1    **Metagenomic Characteristics of Bacterial Response to Petroleum Hydrocarbon**

2    **Contamination in Diverse Environments as Revealed by Functional Taxonomic Strategies**

3    Arghya Mukherjee[1], Bobby Chettri[2], James S. Langpoklakpam[2], Pijush Basak[3@], Aravind

4    Prasad[4‡], Ashis K. Mukherjee[5], Maitree Bhattacharyya[3], Arvind K. Singh[2] and Dhrubajyoti

5    Chattopadhyay[1#*]

6    [1] Department of Biotechnology, University of Calcutta, Kolkata, West Bengal, India

7    [2] Department of Biochemistry, North-Eastern Hill University, Shillong, India.

8    [3] Department of Biochemistry, University of Calcutta, Kolkata, West Bengal, India.

9    [4] Bioinformatics Resources and Applications Facility, Centre for Development of Advance

10   Computing, Pune, India.

11   [5] Department of Molecular Biology and Biotechnology, Tezpur University, Tezpur, India

12   @ Present address: Jagadis Bose National Science Talent Search, Kolkata, West Bengal, India

13   ‡ Present address: Institute of Molecular and Cell Biology, A*STAR, Singapore, Singapore

14   * Present address: Department of Biotechnology, Amity University, Rajarhat, New Town,

15   Kolkata, West Bengal, India

16   # Correspondence and request for materials should be addressed to Dhrubajyoti

17   Chattopadhyay (email: dchattopadhyay@kol.amity.edu).

18   **Abstract**

19   Microbial remediation of oil polluted habitats remains one of the foremost methods for

20   restoration of petroleum hydrocarbon contaminated environments. The development of

21   effective bioremediation strategies however, require an extensive understanding of the

22  resident microbiome of these habitats. Recent developments such as high-throughput

23  sequencing has greatly facilitated the advancement of microbial ecological studies in oil

24  polluted habitats. However, effective interpretation of biological characteristics from these

25  large datasets remains a considerable challenge. In this study, we have implemented

26  recently developed bioinformatic tools for analyzing 65 publicly available 16S rRNA datasets

27  from 12 diverse hydrocarbon polluted habitats to decipher metagenomic characteristics of

28  bacterial communities of the same. We have comprehensively described phylogenetic and

29  functional compositions of these habitats and additionally inferred a multitude of

30  metagenomic features including 255 taxa and 414 functional modules which can be used as

31  biomarkers for effective distinction between the 12 oil polluted sites. We have identified

32  essential metabolic signatures and also showed that significantly over-represented taxa

33  often contribute to either or both, hydrocarbon degradation and additional important

34  functions. Our findings reveal significant differences between hydrocarbon contaminated

35  sites and establishes the importance of endemic factors in addition to petroleum

36  hydrocarbons as driving factors for sculpting hydrocarbon contaminated bacteriomes.

37

38

39

40

41

42

43

**Introduction**

Anthropogenic activities and agents leading to contamination of the environment is one of the major issues that developing and developed industrial societies face today. Petroleum hydrocarbons are the most widespread of these anthropogenic agents and frequently contaminate aquatic and terrestrial ecosystems through releases of hydrocarbon during production, operational use, and transportation. The development, effectiveness and availability of technologies and strategies pose a significant challenge for the remediation, rehabilitation and restoration of these contaminated environments. Many of the technologies developed and in use for the restoration of oil contaminated environments exploit the potential of biological systems, in particular microbial systems, to use these toxic compounds as substrates for growth. Hence, much of the research conducted on bioremediation has concentrated on the capabilities of a single or couple of microbes exhibiting robust and effective growth using petroleum hydrocarbons. However, in the environment, bioremediation is often a complex process involving co-metabolism, cross-induction, inhibition and non-interaction among microbes [1-3], possibly as petroleum hydrocarbons are a mixture of organic pollutants and therefore are used differently by different microbes. These findings, along with others, established bioremediation as a process mediated by a consortium of microbes rather than a few. Thus, characterization of microbial communities of oil contaminated environments could potentially provide guidelines for effective remediation and restoration of such environments.

