Schulz and Schulz, 2005; Prokopenko et al., 2013).

Until recently, these organisms were classified into three genera, *Beggiatoa*, *Thioploca*, and *Thiomargarita*, on the basis of their morphological features (Salman et al., 2011). Although recent studies have revealed that these microorganisms are much more diverse than previously thought (Salman et al., 2013), all known nitrate-storing sulfur oxidizers belong to a particular lineage within the class *Gammaproteobacteria*. Using a large set of 16S rRNA gene sequences, the bacteria in this lineage were reclassified in 2011, and the revised family *Beggiatoaceae* was proposed to encompass all of these bacteria (Salman et al., 2011). The extended family contains three genera along with nine candidate genera, two of which were proposed after the reclassification (Salman et al., 2013). From this family, only a few strains of the genus *Beggiatoa* have been isolated in pure culture (Salman et al., 2013).

As a cultured representative of this family, *Beggiatoa alba* B18LD has been subjected to whole-genome sequencing, and its draft genome sequence is now available in public databases; however, this strain cannot accumulate nitrate. The draft genome sequences of nitrate-storing sulfur oxidizers have been obtained for *Candidatus Isobeggiatoa* and *Candidatus*

Parabeggiatoa, both of which are from coastal marine sediments (Mußmann et al., 2007), and
80 for *Candidatus* Maribeggiatoa, which is from a deep sea sediment that is influenced by
hydrothermal fluid (MacGregor et al., 2013a, 2013b). These genome sequencing projects all
employed whole-genome multiple displacement amplification (MDA) to obtain sufficient
amounts of DNA for sequencing from single bacterial filaments that are expected to consist of
clonal cells. The single-filament approach may be effective for coping with genetic diversities
85 among the morphologically indistinguishable organisms inhabiting the same sediment, but risks
generating chimeric sequences during the process of amplification. Presumably due to the
presence of such chimeric sequences and/or other difficulties (e.g., short reads and repeat
regions in the genomes), sequence assembly in these studies was not fully successful, as
illustrated by large numbers of contigs (>800 contigs). Although these draft genome sequences
90 have provided valuable insights into the physiology and evolution of this group of
microorganisms, the availability of complete genome sequences of large sulfur-oxidizing
bacteria is highly desirable.

Members of the genus *Thioploca* are gliding filamentous bacteria that have a common
sheath surrounding the trichomes. The first description of the genus was made in the early
95 twentieth century, with the type species of *Thioploca schmidlei* obtained from freshwater lake
sediment (Lauterborn, 1907). Marine species with much larger cell sizes were once included in
this genus, but they have been reclassified within a candidate genus, *Candidatus* *Marithioploca*
(Salman et al., 2011). *T. ingrlica* was described as the second species of the genus (Wislouch,
1912) and was listed in the Approved Lists of Bacterial Names after a temporary loss of status
100 as valid name (Maier, 1984). It has been retained in this genus, after the reclassification of the
family *Beggiatoaceae* (Salman et al., 2011), as the sole species whose 16S rRNA gene

sequences are available. Morphologically, *T. ingrlica* was defined as a *Thioploca* species having a trichome 2.0–4.5 µm in diameter (Wislouch, 1912; Maier, 1984). Organisms that fit this description have been found in freshwater and brackish sediments of various localities, and the
105 placement of these organisms in the same species has been generally supported by their 16S rRNA gene sequences (Maier and Murray, 1965; Dermott and Legner, 2002; Høgslund et al., 2010; Kojima et al., 2003; Kojima et al., 2006; Salman et al., 2011; Nemoto et al., 2011; Nemoto et al., 2012; Nishino et al., 1998; Zenskaya et al., 2001; Zenskaya et al., 2009). Nitrate accumulation by *T. ingrlica* was reported in previous studies, though the intracellular
110 concentrations were much lower than in relatives with large vacuoles (Høgslund et al., 2010; Kojima et al., 2007; Zenskaya et al., 2001).

In this study, the complete genome sequence of *T. ingrlica* was reconstituted from metagenomic sequences, uncovering the metabolic and genetic characteristics of this organism. In addition, proteomic analysis was performed to investigate the physiology of *Thioploca* in
115 lake sediments.

Materials and methods

Sampling

Thioploca samples were obtained at a site near the north shore of Lake Okotanpe,
120 approximately 200 m from the site of our previous study (Nemoto et al., 2011). Sediment samples were obtained with an Ekman–Birge grab sampler and were immediately sieved at the site with a 0.25-mm mesh in lake water. The materials retained on the mesh were transferred to lake water and kept at 4 °C in the dark until processing in the laboratory. Upon returning to the laboratory, *Thioploca* filaments were individually removed with forceps from the materials

125 collected by sieving and were then repeatedly washed with filter-sterilized lake water. The
washed filaments were stored at -30 °C for DNA extraction (the samples obtained in 2013) or at
-80 °C for protein extraction (in 2012) or were immediately subjected to protein extraction (in
2011).

130 Sequencing and genome assembly

Genomic DNA was prepared from the washed *Thioploca* filaments using the Wizard
Genomic DNA Purification Kit (Promega). A PCR-free paired-end library was constructed
using the TruSeq DNA PCR-free Sample Prep Kit (Illumina) and sequenced on the Illumina
MiSeq platform (2x300 bp). After trimming low-quality bases (quality score < 25) and adapter
135 sequences, sequence assembly using varying numbers of reads (5,000,000, 4,000,000, 3,000,000,
2,000,000, or 1,000,000 reads) was performed using the Platanus assembler (Kajitani et al.,
2014) with the default setting to determine the optimal number of input sequences for assembly.
Based on the results of this analysis, we used 3,000,000 reads (1,500,000 sequence pairs) for
assembly and obtained 38 scaffolds (>1 kb). Of these, 24 were larger than 5 kb. As 23 out of the
140 24 scaffolds showed similar ranges of GC content and sequence coverage, a single circular
super-scaffold was manually constructed from the 23 scaffolds by PCR. All gaps within and
between scaffolds were closed by sequencing gap-spanning PCR products using an ABI3130xl
DNA sequencer (Applied Biosystems). Because one scaffold (28.75 kb in size) was derived
from a plasmid-like circular DNA molecule, gaps in this scaffold were also closed in the same
145 way. The BWA program (Li and Durbin, 2009) was used for mapping analysis of the Illumina
reads to the reconstituted genome. The reconstituted sequences have been deposited in the
DDBJ/NCBI/EMBL databases (accession no., AP014633).

Annotation

150 Protein-coding sequences (CDSs) were predicted using MetaGeneAnnotator (Noguchi
et al., 2008). The tRNA and rRNA genes were identified using tRNAScan-SE (Lowe and Eddy,
1997) and RNAmmer (Lagesen et al., 2007), respectively. The ORF extraction function in the In
Silico Molecular Cloning Genomic Edition ver. 5.2.65 software (In Silico Biology, Inc.,
Yokohama, Japan) was used for additional gene prediction. Functional annotation of CDSs was
155 performed on the basis of the results of BLASTP searches (Altschul et al., 1997) against the
NCBI non-redundant database, the KEGG database (Kanehisa et al., 2014), and the NCBI
Clusters of Orthologous Groups (COG). CDSs were annotated as hypothetical protein-coding
genes when the CDSs met any of the following three criteria in the top hit of the BLASTP
analysis: (i) E-value greater than $1e-8$, (ii) length coverage less than 30% against query
160 sequence, or (iii) sequence identity less than 60%.

