1    **Gene expression in the deep biosphere**

2

3    William D. Orsi[1*], Virginia P. Edgcomb[1], Glenn D. Christman[2], and Jennifer F. Biddle[2]

4

5    Affiliations:

6    [1] Department of Geology and Geophysics, Woods Hole Oceanographic Institution

7    [2] College of Earth, Ocean, and Environment, University of Delaware

8

9    Keywords: marine subsurface, mRNA, metatranscriptomics, deep biosphere, Illumina

10    sequencing, marine sediment, subseafloor, ODP Leg 201, gene expression

11    Running title:  Gene expression in the deep biosphere.

12    *To whom correspondence should be addressed

13

14    **Scientific ocean drilling has revealed a deep biosphere of widespread microbial life in sub-**

15    **seafloor sediment.  Microbial metabolism in the marine subsurface likely plays an**

16    **important role in global biogeochemical cycles[1-3] but deep biosphere activities are not well**

17    **understood[1].  Here, we describe and analyze the first subseafloor metatranscriptomes from**

18    **anaerobic Peru Margin sediment up to 159 meters below seafloor (mbsf) represented by**

19    **over 1 billion cDNA sequence reads. Anaerobic metabolism of amino acids, carbohydrates,**

20    **and lipids appear to be dominant metabolic processes, and profiles of dissimilatory sulfite**

21    **reductase (*Dsr*) transcripts are consistent with porewater sulfate concentration profiles[1].**

22    **Moreover, transcripts involved in cell division increase as a function of microbial cell**

23    **concentration, indicating that increases in subseafloor microbial abundance are a function**

**of cell division across all three domains of life. These data support calculations[1] and**

**models[4] of subseafloor microbial metabolism and represent the first holistic picture of deep**

**biosphere activities.**

Abundant microbial cells[5, 6] exist in sub-seafloor (>1.5 mbsf) sediment and represent a

significant portion of Earth's biomass[7, 8]. Marine sediment contains Earth's largest pool of

organic carbon, which may be the primary energy source for subsurface microbes[1, 2, 9, 10, 11]. A

model recently suggested biomass turnover rates on the order of thousands of years in the marine

subsurface that are hypothesized to have an impact on global biogeochemical cycling over

geological timescales[4]. Logistical sampling constraints, the complex sediment matrix composed

of organic material and minerals, and low metabolic rates[3,4], have hindered directed testing of

microbial activities at the molecular level in this environment. A better understanding of deep

biosphere activities will help define the deep biosphere's role in global biogeochemical cycles[12].

We optimized an mRNA extraction and amplification protocol for subseafloor sediment,

and combined this with high-throughput sequencing to report the first dataset on microbial gene

expression in the marine subsurface, demonstrating that despite the extremely low metabolic

rates[1, 4], mRNA-based investigations of the deep biosphere are possible and informative. We

use the gene expression data to reconstruct active community metabolism and results support

calculations[1] and models[4] of sub-seafloor microbial activities. The Peru Margin (Ocean Drilling

Program Leg 201, Site 1229D) was analyzed because a wealth of biogeochemical data exists for

this site [e.g. 1, 4,, 6, 9, 10] that exhibits peaks of cell abundance, and profiles of sulfate and methane

suggestive of microbial activity[1] (Figure 1).

Picogram quantities of total RNA were extracted from 25 grams of Peru Margin sediment

from six depths (5, 30, 50, 70, 91, 159 mbsf), consistent with basal levels of microbial activity

47    predicted for this environment[3, 4]. Illumina® sequencing of total cDNA produced over 1 billion

48    reads, with 50% to 85% of reads mapping to open reading frames that were assigned a functional

49    annotation (Table S1).

50         The dominance of transcripts from Firmicutes, Actinobacteria, Alphaproteobacteria, and

51    Gammaproteobacteria (Fig. S1) is consistent with previous cultivation-based, metagenomic, and

52    phylogenetic surveys from Peru Margin subsurface sediment[1, 5, 13, 14], and suggests these to be

53    some of the most active microbial groups. The abundance of gammaproteobacterial transcripts

54    (Fig. S1) suggests that they are likely the most active microbial group in the deeper, anoxic,

55    subseafloor sediment at this site.  Fungal transcripts were also present in every sample ranging in

56    representation from 3% at 70 mbsf to 20% at 5 mbsf.  Archaea and Chloroflexi are present in

57    noticeably low abundance, despite their previous detection at this site[6, 13, 15], suggesting that our

58    approach might miss organisms with lower mRNA expression levels. As such, interpretations of

59    relative abundances should be treated cautiously[16]. Changes in pressure and temperature may

60    have altered gene expression during sampling.   However, low representation of heat shock

61    proteins (a proxy for physiological stress response[17]) in protein coding reads ($< 10^{-5}$ %) suggests

62    the physiological state of most microbes was not significantly altered during sample retrieval and

63    storage.

