



# ERNEST ORLANDO LAWRENCE BERKELEY NATIONAL LABORATORY

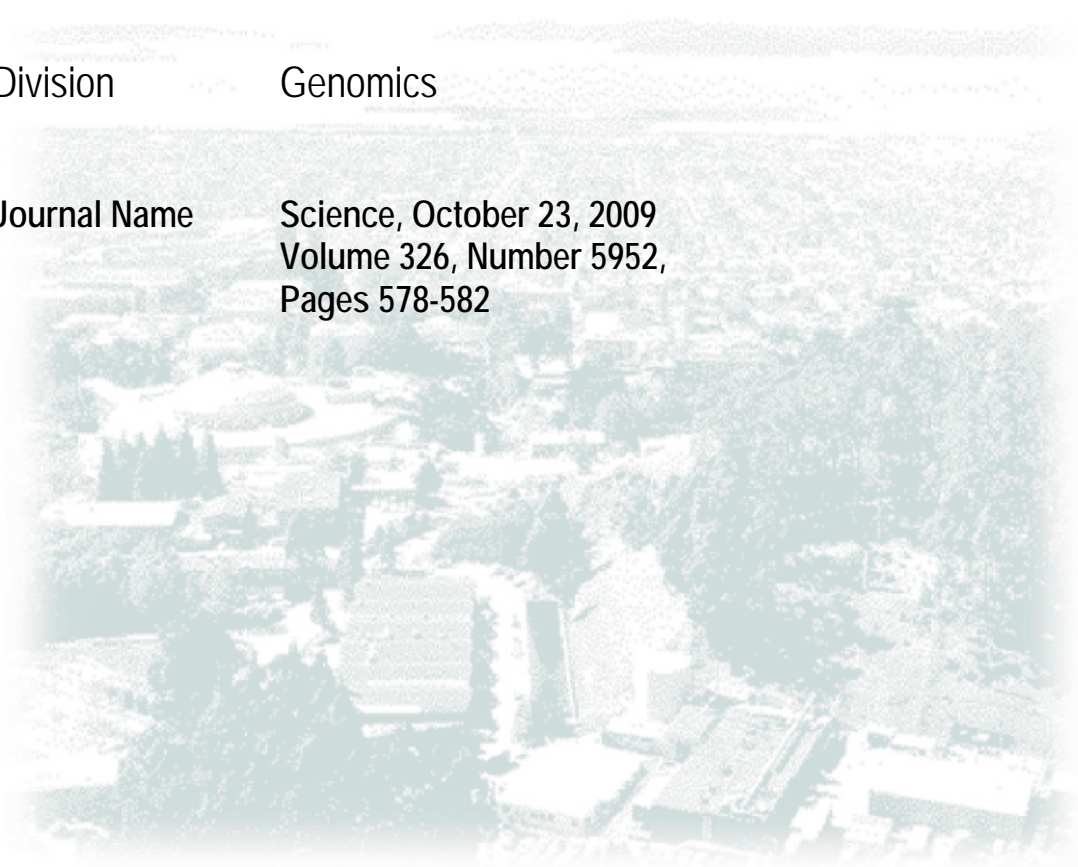
---

Title           Metagenome of a Versatile  
Chemolithoautotroph from  
Expanding Oceanic Dead  
Zones

Author(s)       David A. Walsh, Elena Zaikova,  
Charles L. Howes, Young Song,  
Jody Wright, Susannah G. Tringe,  
Philippe D. Tortell, Steven J. Hallam

Division         Genomics

Journal Name    Science, October 23, 2009  
Volume 326, Number 5952,  
Pages 578-582



1 **Metagenome of a versatile chemolithoautotroph from expanding oceanic dead zones**

2  
3 Summary: Time-resolved metagenomic approaches are used to describe carbon and  
4 energy metabolism of an ecologically relevant microbe from expanding oceanic dead  
5 zones mediating carbon sequestration, sulfur-detoxification and biological nitrogen loss.  
6

7 David A. Walsh<sup>1</sup>, Elena Zaikova<sup>1</sup>, Charles L. Howes<sup>1</sup>, Young Song<sup>1</sup>, Jody Wright<sup>1</sup>,  
8 Susannah G. Tringe<sup>5</sup>, Philippe D. Tortell<sup>2,3</sup>, Steven J. Hallam<sup>1,4\*</sup>  
9

10 <sup>1</sup>Department of Microbiology & Immunology, University of British Columbia

11 <sup>2</sup>Department of Earth & Ocean Sciences, University of British Columbia

12 <sup>3</sup>Department of Botany, University of British Columbia

13 <sup>4</sup>UBC Graduate Program in Bioinformatics

14 <sup>5</sup>Department of Energy Joint Genome Institute  
15

16 \* To whom correspondence should be addressed:

17 University of British Columbia, Department of Microbiology & Immunology, Life  
18 Sciences Institute, 2552-2350 Health Sciences Mall, Vancouver British Columbia  
19 V6T1Z3 Canada

20 Office: (604) 827-3420 FAX: (604) 822-6041 e-mail: shallam@interchange.ubc.ca  
21

22 Running Title:

23 *Versatile energy metabolism in a ubiquitous oxygen minimum zone microbe*

1 **Abstract**

2

3 Oxygen minimum zones (OMZs), also known as oceanic “dead zones”, are widespread  
4 oceanographic features currently expanding due to global warming and coastal  
5 eutrophication. Although inhospitable to metazoan life, OMZs support a thriving but  
6 cryptic microbiota whose combined metabolic activity is intimately connected to nutrient  
7 and trace gas cycling within the global ocean. Here we report time-resolved metagenomic  
8 analyses of a ubiquitous and abundant but uncultivated OMZ microbe (SUP05) closely  
9 related to chemoautotrophic gill symbionts of deep-sea clams and mussels. The SUP05  
10 metagenome harbors a versatile repertoire of genes mediating autotrophic carbon  
11 assimilation, sulfur-oxidation and nitrate respiration responsive to a wide range of water  
12 column redox states. Thus, SUP05 plays integral roles in shaping nutrient and energy  
13 flow within oxygen-deficient oceanic waters via carbon sequestration, sulfide  
14 detoxification and biological nitrogen loss with important implications for marine  
15 productivity and atmospheric greenhouse control.

16

17

18

19

20

21

22

23

24

1 Dissolved oxygen (O<sub>2</sub>) concentration is a critical determinant of marine ecosystem  
2 function and food web structure. Water column O<sub>2</sub> deficit results in habitat compression  
3 and reduced productivity for aerobic respiring organisms with concomitant expansion of  
4 chemolithoautotrophic metabolism (1) manifested by biological nitrogen loss and the  
5 production of climate active trace gases (2, 3). Therefore, OMZ expansion and  
6 intensification (3-5) represents a global ecological phenomenon with potentially  
7 deleterious feedback and forcing effects (1, 6). In order to understand, respond to, or  
8 mitigate these transitions, studies monitoring and modeling dynamics and systems  
9 metabolism of OMZ microbiota in relation to physical and chemical oceanographic  
10 parameters are imperative.

11 Extensive OMZs are found throughout the Eastern North Pacific (ENP), Eastern  
12 South Pacific (ESP), Northern Indian Ocean, and Southwest African shelf waters (*i.e.*  
13 Namibian upwelling system) (3). Although the extent of oxygen-deficiency varies within  
14 and between oceanographic provinces, taxonomic surveys have revealed conserved  
15 patterns of microbial community composition (7-11). In all sites thus far examined, small  
16 subunit ribosomal RNA (SSU rRNA) gene libraries are enriched with sequences  
17 affiliated with chemoautotrophic gill symbionts of deep-sea clams and mussels (EOSA-1  
18 in the ESP and GSO in African shelf waters) (9-11). Phylogenetic analyses indicate that  
19 the GSO/EOSA-1 complex consists of two closely related, co-occurring and uncultivated  
20 lineages, ARCTIC96BD-19 (12) and SUP05 (13) with the latter encompassing the  
21 symbionts (Fig. 1A-B). From an ecophysiological perspective, blooming SUP05  
22 populations have recently been implicated in chemolithotrophic sulfide (H<sub>2</sub>S) oxidation  
23 coupled to nitrate (NO<sub>3</sub><sup>-</sup>) reduction in the Namibian upwelling system (10). Both,

1 ARCTIC96BD-19 and SUP05 populations are also prevalent in non-sulfidic waters of the  
2 ENP and ESP, suggesting alternative or amended modes of energy metabolism. Given the  
3 likely importance of ARCTIC96BD-19 and SUP05 to carbon, nitrogen and sulfur cycling  
4 in marine OMZs, and their largely unexplored metabolic capabilities, a deeper  
5 understanding of both lineages is needed to constrain their respective ecological and  
6 biogeochemical roles.

7         Saanich Inlet, British Columbia is a seasonally anoxic fjord characterized by an  
8 annual cycle of stratification and deep water renewal (14) that is associated with large  
9 water column redox gradients and high rates of trace gas production and consumption  
10 (Fig. S1A-C). Time resolved studies identified pelagic SUP05 as a dynamic and  
11 numerically abundant denizen of the Saanich Inlet water column, representing up to 37%  
12 of total bacteria (Table S2) (11). Closer examination of SUP05 SSU rRNA gene copy  
13 number during onset and progression of seasonal stratification revealed blooming  
14 populations below the oxycline, reaching up to  $4.75 \times 10^5$  copies  $\text{ml}^{-1}$  in regions of  $\text{H}_2\text{S}$  and  
15  $\text{NO}_3^-$  depletion (Fig. S2). Further high-resolution SSU rRNA gene surveys revealed two  
16 SUP05 phylotypes, SI-1 and SI-2 (Fig. 1A and 2A), differing by ~4% nucleotide identity.  
17 While SI-1 dominated suboxic waters throughout the year, SI-2 was less common, and  
18 transiently increased during deep water renewal events, alluding to the presence of  
19 ecologically differentiated populations (Fig. 2A). Given these observations we reasoned  
20 that Saanich Inlet would provide a natural enrichment amenable to environmental  
21 genomic (*i.e.* metagenomic) assembly and metabolic pathway reconstruction of pelagic  
22 SUP05 populations. To this end, we analyzed sixteen bi-directionally end sequenced

1 fosmid libraries constructed from environmental DNA samples spanning oxic to anoxic  
2 waters over the seasonal stratification and deep water renewal cycle (Fig. S1D).

3 Fosmid end sequences were initially screened for putative SUP05 genotypes by  
4 fragment recruitment to closely related symbiont reference genomes (*15, 16*) revealing  
5 extensive coverage within oxygen-deficient sampling intervals (Fig. S3). To reconstruct  
6 and identify SUP05-specific scaffolds, paired-end sequences were assembled and binned  
7 based on shared sequence similarity to symbiont reference genomes and analysis of  
8 intrinsic oligonucleotide composition patterns (Fig. S4A-B). Nineteen scaffolds  
9 encompassing 1.16 million base pairs of SUP05 DNA, herein referred to as the SUP05  
10 metagenome, were identified and taxonomically verified (Table S3). See Table S1 for  
11 additional breakdown of important assembly features. Although derived from a  
12 heterogeneous population, average polymorphism within the SUP05 metagenome was  
13 0.4%, indicating assembly of closely related sympatric donor genotypes. Consistent with  
14 this observation, a single SSU-LSU rRNA operon affiliated with the SI-I phylotype was  
15 identified (Fig. S4). Closer examination of individual fosmid library contributions to the  
16 assembly revealed a majority of paired-end sequences within SUP05 scaffolds were  
17 derived from samples exhibiting elevated SI-1 phylotype abundance, further supporting  
18 scaffold assignments (Fig. 2B). Overall coverage efficiency was assessed based on  
19 recovery of 83 of 93 information processing genes, including 31 of 32 canonically  
20 conserved single copy genes (*17*) (Table S3) suggesting near complete recovery of the  
21 core genome, in turn facilitating downstream efforts to reconstruct and interpret  
22 metabolic pathways.

1           The comparative architecture of SUP05 and symbiont reference genomes revealed  
2 patterns of genetic relatedness consistent with derivation from a common free-living  
3 ancestor. At the same time, differences in gene content and organization provided insight  
4 into the process of genome divergence and niche partitioning. Approximately 20% of  
5 predicted gene content shared between symbiont reference genomes was absent from the  
6 SUP05 metagenome. Genes from this set typically clustered together (Fig. 3), and could  
7 represent genomic features mediating symbiont-host interactions or adaptations to  
8 hydrothermal vent or cold seep settings (18). Of the 861 genes shared between symbiont  
9 reference genomes, 80% were also conserved in the SUP05 metagenome (Fig. 3 and S5-  
10 6), indicating significant metabolic overlap. Many of these genes are predicted to mediate  
11 informational processing steps, particularly translation, although a significant fraction  
12 function in carbon, sulfur, amino acid, and coenzyme metabolism (Fig. S6).  
13 Approximately 35% of the gene content predicted in the SUP05 metagenome was unique,  
14 reflecting characteristic differences between pelagic and symbiotic modes of existence  
15 (Fig. 3 and S5-6). These included genes implicated in DNA uptake and repair,  
16 denitrification and adaptive or stress responses (See supporting online material for more  
17 details).

18           Given the potential importance of pelagic SUP05 populations on the ecology and  
19 biogeochemistry of oxygen-deficient oceanic waters, examination of carbon and energy  
20 metabolism within the SUP05 metagenome is of particular interest. Similar to symbiotic  
21 counterparts, the SUP05 metagenome harbors genes mediating the Calvin-Benson-  
22 Bassham (CBB) cycle for autotrophic carbon assimilation, including a single form II  
23 ribulose 1,5-bisphosphate carboxylase-oxygenase (RubisCO) gene, implicating SUP05 in

1 chemosynthetic carbon fixation within OMZs (19). In addition, a gene encoding  $\beta$ -class  
2 carbonic anhydrase, encoding a potential CO<sub>2</sub> concentrating mechanism, was also  
3 identified. A complete repertoire of genes mediating the conversion of fixed carbon to  
4 hexose and ribose sugars via gluconeogenesis and the non-oxidative branch of the  
5 pentose phosphate pathway was identified along with the majority of tricarboxylic acid  
6 (TCA) cycle components (Fig. S6). However, genes mediating the interconversion of  
7 succinyl-CoA and 2-oxoglutarate were not recovered (Fig. S6), indicating the potential  
8 for obligate autotrophy as posited earlier for clam symbionts (15, 16).