Genomic comparison and phylogenetic analyses

The genome sequences of *Thioploca* relatives were obtained from the following sites.
Beggiatoa sp. SS (*Candidatus* Parabeggiatoa, henceforth “SS”) and *Beggiatoa* sp. PS
165 (*Candidatus* Isobeggiatoa, “PS”) were taken from the NCBI Bacteria draft genome FTP site
(ftp://ftp.ncbi.nlm.nih.gov/genomes/Bacteria_DRAFT); *Beggiatoa* sp. Orange Guaymas
(*Candidatus* Maribeggiatoa, “Orange Guaymas”) and *Beggiatoa alba* B18LD were taken from
DOE-JGI IMG (<https://img.jgi.doe.gov/w>); and *Thioploca araucae* Tha-CCL (*Candidatus*
Marithioploca, “Tha-CCL”) was taken from JCVI
170 (<https://moore.jcvi.org/moore/SingleOrganism.do?speciesTag=THACCL>). Functional

annotation of CDSs was performed as described above. To search for homologous sequences, all-against-all BLASTP searches (E-value < 0.001) were performed for all of the CDSs from the six strains (*T. ingrica* and relatives). On the basis of the BLASTP results, the CDSs were clustered by applying OrthoMCL (Li et al., 2003) with default parameters.

175 Homologous gene clusters consisting of members from two or more strains were used to estimate the gene content phylogeny. The Euclidean distances among strains were calculated from the homolog content matrix among these strains, and complete-linkage hierarchical clustering was then conducted using the `hclust` function in the R software (<http://www.r-project.org/>).

180 A maximum-likelihood phylogenetic tree for *T. ingrica* and its relatives was built upon a concatenated protein sequence alignment of the universal single-copy genes (USCGs) and the *gyrB* gene. The genome of SS was excluded from the analysis because the quality of the draft genome was too low, whereas the sequence of *Thioalkalivibrio nitratreducens* DSM 14787 was included in the analysis. Among the 35 USCGs originally proposed (Raes et al., 2007), four are
185 duplicated in the genomes of Tha-CCL and *T. ingrica* (*rpsG* and *rpsS* in Tha-CCL, *rplM* and *rpsI* in *T. ingrica*). Therefore, these four genes were excluded from the analysis. Because the *valS* and *gyrB* genes of Tha-CCL were both fragmented into two different contigs, they were manually merged. Each homolog cluster was separately aligned using MAFFT software version 7.130b (Kato and Standley, 2013) with default parameters, and then the 32 alignments (31
190 USCGs and *gyrB*) were concatenated. Because some homologs were not found in the draft genomes of PS and Tha-CCL (*gyrB* and *yehF* in PS, *pheS* and *rplK* in Tha-CCL), the sequences of these genes were replaced with gaps in all positions. From the alignment of the concatenated sequences, a maximum likelihood tree was constructed using MEGA5 (Tamura et al., 2011).

The Whelan and Goldman + Frequency model (Whelan et al., 2001) was selected based on the
195 result of a likelihood ratio test. An initial tree for the heuristic search was obtained by applying
the Neighbor-Join and BioNJ algorithms to a matrix of pairwise distances estimated using a JTT
model and then selecting the topology with a superior log likelihood value. A discrete gamma
distribution was used to model evolutionary rate differences among sites (five categories (+G,
parameter = 0.6269)). All positions with more than 50% alignment gaps were eliminated.

200

Protein analysis

Proteomic analyses were conducted using fresh or stored (at -80 °C for 120 days)
samples obtained in 2011 and 2012, respectively. Protein extraction, sodium dodecyl sulfate
polyacrylamide gel electrophoresis (SDS-PAGE), in-gel trypsin digestion, and nano-liquid
205 chromatography tandem mass spectrometry (nano-LC-MS/MS) were performed as previously
described (Watanabe et al., 2012). The reconstituted chromosome sequence was used to
generate the database for protein identification with the Mascot search program ver. 2.4.0
(MS/MS Ion Search; Matrix Science, MA, USA). Search parameters were set as follows: tryptic
digest with a maximum of two missed cleavages; fixed modifications, carbamidomethyl
210 cysteine; variable modifications, methionine oxidation; peptide masses, monoisotopic; positive
charge (+1, +2, +3) of peptide; and mass tolerance of 1.2 Da for precursor ions and 0.8 Da for
product ions. The false discovery rate was estimated using an automatic decoy search against a
randomized database with a significance threshold of $p < 0.05$. Protein detection was judged as
positive when two or more different peptides were detected, and the exponentially modified
215 protein abundance index (emPAI) was calculated as previously described (Ishihama et al., 2005).
The normalized protein content (PC) value was calculated as a percentage of each emPAI in the

summation of all identified proteins.

Results and Discussion

220 Reconstitution of the *T. ingrica* genome

To reconstitute the *T. ingrica* genome, we used the washed *T. ingrica* filaments to generate more than 100,000,000 paired-end metagenomic sequences using the Illumina MiSeq platform. Because low-redundancy sequences derived from minor contaminating bacteria disturbed sequence assembly, 3,000,000 reads were used as input sequences, in accordance with
225 the results of an optimization procedure (see Materials and Methods). As a result, 24 scaffolds larger than 5 kb were constructed, all of which showed a similar range of GC content (from 39.6% to 44.3%, mostly approximately 41%). All but one of the scaffolds showed a similar sequence coverage (from 104x to 116x), suggesting that these scaffolds were derived from a single dominant species. Furthermore, all 16S rRNA gene reads belonging to these scaffolds
230 were identical to the published 16S rRNA sequences of *T. ingrica* samples from Lake Ogawara and Lake Okotanpe (Kojima et al., 2006; Nemoto et al., 2011). We manually constructed a single circular super-scaffold from the 23 scaffolds with similar sequence coverage and, finally, reconstituted a complete circular chromosome by closing all of the gaps within and between the scaffolds. We searched the reconstituted chromosome for the USCGs (Raes et al., 2007) and
235 identified a single homolog for each of the USCGs, except for *rpIM* and *rpsI*, which had been duplicated. This finding confirmed successful reconstitution of the *T. ingrica* genome.

A scaffold that was not included in the reconstituted chromosome showed a GC content (40.3%) similar to the chromosome, but its sequence coverage (381x) was much higher. Sequence analysis revealed that this scaffold encoded a putative replication protein, an integrase,

240 a lytic transglycosidase, and a response regulator (all other putative genes encoded hypothetical proteins), suggesting that the scaffold represented a plasmid-like DNA molecule. In fact, using PCR examination and gap filling, we reconstituted a single circular sequence (28,669 bp) from this scaffold. Although this plasmid-like element encoded an integrase, we did not find any evidence that it had been integrated into the chromosome of *T. ingrica*. The similarity in GC
245 content suggested that the plasmid-like element may be associated with *T. ingrica*, but this idea needs to be experimentally confirmed.

In a mapping analysis of Illumina reads to the reconstituted sequences, 77.1% and 1.7% of the reads used for assembly mapped to the *T. ingrica* chromosome and the plasmid-like element, respectively. This finding indicated that the DNA from the washed *T. ingrica* filaments
250 represented ca. 79% of the DNA preparation used for genome reconstitution.