64         Dissimilatory sulfate reduction may represent a major form of microbial metabolism and

65    energy production in the sub-seafloor[1, 2, 18] and is indicated by porewater sulfate concentrations at

66    Site 1229[1] (Fig. 1).  Representation of *Dsr* transcripts was highest in sediment with sulfate

67    profiles suggestive of biogenic sulfate reduction (Fig. 1) and supports biogeochemical evidence

68    for sulfate reduction at this site[1, 4].  Surprisingly, transcripts coding for dissimilatory nitrate

69    reductases (*Nar*) were represented throughout the sediment column, despite no measureable

70    nitrate.  The origin of nitrate as a substrate in this sediment is unknown, but could potentially be

71    produced as a by-product of anaerobic ammonium oxidation.  Once produced, nitrate would

72    likely not accumulate to measurable concentrations given the higher free energy yield of nitrate

73    as electron acceptor compared to the dominant electron acceptors in this environment, sulfate

74    and iron.  Nitrate reduction appears to be performed predominantly by Alphaproteobacteria and

75    Betatproteobacteria at most depths (Fig. 1) and the resulting nitrite is likely reduced by Fungi,

76    Gammaproteobacteria, and Firmicutes (Fig. S3). In contrast, Deltaproteobacteria and Firmicutes

77    are the dominant groups expressing *Dsr* transcripts at 5 and 30 mbsf, and Gammaproteobacteria

78    were the only group with detectable *Dsr* transcripts at deeper depths (Fig. 1).  Expression of *Dsr*

79    transcripts from a methanogenic lineage (Fig 1) in the deep biosphere supports the evidence that

80    anaerobic oxidation of methane (AOM) may not be an obligate syntrophic process[19].

81         Gene expression from methanogenic lineages was found, including Methanosarcinales,

82    which contain the anaerobic methane-oxidizing group ANME-2[20] (Fig. S4).  However, we did

83    not detect any transcripts coding for methyl-coenzyme reductase M (*mcrA*), arguably the best

84    diagnostic enzyme for AOM and methanogenesis.  This could be explained by low levels of

85    archaeal mRNA expression and a masking of *mcrA* gene expression by archaeal housekeeping

86    genes.  As a DNA-based study detected *mcrA* genes from this site[21], this explanation seems

87    likely.  Consistent with DNA-based observations from other sites[20], gene expression from

88    methanogens was detected in the sulfate reduction zones (Fig. S4).  Methylotrophic

89    methanogenesis has been documented in shallow sediment sulfate reduction zones that contain

90    non-competitive substrates such as trimethylamine[22, 23].  Our detection of trimethylamine

91    methyltransferase transcripts from Methanosarcinales and Methanobacteriales (Fig S4) suggests

92    that this process occurs in the deep subseafloor and support previous suggestions of biogenic

93    methane at this site[1]. While Crenarchaeota have been suggested to be dominant at this site[6, 13, 15],

94    they are a minority contribution to the metatranscriptome (Fig S1) even with incorporating new,

95    partially completed, single cell genomes from shallow sediments[24] (Table S2). One explanation

96    is that Crenarchaeota may have relatively low levels of mRNA expression in the deep biosphere.

97          A model suggests turnover of microbial biomass in this environment[4], but at the

98    extremely low metabolic rates proposed it is unknown whether growth yield leads to cell division

99    or to biomass turnover without division[4, 25]. Representation of transcripts involved in cell

100   division (Table S3) increases at sulfate methane transition zones (SMTZs) where cell abundances

101   increase by an order of magnitude ($p = 0.03$, Figs 1, S5). Our data suggest that the portion of the

102   vegetative population that is actively dividing is largest in the SMTZs, and that observed peaks

103   in cell counts at SMTZs are a result of *in situ* cell division. Cell division transcripts from all three

104   domains of life strongly indicate a diversity of actively dividing cells in deeply buried sediment,

105   including Fungi. The dominance of transcripts involved in amino acid metabolism (Fig. 2) and

106   coding for peptidases (Fig. S6) support a recent model of amino acid turnover in the deep

107   biosphere[4] and evidence for peptidase activity in shallow marine sediments[24].

108          Microbial motility has been proposed for deep sediment[5], however, calculations of mean

109   metabolic rates suggest that flagellar motility may not be possible in the deep biosphere[26]. We

110   detected expressed ORFs involved in in flagellar, gliding, and twitching based motility (Table

111   S3) up to 159 mbsf (Fig. 3) and the abundance of these categories decreases with decreasing

112   sediment porosity ($p = 0.01$, Fig. 3), indicating that microbial motility is related to the space

113   available for movement. The evidence for motility presented here implies that metabolic rates are

114   not equal across all cells in the deep biosphere and that some cells may be significantly more

115   metabolically active than others.  The offset in taxonomic assignment of motility reads (Fig. S7)

116   relative to total mRNA reads (Fig. S1) is suggestive of such differences.

117        DNA repair may represent a mechanism by which microbes in the deep biosphere are

118   able to cope with the slow degradation of DNA over geological timescales due to spontaneous

119   chemical or radiolytic reactions in the subseafloor[25, 26].    The representation of DNA repair

120   transcripts involved in nucleotide excision and mismatch repair (Table S3) increases linearly

121   with sediment depth (p = 0.004, Fig. 3).  This suggests DNA repair is a survival mechanism for

122   microbial populations in ancient sediment and supports the suggestion that dormancy may not be

123   a feasible survival strategy for the deep biosphere, because it does not completely arrest the slow

124   degradation of DNA[25, 26].