Until recently, it was only possible to study a handful of microorganisms of interest isolated from source materials (as blood, soil, water or air), given the restrictions of the composition of culture media which cannot reflect and mimic the dynamic nutrient fluxes of the source environment. Indeed, only 1% of microorganisms were found to be cultivable

68    using a set of media from the highly characterized soil rhizosphere [4]. The advent of high

69    throughput massively-parallel sequencing methods has however, allowed us to investigate

70    the entire complement of organisms inhabiting a certain environment. These next-

71    generation sequencing methods (NGS) include a variety of methods to holistically study any

72    biological system such as amplicon sequencing (for variant identification and phylogenetic

73    surveys), whole genome shotgun sequencing (single organism genome and metagenomes)

74    and RNA-Seq (transcriptomes, metatranscriptomes and identification of non-regular RNAs).

75    These powerful methods have ushered in rapid advances in bioinformatics approaches

76    leading to development of software capable of handling huge amounts of data and offering

77    meaningful biological interpretations of the same. Although a technological breakthrough in

78    modern science, a number of NGS methods as metagenomic and transcriptomic/

79    metatranscriptomic sequencing are still expensive and hence, most studies on ecological

80    processes on bioremediation report marker surveys as 16S rRNA gene amplicon sequencing

81    when dealing with a large number of samples. Thus, in general, most of these studies

82    concentrated on interpretations from microbial community composition but inferred poorly

83    regarding functional and metabolic properties of the same.

84         Recently, with the implementation of the Human Microbiome Project (HMP),

85    bioinformatic advancements have been furthered through the development of powerful

86    new computational tools for effective interpretation and visualization of taxonomic and

87    functional composition of microbial communities [5,6]. These tools have obvious applications

88    for the analysis of huge amounts of microbial genomic/amplicon/transcriptomic data

89    collected from other sources such as soil, water and so on. Some particularly interesting

90    computational tools allow to explain the complex mutual interactions and heterogeneity

91   inherent in microbial communities through network-based correlation analyses [7], prediction

92   of metagenomic biomarkers [8] and prediction of metagenomes from 16S rRNA data [9].

93       It is well understood that depending on the environment, the method of

94   bioremediation will vary. However, essential information required for development of these

95   technologies include the response of microbes to petroleum hydrocarbons and their

96   dynamics with the immediate environment. Unfortunately, despite the large amount of

97   work done on microbial community composition across a myriad of oil contaminated

98   environments, mainly through 16S rRNA amplicon sequencing, no attempt has been made

99   to find differential metagenomic signatures among these studies. In the present study, we

100  have aimed to investigate the taxonomic and functional characteristics of diverse oil

101  contaminated environments using recent bioinformatics tools through an evolving pipeline

102  to process metagenomic data. Due to the current paucity of metagenomic datasets for this

103  kind of study and for the large availability of them, we used 61 publicly available 16S rRNA

104  datasets and 4 from this study as inputs for our analysis. Consequently, metagenomic level

105  characteristics of bacterial composition and metabolic potential were comprehensively

106  deduced for oil contaminated soils in north-east India along with 11 other petroleum

107  hydrocarbon contaminated habitats. We inferred an array of differentially abundant

108  taxonomic and functional features which may be used as biomarkers for successful

109  distinction of different oil contaminated habitats as well as for monitoring of bioremediation

110  efforts in the same. Additionally, we deduced important metabolic pathways for all

111  contaminated environments. Evaluation of correlation between taxa and functional

112  orthologs was also carried out along with estimation of metagenomic contributions to

113  hydrocarbon degradation to detect taxa responsible for critical functions in oil polluted

114  habitats. Furthermore, a network of bacterial interaction patterns was inferred to deduce

115 complex co-occurrence and co-exclusion relationships in these environments. We found

116 that phylogenetic and functional composition oil contaminated bacteriomes were

117 significantly different to each other and greatly influenced by immediate environmental

118 factors along with petroleum hydrocarbon contamination. Our investigation provides novel

119 and valuable insights into the differential nature of various oil polluted habitats and

120 hopefully improves upon previous understanding of these environments.

121

122 **Materials and Methods**

123 Ethics statement

124

125 No specific permits were required for the samplings carried out at the fields described in

126 Noonmati or Barhola in Assam, India. The study sites are not privately owned or protected

127 in any terms. Also, the field work did not involve any protected or endangered species.

128

129 Collection of soil samples

130

131 Oil contaminated soil samples were collected from Noonmati Oil Refinery in Guwahati and

132 from oilfields in Barhola, both in Assam, India. Soil samples were collected in triplicate from

133 both sites from the surface (0-10 cm) and beneath (20-40 cm) using sterile equipment. The

134 soil samples were transported to the laboratory on ice within 48-72 hours. In the laboratory,

135 replicates for each sample was mixed and homogenized to uniformity to form four

136 composite samples prior to isolation of soil DNA.