General features of the *T. ingrica* genome

The general features of the *T. ingrica* chromosome are shown in Table 1 and Figure 1. Nucleotide position 1 of the chromosome was arbitrarily determined, as the replication origin
255 could not be identified. The genome lacks an identifiable *dnaA* gene, and the GC skew fluctuates irregularly throughout the chromosome. The lack of *dnaA* may be a characteristic shared by nitrate-storing sulfur oxidizers, while *dnaA* was identified in *B. alba* B18LD. In all available draft genomes of nitrate-storing sulfur oxidizers, no *dnaA* gene was found by BLASTP or BLASTN searches when using the sequence of *B. alba* as a query.

260 Approximately 48% of the predicted genes of *T. ingrica* were classified into COG functional categories (Figure 1, Supplementary Table S1). No significant differences were identified among *T. ingrica*, Orange Guaymas, and *B. alba* B18LD in their distributions of

genes into COG functional categories (Supplementary Table S1).

265 Genomic comparison and phylogenetic analyses

Homolog distributions among the three strains (*T. ingrica*, Orange Guaymas, and *B. alba* B18LD) are summarized in Supplementary Figure 1 and Supplementary Table S2. The number of genes shared by the three strains suggests that the number of core genes in the family *Beggiatoaceae* is fewer than 1,500.

270 The phylogenetic trees constructed based on the sequences of USCGs or on gene content are shown in Figure 2. Both trees located *T. ingrica* at a position between its marine relatives and *B. alba*. This result is consistent with the phylogenetic position of *T. ingrica* in the family *Beggiatoaceae*, which was inferred by 16S rRNA sequence analysis (Salman et al., 2011).

275

Proteomic analysis

Metagenomic analysis indicated that 79% of the DNA in the samples originated from *T. ingrica*. Considering that *T. ingrica* has a very large cell size and a normal genome size, *T. ingrica* cells may have a large protein/DNA ratio. If so, protein extracts obtained from the washed *T. ingrica* filaments may contain only a trace amount of proteins from microorganisms contaminating the filament sample. From two samples obtained in different years, we detected 864 and 826 proteins; among these, 563 were detected in both samples (Supplementary Tables S3 and S4). Among the 54 ribosomal proteins encoded in the genome, only 21 were detected in both samples, and 20 were not detected in either sample. This result indicated that only a limited portion of the *T. ingrica* proteome was detected and that a considerable portion of the proteins

285

might have been missed by this analysis. However, the result also suggested that the proteins detected in this analysis were contained in the samples at relatively high concentrations. Therefore, the detected proteins are most likely *T. ingrica* proteins, especially those with high relative abundances. Protein abundance was evaluated based on emPAI, and the relative abundances of major proteins are shown in Figure 3. It should be noted that the samples subjected to the protein analysis were mixtures of cells positioned in different layers of sediment. Therefore, the detected proteins might derive from cells experiencing different environmental conditions. For instance, key enzymes for respiration with each of three electron acceptors (oxygen, nitrate, and dimethyl sulfoxide) were all detected in both samples, but there is no way to know whether these proteins originated from the same cell or from cells at different positions in the redox gradient. Nevertheless, the detected proteins provide some insight into the ecophysiology of *T. ingrica*, as described below.

Sulfur oxidation

Although an array of circumstantial evidence has suggested sulfur oxidation by *Thioploca*, there has been no solid evidence to indicate that *T. ingrica* is truly a sulfur oxidizer. In the *T. ingrica* genome, genes involved in the proposed sulfur oxidation pathway were identified, and their gene products were detected in the protein analysis (Figures 3 and 4). The sulfur oxidation pathway of *T. ingrica* is basically the same as that proposed for its closest relatives, but it is notable for being the first case to confirm the presence of a full set of genes. Further, the sequences of all of the genes are now fully available.

As shown in Figures 3 and 4, it appears that sulfide oxidation by *T. ingrica* is mediated by the flavocytochrome *c*/sulfide dehydrogenase encoded by the *fccAB* (*soxEF*) genes

(THII_1691-1692). The gene for another enzyme mediating sulfide oxidation, *sqr*, which
310 encodes a sulfide-quinone oxidoreductase, was also identified in the *T. ingrica* genome
(THII_0779).

It appears that *T. ingrica* oxidizes elemental sulfur to sulfite via the dissimilatory
sulfite reductase (DSR) system. *T. ingrica* possesses a full set of genes for a DSR system,
dsrABEFHCMKLJOP (THII_0611-0622), *dsrNR* (THII_3808-3809), and *dsrS* (THII_1535). In
315 general, genes for the DSR system are mutually exclusive with the *soxCD* genes, with an
exception in an environmental fosmid sequence (Lenk et al., 2012). Accordingly, genes
corresponding to *soxCD* were not found in the *T. ingrica* genome.

For sulfite oxidation to sulfate, *T. ingrica* contains the *aprBA* genes, encoding an
adenosine-5'-phosphosulfate (APS) reductase (THII_3105-3106); the *sat* gene, encoding a
320 sulfate adenylyltransferase (THII_1057); and the *hdrAACB* genes (THII_2277-2280). The
proteins encoded by the *hdr* genes are thought to be components of the Qmo complex, which
interacts with the APS reductase (Dahl et al., 2013). Sulfite oxidation by these enzymes occurs
in the cytoplasm, where sulfite is generated by DsrAB (Pott and Dahl, 1998). The *T. ingrica*
genome also encodes an enzyme that mediates direct sulfite oxidation to sulfate (the *sorAB*
325 genes; THII_2320-2321), but this sulfite oxidase usually works outside of the cytoplasm
(Kappler et al., 2000; Kappler and Bailey, 2005).

The *T. ingrica* genome sequence further suggested that this microorganism is capable
of thiosulfate oxidation by a system referred to as the Sox system (Friedrich et al., 2001). Of the
three core components (SoxAX, SoxYZ, and SoxB) that are thought to be indispensable for
330 thiosulfate oxidation by this system (Hansen et al., 2006), two are encoded in the *soxYZB* gene
cluster (THII_1577-1579) in *T. ingrica* genome. However, genes encoding SoxA and SoxX

were not identified. These genes are usually located in a *soxXYZAB* gene cluster (Kappler and Maher, 2013). In *T. ingrica*, the function of the SoxAX complex might be substituted by a single protein, encoded by a gene located adjacent to the *soxYZB* cluster (THII_1580). The protein
335 encoded by this gene is closely related to a protein designated as “SoxXA” in the PS genome. Although these proteins are larger than conventional SoxA proteins, their sequences exhibit a partial similarity to SoxA.

The lifestyle of *T. ingrica* as a sulfur oxidizer, which was deduced from the genome sequence (Figure 4), is supported by the fact that most of the enzymes described above were
340 among the most abundant proteins identified in the proteomic analysis of the samples directly collected from lake sediment (Figure 3).