125        Fungal metabolic transcripts confirm previous suggestions of living fungi in the

126   subseafloor[9, 13, 27], and are the first direct evidence for active fungal metabolism in the deep

127   biosphere.   Five percent of transcripts involved in carbohydrate, amino acid, and lipid

128   metabolism were assigned to Fungi, suggesting that Fungi play an overlooked role in organic

129   carbon turnover in sub-seafloor sediment (Fig. 2).   Fungal expression of transcripts coding for

130   hydrolases involved in protein, carbohydrate, and lipid degradation (Fig. S6) indicates they

131   degrade a variety of organic carbon substrates in deep subseafloor sediment.

132        Microbial expression of antibiotic defense mechanisms, polyketide synthases, and non-

133   ribosomal proteins was detected (Fig. S8).  Polyketide synthases and non-ribosomal proteins are

134   involved in the biosynthesis of natural products (*e.g.* antibiotics, immunosuppressants,

135   antifungals) of clinical and industrial importance. These findings warrant further investigation

136   into potentially novel secondary metabolites produced by the deep biosphere, and support the

137     hypothesis that the deep biosphere may represent a "seed bank" of novel biotechnological and

138     biomedical innovation[28].

139          A comparison of the metatranscriptomic data to existing metagenomic datasets from this

140     site[13, 29] reveals an increased representation of key metabolic and cell cycle functional genes in

141     the metatranscriptome including those involved in DNA repair, DNA replication, transcription,

142     amino acid biosynthesis, and lipid biosynthesis (Fig. 4). The significant difference between

143     mRNA and metagenome samples with similar biogeochemical profiles (upper SMTZ and 50

144     mbsf: 5/12 samples) suggests these to be some of the more active processes.  Although not a

145     primary group in the overall annotations, activity of Archaea in the deep biosphere is highlighted

146     by archaeal ATPase and DNA polymerase transcripts that are overrepresented in the

147     metatranscriptomes relative to metagenomes ($p < 0.0005$).   An analysis of similarity test

148     (ANOSIM) indicates that the gene expression approach captures a significantly different picture

149     of microbial activities compared to DNA based data ($p=0.001$, Fig. S9).   As deep biosphere

150     studies move forward, joint investigation of both nucleic acid pools are needed for full

151     interpretation of metabolic activity and potential.

152          Metatranscriptomic analysis enables a refined view of deep biosphere activities.

153     Microbial activity in deeply buried marine sediment is important because the collective activities

154     of subsurface microbiota directly influences whether elements such as carbon are sequestered for

155     millions of years in sediment or returned to the ocean, impacting food webs and climate[12]. Our

156     data suggest the latter is mediated by diverse metabolic activities across all three domains of life

157     in the sub-seafloor.

158

159     **METHODS SUMMARY**

160    **Sample collection.**  Subsurface sediment samples from the continental shelf of Peru, Ocean

161    Drilling Program (ODP) Site 1229D (77º 57.4590' W, 10º 58.5721' S), were obtained during

162    ODP Leg 201 on March 6[th], 2002.

163    **RNA extraction, purification, and amplification.**  RNA was extracted from 25 g subseafloor

164    sediment according to the protocol described by Orsi *et al* [26] using the FastRNA Pro Soil-Direct

165    Kit[®] (MP Biomedicals, Solon, OH).  In addition to the manufacturers instructions, physical and

166    chemical adjustments to the sample were used to increase RNA yield and purity (see

167    Supplemental Methods).  DNA was removed using the Turbo DNA-free[®] kit (Life Technologies,

168    Grand Island, NY), increasing the incubation time to 1 hour to ensure rigorous DNA removal.

169    The MEGA-Clear[®] RNA Purification Kit (Life Technologies, Grand Island, NY) was used to

170    further purify the RNA.    Removal of contaminating DNA in RNA extracts was confirmed by

171    the absence of visible amplification of SSU rRNA genes after 35 cycles of PCR using the RNA

172    extracts as template.  Total RNA was used as template for cDNA amplification using the Ovation

173    RNA-Seq v2 System[®] (NuGEN technologies).

174    **Bioinformatic    analyses.**        Quality    control    was    performed    using    FastQC

175    (http://www.bioinformatics.babraham.ac.uk/projects/fastqc/).  Read assembly and mapping were

176    performed in CLC Genomics Workbench 5.0 (CLC Bio Inc.). The Rapid Analysis of Multiple

177    Metagenomes with a Clustering and Annotation Pipeline (RAMMCAP) available through

178    CAMERA (http://camera.calit2.net/) was used to annotate contigs against COG and Pfam

179    databases. Heatmaps and statistical tests were performed in R (http://www.r-project.org/) using

180    the vegan (http://vegan.r-forge.r-project.org/) and matR (metagenomics.anl.gov) packages.

181    Taxonomic assignments of contigs were performed using PhymmBL[30] with addition of fungal

182    genomes available in the NCBI RefSeq and JGI databases and four partial single cell archaeal

183    genomes from a shallow sediment site[24].

184

185

186    **REFERENCES**

187    1.    D'Hondt, S. *et al.* Distributions of microbial activities in deep subseafloor sediments.

188          *Science* **306**, 2216-2221 (2004).

189    2.    Schrenk, M. O., Huber, J. A., & Edwards, K. J.  Microbial provinces in the subseafloor.

190          *Ann Rev Mar Sci* **2**, 279-304 (2010).