9 From the standpoint of energy metabolism, the SUP05 metagenome harbors a  
10 diverse repertoire of genes mediating chemolithotrophic oxidation of reduced sulfur  
11 compounds. Genes encoding flavocytochrome *c*/sulfide dehydrogenase (*fccAB*) unique to  
12 the SUP05 metagenome, and sulfide quinone oxidoreductase (*sqr*) conserved between  
13 pelagic SUP05 and symbiont reference genomes mediating the oxidation of H<sub>2</sub>S to  
14 elemental sulfur (S<sup>0</sup>) were identified (Fig. S6). The presence of two enzymatic complexes  
15 may facilitate sulfur-based energy metabolism under variable sulfide regimes. For  
16 instance, FccAB is thought to be functionally significant at low sulfide concentrations  
17 (20). In addition, sirohaem dissimilatory sulfite reductase subunits (*dsrAB*), APS  
18 reductase (*apr*), ATP sulfurylase (*sat*) mediating the complete oxidation of S<sup>0</sup> to sulfate,  
19 and the Sox pathway (*soxABXYZ*) for thiosulfate (S<sub>2</sub>O<sub>3</sub><sup>2-</sup>) oxidation (Fig. S6) (21) were  
20 also conserved between pelagic SUP05 and symbiont reference genomes. The capacity to  
21 obtain electrons from S<sub>2</sub>O<sub>3</sub><sup>2-</sup> may be of considerable ecological relevance given that  
22 mixing of sulfidic and oxygenated water masses results in S<sub>2</sub>O<sub>3</sub><sup>2-</sup> accumulation due to the  
23 chemical oxidation of H<sub>2</sub>S (22, 23). Moreover, the apparent absence of *soxCD* sulfur



1 dehydrogenase genes suggests the capacity to store  $S^0$ , which can be subsequently  
2 oxidized via the reverse DSR pathway thereby provisioning SUP05 in the absence of  
3 ambient reductant (24). Indeed, *soxCD* homologues are also absent from symbiont  
4 reference genomes, and sulfur globule formation has been associated with a subset of  
5 clam symbionts (25).

6 Although symbiont reference genomes harbor multiple aerobic respiratory  
7 complexes (15, 16), none were recovered in the SUP05 metagenome consistent with a  
8 facultative or strictly anaerobic lifestyle. Indeed, all the enzymatic machinery needed to  
9 reduce  $NO_3^-$  to the greenhouse gas nitrous oxide ( $N_2O$ ) including membrane-bound  
10 (*narKK<sub>2</sub>GHJI*) and periplasmic (*napFBAHGD*) dissimilatory nitrate reductases  
11 potentially operating under high and low  $NO_3^-$  conditions, respectively (26, 27), copper-  
12 containing nitrite reductase (*nirK*), and  $N_2O$  forming nitric oxide reductase (*norCB*) (Fig.  
13 4 and S6) were identified, mechanistically implicating pelagic SUP05 in biological  
14 nitrogen loss from oxygen-deficient oceanic waters. Moreover, the genomic co-  
15 localization of sulfur oxidation and denitrification genes suggests a highly integrated and  
16 redox-sensitive energy metabolism (Fig. 4). For example, the genomic proximity of  
17 Crp/Fnr transcriptional regulators with *fcc* and *nap* gene clusters (Fig. 4) may indicate  
18 coordinated gene expression in response to changing redox status (20, 26, 28). In both the  
19 Namibian upwelling system and Saanich Inlet, blooming SUP05 populations occur in  
20 regions of  $H_2S$  and  $NO_3^-$  depletion (Fig. S2) where coexpression of *nap* and *fcc* genes  
21 clusters may become critical for survival as energetic substrates become limiting.

22 Curiously, we identified more than ten putative toxin-antitoxin (TA) modules  
23 unique to the SUP05 metagenome, indicating a highly regulated stress response (Table

1 S4). TA modules consist of a stable toxin and a labile antitoxin and are commonly  
2 associated with environmental bacteria where they control induction of reversible cellular  
3 stasis (29). Of specific interest is a TA module of the RelE superfamily identified within  
4 an operon encoding molybdopterin-guanine dinucleotide synthase (*mobA*) (Fig. 4). The  
5 product of MobA, molybdopterin-guanine dinucleotide (MGD), is an essential cofactor  
6 for all described classes of nitrate reductase (30) and therefore *mobA* expression is  
7 integral to denitrification in pelagic SUP05 populations. Severe  $\text{NO}_3^-$  limitation could  
8 limit *mobA* expression leading to co-repression of the embedded TA module. This would  
9 result in activation of the RelE toxin, through degradation of the labile antitoxin, and  
10 induction of cellular stasis (31). In this regard, the integration of a TA system into a  
11 denitrification regulon may allow SUP05 to persist during periods of extreme  $\text{NO}_3^-$   
12 limitation, analogous to other forms of nutritional stress response (*e.g.* amino acid  
13 starvation in *E. coli*).

14 As the number of studies surveying OMZ community structure increases, the  
15 ubiquity and abundance of the SUP05 lineage becomes ever more apparent. Analysis of  
16 the SUP05 metagenome, and the water column disposition of pelagic SUP05 with respect  
17 to  $\text{H}_2\text{S}$  and  $\text{NO}_3^-$  gradients, resolves a chemolithoautotrophic metabolism based on  
18 oxidation of reduced sulfur compounds with  $\text{NO}_3^-$  through multiple and highly regulated  
19 bioenergetic routes. Paradoxically, as “dark” primary producers, blooming SUP05  
20 populations have the potential to sequester large amounts of  $\text{CO}_2$  while simultaneously  
21 producing  $\text{N}_2\text{O}$  via  $\text{NO}_3^-$  respiration. Therefore the SUP05 metagenome provides a  
22 functional template for analysis of gene expression in relation to climatologically relevant  
23 biogeochemical transformations within oxygen-deficient oceanic waters. We anticipate

1 that this information will become an essential resource for comparative analysis of  
2 ecotype diversification within the SUP05 and ARCTIC96BD-19 lineages aiding in the  
3 development of monitoring tools to assess and model microbial community responses to  
4 OMZ expansion and intensification.

5

## 6 **References**

7

- 8 1. R. J. Diaz, R. Rosenberg, *Science* **321**, 926 (2008).
- 9 2. K. R. Arrigo, *Nature* **437**, 349 (2005).
- 10 3. A. Paulmier, D. Ruiz-Pino, *Progress in Oceanography* **In press**, (2008).
- 11 4. F. A. Whitney, H. J. Freeland, M. Robert, *Prog. in Oceanog.* **75**, 179 (2007).
- 12 5. L. Stramma, G. C. Johnson, J. Sprintall, V. Mohrholz, *Science* **320**, 655 (2008).
- 13 6. P. G. Brewer, E. T. Peltzer, *Science* **324**, 347 (2009).
- 14 7. B. M. Fuchs, D. Woebken, M. V. Zubkov, P. Burkill, R. Amann, *Aquat Microb*  
15 *Ecol* **39**, 145 (2005).
- 16 8. D. Woebken, B. A. Fuchs, M. A. A. Kuypers, R. Amann, *Applied and*  
17 *Environmental Microbiology* **73**, 4648 (2007).
- 18 9. H. Stevens, O. Ulloa, *Environ Microbiol* **10**, 1244 (2008).
- 19 10. G. Lavik *et al.*, *Nature* **457**, 581 (2009).
- 20 11. E. Zaikova *et al.*, *Environ Microbiol* (submitted).
- 21 12. N. Bano, J. T. Hollibaugh, *Appl Environ Microbiol* **68**, 505 (2002).
- 22 13. M. Sunamura, Y. Higashi, C. Miyako, J. Ishibashi, A. Maruyama, *Appl Environ*  
23 *Microbiol* **70**, 1190 (2004).

- 1 14. J. J. Anderson, A. H. Devol, *Estuarine and Coastal Marine Science* **1**, 1 (1973).
- 2 15. H. Kuwahara *et al.*, *Curr Biol* **17**, 881 (2007).
- 3 16. I. L. G. Newton *et al.*, *Science* **315**, 998 (2007).
- 4 17. F. D. Ciccarelli *et al.*, *Science* **311**, 1283 (2006).
- 5 18. I. L. Newton, P. R. Girguis, C. M. Cavanaugh, *BMC Genomics* **9**, 585 (2008).
- 6 19. G. Jost, M. V. Zubkov, E. Yakushev, M. Labrenz, K. Jurgens, *Limnology and*  
7 *Oceanography* **53**, 14 (2008).
- 8 20. M. Mussmann *et al.*, *PLoS Biol* **5**, e230 (2007).
- 9 21. C. G. Friedrich, F. Bardischewsky, D. Rother, A. Quentmeier, J. Fischer, *Curr*  
10 *Opin Microbiol* **8**, 253 (2005).
- 11 22. J. Zopfi, T. G. Ferdelman, B. B. Jorgensen, A. Teske, B. Thamdrup, *Marine*  
12 *Chemistry* **74**, 29 (2001).
- 13 23. R. B. Cardoso *et al.*, *Biotechnology and Bioengineering* **95**, 1148 (2006).
- 14 24. C. Dahl *et al.*, *J Bacteriol* **187**, 1392 (2005).
- 15 25. R. D. Vetter, *Marine Biology* **88**, 33 (1985).
- 16 26. V. Stewart, Y. Lu, A. J. Darwin, *J Bacteriol* **184**, 1314 (2002).
- 17 27. H. Wang, C. P. Tseng, R. P. Gunsalus, *J Bacteriol* **181**, 5303 (1999).
- 18 28. S. Spiro, J. R. Guest, *FEMS Microbiol Rev* **6**, 399 (1990).
- 19 29. D. P. Pandey, K. Gerdes, *Nucleic Acids Res* **33**, 966 (2005).
- 20 30. D. J. Richardson, B. C. Berks, D. A. Russell, S. Spiro, C. J. Taylor, *Cell Mol Life*  
21 *Sci* **58**, 165 (2001).
- 22 31. S. K. Christensen, M. Mikkelsen, K. Pedersen, K. Gerdes, *Proc Natl Acad Sci U S*  
23 *A* **98**, 14328 (2001).

32. This work was performed under the auspices of the U.S. Department of Energy's Office of Science, Biological, and Environmental Research Program and by the University of California, Lawrence Berkeley National Laboratory, Lawrence Livermore National Laboratory, under contract no. DE-AC02-05CH11231, Lawrence Livermore National Laboratory under contract no. DE-AC52-07NA27344, and Los Alamos National Laboratory under contract no. DE-AC02-06NA25396. This work was also supported by grants from the Natural Sciences and Engineering Research Council (NSERC) of Canada 328256-07 and STPSC 356988, Canada Foundation for Innovation (CFI) 17444; Canadian Institute for Advanced Research (CIFAR), and the Center for Bioinorganic Chemistry (CEBIC). D.A.W. was supported by NSERC, Killam Trust, and the Tula Foundation-funded Centre for Microbial Diversity and Evolution (CMDE). We thank M. Robert (Institute of Ocean Sciences, Sidney, BC, Canada), C. Payne, L. Pakhomova, and J. Granger (UBC) for help in sampling and chemical analyses and the captains and crews of the *CCGS John P. Tulley* and *HMS John Strickland* for logistical support. We thank the Joint Genome Institute, including K. Barry, S. Pitluck, and E. Kirton, for technical assistance and A. Page, K. Mitchell, and S. Lee in the Hallam laboratory for reading the manuscript. This metagenome project has been deposited at the DNA DataBank of Japan and European Molecular Biology Laboratory, and GenBank, under the project accession ACSG00000000. The version described in this paper is the first version, ACSG01000000. SSU rRNA gene sequences were deposited at GenBank under the accession numbers GQ345343-GQ351265, and fosmid sequences were deposited under the accession numbers GQ351266 to GQ351269 and GQ369726.

1

2 **Figure Titles**

3

4 **Fig. 1** (A) Phylogenetic tree of ARCTIC96BD-19 and SUP05 lineages based on  
5 comparative SSU rRNA gene analysis. The tree was inferred using maximum likelihood  
6 implemented in PHYML (B) Relative abundance of ARCTIC96BD-19 and SUP05 SSU  
7 rRNA sequences recovered from Saanich Inlet (SI), eastern North Pacific (ENP) (this  
8 study), eastern South Pacific (ESP) (9), and southwest African shelf waters (Namibia)  
9 (10).

10

11 **Fig. 2** (A) Phylotype abundance of SUP05 SI-1 (black circles) and SI-2 (white circles)  
12 based on recovery of SSU rDNA sequences in PCR generated clone libraries (B) mean  
13 depth of coverage of the SUP05 metagenome plotted over the Saanich Inlet nitrate profile  
14 during the 2006-2007 season. Sample depths and dates are noted on the axes.

15

16 **Fig. 3** Gene content comparison between SUP05 metagenome and symbiont reference  
17 genomes. Nested circles from outermost to innermost represent the following. (i and ii)  
18 COG functional predictions on the forward and reverse strands of the *R. magnifica*  
19 reference genome. (Fig. S7 for color designation), (iii) Conservation of gene content (iv)  
20 Genes conserved in symbionts, but absent from the SUP05 metagenome. Inset Venn  
21 diagram depicts predicted gene distribution among SUP05 metagenome and symbiont  
22 reference genomes. Values correspond to the number of shared genes among overlapping  
23 genomes, using each genome as the original query. The dotted line represents the open-

1 genome configuration of the SUP05 metagenome. \*The discrepancy in core size when  
2 SUP05 metagenome is employed as query (774) compared to symbionts (~683) reflects  
3 gene content redundancy in the metagenome assembly.