Nitrate respiration

A full set of genes required for respiratory reduction of nitrate to N₂ was identified in
345 the *T. ingrica* genome (Figure 4). *T. ingrica* contains a *narGHJI* gene cluster (THII_1673-1676), encoding a membrane-bound nitrate reductase that mediates nitrate reduction to nitrite. In contrast to its closest relatives, the *napAB* genes that encode a periplasmic nitrate reductase were not detected in the *T. ingrica* genome. *T. ingrica* appears to generate N₂ as the end-product of respiration, as it possesses *nirS*, encoding a nitrite reductase (THII_2875); *norCBQ*, encoding
350 a nitric oxide reductase (THII_0334-0336); and *nosZ*, encoding a nitrous oxide reductase (THII_2884). The *norD* gene (THII_0363), which encodes a nitric oxide reductase activation protein, was also identified in the genome. Among these, the *nosZ* gene, which is directly responsible for the generation of N₂, has never been found in the genomes of close relatives of *T. ingrica*. The *nrfA* gene, which encodes an enzyme for dissimilatory nitrite reduction to

355 ammonium, was not found in *T. ingrica*.

The active reduction of nitrate to N₂ by *Thioploca* in sediments has gained strong support from proteomic analysis; enzymes involved in all steps of the successive reduction (Figure 4) were detected, and many of these were major components in the samples (Figure 3). Recently, a large contribution by *Candidatus* Marithioploca to nitrogen loss in anoxic marine
360 sediments was suggested by Prokopenko et al. (2013). This hypothesis was premised on nitrate reduction to ammonium by *Candidatus* Marithioploca and on the cooperative involvement of anaerobic ammonia-oxidizing bacteria. Interaction with ammonia oxidizers was also suggested for sulfur oxidizers inhabiting hydrothermal sediments, but the retention of nitrogen is assumed to take place in this case (Winkel et al., 2013). *T. ingrica* may contribute to nitrogen loss in
365 freshwater lake environments in a different way from that of its marine counterparts.

Nitrogen assimilation

Nitrogen (N₂) fixation by *Beggiatoa* species has been previously reported (Nelson et al., 1982), and genes involved in diazotrophy were identified in the *Beggiatoa alba* B18LD
370 genome. In the *T. ingrica* genome, *nifHDKT* (THII_1806-1809), *nifY* (THII_1811), *nifENX* (THII_1813-1815), *nifZW* (THII_3085-3086), *nifV* (THII_3088), *nifQ* (THII_3125), *nifB* (THII_3130), *nifM* (THII_3738), and two copies of *nifA* (THII_2556, THII_2572) were identified as genes putatively involved in N₂ fixation and its regulation. Among the products of these *nif* genes, the regulatory protein NifA encoded by THII_2556 was detected in the
375 proteomic analysis of both of the samples obtained in 2011 and 2012 (Supplementary Table S4). The others were not detected in the protein analysis, except for NifY, which was detected only in the sample from 2011.

In the sediment of Lake Okotanpe, nitrate and ammonium are both available (Nemoto et al., 2011); thus, N₂ fixation would not be advantageous because of the energetically high cost required for the process. The *T. ingrica* genome encodes both an ammonium transporter (THII_2433) and an assimilatory nitrate reductase, although these proteins were not detected in the protein analysis. *T. ingrica* may also use organic compounds as nitrogen sources, as several subunits for amino acid transporters were detected in the proteomic analysis. Taken together, multiple and flexible nitrogen assimilation pathways are deduced for *T. ingrica*.

385

Nitrate storage

The capacity to store a large amount of nitrate within cells is a unique property of *Thioploca* and closely related sulfur oxidizers. Based on genome sequence information for large marine sulfur oxidizers, a hypothetical model for nitrate accumulation in the vacuole was proposed in analogy to that in plants (Mußmann et al., 2007). The model includes two processes: the electrochemical gradient formed by proton pumping and the NO₃⁻/H⁺ antiporter that depends on the gradient. When the model was proposed, three enzymes were taken into consideration as candidate enzymes responsible for the electrochemical gradient. All three enzymes, namely, a vacuolar-type ATPase, an H⁺-pyrophosphatase (THII_0754), and a Ca²⁺-translocating ATPase (THII_0386), are encoded in the *T. ingrica* genome. The latter two are also present in the *B. alba* genome, which lacks nitrate-storing capacity. In the originally proposed model, the vacuolar-type ATPase and H⁺-pyrophosphatase were assumed to generate a proton motive force, but experiments using *Candidatus Allobeggiatoa halophila* revealed that these enzymes work in the reverse direction by consuming the proton gradient to generate ATP and pyrophosphate (Beutler et al., 2012). The originally proposed model also predicted that the

400

NO₃⁻/H⁺ antiporter is another key protein directly responsible for nitrate accumulation. Whereas BgP0076 and BgP4800 in the genome of PS were assumed to encode NO₃⁻/H⁺ antiporters, corresponding genes were not identified in *T. ingrica*. This finding indicates that the previously proposed model cannot be fully applicable to *T. ingrica*. Among the proteins mentioned above,
405 an H⁺-pyrophosphatase and a Ca²⁺-translocating ATPase were both detected in the proteomic analyses of filament samples collected both in 2011 and 2012 (Supplementary Table S4).

The molecular mechanism of nitrate accumulation in *T. ingrica* and relatives is still not fully understood, and further studies are necessary. The complete genome sequence information of *T. ingrica* obtained in this work will facilitate such studies.

410

Carbon metabolism

Many genes for the enzymes constituting the tricarboxylic acid (TCA) cycle have been identified in other members of the family *Beggiatoaceae*, but some of the genes are missing from their draft genomes (MacGregor et al., 2013a). In the genome of *T. ingrica*, a full set of
415 genes for TCA cycle enzymes was identified. Most of these enzymes were also detected in the protein analysis, confirming that the cycle was actually operating in *T. ingrica* cells. Similarly, genes encoding enzymes for the glycolytic pathway were also identified in the reconstituted genome, and the enzymes mediating all pathway steps were detected in the protein analysis (Supplementary Table S4).

420

In previous studies, acetate assimilation by *T. ingrica* was demonstrated by microautoradiography (Kojima et al., 2007; Høgslund et al., 2010). The gene for acetate uptake, *actP*, which encodes a cation/acetate symporter (THII_3816), was identified in the genome, as was the *acs* gene, encoding an acetyl-CoA synthetase (THII_1554) that converts acetate into

acetyl-CoA. In addition to the TCA cycle, enzymes in the glyoxylate cycle (isocitrate lyase and
425 malate synthase, encoded by *aceA* and *aceB*, respectively; THII_0533-0534) are also encoded
by the genome; thus, acetyl-CoA can probably be used both for dissimilatory and assimilatory
carbon metabolism.

Inorganic carbon fixation by *T. ingrica* has also been demonstrated in a previous study
(Høgslund et al., 2010), and the key enzymatic activities of the reductive pentose phosphate
430 cycle were detected in some strains of *Beggiatoaceae* (McHatton et al., 1996). In the *T. ingrica*
genome, the *rbcLS* genes (THII_3311-3312), which encode the large and small chains of form I
ribulose biphosphate carboxylase (RuBisCO), were identified, in addition to other 8 genes for
the enzymes of Calvin-Benson-Bassham cycle. The gene for form II RuBisCO was not found in
the genome.

435 As described above, the protein analysis suggested that *T. ingrica* cells gain energy
from sulfur oxidation. The detection of other key enzymes, namely, the gene products of *rbcLS*,
actP and *acs*, in both protein samples further supported the notion that acetate and bicarbonate
are serving as carbon sources for *T. ingrica* in lake sediment. Methylophony in *Beggiatoa alba*
has been demonstrated, and genes for key enzymes have been identified (Jewell et al., 2008),
440 but such genes were not found in the *T. ingrica* genome.