191    3.    Jorgensen, B. B. & D'Hondt, S. A starving majority deep beneath the seafloor. *Science*

192          **314**, 932-934 (2006).

193    4.    Lomstein, B. A., Langerhuus, A. T., D'Hondt, S., Jorgensen, B. B., & Spivack, A. J.

194          Endospore abundance, microbial growth and necromass turnover in deep sub-seafloor

195          sediment. *Nature* **484**, 101-104 (2012).

196    5.    Parkes, J., Cragg, B., & Wellsbury, P. Recent studies on bacterial populations and

197          processes in subseafloor sediments: A review. *Hydrogeol J* **8**, 11-28 (2000).

198    6.    Biddle, J. F. *et al.* Heterotrophic Archaea dominate sedimentary subsurface ecosystems

199          off Peru. *Proc Natl Acad Sci USA* **103**, 3846-3851 (2006).

200    7.    Kallmeyer, J., Pockalny, R., Adhikari, R., Smith, D. C., & D'Hondt, S.  Global

201          distributions of microbial abundance and biomass in subseafloor sediment. *Proc Natl*

202          *Acad Sci U S A* **109**, 16213-16216  (2012).

203    8.    Whitman, W. B., Coleman, D. C., & Wiebe, W. J.  Prokaryotes: The unseen majority.

204          *Proc Natl Acad Sci U S A* **95**, 6578-6583 (1998).

205    9.    Biddle, J. F., House, C. H., & J. E. Brenchley.  Microbial stratification in deeply buried
206         marine sediment reflects changes in sulfate/methane profiles. *Geobiology* **3**, 287-295
207         (2005).

208    10.   D'Hondt, S. *et al.* Subseafloor sedimentary life in the South Pacific Gyre.  *Proc Natl*
209         *Acad Sci U S A* **106**, 11651-11656 (2009).

210    11.   D'Hondt, S., Rutherford, S., & Spivack, A. J.  Metabolic activity of subsurface life in
211         deep-sea sediments. *Science* **295**, 2067-2070 (2002).

212    12.   Hinrichs, K. U. & Inagaki, F.  Downsizing the deep biosphere.  *Science* **338**, 204-205
213         (2012).

214    13.   Biddle, J. F., Fitz-Gibbon, S., Schuster, S. C., Brenchley, J. E., & House, C. H.
215         Metagenomic signatures of the Peru Margin subseafloor biosphere show a genetically
216         distinct environment.  *Proc Natl Acad Sci USA* **105**, 10583-10588 (2008).

217    14.   Teske, A. In: *Proceedings of the Ocean Drilling Program, Volume 201, Scientific Results.*
218         B. B. Jørgensen *et al.*, Eds. (ODP, College Station, TX, 2006), Chapter 2, pp. 1–19
219         (http://www-odp.tamu.edu/publications/201_SR/120/120.htm).

220    15.   Lipp, J. *et al*.  Significant contribution of Archaea to extant biomass in marine subsurface
221         sediments. *Nature* **454**, 991-994 (2008).

222    16.   Moran, M. A. *et al*.  Sizing up metatranscriptomics. *ISME J* **7**, 237-243 (2012).

223    17.   Gao, H. *et al*. Global transcriptome analysis of the heat schock response of Shewanella
224         oneidensis.  *J Bacteriol* **22**:7766-7803 (2004).

225    18.   Jørgensen, B. B., D'Hondt, S., & Miller, D. J.  In: *Proceedings of the Ocean Drilling*
226         *Program, Volume 201, Scientific Results,* (ODP, College Station, TX, 2006), pp. 1–45
227         (www-odp.tamu.edu/publications/201_SR/201sr.htm).

228  19.  Milucka, J. *et al*.  Zero valent sulphur is a key intermediate in marine methane oxidation.

229       *Nature* **491**, 541-546 (2012).

230  20.  Lever, M. Functional gene surveys from ocean drilling expeditions - a review and

231       perspective.  *FEMS Microbiol Eco* **84**, 1-23 (2013).

232  21.  Webster, G., Parkes, R. J., Cragg, B. A., Newberry, C.J., Weightman, A. J., Fry, J. C.

233       Prokaryotic community composition and biogeochemical processes in deep subseafloor

234       sediments from the Peru Margin.  *FEMS Microbiol Ecol*  **58**, 65-85 (2006).

235  22.  Oremland, R. S. & Polcin, S. Methanogenesis and Sulfate Reduction: Competitive and

236       noncompetitive substrates in estuarine sediments. *Appl  Environ Microbiol* **44**, 1270-1276

237       (1982).

238  23.  Valentine, D. L. Emerging topics in marine methane biogeochemistry.  *Annu Rev Mar Sci*

239       **3**:147-171 (2011).

240  24.  Lloyd, K. *et al*.  Predominant archaea in marine sediments degrade detrital proteins.

241       *Nature* **496**, 215-218 (2013).

242  25.  Jorgensen, B. B. Deep subseafloor microbial cells on physiological standby.  *Proc Natl*

243       *Acad Sci USA* **108**, 18193-18194 (2011).

244  26.  Hoehler, T. M & Jorgensen, B. B. Microbial life under extreme energy limitation.  *Nat*

245       *Rev Microbiol* **11**, 83-94 (2013).