4

5 **Fig. 4** Alignment of an ungapped region of a SUP05 metagenomic scaffold, encoding  
6 genes involved in nitrate and sulfur metabolism, with the corresponding genomic regions  
7 of symbiont reference genomes. The height of red bars corresponds to nucleotide  
8 similarity over conserved genomic regions. Proper scaffold assembly across this region  
9 was verified by full-length sequencing of two overlapping fosmids.

10

11

12

13

14

15

16

17

18

19

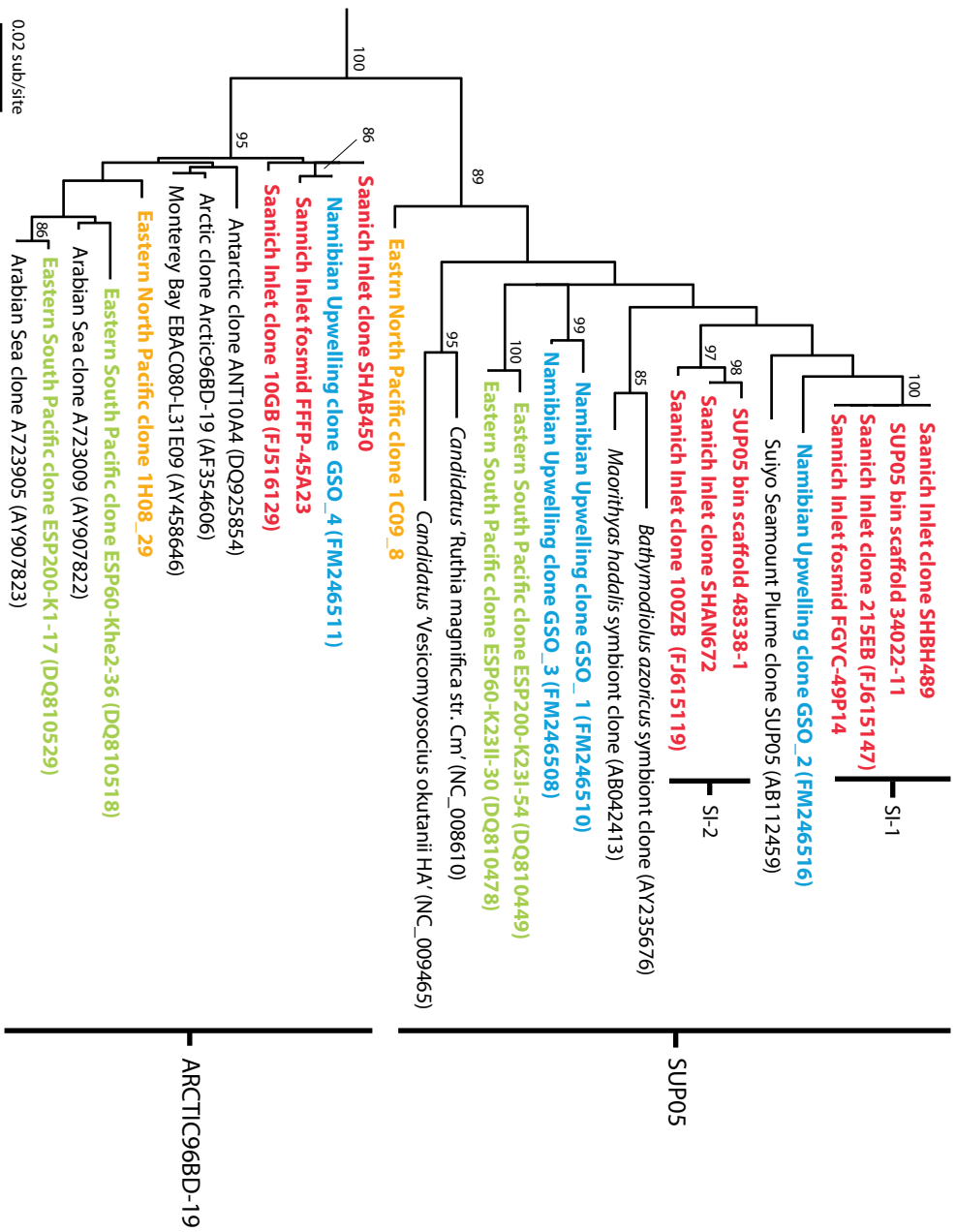
20

21

22

23

A



B

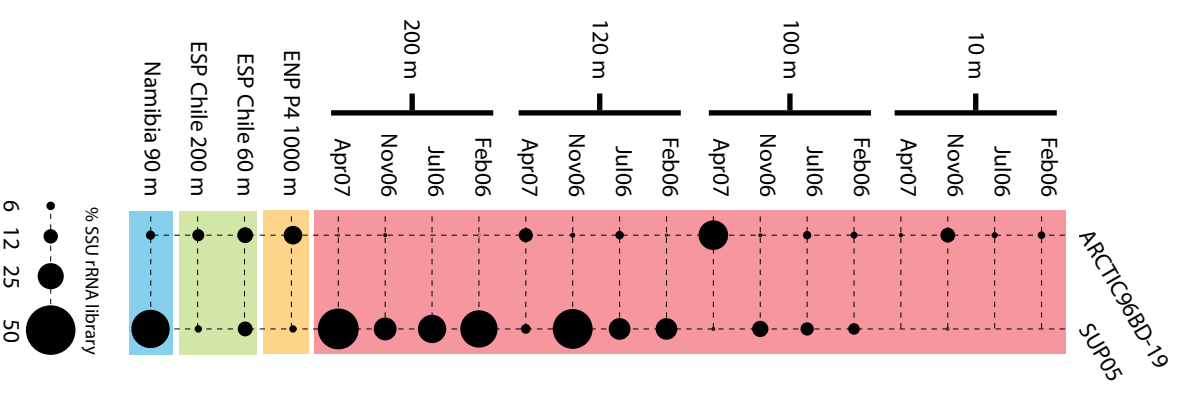


Fig. 1



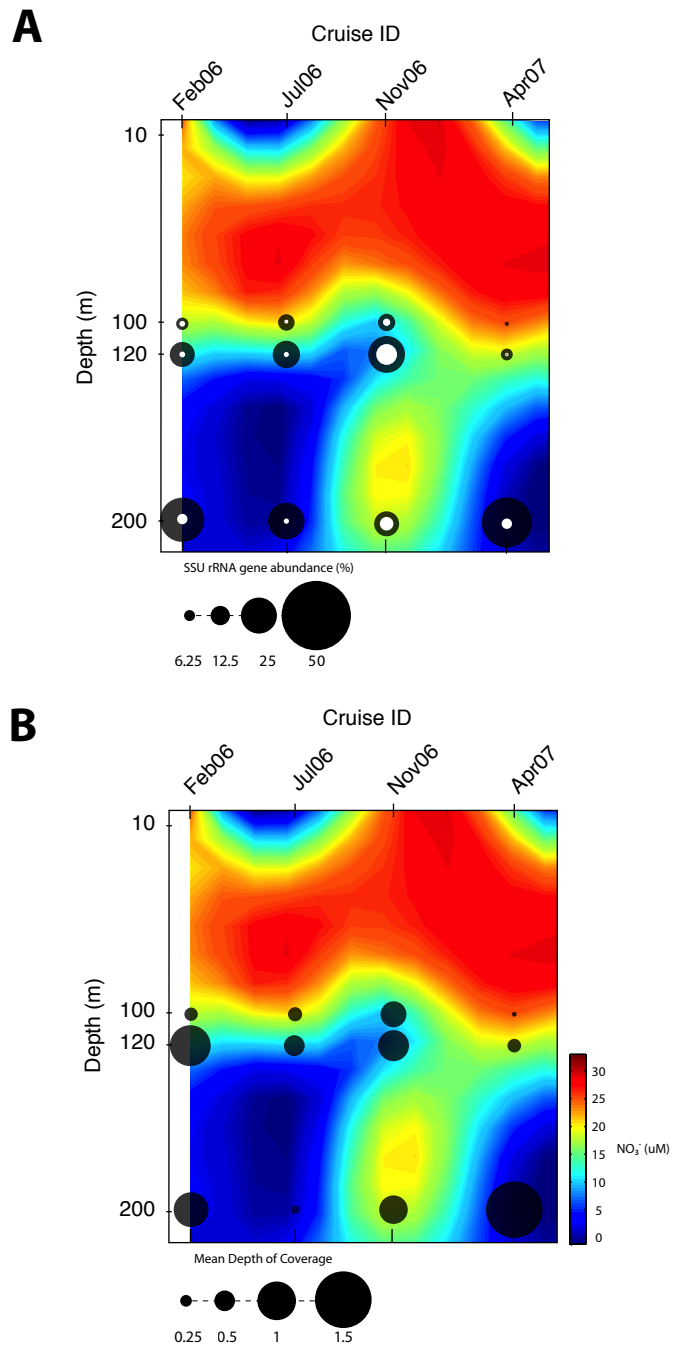


Fig. 2

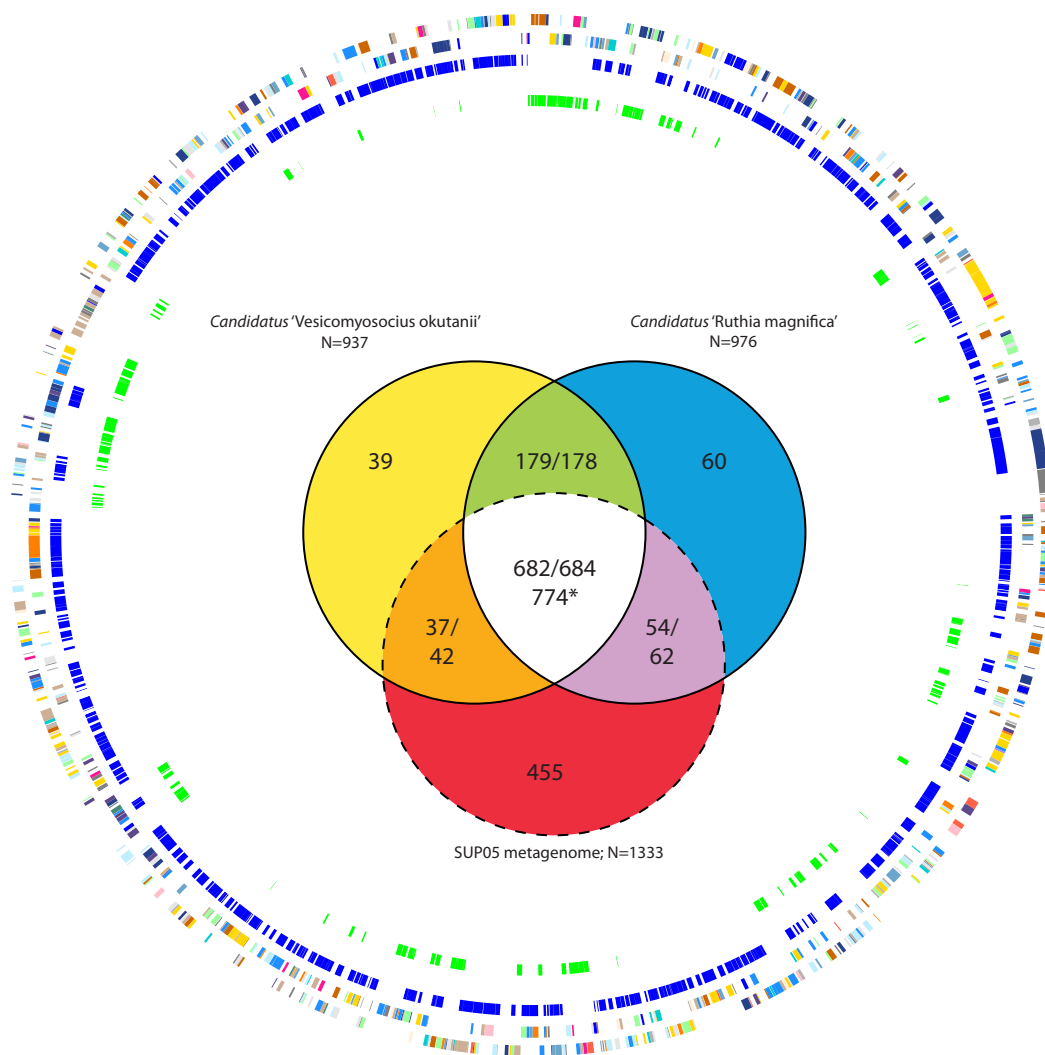


Fig. 3

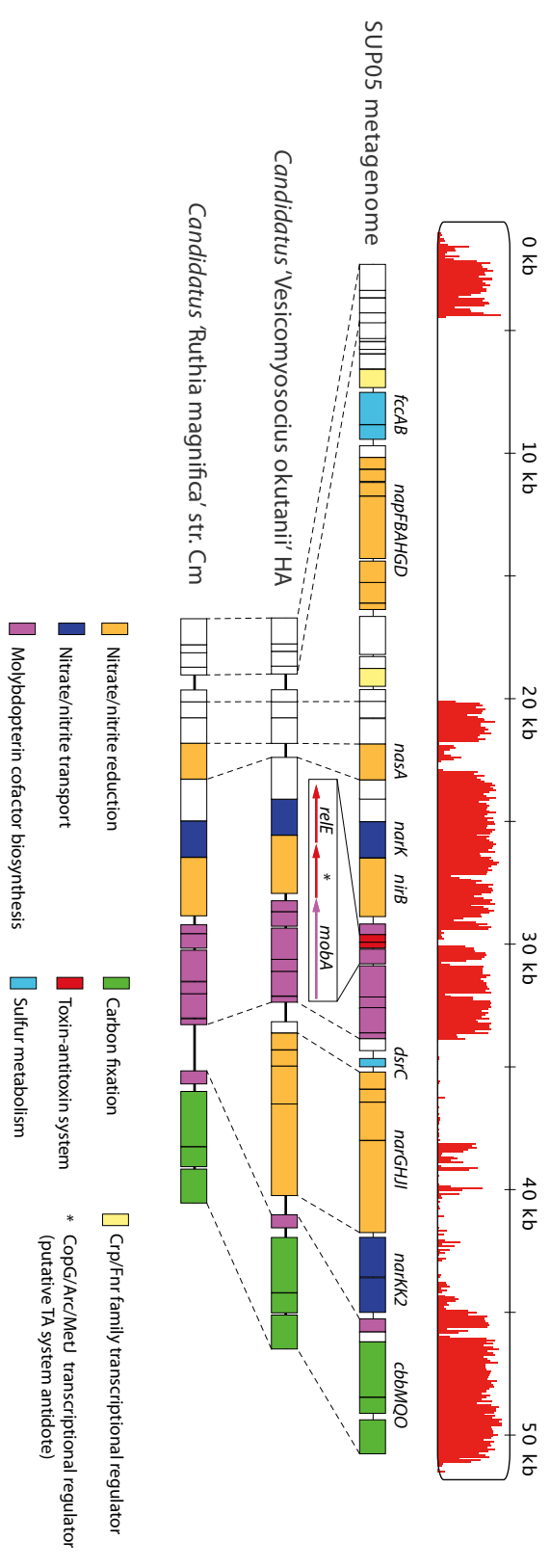


Fig. 4

1 **Supporting Online Material for**

2  
3 **Metagenome of a versatile chemolithoautotroph from expanding oceanic dead zones**