Oxygen and dimethyl sulfoxide (DMSO) respiration

T. ingrica most likely has the capacity to use oxygen as a terminal electron acceptor
for respiration, as it contains the *ccoNOQP* gene cluster, which encodes a high-affinity
445 cytochrome c bb3-oxidase (THII_1615-1617), and as the proteomic analysis detected the CcoN,
CcoO, and CcoP proteins (Supplementary Table S4). In the PS genome, the genes for the

low-affinity cytochrome c aa3-oxidase were also found, but their counterparts were not identified in the *T. ingrica* genome. Notably, all of the subunits of DMSO reductase (THII_3261-3263) were detected in both samples subjected to protein analysis, suggesting that

450 DMSO is also utilized by *T. ingrica* for respiration.

Phosphorus metabolism

Accumulation of phosphorus in the form of polyphosphate has been reported for some bacteria of the family *Beggiatoaceae* (Brock et al., 2012; Brock and Schulz-Vogt, 2011; Schulz and Schulz, 2005). In previous studies, *T. ingrica* samples were subjected to toluidine blue staining to visualize intracellular polyphosphate granules, but only negative results were obtained (Høgslund et al., 2010; Nemoto et al., 2011). In the *T. ingrica* genome, the gene for polyphosphate kinase (*ppk*) was identified (THII_2967). Polyphosphate is a multifunctional molecule, and the presence of this gene may not yield the potential to form polyphosphate granules; however, the *ppgK* gene, encoding a polyphosphate glucokinase, was also identified in 460 the genome (THII_0714). Previously, the phytase-encoding gene found in PS was discussed in relation to the acquisition of inorganic phosphate (Mußmann et al., 2007). This gene is also present in the *T. ingrica* (THII_2499) and *B. alba* genomes. In the protein analysis, polyphosphate kinase, polyphosphate glucokinase, and phytase were all repeatedly detected.

465

Osmoregulation by the glycine betaine/proline transporter

Thioploca is unique in habitat preference among nitrate-storing sulfur oxidizers. In contrast to its relatives living in marine sediments, *T. ingrica* inhabits freshwater and brackish environments. To adapt to habitats of differing salinity, osmoregulation systems should play

470 important roles. *T. ingrica* has a full set of genes for the ProU glycine betaine/proline transport system. The system consists of the ATP-binding protein ProV (THII_1898), the permease protein ProW (THII_1897), and the substrate-binding protein ProX (THII_1896). The ProU system is the transport system for various osmolytes, such as glycine betaine, proline, and proline betaine, in Gram-negative bacteria (Sleator and Collin, 2001). Some bacteria harboring
475 the ProU system can cope with increased environmental osmolarity by accumulating glycine betaine (Lucht and Bremer, 1994). A full set of genes for this system is also present in the genome of Orange Guaymas, obtained from marine sediment under the influence of hydrothermal fluid in the Guaymas Basin (MacGregor et al. 2013a). It was suggested that the habitat of Orange Guaymas is exposed to large fluctuation of temperature (McKay et al. 2012),
480 presumably because of changes in supply of the hot water. The osmoregulation systems may also be effective to adapt to frequent changes in composition of surrounding water, brought about by fluctuating mixing ratio of hydrothermal fluid.

Conclusion

485 The complete genome sequence of *T. ingrica* was successfully reconstituted from metagenomic sequences, thus representing the first complete genome sequence of a nitrate-storing sulfur oxidizer. We found that *T. ingrica* possesses all of the genes required for complete denitrification. The presence of the denitrification system in *T. ingrica* was further confirmed by protein analysis, as were several other physiologically important functions, such
490 as sulfur oxidation and inorganic carbon fixation. The complete genome sequence of *T. ingrica* will be a valuable genetic basis for a wide range of future studies on nitrate-storing sulfur oxidizers, which are important constituents of aquatic ecosystems.

Acknowledgements

495 This work was supported by a KAKENHI for Innovative Areas “Genome Science”
(No. 221S0002) from the Ministry of Education, Culture, Sports, Science and Technology
(MEXT) of Japan. This study was also supported by JSPS KAKENHI Grant Number 22370005.

Conflict of interest

500 The authors declare no conflict of interest.

Supplementary information is available at ISMEJ's website.

References

- 505 Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W *et al.* (1997) Gapped
BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids
Res.* 25:3389–3402.
- Beutler M, Milucka J, Hinck S, Schreiber F, Brock J, Mussmann M *et al.* (2012) Vacuolar
510 respiration of nitrate coupled to energy conservation in filamentous *Beggiatoaceae*. *Environ
Microbiol* 14:2911–2919
- Brock J, Rhiel E, Beutler M, Salman V, Schulz-Vogt HN. (2012) Unusual polyphosphate
inclusions observed in a marine *Beggiatoa* strain. *Antonie van Leeuwenhoek.* 101:347–357
- 515 Brock J, Schulz-Vogt HN (2011) Sulfide induces phosphate release from polyphosphate in

cultures of a marine *Beggiatoa* strain. ISME J 5:497–506

520 Dahl C, Franz B, Hensen D, Kesselheim A, Zigann R (2013) Sulfite oxidation in the purple sulfur bacterium *Allochromatium vinosum*: identification of SoeABC as a major player and relevance of SoxYZ in the process. Microbiology 159: 2626–2638

525 Dermott R, Legner M (2002) Dense mat-forming bacterium *Thioploca ingrica* (*Beggiatoaceae*) in eastern Lake Ontario: implications to the benthic food web. J Great Lakes Res 28:688–697

Friedrich CG, Rother D, Bardischewsky F, Quentmeier A, Fischer J. (2001) Oxidation of reduced inorganic sulfur compounds by bacteria: emergence of a common mechanism? Appl Environ Microbiol 67: 2873–2882

530 Hensen D, Sperling D, Trüper HG, Brune DC, Dahl C. (2006). Thiosulphate oxidation in the phototrophic sulphur bacterium *Allochromatium vinosum*. Mol Microbiol 62: 794–810.

535 Høgslund S, Nielsen JL, Nielsen LP (2010) Distribution, ecology and molecular identification of *Thioploca* from Danish brackish water sediments. FEMS Microbiol Ecol 73:110–120

Ishihama Y, Oda Y, Tabata T, Sato T, Nagasu T, Rappsilber J *et al.* (2005) Exponentially modified protein abundance index (emPAI) for estimation of absolute protein amount in proteomics by the number of sequenced peptides per protein. Mol Cell Proteomics 4:1265–1272.

540

Jewell T, Huston SL, Nelson DC. (2008) Methylotrophy in freshwater *Beggiatoa alba* strains. *Appl Environ Microbiol* 74: 5575–5578

Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M *et al.* (2014) Efficient de
545 novo assembly of highly heterozygous genomes from whole-genome shotgun short reads. *Genome Res* 24:1384–1395.

Kanehisa M, Goto S, Sato Y, Kawashima M, Furumichi M, Tanabe M (2014) Data, information, knowledge and principle: back to metabolism in KEGG. *Nucleic Acids Res* 42:D199–205.