246  27.  Orsi, W., Biddle, J., Edgcomb.  Deep sequencing of subseafloor eukaryotic rRNA reveals

247       active Fungi across multiple subsurface provinces. *PLoS ONE* **8**, e56335 (2013).

248  28.  Parkes, R. J. & Wellsbury P. in  *Microbial Diversity and Bioprospecting.*  (ed Bull,  A.T.)

249       120-129  (ASM Press, 2004).

250  29.  Martino, A. J. *et al*. Novel degenerate PCR method for whole-genome amplification

251 applied to Peru Margin (ODP Leg 201) subsurface samples. *Front Microbiol* **3**, 17

252 (2012).

253 30. Brady, A. & Salzberg, S. L. Phymm and PhymmBL: Metagenomic phylogenetic

254 classification with interpolated markov models. *Nature Methods* **6**, 673-676 (2009).

255 -

256 **Supplementary Information** is linked to the online version of the paper at

257 www.nature.com/nature.

266 **Author Contributions** W.O. performed experiments, analyzed data, and wrote the paper; W.O.,

267 J.B., and V.E. designed experiments and developed ideas. W.O. and G.C. developed analytical

268 tools. All authors participated in data interpretation and provided editorial comments on the

269 manuscript.

270 **Author Information** Data has been deposited in the NCBI Short Read Archive under accession

271 number SRA058813 and in MG RAST (metagenomics.anl.gov) under accession numbers

272 4515478.3, 4515477.3, 4515476.3, 4510337.3, 4510336.3, and 4510335.3. Reprints and

273 permission information is available at www.nature.com/reprints. The authors declare no

274  competing interests.  Correspondence and requests for materials should be addressed to W.O.

275  (william.orsi@gmail.com).

276

277

278
279  **Figure and Table legends**

280

281  **Figure 1:** Biogeochemical and gene expression profiles of the deep biosphere from Peru Margin

282  sediments, IODP site 1229D.     **(a)** Cell abundance, sulfate concentrations, and methane

283  concentrations, dotted lines indicate the SMTZs.  Values were taken from the Ocean Drilling

284  Program Janus Database (http://www-odp.tamu.edu/database/).  **(b)** Proportion of cell division

285  transcripts within the cluster of orthologous genes (COG) class D (cell cycle control/cell

286  division/chromosome partitioning, n = 30.22 million reads), see Table S3 for description of cell

287  division proteins. The proportion of **(c)** *Dsr* and **(d)** *Nar* transcripts relative to total transcripts

288  involved in energy production (COG class C, n = 92.33 million reads).  See Figure S2 for

289  number of sequences and ORFs used in each comparison, and *E*-values for hits in the COG

290  database.

291

292  **Figure 2:** The proportion of reads mapping to ORFs assigned to amino acid, lipid, and

293  carbohydrate metabolism (eleven most dominant taxa shown).  Note the relative abundance of

294  amino acid metabolism (both anabolic and catabolic) relative to lipid and carbohydrate

295  metabolism across all depths.  See Figure S2 for the number of sequences and ORFs used in each

296  comparison, and *E*-values for hits in the COG database.

297

298 **Figure 3:** Transcripts involved in cell motility and DNA repair. **(a)** The percentage of reads

299 mapping to ORFs coding for proteins involved in different modes of cellular motility (see Table

300 S3 for descriptions). **(b)** A correlation of cell motility transcripts versus sediment porosity ($R^2 =$

301 0.8, p = 0.01) and 95% prediction interval (red dotted lines). **(c)** The percentage of reads

302 mapping to ORFs involved in DNA repair (only eleven most dominant taxa are shown, see Table

303 S3 for descriptions). **(d)** A correlation of DNA repair transcripts versus sediment depth ($R^2 = 0.9$,

304 p = 0.004) and 95% prediction interval (red dotted lines). See Figure S2 for the number of

305 sequences and ORFs used in each comparison and *E*-values for ORF hits in COG database.

306

307 **Figure 4:** A comparison of gene expression data to existing metagenomic studies[13, 29] from

308 IODP site 1229. Functional genes significantly (Kruskal-Wallis test, p < 0.0005)

309 overrepresented in the metatranscriptomic samples relative to metagenomic data include DNA

310 repair and replication transcripts, RNA polymerase, and archaeal ATPase and DNA polymerase

311 transcripts. The dendrogram represents a UPGMA hierarchical clustering analysis (Manhattan

312 distance) of significantly overrepresented mRNA transcripts, note the complete separation of

313 mRNA samples from DNA samples.

314

315 **Methods**

316 **Sample collection and storage** Subsurface sediment samples from the continental shelf of Peru,

317 Ocean Drilling Program (ODP) Site 1229D (77º 57.4590' W, 10º 58.5721' S), were obtained

318 during ODP Leg 201 on March 6[th], 2002. Careful precautions were taken during sampling to

319 avoid contamination during the sampling process. For IODP cores, contamination tests were

320 performed using Perfluorocarbon tracers and fluorescent microspheres (for more information see

321    http://www-odp.tamu.edu/publications/201_IR).  Sediment samples were immediately frozen at -

322    80 ºC after sampling and stored at -80 ºC until used for mRNA extractions in this study (10 year

323    storage time at -80 ºC).