4 David A. Walsh<sup>1</sup>, Elena Zaikova<sup>1</sup>, Charles L. Howes<sup>1</sup>, Young Song<sup>1</sup>, Jody Wright<sup>1</sup>, Susannah G.  
5 Tringe<sup>5</sup>, Philippe D. Tortell<sup>2,3</sup>, Steven J. Hallam<sup>1</sup>,  
6

7 **The Supporting Online Material includes:**

8 Materials and Methods  
9 Supplementary Text  
10 Supplementary Tables S1 – S4  
11 Supplementary Figures S1 – S7  
12

13 **Materials and Methods**

14 **Sample collection and processing.** Sample collection, O<sub>2</sub> and NO<sub>3</sub><sup>-</sup> measurements, microbial  
15 biomass concentration, environmental DNA extraction, and quantitative PCR analysis of bacteria  
16 and SUP05 SSU rRNA gene sequences were performed as previously described (1). Samples for  
17 H<sub>2</sub>S measurement were fixed with 2% (final concentration) zinc acetate and analyzed by  
18 spectrophotometer as described in (2).  
19

20 **Saanich Inlet SSU rRNA gene library construction and analysis.** For each of sixteen Saanich  
21 Inlet environmental DNA samples (see Fig S1), clone libraries were constructed and sequenced  
22 from pooled PCR products amplified using the bacterial-specific forward primer 27F (5'-  
23 AGAGTTTGATCCTGGCTCAG-3') and reverse primer 1391R (5'-  
24 GACGGGCGGTGWGTRCA-3') as described at <http://jgi.doe.gov>. Chimeric sequences,  
25 identified using Bellerophon (3) and Mallard (4), were excluded from the sequence dataset. A

1 total of 5,753 partial bacterial SSU rRNA clones passed the quality and chimera check. SSU  
2 rRNA gene sequences were aligned using the SILVA aligner (5), imported into an ARB database  
3 (6) with the same alignment, and 1,067 SUP05 sequences were identified based on their position  
4 in the SILVA reference tree. The 1,067 SUP05 sequences were extracted from ARB, realigned  
5 with MUSCLE (7) and two sequence clusters (i.e. the SI-1 and SI-2 phylotypes) were defined  
6 after assignment of sequences to operational taxonomic units (OTUs) by applying DOTUR (8)  
7 analysis to an uncorrected distance matrix generated with dnadist (9). The frequency distribution  
8 of SI-1 and SI-2 phylotypes across libraries was calculated by dividing phylotype abundance by  
9 the total number of library clones. ARCTIC96-BD-19 clones were identified in a similar manner.  
10 The phylogenetic tree (Figure 1a) was inferred from a manually refined MUSCLE alignment of  
11 partial sequences by PHYML, using an HKY + 4 $\Gamma$  + I model of evolution and estimated values  
12 for the  $\alpha$  parameter of the  $\Gamma$  distribution, the proportion of invariable sites, and the  
13 transition/transversion (10). The confidence of each node was determined by assembling a  
14 consensus tree of 100 bootstrap replicates.

15

16 **Eastern North Pacific SSU rRNA gene analysis.** An environmental DNA sample extracted  
17 from microbial biomass collected from within the OMZ (1000 m) at Station P4 (48° 39.0' N,  
18 126° 40.0' W) of the Line P oceanographic time series ([http://www.pac.dfo-](http://www.pac.dfo-mpo.gc.ca/sci/osap/projects/linepdata/default_e.htm)  
19 [mpo.gc.ca/sci/osap/projects/linepdata/default\\_e.htm](http://www.pac.dfo-mpo.gc.ca/sci/osap/projects/linepdata/default_e.htm)) was used to construct a bacterial SSU  
20 rRNA gene clone library that was then sequenced as previously described (1). In total, 170  
21 nonchimeric sequences were generated and analyzed as describe above.

22

1 **Fosmid library construction and sequencing.** Sixteen fosmid libraries (7,680 clones/libraries)  
2 were constructed from DNA samples collected from four depths during four cruises over the  
3 2006-07 seasonal stratification cycle (see Fig. S1). Prior to cloning, ~5 µg of environmental  
4 DNA was further purified on a CsCl density gradient as previously described (11). Fosmid  
5 libraries were prepared using the CopyControl Fosmid Library Production Kit (Epicentre).  
6 Briefly, ~1 µg of CsCl-purified DNA was blunt end repaired and separated on a 1% low melt  
7 agarose pulse-field gel O/N at 6 V/cm. The 40-50 kb fragment range was excised and gel  
8 purified using agarase, followed by concentration using an Amicon Ultracel 10K filter device.  
9 DNA was ligated into the pCC1fos vector, packaged using the MaxPlax lambda packaging  
10 extract, and used to transfect TransforMax EPI300 *E. coli* cells. Transfected cells were plated on  
11 selective agar and fosmid clones picked using the Q-Pix robotic colony picker and grown in  
12 selective media for DNA sequencing. Bidirectional end sequencing of fosmids was performed  
13 with standard M13 -28 or -40 primers and the BigDye sequencing kit (Applied Biosystems). The  
14 reactions were purified by a magnetic bead protocol and run on an ABI PRISM 3730 (Applied  
15 Biosystems) capillary DNA sequencer (for research protocols, see <http://jgi.doe.gov>).

16  
17 **SUP05-focused assembly of fosmid end sequence data.** All 243,264 fosmid end sequences (see  
18 Table S1) were initially assembled with Phrap using parameters (minmatch 30 maxmatch 55  
19 minscore 55 max\_subclone\_size 50000 revise\_greedy vector\_bound 20) reported in (12). Phrap-  
20 generated contigs were linked into “ambiguous” scaffolds (*i.e.* allowed contig overlap) with  
21 Bambus (13), using the default settings, based on fosmid paired-end information. The initial  
22 assembly was comprised of 31,766 scaffolds with a total length (excluding gaps) and span of  
23 35.6 Mb and 85.2 Mb respectively, and included 192 scaffolds  $\geq$  5 kb in length. Of the 192

1 scaffolds, those that exhibited highest similarity to the genomes of either of two symbiont  
2 reference species, *Candidatus* ‘Ruthia magnifica’ or *Candidatus* ‘Vesicomysocius okutanii’,  
3 upon a BLASTN search of the NCBI refseq database were flagged as potential SUP05 scaffolds  
4 and targeted for further assembly as follows. Overlapping contigs within these ambiguous  
5 scaffolds were identified in Sequencher using a 90% identity cut-off and 100 bp minimum  
6 overlap. Those contigs that cross-assembled in the same order and orientation described by  
7 Bambus were then reassembled using the miniassembly tool in consed (14) with the following  
8 parameters: minscore 55 minmatch 30 forcelevel 10. Addition of the forcelevel parameter  
9 decreases assembly stringency and permits assembly between polymorphic reads, leading to the  
10 collapse of short non-polymorphic contigs into longer polymorphic scaffold. These next  
11 generation contigs were then re-ordered and oriented into ambiguous scaffolds with Bambus and  
12 the whole process was repeated until contig extension no longer occurred. Upon final scaffolding  
13 of contigs, the untangle script included with Bambus was used to break ambiguous scaffold into  
14 single linear scaffolds (*i.e.* disallowed contig overlap). In addition, the stringency of the final  
15 assembly was increased by requiring a minimum of four read pairs to link contig into scaffolds.  
16 The final assembly was comprised of 33,630 scaffolds with a total length (excluding gaps) and  
17 span of 34.9 Mb and 38.8 Mb respectively, and included 86 scaffolds  $\geq$  5 kb in length (Figure  
18 S4a). Within the assembly, 5,703 polymorphic sites, defined by  $\geq$  2 reads exhibiting high quality  
19 discrepancies with the scaffold consensus sequence, were identified.

20

21 **Identification of putative SUP05 scaffolds.** Predicted protein-encoding genes from all scaffolds  
22  $\geq$  5 kb were searched against the two clam symbiont genomes using BLASTP. A subset of  
23 scaffold enriched in SUP05 was identified using the following criterion:  $\geq$  3 syntenic open

1 reading frames (ORFs) of high sequence similarity (blast score ratio  $\geq 0.6$ ) with either symbiont  
2 reference genome. The initial bin was comprised of 38 scaffolds (total length and span of 1.36  
3 Mb and 2.54 Mb, respectively), contained two rRNA operons, and the corresponding SSU rRNA  
4 genes were SI-1 and SI-2 phylotypes, respectively (Figure 1a). Principal component analysis  
5 (PCA) of intrinsic DNA signatures was then used to explore the oligonucleotide frequencies of  
6 the 38 scaffolds. The frequencies of all tri-, tetra-, and penta-nucleotides were summed for each  
7 scaffold and their over- and under- representations were evaluated against their expected  
8 frequencies using a maximal-order Markov model as previously described (15). In addition to  
9 scaffolds, four full length overlapping fosmid (~40 kb in length each) known to originate from  
10 SUP05 based on an SI-1 SSU rRNA gene sequence were included in order to identify the  
11 expected oligonucleotide patterns of SUP05 sequences. The resulting Z scores (15) were  
12 normalized by length, imported into the statistical package PC-ORD 5.10 (16) and subjected to  
13 PCA using a Correlation Cross-products matrix. Visualization of the first two axes of the PCA  
14 analysis revealed many scaffolds that were uniquely positioned in the ordination space, but also  
15 revealed a single cluster consisting of 19 scaffolds that exhibited very similar nucleotide pattern  
16 compositions (Figure S4b, red circles). We interpret this subset of related scaffolds as arising  
17 from the same closely related population of SUP05 cells in Saanich. In support of this  
18 interpretation, the SI-I phylotype was observed in this scaffold subset (Figure 1a). Nucleotide  
19 composition patterns for the four fully sequenced and overlapping fosmids (spanning ~140 Kb)  
20 affiliated with SI-1 on the basis of SSU rRNA gene linkage were indistinguishable from these  
21 scaffolds further reinforcing the accuracy and resolution of our binning method (Figure S4b,  
22 green circles). The remaining scaffolds (Figure S4b, blue circles), may represent less abundant  
23 SUP05 genotypes, as evident from the presence of an SI-2 SSU rRNA gene sequence, or



1 unrelated bacterial taxa sharing genomic similarity with clam symbionts. As their taxonomic  
2 identity is questionable, they were removed prior to further analysis and the remaining scaffolds  
3 are herein referred to as the SUP05 metagenome (Table 1). An automated phylogenetic approach  
4 using Phylogenie (17) was used on a set of 76 typically conserved genes, to assess allelic  
5 variation and to verify the absence of non-SUP05 scaffolds in the SUP05 metagenome (see  
6 Table S2).

7  
8 **SUP05 metagenome annotation and comparative genome analysis.** Scaffolds were annotated  
9 using the FGENESB gene calling pipeline from Softberry ([www.softberry.com/berry.phtml](http://www.softberry.com/berry.phtml),  
10 Mount Kisco, NY). Ribosomal RNA and tRNA genes were identified using BLASTN and  
11 tRNA-Scan (18). Functional annotation and classification of the predicted proteome was  
12 performed by using BLASTP homology searches against COG, KEGG, and NCBI nr public  
13 database, and domain analysis with InterProScan (<http://www.ebi.ac.uk/Tools/InterProScan>)  
14 (19). Metabolic pathways were constructed based on KEGG and MetaCyc (20). Comparison of  
15 gene content between the SUP05 metagenome and the reference symbiont genomes was  
16 performed by BLAST Score Ratio (BSR) analysis with a BSR cut off = 0.4 (21). Mapping of  
17 shared gene content and COG categories onto the *R. magnifica* reference genome was performed  
18 with GenomeViz ([www.uniklinikum-giessen.de/genome](http://www.uniklinikum-giessen.de/genome))(22). Multi genome alignment of  
19 SUP05 scaffolds and symbiont genomes was performed with Mauve  
20 (<http://asap.ahabs.wisc.edu/mauve/>) (23).