550

Kappler U, Bailey S. (2005) Molecular basis of intramolecular electron transfer in sulfite-oxidizing enzymes is revealed by high resolution structure of a heterodimeric complex of the catalytic molybdopterin subunit and a *c*-type cytochrome subunit. *J Biol Chem* 280:24999–25007.

555

Kappler U, Bennett B, Rethmeier J, Schwarz G, Deutzmann R, McEwan AG *et al.* (2000) Sulfite:cytochrome *c* oxidoreductase from *Thiobacillus novellus*. Purification, characterization, and molecular biology of a heterodimeric member of the sulfite oxidase family. *J Biol Chem* 275:13202–13212.

560

Kappler U, Maher MJ (2013) The bacterial SoxAX cytochromes. *Cell Mol Life Sci.* 70:977–992

Katoh K, Standley DM (2013) MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol* 30:772–780

565

Kojima H, Koizumi Y, Fukui M (2006) Community structure of bacteria associated with sheaths of freshwater and brackish *Thioploca* species. *Microb Ecol* 52:765–773

Kojima H, Nakajima T, Fukui M (2007) Carbon source utilization and accumulation of
570 respiration-related substances by freshwater *Thioploca* species. *FEMS Microbiol Ecol* 59:23–31

Kojima H, Teske A, Fukui M (2003) Morphological and phylogenetic characterizations of freshwater *Thioploca* species from Lake Biwa, Japan, and Lake Constance, Germany. *Appl Environ Microbiol* 69:390–398

575

Lagesen K, Hallin PF, Rødland E, Stærfeldt HH, Rognes T, Ussery DW (2007) RNAmmer: consistent annotation of rRNA genes in genomic sequences. *Nucleic Acids Res* 35:3100–3108

Lauterborn R (1907) Eine neue Gattung der Schwefelbakterien (*Thioploca schmidlei* nov. gen.
580 nov. spec.). *Ber Dtsch Bot Ges* 25:238–242

Lenk S, Moraru C, Hahnke S, Arnds J, Richter M, Kube M *et al.* (2012) *Roseobacter* clade bacteria are abundant in coastal sediments and encode a novel combination of sulfur oxidation genes. *ISME J.* 6:2178–2187

585

Li H, Durbin R (2009) Fast and accurate short read alignment with Burrows-Wheeler Transform. *Bioinformatics* 25:1754–1760.

Li L, Stoeckert CJ Jr, Roos DS (2003) OrthoMCL: identification of ortholog groups for
590 eukaryotic genomes. *Genome Res* 13:2178–2189.

Lowe TM, Eddy SR (1997) tRNAscan-SE: a program for improved detection of transfer RNA genes in genomic sequence. *Nucleic Acids Res* 25:955–964

595 Lucht JM, Bremer E (1994) Adaptation of *Escherichia coli* to high osmolarity environments: Osmoregulation of the high-affinity glycine betaine transport system ProU. *FEMS Microbiol Rev* 14:3–20

MacGregor BJ, Biddle JF, Harbort C, Matthysse AG, Teske A (2013a) Sulfide oxidation, nitrate
600 respiration, carbon acquisition, and electron transport pathways suggested by the draft genome of a single orange Guaymas Basin *Beggiatoa* (*Cand. Maribeggiatoa*) sp. filament. *Mar Genomics* 11:53–55.

MacGregor BJ, Biddle JF, Teske A (2013b) Mobile elements in a single-filament orange
605 Guaymas Basin *Beggiatoa* sp. genome: evidence for genetic exchange with cyanobacteria. *Appl Environ Microbiol* 79:3974–3985

Maier SH, Murray GE (1965) The fine structure of *Thioploca ingrica* and a comparison with

Beggiatoa. Can J Microbiol 11:645–663

610

Maier S (1984) Description of *Thioploca ingraca* sp. nov., nom. rev. Int J Syst Bacteriol 34: 344–345

615 McHatton SC, Barry JP, Jannasch HW, Nelson DC (1996) High nitrate concentrations in vacuolate, autotrophic marine *Beggiatoa* spp. Appl Environ Microbiol 62:954–958

McKay LJ, MacGregor BJ, Biddle JF, Albert DB, Mendlovitz HP, Hoer DR *et al.* (2012) Spatial heterogeneity and underlying geochemistry of phylogenetically diverse orange and white *Beggiatoa* mats in Guaymas Basin hydrothermal sediments. Deep-Sea Res. I 67: 21–31.

620

Mußmann M, Hu FZ, Richter M, de Beer D, Preisler A, Jørgensen BB *et al.* (2007) Insights into the genome of large sulfur bacteria revealed by analysis of single filaments. PLoS Biol 5:1923–1937

625 Nelson DC, Waterbury John B, Jannasch Holger W. (1982). Nitrogen fixation and nitrate utilization by marine and freshwater *Beggiatoa*. Arch Microbiol 133: 172–177.

Nemoto F, Kojima H, Fukui M (2011) Diversity of freshwater *Thioploca* species and their specific association with filamentous bacteria of the phylum *Chloroflexi*. Microb Ecol

630 62:753–764

Nemoto F, Kojima H, Ohtaka A, Fukui M (2012) Filamentous sulfur-oxidizing bacteria of the genus *Thioploca* from Lake Tonle Sap in Cambodia. *Aquat Microb Ecol* 66: 295–300

635 Nishino M, Fukui M, Nakajima T (1998) Dense mats of *Thioploca*, gliding filamentous sulfur-oxidizing bacteria in Lake Biwa, central Japan. *Water Res* 32:953–957

Noguchi H, Taniguchi T, Itoh T (2008) MetaGeneAnnotator: detecting species-specific patterns of ribosomal binding site for precise gene prediction in anonymous prokaryotic and phage

640 genomes. *DNA Res* 15:387–396

Pott AS, Dahl C (1998) Sirohaem sulfite reductase and other proteins encoded by genes at the *dsr* locus of *Chromatium vinosum* are involved in the oxidation of intracellular sulfur. *Microbiology* 144:1881–1894

645

Prokopenko MG, Hirst MB, De Brabandere L, Lawrence DJ, Berelson WM, Granger J *et al.* (2013) Nitrogen losses in anoxic marine sediments driven by *Thioploca*–anammox bacterial consortia. *Nature* 500:194–198

650 Raes J, Korb J, Lercher MJ, von Mering C, Bork P (2007) Prediction of effective genome size in metagenomic samples. *Genome Biol* 8:R10.

Salman V, Amann R, Girth AC, Polerecky L, Bailey JV, Høglund S *et al.* (2011) A single-cell sequencing approach to the classification of large, vacuolated sulfur bacteria. *System Appl*

655 Microbiol 34:243–259

Salman V, Bailey JV, Teske A (2013) Phylogenetic and morphologic complexity of giant sulphur bacteria. *Antonie van Leeuwenhoek* 104:169–186

660 Schulz HN, Schulz HD (2005) Large sulfur bacteria and the formation of phosphorite. *Science* 307:416–418

Sleator RD, Colin H (2001) Bacterial osmoadaptation: the role of osmolytes in bacterial stress and virulence. *FEMS Microbiol Rev* 26:49–71

665

Tamura K, Peterson D, Peterson N, Stecher G, Nei M, Kumar S (2011) MEGA5: Molecular Evolutionary Genetics Analysis using Maximum Likelihood, Evolutionary Distance, and Maximum Parsimony Methods. *Mol Biol Evol* 28:2731–2739.