324    **RNA extraction and purification** Extraction of subseafloor RNA was performed according to

325    the protocol of Orsi *et al* [26].  To summarize, RNA was extracted from 25 grams of sediment

326    using the FastRNA Pro Soil-Direct Kit[®] (MP Biomedicals, Solon, OH).  It was necessary to scale

327    up the volume of sediment that is typically extracted with the kit (~0.5 grams) due to the low

328    biomass inherent to marine subsurface samples. All tubes, tips, and disposables used were

329    certified RNAse free and all extraction procedures were performed in a laminar flow hood to

330    reduce aerosol contamination by bacterial and fungal cells/spores. Five 15ml Lysing Matrix E[®]

331    tubes (MP Biomedicals, Solon, OH) were filled with 5 g sediment and 5 ml of Soil Lysis

332    Solution[®] (MP Biomedicals, Solon, OH).  Tubes were vortexed to suspend the sediment and Soil

333    Lysis Solution[®] was added to the tube leaving 1 ml of headspace. Tubes were then homogenized

334    for 60 seconds on the FastPrep-24 homogenizer[®] (MP Biomedicals, Solon, OH) with a setting of

335    4.5.  Contents were pooled into two 50ml tubes and centrifuged for 30 minutes at 4,000 RPM

336    (3220 x g) at room temperature (RT). Supernatants were combined in a new 50ml tube and 1/10

337    volume of 2M Sodium Acetate (pH 4.0) was added.  An equal volume of phenol-chloroform (pH

338    6.5) was added and vortexed for 30 seconds, incubated for 5 minutes at room temperature, and

339    spun at 4000 RPM (3220 x g) for 20 minutes at 4 ºC.  The aqueous phase was transferred to a

340    new 50ml tube.  Nucleic acids were precipitated by adding 2.5 and 1/10 volumes 100% ethanol

341    and 3M Sodium Acetate, respectively, and incubating overnight at -80 ºC.   The next day, tubes

342    were spun at 4000 RPM (3220 x g) for 60 minutes at 4 ºC and the supernatant removed.   Pellets

343    were washed with 70% ethanol, spun for 15 minutes at 4 ºC, and air-dried.  Dried pellets were

344  resuspended with 0.25 ml RNAse-free sterile water and combined into a new 1.5ml tube. 1/10

345  volume of 2M Sodium Acetate (pH 4.0) and an equal volume of phenol:chloroform (pH 6.5)

346  were added, vortexed for 1 minute, and incubated for 5 minutes at RT.  This was necessary to

347  remove residual organic material (*i.e.* humic acids) resulting from the rather large

348  pellet/precipitate.  After centrifuging at 14,000 RPM (20,817 x g) for 10 minutes at 4 °C, the top

349  phase was removed into a new 1.5ml tube. 0.7 volumes of 100% isopropanol was added and

350  incubated for 1 hour at -20 °C (to precipitate nucleic acids). Tubes were then centrifuged for 20

351  minutes at 14,000 (20,817 x g) RPM at 4 °C and the supernatant removed.  Pellets were washed

352  with 70% ethanol and centrifuged at 14,000 RPM (20,817 x g) for 5 minutes at 4 °C.  After

353  removing ethanol and air-drying, pellets were resuspended in 0.2 ml of RNAse free sterile water.

354  DNA was removed using the Turbo DNA-free® kit (Life Technologies, Grand Island, NY),

355  increasing the incubation time to 1 hour to ensure rigorous DNA removal.  After this step,

356  samples were taken through the protocol supplied with the FastRNA Pro Soil-Direct kit® to the

357  end (starting at the RNA Matrix® and RNA Slurry® addition step), including the column

358  purification step to remove residual humic acids (see FastRNA Pro Soil-Direct Kit® manual).

359  Extraction blanks were performed (adding sterile water instead of sample) to ensure that

360  aerosolized contaminants did not enter sample and reagent tubes during the extraction process.

361  Absence of DNA and RNA contamination was confirmed by no visible amplification of small

362  subunit (SSU) rRNA and rRNA genes from extraction blanks after 35 cycles of PCR and RT-

363  PCR.

364         After RNA extraction, used the MEGA-Clear® RNA Purification Kit (Life Technologies,

365  Grand Island, NY) to purify the RNA.  This kit removes short RNA fragments (mostly produced

366  during the extraction protocol) and residual inhibitors (*i.e*. humics).  We followed the protocol all

367     the way through the optional precipitation/concentration step, resuspending the RNA pellet in 10

368     microliters of RNAse free sterile water.  Prior to cDNA amplification, the removal of

369     contaminating DNA in RNA extracts was confirmed by the absence of visible amplification of

370     SSU rRNA genes after 35 cycles of PCR using the RNA extracts as template.

371     **cDNA amplification and Illumina sequencing** Five microliters of purified RNA was used as

372     template for whole cDNA amplification using the Ovation RNA-Seq v2 System[®] (NuGEN

373     technologies,  http://www.nugeninc.com/nugen/index.cfm/products/cs/ngs/rna-seq-v2/).  We

374     followed the manufacturers instructions for cDNA amplification, and the resulting quantity of

375     cDNA was checked on a Nanodrop (Thermo Scientific) and Fluorometer (Qubit 2.0, Life

376     Technologies).  Quality of the amplified cDNA was checked on a Bioanalyzer (Agilent

377     Biotechnologies) prior to Illumina[®] sequencing.  Illumina[®] library preparation and paired-end

378     sequencing was performed at the University of Delaware Sequencing and Genotyping Center

379     (Delaware Biotechnology Institute, Newark DE).