21

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23

## Supplementary Text

### **Quantification of total bacterial and SUP05 SSU rRNA gene sequences during the 2008**

### **seasonal stratification cycle in Saanich Inlet.**

To assess SUP05 distribution with respect to vertical profiles of oxygen (O<sub>2</sub>), nitrate (NO<sub>3</sub><sup>-</sup>), and sulfide (H<sub>2</sub>S), we quantified SUP05 SSU rRNA gene sequences at three time intervals representing the main stages of the 2008 seasonal stratification cycle (Fig. S2). On April 9<sup>th</sup>, an extensive suboxic zone, defined by the absence of detectable O<sub>2</sub> or H<sub>2</sub>S, was observed between ~135 and 185 m, NO<sub>3</sub><sup>-</sup> was depleted below 185 m and H<sub>2</sub>S was detected at 200 m. SUP05 was detected at  $\geq 1.4 \times 10^5$  copies ml<sup>-1</sup> below 150 m and peaked sharply at the base of the suboxic zone, coinciding with the point of NO<sub>3</sub><sup>-</sup> depletion. By August 11<sup>th</sup>, the suboxic zone had narrowed to within ~150 and 165 m and sulfidic water was detected at 185 m. SUP05 increased in abundance from a depth of 100 to 150 m and then remained at  $\sim 2.5 \times 10^5$  copies ml<sup>-1</sup> through the suboxic zone and into the sulfidic waters below. Although the absolute abundance of SUP05 decreased at 185 m from April 9<sup>th</sup> to August 11<sup>th</sup>, as a proportion of the bacterial community SUP05 abundance increased from 20 to 30% of total bacteria. On October 15<sup>th</sup>, deep water renewal was underway and a large suboxic zone between 90 and 135 m was observed that was bracketed by oxygenated water above and below. H<sub>2</sub>S was not detected at any depths, while NO<sub>3</sub><sup>-</sup> was between 20-30 uM above and below the suboxic zone and decreased rapidly within the suboxic zone, reaching the detection limit at 135 m. Overall, SUP05 abundance was lower during deepwater renewal than during the earlier stages of stratification however the SUP05 SSU rRNA gene sequences were still detected at all depths below 100 m and ranged between  $\sim 3.0 \times 10^4$  and  $1.5 \times 10^5$  copies ml<sup>-1</sup>. The highest bacterial

1 abundance was observed within the suboxic zone during renewal ( $4 \times 10^6$  copies  $\text{ml}^{-1}$ ), but SUP05  
2 was only a minor component of this deep water fall bloom as it was consistently below 10% of  
3 total detected bacteria. The abundance of SUP05 within water characterized by very little to no  
4  $\text{O}_2$  situated above sulfidic water suggests reliance on reduced sulfur compounds for energy  
5 conservation and also suggests the ability to do so anaerobically, most likely through the  
6 reduction of  $\text{NO}_3^-$ , in support of previous findings (24).

7  
8 **Alignment of unassembled fosmid end reads to symbiont reference genomes.** We  
9 investigated the overall genome sequence conservation amongst the Saanich Inlet SUP05  
10 populations and the symbiont reference genomes using fragment recruitment plots generated  
11 using nucmer (25). Nucmer plots comparing Saanich Inlet metagenomic data to the *R. magnifica*  
12 genome were constructed with the following parameters and cut-offs: breaklength = 1000,  
13 maximum gap length = 200, and minimum match length = 10. Alignment of all the unassembled  
14 fosmid end reads to the *R. magnifica* genome resulted in greater than 6000 reads (3% of total),  
15 averaging 79% nucleotide identity, aligning to the reference genome (Fig. S3). Similar results  
16 were obtained for *V. okutanii* as the reference, as the two symbionts have highly conserved  
17 genome architecture (26). The distribution of aligned reads over the length of the reference  
18 genome varied considerably between different Saanich Inlet samples, yet a similar pattern to the  
19 SUP05 SSU rRNA gene distribution was observed. Coverage was greatest in samples from 200  
20 m collected in the winter and spring, however significant coverage was detected below 100 m  
21 depth throughout the year. Such high coverage is striking and suggests significant conservation  
22 of gene content between the symbiont genomes and pelagic SUP05. Moreover, only 28% of the  
23 aligning sequences were derived from fosmid mate pairs, suggesting a large proportion of genes

1 present in the Saanich Inlet SUP05 population are either absent or poorly conserved within the  
2 reference genome.

3  
4 **Insights into symbiosis.** Until now, the closest genome-sequenced relative of the clam  
5 symbionts was the sulfur-oxidizing chemoautotroph *Thiomicrospira crunogena* (27). Although  
6 *T. crunogena* has proved useful in comparative investigation of symbiont gene content (28), the  
7 availability of metagenomic data from a pelagic member of the symbiont lineage allows for  
8 further study of the differences between symbiont and pelagic members of SUP05. The largest  
9 category of genes specific to the SUP05 metagenome were those involved in DNA replication  
10 recombination and repair, and have likely been lost from symbionts as is typical for obligate  
11 intracellular bacteria (29). Moreover, ribonucleotide biosynthesis in pelagic SUP05 can proceed  
12 via the alpha and beta subunits of ribonucleotide reductase (*nrdAB*) that is conserved with the  
13 symbionts. However, the SUP05 metagenome specifically encodes a second, oxygen sensitive  
14 reductase (*nrdD*), further supporting the facultative or strict anaerobic lifestyle of the pelagic  
15 SUP05 population in Saanich Inlet. Other genes involved in DNA metabolism specific to the  
16 SUP05 metagenome include DNA internalization and competence genes as well as a Type II  
17 secretion/type IV pilus system that may be involved in DNA uptake or perhaps protein secretion  
18 or twitching motility (30). The SUP05 metagenome also has many genes involved in inorganic  
19 ion transport and metabolism, reflecting the variable external environment of pelagic SUP05.  
20 These include high and low affinity sulfate transporter, nitrate/nitrite transporters, TonB  
21 dependent receptors, ferrous iron transport proteins, bacterioferritin, and alkylphosphonate  
22 uptake proteins.

1           Recently, a comparative analysis of two symbiont genomes lead to the hypothetical  
2 reconstruction of the last common symbiotic ancestor's (LCSA's) carbon and energy metabolism  
3 (28). By in large, our analysis of the SUP05 metagenome supports the findings that LCSA had a  
4 complex sulfur metabolism including both the Sox system and reverse DSR pathways of sulfur  
5 compound oxidation, fixed CO<sub>2</sub> via the CBB cycle, and had the ability to respire nitrate  
6 anaerobically via a membrane bound dissimilatory nitrate reductase. In addition, it was  
7 previously determined that the larger *R. magnifica* genome had many genes involved in the  
8 biosynthesis of polysaccharides and peptidoglycan that were not present in the slightly smaller *V.*  
9 *okutanii* genome. Many of these genes were also absent from the SUP05 metagenome (Figure  
10 S5) and hence appear specific to the *R magnifica* symbiont genome where they may play an  
11 interactive role with the clam host. A complete functional breakdown of genes shared among and  
12 between symbionts and the SUP05 metagenome is available from the authors upon request.  
13

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23

**Supplementary Figure Legends**

**Fig. S1** The seasonal stratification cycle in Saanich Inlet. (A) Oxygen (O<sub>2</sub>) (B) nitrate (NO<sub>3</sub><sup>-</sup>) and (C) sulfide (H<sub>2</sub>S) concentrations plotted from February 2006 to December 2008. H<sub>2</sub>S data collection began in April 07. The inset at the bottom left of the Fig. (D) outlines the depths and dates of samples employed in environmental genomic analysis

**Fig. S2** Quantification of SUP05 and bacterial SSU rRNA genes during the main stages of the 2008 seasonal stratification and renewal cycle. The grey box denotes the non-sulfidic suboxic zone.

**Fig. S3** Identification of metagenomic reads originating from SUP05 in Saanich Inlet. For each water sample, fosmid end reads were aligned to the *R. magnifica* symbiont reference genome using nucmer (25).

**Fig. S4** Overview of scaffolds identified during the SUP05 assembly and binning process (A) Net scaffold length versus sequence depth of the assembled Saanich Inlet scaffolds, highlighting the set of scaffolds enriched with SUP05. The complete scaffold set comprises 13,628 scaffolds ≥ 1kb comprised two or more reads, with a combined net length of 21.6 Mb. (B) Visualization of the first two components of a principal component analysis of the SUP05-enriched scaffolds and known fully-sequenced SUP05 fosmids, in which Z-scores for all possible 64 tri-mers, 256 tetra-mers, and 1024 penta-mers were calculated with TETRA (31) and normalized by length. Scaffolds binned to SUP05 are in red while scaffolds that were included in the enriched scaffold

1 set but were not binned to SUP05 are in blue. The complete scaffold set is in grey. Fully  
2 sequenced SUP05 fosmids are in green. Scaffolds and fosmids containing SSU rRNA genes are  
3 indicated.

4

5 **Fig. S5** Functional analysis of the genes shared between the SUP05 members based on (A) COG  
6 and (B) KEGG categories. Column colors correspond to the same sectors of the Venn diagram in  
7 Fig. 3. The number of genes is normalized to the nonredundant size of each sector of the Venn  
8 diagram, which is included in the column headers in brackets. Therefore, one can compare the  
9 functional distribution of genes within any one sector (across columns) but not between sectors  
10 (across rows).

11

12 **Fig. S6** Pathways of central carbon, nitrogen, and sulfur metabolism in the SUP05 metagenome.  
13 Genes are color coded by their distribution amongst the three SUP05 representatives as revealed  
14 in the Venn diagram presented in Fig. 3.

15

16 **Fig. S7** Color coding system for COG categories used in Fig. 3.

17

## 1 **References**

2

- 3 1. E. Zaikova *et al.*, *Environ Microbiol*, (submitted).
- 4 2. J. D. Cline, *Limnol. Oceanogr.* 14, 288 (1969).
- 5 3. T. Huber, G. Faulkner, P. Hugenholtz, *Bioinformatics* 20, 2317 (2004).
- 6 4. K. E. Ashelford, N. A. Chuzhanova, J. C. Fry, A. J. Jones, A. J. Weightman, *Appl*  
7 *Environ Microbiol* 72, 5734 (2006).
- 8 5. E. Pruesse *et al.*, *Nucleic Acids Research* 35, 7188 (2007).
- 9 6. W. Ludwig *et al.*, *Nucleic Acids Research* 32, 1363 (2004).
- 10 7. R. C. Edgar, *Nucleic Acids Res* 32, 1792 (2004).
- 11 8. P. D. Schloss, J. Handelsman, *Appl Environ Microbiol* 71, 1501 (2005).
- 12 9. Felsenstein, PHYLIP (Phylogeny Inference Package) version 3 (2004).
- 13 10. S. Guindon, F. Lethiec, P. Duroux, O. Gascuel, *Nucleic Acids Res* 33, W557 (2005).
- 14 11. S. J. Hallam *et al.*, *Science* 305, 1457 (2004).
- 15 12. K. Mavromatis *et al.*, *Nat Methods* 4, 495 (2007).
- 16 13. M. Pop, D. S. Kosack, S. L. Salzberg, *Genome Res* 14, 149 (2004).
- 17 14. D. Gordon, C. Abajian, P. Green, *Genome Res* 8, 195 (1998).
- 18 15. H. Teeling, A. Meyerdierks, M. Bauer, R. Amann, F. O. Glockner, *Environ Microbiol* 6,  
19 938 (2004).
- 20 16. B. McCune, J. B. Grace, D. L. Urban. Mjrm Software Design (2002).
- 21 17. T. Frickey, A. N. Lupas, *Nucleic Acids Res* 32, 5231 (2004).
- 22 18. T. M. Lowe, S. R. Eddy, *Nucleic Acids Research* 25, 955 (1997).
- 23 19. S. Hunter *et al.*, *Nucleic Acids Research* 37, D211 (2009).
- 24 20. R. Caspi *et al.*, *Nucleic Acids Research* 36, D623 (2008).



- 1 21. D. A. Rasko, G. S. Myers, J. Ravel, *BMC Bioinformatics* 6, 2 (2005).
- 2 22. R. Ghai, T. Hain, T. Chakraborty, *Bmc Bioinformatics* 5, (2004).
- 3 23. A. C. E. Darling, B. Mau, F. R. Blattner, N. T. Perna, *Genome Research* 14, 1394 (2004).
- 4 24. G. Lavik *et al.*, *Nature* 457, 581 (2009).
- 5 25. A. L. Delcher, A. Phillippy, J. Carlton, S. L. Salzberg, *Nucleic Acids Research* 30, 2478
- 6 (2002).
- 7 26. H. Kuwahara *et al.*, *Extremophiles* 12, 365 (2008).
- 8 27. K. M. Scott *et al.*, *PLoS Biol* 4, e383 (2006).
- 9 28. I. L. Newton, P. R. Girguis, C. M. Cavanaugh, *BMC Genomics* 9, 585 (2008).
- 10 29. N. A. Moran, *Curr Opin Microbiol* 6, 512 (2003).
- 11 30. F. F. Evans, S. Egan, S. Kjelleberg, *Environ Microbiol* 10, 1101 (2008).
- 12 31. H. Teeling, J. Waldmann, T. Lombardot, M. Bauer, F. O. Glockner, *BMC Bioinformatics*
- 13 5, 163 (2004).
- 14
- 15

Table S1. SUP05 metagenome features

Specifications	
No. of scaffolds, contigs	19, 90
Total length (Mb)	1.16
Total span (Mb)	1.91
Average depth of coverage	7.25
No. of assembled end reads, % of total	10449, 4.3
Average polymorphism (%)	0.4
Average G+C content (%)	40
ORF content	
No. of ORFs	1333
Predicted functional in COGs <sup>1</sup>	1033
Predicted functional in KEGG <sup>1</sup>	771
Conserved hypothetical	216
Hypothetical <sup>2</sup>	74
rRNA genes	
16S-5S-23S operon	1
tRNAs	29