670 Watanabe T, Kojima H, Fukui M (2012) Draft genome sequence of a psychrotolerant sulfur-oxidizing bacterium, *Sulfuricella denitrificans* skB26, and proteomic insights into cold adaptation. *Appl Environ Microbiol* 78:6545–6549

Whelan S, Goldman N (2001) A general empirical model of protein evolution derived from
675 multiple protein families using a maximum-likelihood approach. *Mol Biol Evol* 18:691–699.

Winkel M, de Beer D, Lavik G, Peplies J, Mußmann M (2013) Close association of active

nitrifiers with *Beggiatoa* mats covering deep-sea hydrothermal sediments. *Environ Microbiol* 16:1612–1626

680

Wislouch SM (1912) *Thioploca ingraca* nov. sp. *Berichte der deutschen botanischen Gesellschaft* 30: 470–474

Zemskaya TI, Namsaraev BB, Dul'tseva NM, Khanaeva TA, Golobokova LP, Dubinina GA *et al.* (2001) Ecophysiological characteristics of the mat-forming bacterium *Thioploca* in bottom sediments of the Frolikha Bay, northern Baikal. *Microbiology* 70:335–341

Zemskaya TI, Chernitsyna SM, Dul'tseva NM, Sergeeva VN, Pogodaeva TV, Namsaraev BB (2009) Colorless sulfur bacteria *Thioploca* from different sites in Lake Baikal. *Microbiology* 78:117–124

690

Zopfi J, Kjær T, Nielsen LP, Jørgensen BB (2001) Ecology of *Thioploca* spp.: nitrate and sulfur storage in relation to chemical microgradients and influence of *Thioploca* spp. on the sedimentary nitrogen cycle. *Appl Environ Microbiol* 67:5530–5537

695 **Figure legends**

Figure 1. Circular view of the reconstituted *T. ingrica* chromosome. The outermost two circles show predicted CDSs on the plus and minus strands. The CDSs are color-coded by COG categories (unclassified genes are shown in gray). The third and fourth circles indicate GC content (mean value = 41%) and GC skew, respectively. Genes encoding rRNA (red) and tRNA (gray) are shown in the two innermost circles.

700

Figure 2. Phylogenetic trees of *T. ingrica* and its relatives. (A) A tree constructed using the concatenated sequences of USCGs. The numbers at each node represent the percentage values from 1,000 bootstrap resamplings. (B) A tree constructed using the gene contents of each strain.

705 Figure 3. Overviews of the proteomic analysis results from washed *Thioploca* filaments obtained in 2011 (A) and 2012 (B). The proteins are sorted according to their relative abundances in each sample, and the 150 most abundant proteins are shown. The data for the 500 most abundant proteins are outlined in the small nested boxes. Proteins involved in sulfur oxidation and nitrate respiration are indicated in red and green, respectively. The RuBisCO proteins are indicated by blue bars, and ribosomal proteins are shown in black for comparison.

710

Figure 4. Pathways for sulfur oxidation and dissimilatory nitrate reduction deduced from the *T. ingrica* genome. Proteins in red were detected in both samples subjected to protein analysis, and those in green were detected in one of the samples. Proteins encoded by the genome but not detected in the protein analysis are

shown in blue. Gray dashed lines indicate alternative pathways mediated by the enzymes (shown in gray) that

715 are not encoded by the *T. ingrica* genome but are found in its close relatives.