380     **Quality control and assembly** Quality control of the dataset was performed using FastQC

381     (http://www.bioinformatics.babraham.ac.uk/projects/fastqc/), with a quality score cutoff of 28.

382     Approximately 1 billion paired-end reads that passed quality control were imported into CLC

383     Genomics Workbench 5.0[®] (CLC Bio Inc.) and assembled using the paired-end Illumina

384     assembler.  Contigs were assembled over a range of kmer sizes (20, 50, 60, 64) with a minimum

385     contig size cutoff of 300 nucleotides.  The kmer size of 50 resulted in the highest number of

386     contigs and these contigs were chosen for use in downstream analyses.  To reduce the formation

387     of chimeric assemblies, we used a paired-end sequencing approach and performed assemblies

388     without scaffolding.  Reads were mapped onto the contigs using the read mapping option in CLC

389     Genomics Workbench to retain information on relative abundance of contigs.

390 **Functional annotation of contigs** Contigs were submitted to CAMERA (Community

391 Cyberinfrastructure for Advanced Microbial Ecology Research and Analysis,

392 http://camera.calit2.net/) and assigned to clusters of orthologous gene (COG) families, gene

393 ontologies (GO), and protein families (Pfam), using the Rapid Analysis of Multiple

394 Metagenomes with a Clustering and Annotation Pipeline (RAMMCAP) using the 6 reading

395 frame translation option for open reading frame (ORF) prediction and BLASTn for rRNA

396 identifications. The cutoff criterion *E*-value of $10^{-5}$ was used for BLASTx searches against the

397 COG, Pfam, and TIGRfam databases. For identification of bacterial and archaeal ORFs, the

398 RAMMCAP analyses were performed using the bacterial and archaeal genetic code (-t 11 in

399 advanced options). For identification of fungal ORFs, additional RAMMCAP analyses were

400 performed using the standard genetic code for eukaryotes and the alternative yeast genetic code

401 (-t 1 and –t 12 in advanced options). For comparative analysis of the metatranscriptomes to

402 existing metagenomes from ODP Site 1229D we submitted the metatranscriptomes to MG-

403 RAST (metagenomics.anl.gov), which were annotated according to the standard bioinformatics

404 pipeline (http://blog.metagenomics.anl.gov/mg-rast-for-the-impatient-readme-1st/).

405 **Taxonomic annotation of contigs** Contigs were assigned to high-level taxonomic groups (Class

406 level and above) using PhymmBL[30]. In addition to the default interpolated markov model

407 (IMM) database (that contains only bacterial and archaeal genomes), all fungal genomes

408 available in the NCBI RefSeq database and JGI database, along with several representative

409 protistan and plant genomes were added to the IMM database (using the customGenomicData.pl

410 script available with the PhymmBL download) to facilitate identification of eukaryotic contigs.

411 Cutoffs for annotation accuracy were chosen based on the default recommendations. Taxonomic

412 identifications of contigs made using PhymmBL[30] were integrated with the functional

413    annotations from CAMERA (BLASTx searches against the COG database and HMMer searches

414    against Pfam database) and the read mapping information from assemblies. This was done using

415    several custom PERL scripts that are available from the authors upon request.

416    **Statistical analyses** Analyses of overexpression of expressed genes relative to metagenome

417    samples was performed using the R statistical package (http://www.r-project.org/), with the MG-

418    RAST matR library (metagenomics.anl.gov).   To maintain abundance information, assembled

419    contig sequences from each sample were uploaded to MG RAST with the read mapping

420    abundance added to the fasta headers as specified on the MG RAST website. Statistically

421    significant differences in overexpressed functional genes relative to genes detected in

422    metagenomes were determined by a Kruskal-Wallis test with a p value cutoff of 0.0005.   All

423    rRNA reads were removed from both metagenomic and metatranscriptomic datasets prior to

424    comparison.  Data were normalized in MG RAST with a log based transformation:

425    $$Y_{s,i} = \log_2 (X_{s,i} + 1)$$

426    Where $X_{s,i}$ represents an abundance measure ($i$) in sample ($s$).   Log transformed counts from each

427    sample were then standardized (data centering) according to the following equation:

428    $$Z_{s,i} = [(Y_{s,i} - Y_s)/ \sigma_s)$$

429    Where $Z_{s,i}$ is the standardized abundance of an individual measure $Y_{s,i}$ (log transformed from

430    previous equation).  From each log transformed measure of ($i$) in sample ($s$), the mean of all

431    transformed values ($Y_s$) is subtracted and the difference is divided by the standard deviation ($\sigma_s$)

432    of all log-transformed values for the given sample. After log transformation and standardization,

433    the values for the functional categories within each sample were scaled from 0 (minimum value

434    of all samples) to 1 (maximum value of all samples), which is a uniform scaling that does not

435    affect the relative differences of values within a single sample or between 2 or more samples.