<sup>1</sup>Expectation-value cut off 10E-10

<sup>2</sup>No similarity at an expectation value cut off 10E-5

**Table S2. Quantification of SUP05 SSU rRNA genes in Saanich Inlet and summary of fosmid end sequence data**

Cruise ID	Sample Date (mm/dd/yy)	Sample Depth (m)	O <sub>2</sub> (μM)	NO <sub>3</sub> <sup>-</sup> (μM)	% bacterial SSU rRNA gene library <sup>1</sup>	SUP05 SSU rRNA copies/ml (s.d.) <sup>2</sup>	No. of fosmid end reads	No. of fosmid end reads in SUP05 metagenome (% of total)
Feb06	02/18/06	10	212.0	26.7	0/0	1.22x10 <sup>3</sup> (1.1 x10 <sup>3</sup> )	13339	4 (<1)
Feb06	02/18/06	100	51.1	21.0	7/2	1.66x10 <sup>4</sup> (2.24 x10 <sup>3</sup> )	12709	483 (4)
Feb06	02/18/06	125	05	9.7	17/2	6.24x10 <sup>4</sup> (2.15 x10 <sup>4</sup> )	13860	1631 (12)
Feb06	02/18/06	215	1.4	1.8	31/6	5.48x10 <sup>4</sup> (1.33 x10 <sup>4</sup> )	13453	1400 (10)
Jul06	07/06/06	10	381.6	0.3	0/0	0 (0)	13754	1 (<1)
Jul06	07/06/06	100	22.9	20.4	10/1	5.57x10 <sup>4</sup> (2.32 x10 <sup>4</sup> )	12921	513 (4)
Jul06	07/06/06	120	6.5	20.2	19/2	2.48x10 <sup>5</sup> (7.89 x10 <sup>4</sup> )	13321	835 (6)
Jul06	07/06/06	200	0	0.1	25/2	9.65x10 <sup>4</sup> (2.46x10 <sup>4</sup> )	12443	265 (2)
Nov06	11/14/06	10	249.0	25.5	1/0	0 (0)	12933	1 (<1)
Nov06	11/14/06	100	15.4	13.0	11/4	1.05x10 <sup>5</sup> (3.88x10 <sup>4</sup> )	13151	824 (6)
Nov06	11/14/06	120	9.8	8.9	26/14	7.06x10 <sup>4</sup> (5.74 x10 <sup>3</sup> )	14810	1014 (7)
Nov06	11/14/06	200	54.0	19.8	17/8	1.07x10 <sup>5</sup> (2.07 x10 <sup>4</sup> )	13627	990 (7)
Apr07	04/24/07	10	316.1	18.2	0/0	0 (0)	13811	3 (<1)
Apr07	04/24/07	100	67.3	26.5	1/0	6.35x10 <sup>3</sup> (9 x10 <sup>2</sup> )	13286	70 (<1)
Apr07	04/24/07	120	26.9	20.6	7/1	1.58x10 <sup>4</sup> (2.37 x10 <sup>3</sup> )	14522	390 (3)
Apr07	04/24/07	200	1.1	0.0	36/6	3.70x10 <sup>5</sup> (5.86 x10 <sup>4</sup> )	14066	2026 (15)

<sup>1</sup>Report as SUP05 SI-1/SI-2 phylotypes

<sup>2</sup>Quantification of SUP05 SSU rRNA gene copies during Feb06, Jul06, and Nov06 were first reported in (Zaikova et al.)

Table S3. Identification and phylogenetic identity of typically conserved genes in the SUP05 metagenome<sup>1</sup>

COG family	Conserved genes	No. of genes (% id) <sup>2</sup>	SUP05 metagenome locus tag <sup>3</sup>	<i>R. magnifica</i> accession no.	SUP05 monophyly <sup>4</sup>
<b>Large subunit ribosomal proteins</b>					
<b>COG0080</b>	<b>L11</b>	<b>1</b>	<b>Sup05_0527</b>	<b>YP_904007</b>	<b>+</b>
<b>COG0081</b>	<b>L1</b>	<b>1</b>	<b>Sup05_0526</b>	<b>YP_904006</b>	<b>+</b>
<b>COG0087</b>	<b>L3</b>	<b>1</b>	<b>Sup05_0550</b>	<b>YP_903428</b>	<b>+</b>
COG0088	L4	1	Sup05_0551	YP_903429	+
COG0089	L23	1	Sup05_0552	YP_903430	+
COG0090	L2	1	Sup05_0553	YP_903431	+
<b>COG0091</b>	<b>L22</b>	<b>1</b>	<b>Sup05_0555</b>	<b>YP_903433</b>	<b>+</b>
<b>COG0093</b>	<b>L14</b>	<b>1</b>	<b>Sup05_0559</b>	<b>YP_903438</b>	<b>+</b>
<b>COG0094</b>	<b>L5</b>	<b>1</b>	<b>Sup05_0561</b>	<b>YP_903440</b>	<b>+</b>
<b>COG0097</b>	<b>L6P/L9E</b>	<b>1</b>	<b>Sup05_0277</b>	<b>YP_903443</b>	<b>+</b>
<b>COG0102</b>	<b>L13</b>	<b>1</b>	<b>Sup05_1311</b>	<b>YP_904120</b>	<b>+</b>
<b>COG0197</b>	<b>L16/L10E</b>	<b>1</b>	<b>Sup05_0557</b>	<b>YP_903435</b>	<b>+</b>
COG0198	L24	1	Sup05_0560	YP_903439	+
<b>COG0200</b>	<b>L15</b>	<b>1</b>	<b>Sup05_0281</b>	<b>YP_903447</b>	<b>+</b>
COG0203	L17	2 (99)	Sup05_0288 Sup05_0652	YP_903454	+
COG0211	L27	1	Sup05_1289	YP_904102	+
COG0222	L7/L12	1	Sup05_0524	YP_904004	+
COG0227	L28	2 (100)	Sup05_1278 Sup05_0596	YP_904171	+
COG0230	L34	1	NA	YP_904212	ND
COG0244	L10	1	Sup05_0525	YP_904005	+
COG0254	L31	1	Sup05_1222	YP_904203	+
COG0255	L29	1	NA	YP_903436	ND
<b>COG0256</b>	<b>L18</b>	<b>1</b>	<b>Sup05_0278</b>	<b>YP_903444</b>	<b>+</b>
COG0257	L36	1	NA	YP_903449	ND
COG0261	L21	1	Sup05_1290	YP_904103	+
COG0267	L33	2	NA	YP_904170	ND
COG0291	L35	2	NA	YP_903862	ND
COG0292	L20	1	Sup05_1135	YP_903861	+
COG0333	L32	1	Sup05_0983	YP_903712	+
<b>COG0335</b>	<b>L19</b>	<b>0</b>	<b>No hits found</b>	<b>YP_904059</b>	<b>ND</b>
COG0359	L9	1	Sup05_0377	YP_903879	+
COG1825	L25	1	Sup05_0668	YP_904242	+
COG1841	L30/L7E	1	Sup05_0280	YP_903446	+
<b>Small subunit ribosomal proteins</b>					
<b>COG0048</b>	<b>S12</b>	<b>1</b>	<b>Sup05_0234</b>	<b>YP_903423</b>	<b>+</b>
<b>COG0049</b>	<b>S7</b>	<b>1</b>	<b>Sup05_0235</b>	<b>YP_903424</b>	<b>+</b>
COG0051	S10	2 (100)	Sup05_0549 Sup05_0238	YP_903427	+
<b>COG0052</b>	<b>S2</b>	<b>1</b>	<b>Sup05_1214</b>	<b>YP_904247</b>	<b>+</b>
<b>COG0092</b>	<b>S3</b>	<b>1</b>	<b>Sup05_0556</b>	<b>YP_903434</b>	<b>+</b>
<b>COG0096</b>	<b>S8</b>	<b>1</b>	<b>Sup05_0276</b>	<b>YP_903442</b>	<b>+</b>
<b>COG0098</b>	<b>S5</b>	<b>1</b>	<b>Sup05_0279</b>	<b>YP_903445</b>	<b>+</b>
<b>COG0099</b>	<b>S13</b>	<b>1</b>	<b>Sup05_0283</b>	<b>YP_903450</b>	<b>+</b>
<b>COG0100</b>	<b>S11</b>	<b>1</b>	<b>Sup05_0284</b>	<b>YP_903451</b>	<b>+</b>

<b>COG0103</b>	<b>S9</b>	<b>1</b>	<b>Sup05_1310</b>	<b>YP_904119</b>	<b>+</b>
<b>COG0184</b>	<b>S15P/S13E</b>	<b>1</b>	<b>Sup05_0244</b>	<b>YP_903544</b>	<b>+</b>
COG0185	S19	1	Sup05_0554	YP_903432	+
<b>COG0186</b>	<b>S17</b>	<b>1</b>	<b>Sup05_0558</b>	<b>YP_903437</b>	<b>+</b>
COG0199	S14	2	NA	YP_903441	ND
COG0228	S16	1	Sup05_0593	YP_904167	+
COG0238	S18	1	Sup05_0376	YP_903878	+
COG0268	S20	1	Sup05_1171	YP_903947	+
COG0360	S6	1	Sup05_0375	YP_903877	+
<b>COG0522</b>	<b>S4</b>	<b>1</b>	<b>Sup05_0285</b>	<b>YP_903452</b>	<b>+</b>
COG0539	S1	1	Sup05_1121	YP_903814	+
COG0828	S21	1	Sup05_0653	YP_903455	+
<b>tRNA synthetases</b>					
COG0008	Glutamyl-tRNA synthetase	0	No hits found	YP_904093	ND
COG0008	Glutamyl-tRNA synthetase	1	Sup05_0971	YP_903701	+
COG0013	Alanyl-tRNA synthetase	0	No hits found	YP_904068	ND
COG0016	Phenylalanyl-tRNA synthetase alpha subunit	1	Sup05_1136	YP_903859	+
COG0017	Aspartyl/asparaginyl-tRNA synthetases	0	No hits found	No hits found	
COG0018	Arginyl-tRNA synthetase	1	Sup05_0153	YP_903350	+
COG0060	Isoleucyl-tRNA synthetase	1	Sup05_0784	YP_903587	+
COG0072	Phenylalanyl-tRNA synthetase beta subunit	1	Sup05_1137	YP_903858	+
COG0124	Histidyl-tRNA synthetase	1	Sup05_0384	YP_903884	+
COG0162	Tyrosyl-tRNA synthetase	1	Sup05_0054	YP_903397	+
COG0172	Seryl-tRNA synthetase	1	Sup05_1269	YP_904185	+
COG0173	Aspartyl-tRNA synthetase	1	Sup05_0825	YP_903634	+
COG0180	Tryptophanyl-tRNA synthetase	1	Sup05_0782	YP_903585	+
COG0215	Cysteinyl-tRNA synthetase	2 (91)	Sup05_0188 Sup05_0135	YP_903366	+
COG0441	Threonyl-tRNA synthetase	2 (98)	Sup05_0348 Sup05_1133	YP_903864	+
COG0442	Prolyl-tRNA synthetase	1	Sup05_0877	YP_903595	+
<b>COG0495</b>	<b>Leucyl-tRNA synthetase</b>	<b>0</b>	<b>No hits found</b>	<b>YP_903280</b>	<b>ND</b>
COG0525	Valyl-tRNA synthetase	1	Sup05_0929	YP_903692	+
COG0751	Glycyl-tRNA synthetase, beta subunit	1	Sup05_0457	YP_903923	+
COG0752	Glycyl-tRNA synthetase, alpha subunit	1	Sup05_0477	YP_903939	+
COG1190	Lysyl-tRNA synthetase	2 (99)	Sup05_1123 Sup05_0305	YP_903816	+
	Methionyl-tRNA synthetase	1	Sup05_1115	YP_903792	+
<b>COG0201</b>	<b>Preprotein translocase subunit SecY</b>	<b>1</b>	<b>Sup05_0282</b>	<b>YP_903448</b>	<b>+</b>
COG0341	Preprotein translocase subunit SecF	1	Sup05_1265	YP_904182	+
COG0342	Preprotein translocase subunit SecD	1	Sup05_1267	YP_904183	+
COG0653	Preprotein translocase subunit SecA	0	No hits found	YP_903288	ND
COG0690	Preprotein translocase subunit SecE	0	No hits found	YP_904009	ND
COG0706	Preprotein translocase subunit YidC	1	Sup05_0697	YP_904215	+
COG1314	Preprotein translocase subunit SecG	0	No hits found	No hits found	
COG1862	Preprotein translocase subunit YajC	1	Sup05_1268	YP_904184	+
COG1952	Preprotein translocase subunit SecB	0	No hits found	No hits found	
<b>RNA polymerase subunits</b>					
<b>COG0085</b>	<b>DNA-directed RNA polymerase, beta subunit/140 kD subunit</b>	<b>1</b>	<b>Sup05_0523</b>	<b>YP_904003</b>	<b>+</b>
COG0086	DNA-directed RNA polymerase, beta' subunit/160 kD subunit	2	Sup05_0521 Sup05_0522 (potential split gene)	YP_904002	+

<b>COG0202</b>	<b>DNA-directed RNA polymerase, alpha subunit/40 kD subunit</b>	<b>2</b>	<b>Sup05_0287</b> <b>Sup05_0286</b> (potential split gene)	<b>YP_903453</b>	<b>+</b> <b>+</b>
COG0568	DNA-directed RNA polymerase, sigma 32	0	No hits found	YP_903783	+
COG0568	DNA-directed RNA polymerase, sigma 70	1	Sup05_0106	YP_903375	+
COG1758	DNA-directed RNA polymerase, subunit K/omega	1	NA	YP_903640	ND
<b>COG0012</b>	<b>Predicted GTPase</b>	<b>1</b>	<b>Sup05_0666</b>	<b>YP_904244</b>	<b>+</b>
<b>COG0533</b>	<b>Metal-dependent protease</b>	<b>1</b>	<b>Sup05_0585</b>	<b>YP_904157</b>	<b>+</b>

<sup>1</sup> Genes outlined in bold are the 32 universally conserved single copy genes reported in Ciccarelli et al 2006. Genes missing from the SUP05 metagenome, but present in the *R. magnifica* reference genome are highlighted in green.

<sup>2</sup> For those genes present in multiple copies, the pairwise amino acid sequence similarity is presented in brackets.

<sup>3</sup> NA refers to genes that were present in the SUP05 metagenome but not automatically annotated by fgenesb gene-calling software.

<sup>4</sup> + refers to genes that formed a phylogenetically coherent cluster with the symbiont reference genomes upon automated phylogenetic reconstruction.

Table S4. Toxin-antitoxin systems identified in the SUP05 metagenome.