720

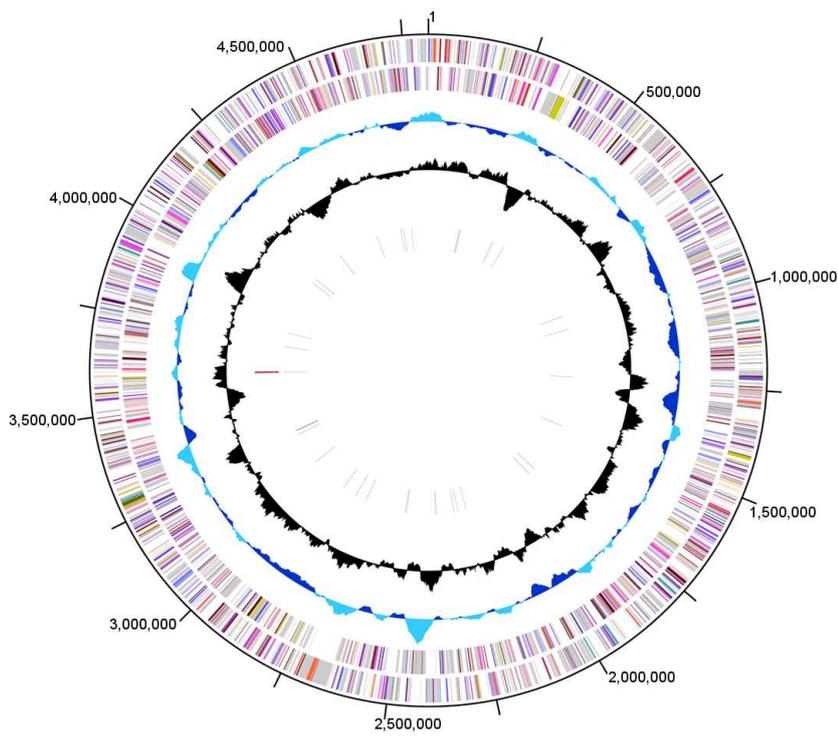


Fig.1

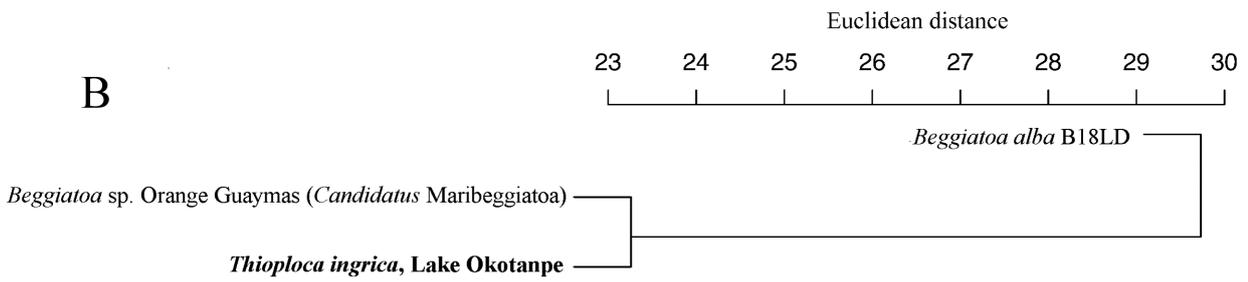
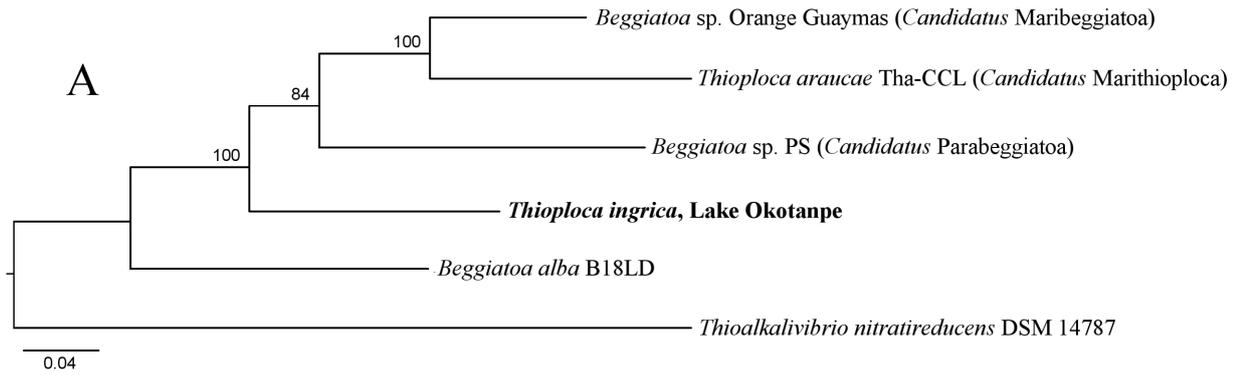
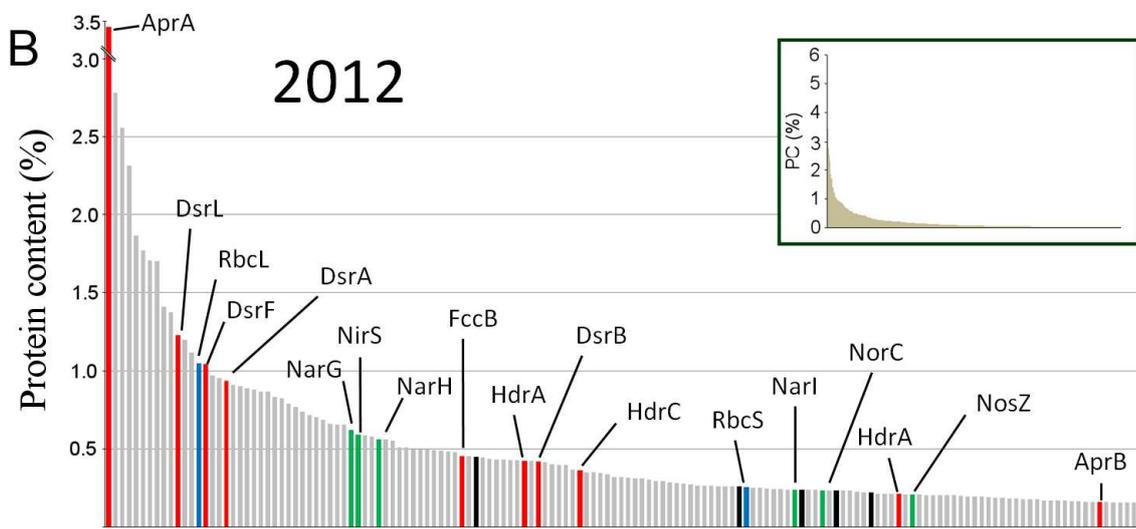
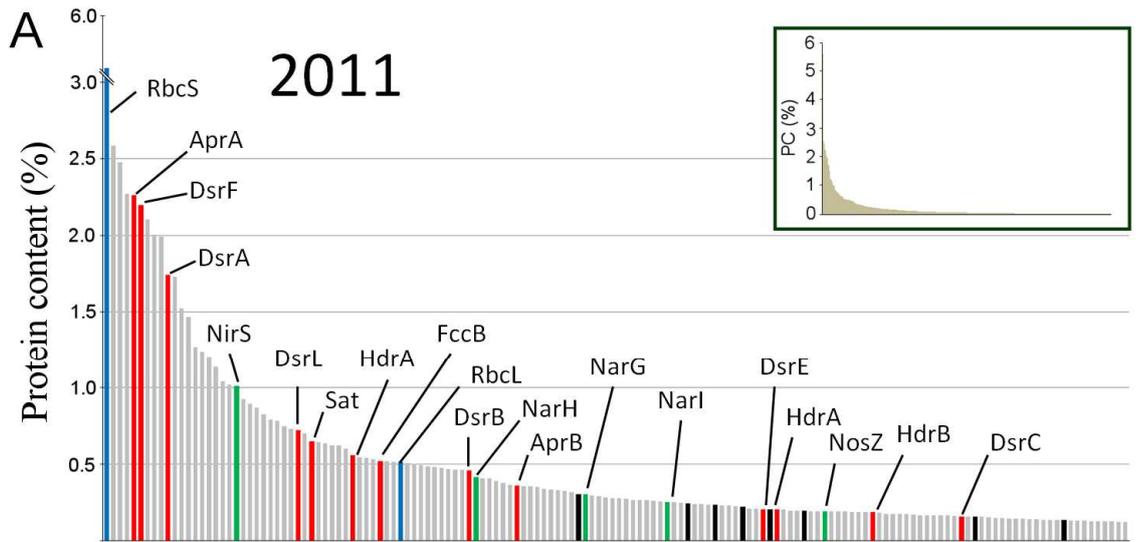


Fig.2



725

Fig.3

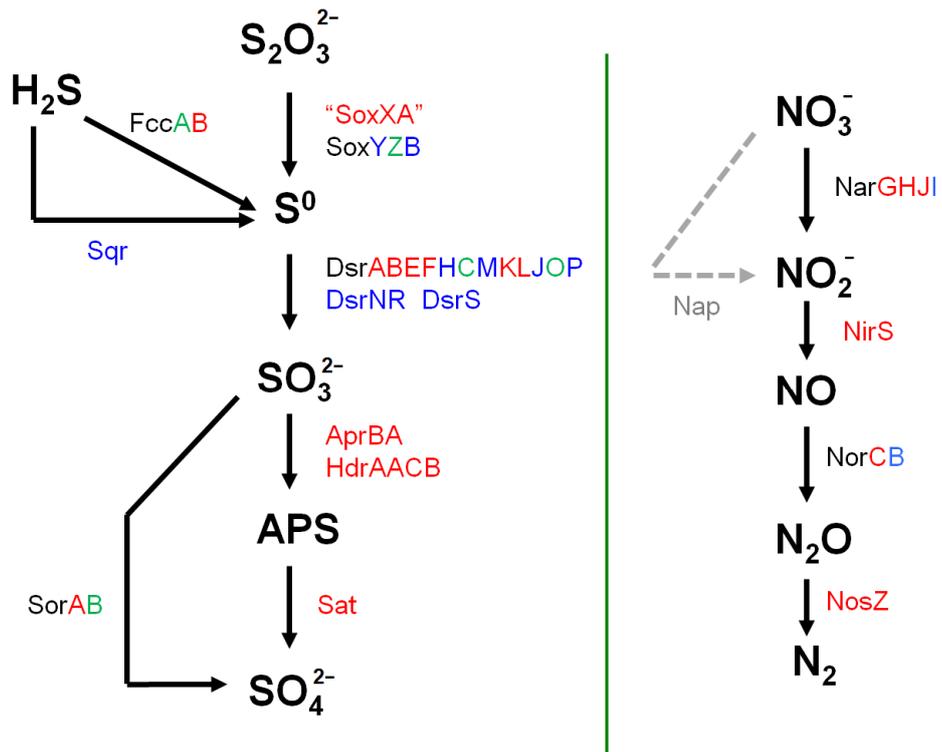


Fig.4