436    This procedure places the value of functional categories (*i.e*. COG categories) from each sample

437    on a scale from 0 to 1 and was used to produce figures (*i.e*. heatmaps or principal component

438    analysis) where the abundance range is on a scale from 0 to 1 (*i.e.* Figure 4).  Normalized data

439    that passed the Kruskal-Wallis test (p value cutoff criterion 0.0005) were used as input for

440    heatmap presentation, UPGMA hierarchical clustering, and principal component analysis in R,

441    using the matR package (metagenomics.anl.gov).  Analysis of similarity (ANOSIM) analyses

442    were performed on the normalized data in R, using the vegan package (http://vegan.r-forge.r-

443    project.org/).  ANOSIM was performed with 999 permutations using a Bray-Curtis distance

444    metric.  Correlations of gene expression data with geochemical and geophysical metadata were

445    performed using the lm and predict commands in R, which are used to fit linear models to

446    relationships between two different variables.  The data for these analyses were normalized in

447    the same fashion as Figures 1, 2, 3, S3, S4, S5, S6 and S8 (*i.e*. the relative abundance, per sample,

448    of transcripts mapping to ORFs that were annotated to each functional COG category).

**a.**

Depth (mbsf)

Cell concentration (log₁₀ cm³) — red
SO₄⁻² (mM) — yellow
CH₄ (mM) — blue

$Cell\ concentration\ (log_{10}\ cm^3)$
$SO_4^{-2}\ (mM)$
$CH_4\ (mM)$

**b.** Cell division:COG class D

**c.** *Dsr*:Energy production

**d.** *Nar*:Energy production

Legend:
- Crenarchaeota
- Actinobacteria
- Synergistes
- Alphaproteobacteria
- Betaproteobacteria
- Fungi
- Methanomicrobiales
- Deltaproteobacteria
- Gammaproteobacteria
- Euryarchaeota
- Firmicutes

Amino acid transport and metabolism (COG class E)
- 5 mbsf
- 30 mbsf
- 50 mbsf
- 70 mbsf
- 91 mbsf
- 159 mbsf

Lipid transport and metabolism (COG class I)
- 5 mbsf
- 30 mbsf
- 50 mbsf
- 70 mbsf
- 91 mbsf
- 159 mbsf

Carbohydrate transport and metabolism (COG class G)
- 5 mbsf
- 30 mbsf
- 50 mbsf
- 70 mbsf
- 91 mbsf
- 159 mbsf

Percent of reads annotated against the COG database

Legend:
- Actinobacteria
- Bacteroidetes
- Firmicutes
- Fungi
- Alphaproteobacteria
- Betaproteobacteria
- Deltaproteobacteria
- Gammaproteobacteria
- Crenarchaeota
- Synergistes
- Euryarchaeota

**a.**

Sediment depth

- 5 mbsf
- 30 mbsf
- 50 mbsf
- 70 mbsf
- 91 mbsf
- 159 mbsf

Percent of reads annotated against the COG database

| ■ Flagellar | ■ Gliding | ■ Twitching |

**b.**

Sediment porosity (%)

Percent of annotated reads

**c.**

Sediment depth

- 5 mbsf
- 30 mbsf
- 50 mbsf
- 70 mbsf
- 91 mbsf
- 159 mbsf

Percent of reads annotated against the COG database

■ Actinobacteria    ■ Bacteroidetes    ■ Firmicutes
■ Fungi    ■ Alphaproteobacteria    ■ Betaproteobacteria
■ Deltaproteobacteria    ■ Gammaproteobacteria    ■ Crenarchaeota
■ Euryarchaeota    ■ Nitrospirae

**d.**

Sediment depth (mbsf)

Percent of annotated reads

Read abundance

0 ■■■■■■ 1

| | |
|---|---|
| DNA segregation *FtsK* (COG1674) | |
| DNA recombinase *XerD* (COG4974) | |
| RNA Polymerase (COG0568) | |
| ATP synthase (COG0056) | |
| Pyrophosphohydrolases/synthetases (COG0317) | |
| Riboflavin biosynthesis (COG0108) | |
| Ornithine aminotransferase (COG4992) | |
| Glutamate dehydrogenase (COG0334) | |
| DNA repair (COG1197) | |
| Adenylosuccinate synthase (COG0104) | |
| DNA/RNA helicases (COG0553) | |
| ATP synthase (COG0055) | |
| Malate dehydrogenase (COG0281) | |
| Urocanate hydratase (COG2987) | |
| Pyruvate/oxaloacetate carboxyltransferase (COG5016) | |
| Inorganic pyrophosphatase (COG3808) | |
| Flavoproteins (COG0426) | |
| Archaeal ATPase (COG1155) | |
| Archaeal ATPase (COG1156) | |
| Lipid biosynthesis (COG1260) | |
| DNA exinuclease (COG0178) | |
| Protein transport (COG4608) | |
| DNA gyrase (COG0187) | |
| Amino acid biosynthesis (COG0458) | |
| DNA excision repair (COG0556) | |
| RNA polymerase (COG0086) | |
| RNA polymerase (COG0085) | |
| NADH:ubiquinone oxidoreductase (COG1894) | |
| Ferredoxin *NapF* (COG1145) | |
| Translation initiation factor (COG5257) | |
| Archaeal DNA polymerase (COG1933) | |

| 5 | 30 | 50 | 159 | 91 | 70 | | 1 | 16 | 32 | 32 | 50 | 1 | Depth (mbsf) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | mRNA | | | | | | | DNA | | | | |