Toxin family <sup>1</sup>	Locus tag	Putative antidote locus tag	Genes in operon with TA system
<i>relE/parE</i>	Sup05_0405	Sup05_0404	MGD synthase
	Sup05_0483	Sup05_0482	-
	Sup05_0538	-	-
	Sup05_0866	-	-
	Sup05_0615	-	-
	Sup05_0057	Sup05_0056	-
	Sup05_0104	Sup05_0103	PRPP synthetase
	Sup05_0829	Sup05_0830	-
	Sup05_0293	-	-
	Sup05_1093	Sup05_1092	Selenophosphate synthase
	<i>hicB</i>	Sup05_1327	Sup05_1326
Sup05_0962		Sup05_0961	Integrases, transposases

<sup>1</sup>Family assignment is based on conserved domain search at NCBI

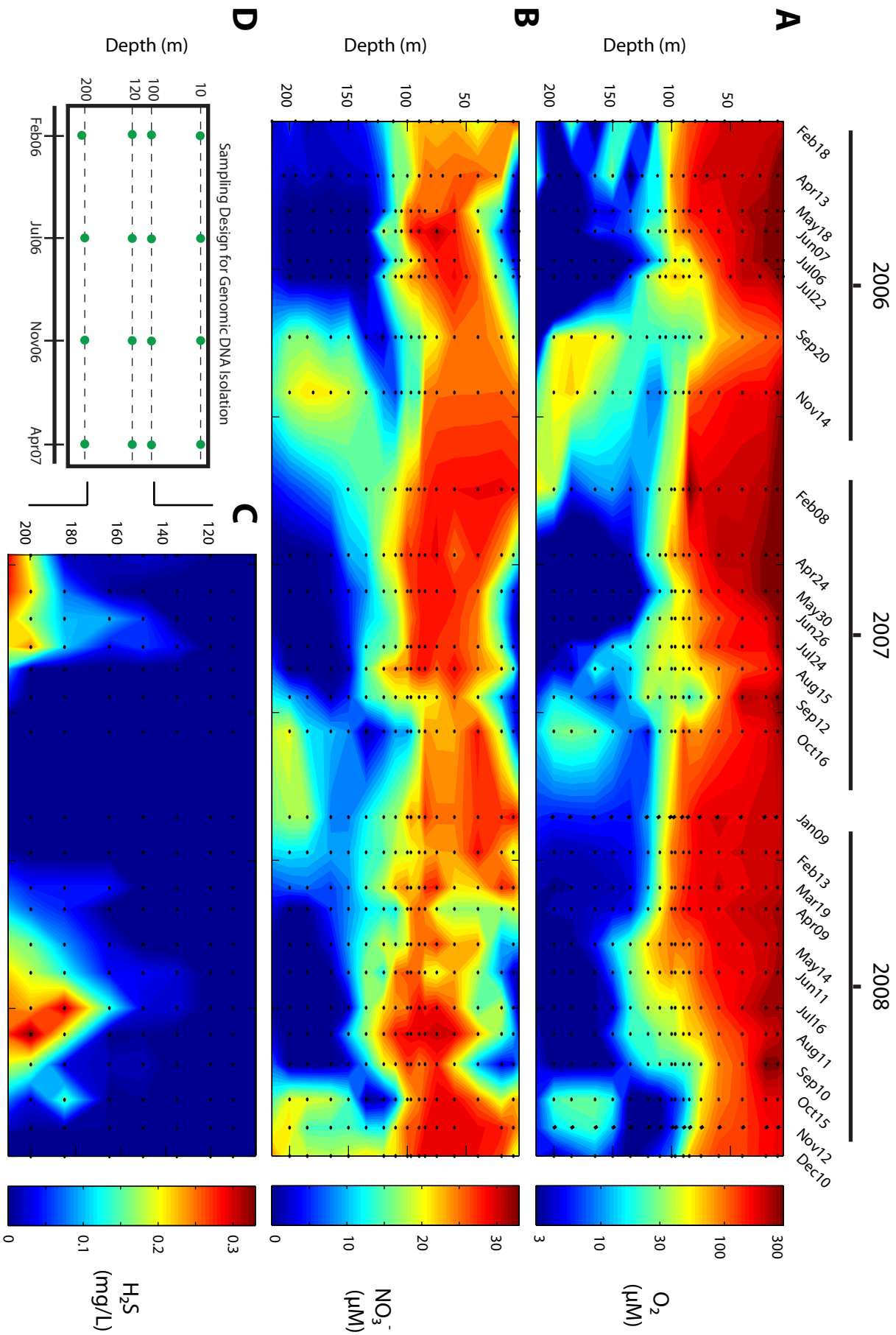


Fig. S1



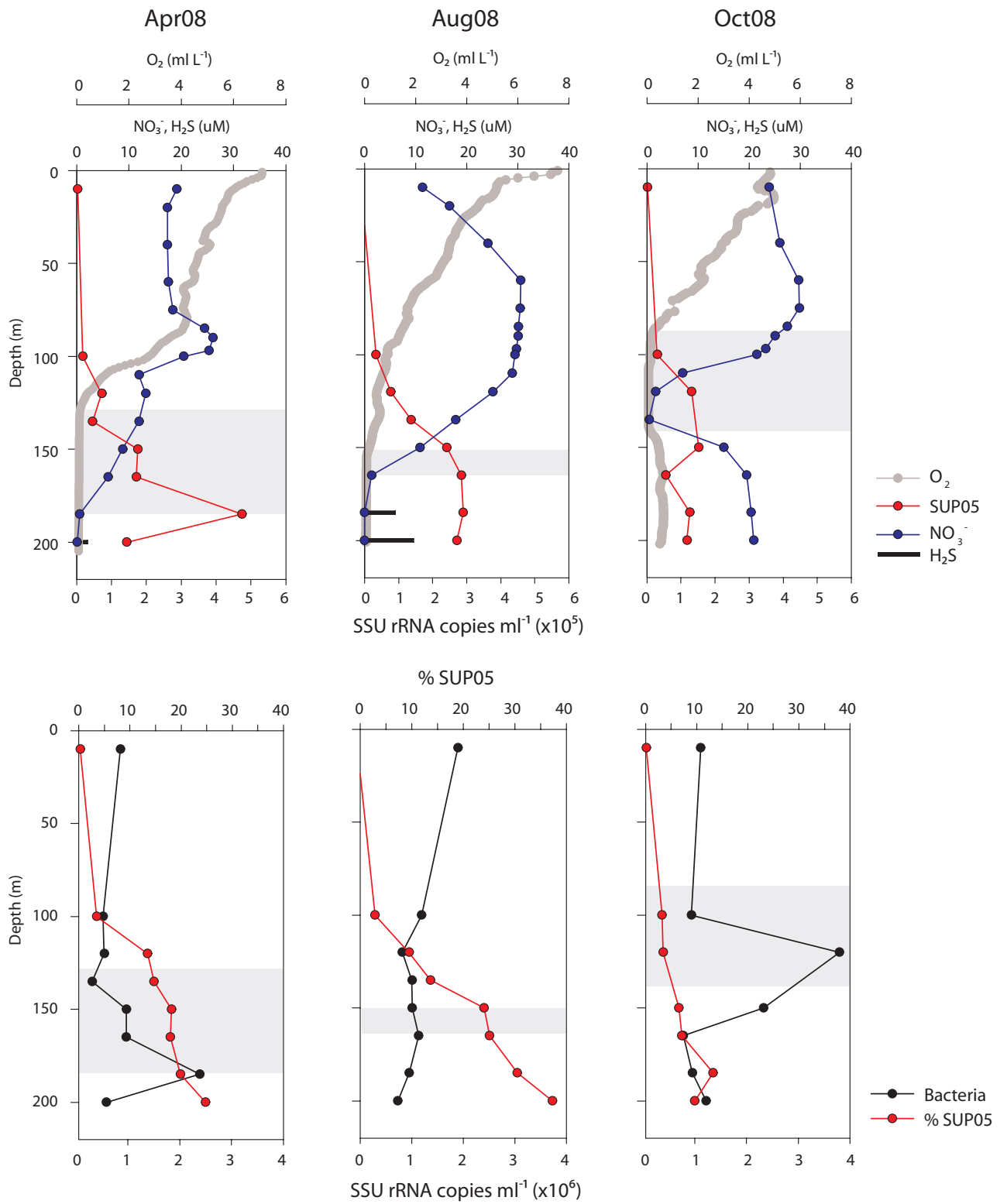


Fig. S2

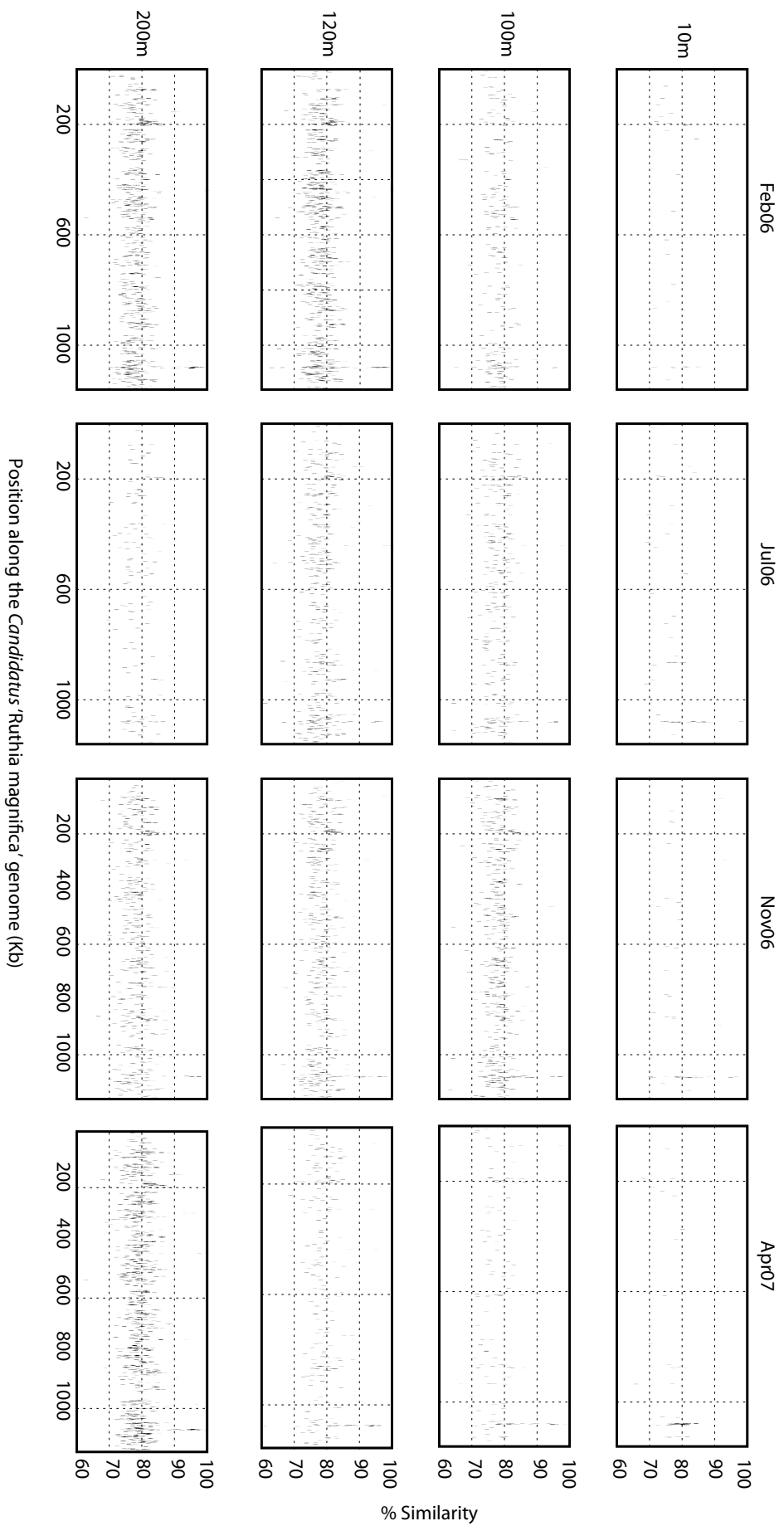


Fig. S3

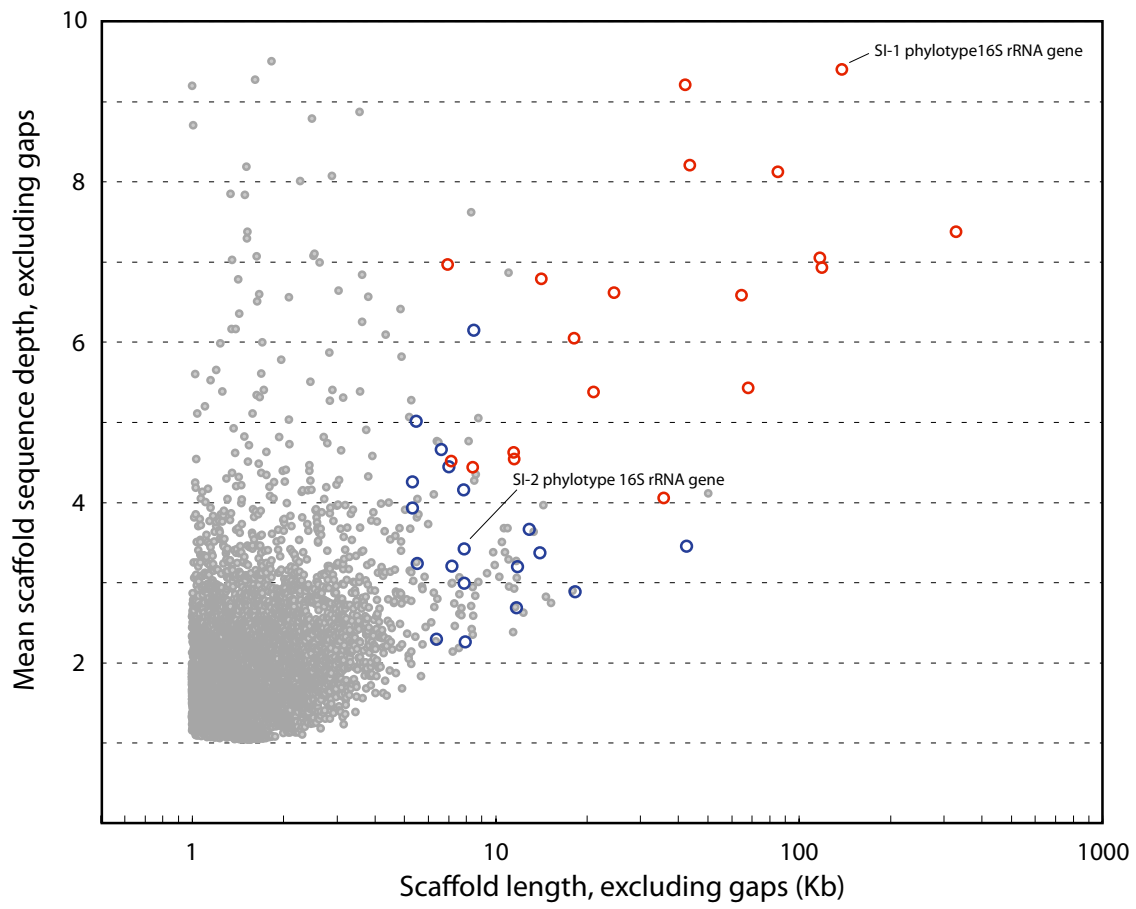
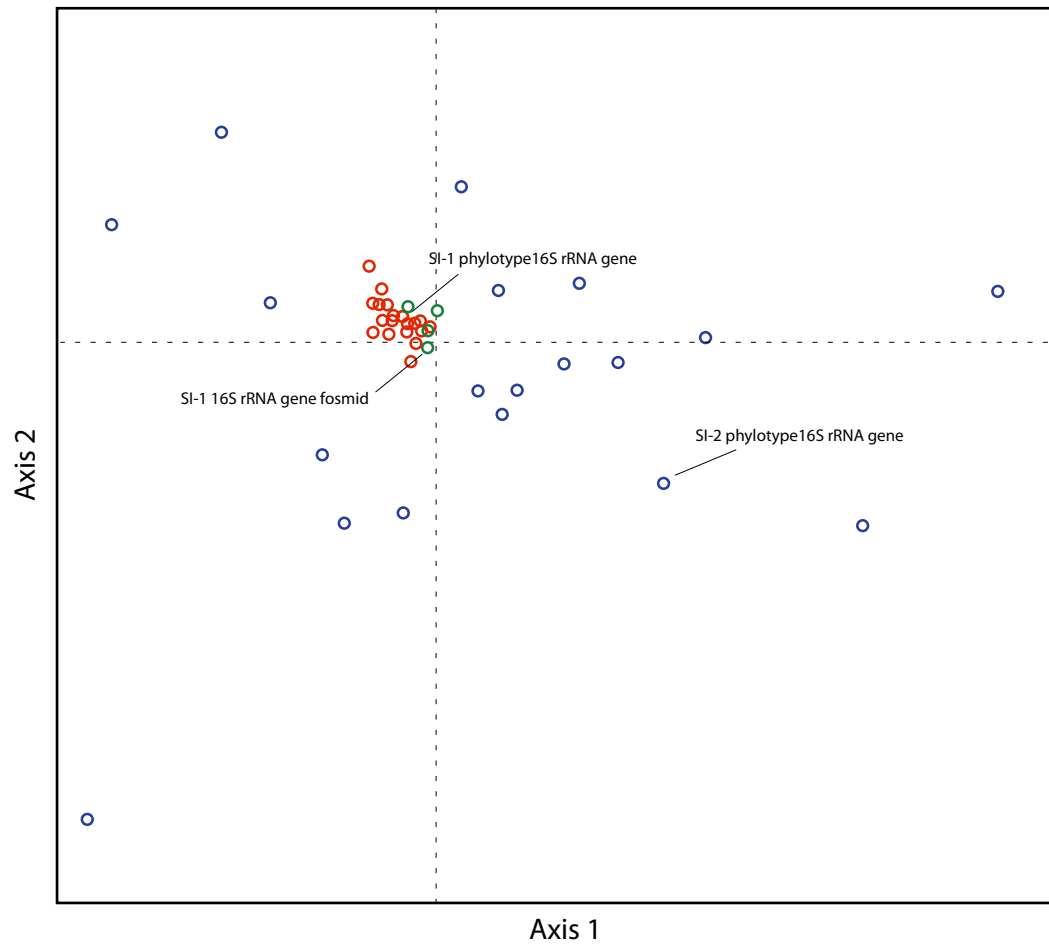
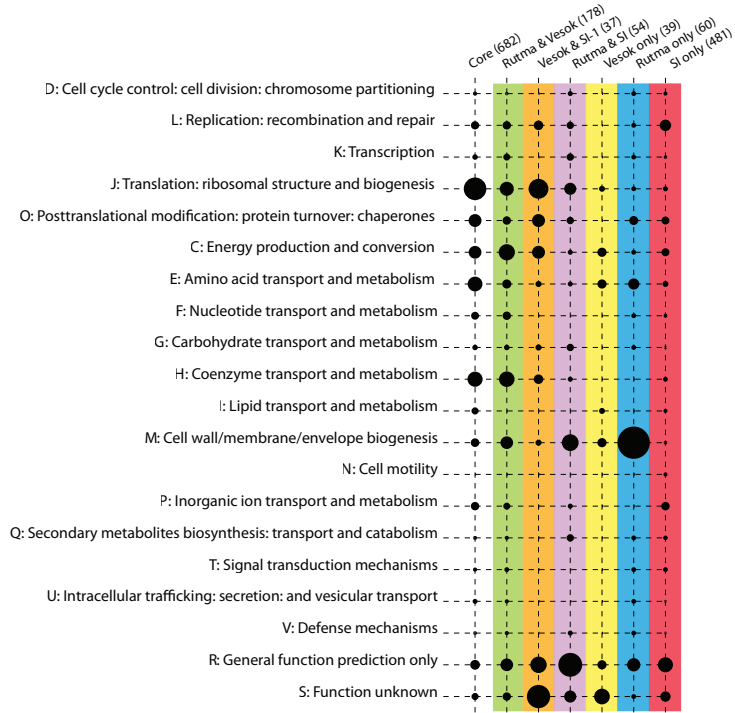
**A****B**

Fig. S4

**A**



**B**

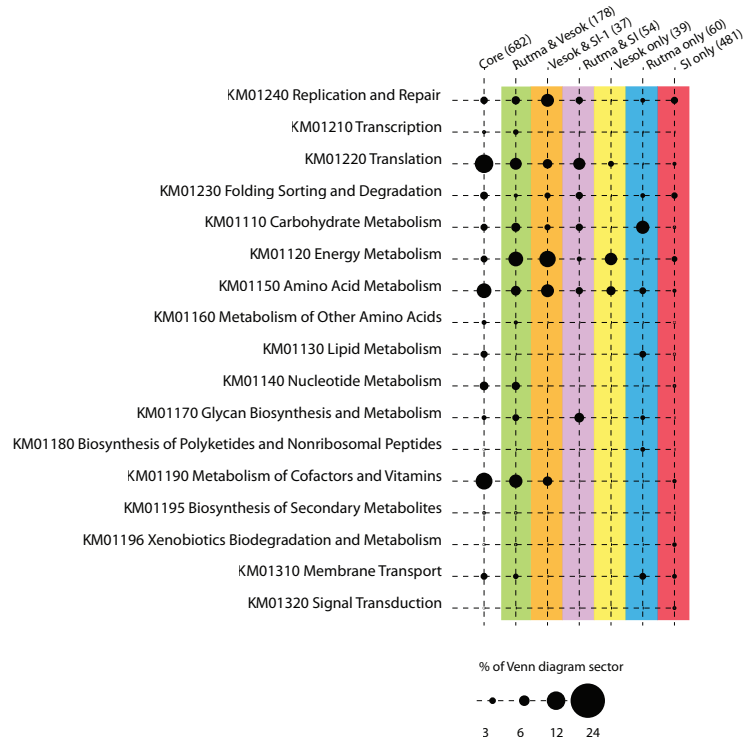
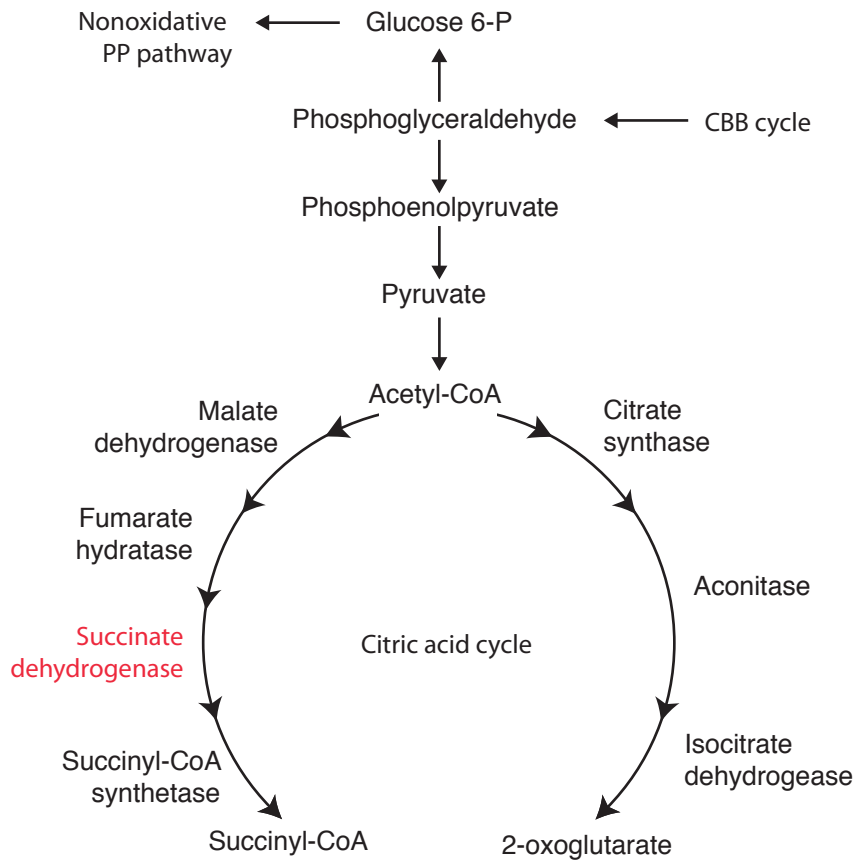
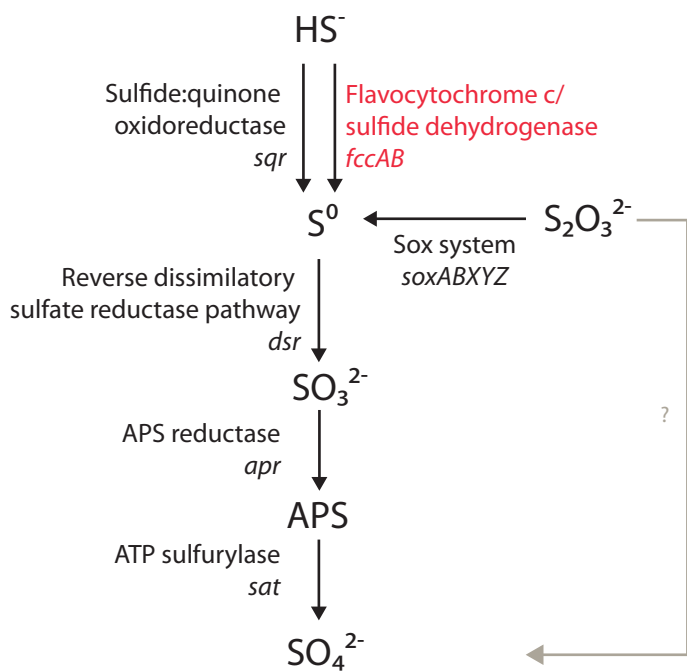


Fig. S5

# Central carbon metabolism



## Sulfur oxidation



## Denitrification

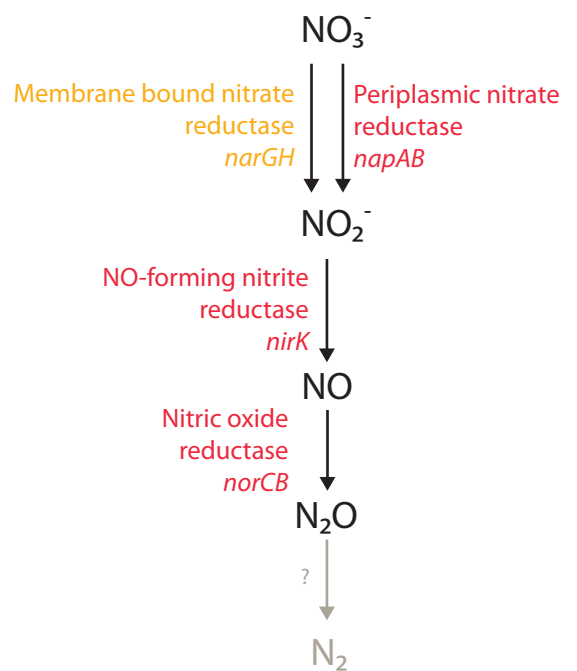



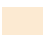


















Fig. S6

-  J Translation, ribosomal structure and biogenesis
-  K Transcription
-  L DNA replication, recombination and repair
-  D Cell division and chromosome partitioning
-  T Signal transduction mechanisms
-  M Cell envelope biogenesis, outer membrane
-  N Cell motility and secretion
-  U Intracellular trafficking and secretion
-  O Posttranslational modification, turnover, chaperones
-  C Energy production and conversion
-  G Carbohydrate transport and metabolism
-  E Amino acid transport and metabolism
-  F Nucleotide transport and metabolism
-  H Coenzyme metabolism
-  I Lipid metabolism
-  P Inorganic ion transport and metabolism
-  Q Secondary metabolites biosynthesis, transport, catabolism
-  R General function prediction only
-  S Function unknown
-  - Not in COGs