137

138 Extraction of total soil DNA

139

140    Soil DNA was isolated using the PowerMax Soil DNA Isolation Kit (MoBio Labs, Carlsbad, CA,

141    USA) according to the manufacturer's protocol. Integrity of isolated soil DNA was then

142    checked in a 0.8% agarose gel, while a NanoDrop 2000 Spectrophotometer (Thermo

143    Scientific, Wilmington, DE, USA) was used to obtain qualitative (260:280 and 260:230 ratios)

144    and quantitative estimates.

145

146    16S rRNA PCR Amplification and Pyrosequencing

147

148    The V1-V3 region of the bacterial 16S rRNA gene was amplified using specially designed

149    fusion primers, consisting of both template derived sequence and a specific sequence as

150    recommended by the "Sequencing Technical Bulletin No. 013-2009" (454 Life Sciences,

151    Branford, CT, USA). The Fusion forward primer (5'-3') consisted of a Roche A adapter

152    (CCATCTCATCCCTGCGTGTCTCCGAC), a key sequence (TCAG), a 10 bp Multiplex Identifier

153    (MID) sequence and a template specific sequence (GAGTTTGATCMTGGCTCAG) derived from

154    the universal 16S rRNA primer 27F [10]. The Fusion reverse primer consisted of a Roche B

155    adapter (CCTATCCCCTGTGTGCCTTGGCAGTC), a key sequence (TCAG) and a template specific

156    sequence (ATTACCGCGGCTGCTGG) derived from the universal 16S rRNA gene primer 541R

157    [11]. V1-V3 region of the 16S rRNA gene was amplified through PCR in a 25 µl reaction mix

158    containing 2.5 µl Fast Start Buffer (10X), 1 µl each of Fusion forward primer (10 µM) and

159    Fusion reverse primer (10 µM), 0.5 µl dNTPs (10 mM), 0.5 µl Fast Start Taq Polymerase

160    (5U/µl Fast Start High Fidelity PCR System, Roche) with the rest made up with water. The

161    PCR was carried out in a Veriti Thermal Cycler (Thermo Scientific, Wilmington, DE, USA)

162    under the following conditions: initial denaturation at 94°C for 3 minutes, followed by 25

163    cycles of denaturation at 94 °C for 15 secs, annealing at 58 °C for 45 secs and elongation at

164    72 °C for 1 min, with a final elongation at 72 °C for 10 mins. Thereafter samples were stored

165    at -20 °C, if required.

166

167    Sequence processing and taxonomic analysis of 16S rRNA data

168

169    Raw 16S rRNA sequencing reads were checked for quality using FastQC [12] and subsequently

170    processed using mothur [13], which included trimming of adapters, keys, MIDs and primers

171    from the raw sequences. Sequences were further filtered for quality using the following

172    non-default parameters: maxhomop = 6, maxambig = 0, maxlength = 575, minlength = 200,

173    qwindowaverage = 30, bdiffs = 1, pdiffs = 2, and tdiffs = 2. Filtered high quality sequences

174    were then aligned to the mothur implementation of the SILVA database and trimmed for

175    the alignment region. Chimeric sequences were then removed from the datasets using the

176    mothur implementation of Uchime [14]. Filtered sequences were then taxonomically classified

177    using the May 2013 release of the Greengenes database [15] and contaminating archaeal,

178    eukaryal, mitochondrial and chloroplast sequences or sequences classified as unknown were

179    removed from further analysis. mothur was further used to compute coverage, boneh index

180    (for additional 500 sequences), observed species richness, and alpha diversity metrics

181    through the estimation OTUs at 0.03 level of phylogenetic divergence. OTUs were further

182    classified using the taxonomy file generated in the steps before. Taxonomically classified

183    OTUs were converted to number of sequences and  visualized as circular cladograms using

184    the standalone graphical tool GraPhlan v0.95 [16].

185

186    Collection and quality filtering of 16S rRNA datasets from oil contaminated environments

8

187

188    Sixty-one 16S rRNA datasets on oil degradation studies from 11 different environments

189    collected from publicly available resources along with four samples from this study were

190    used for the present study (Table 1, Supplementary Table S1). These included four datasets

191    representing upper soil layers of the Tundra biome (Tu), four from subsurface layers of the

192    Tundra biome (Tb), four from the permafrost layers of the Tundra (Tp), nine from surface

193    soil of Chinese oil refineries (C), twelve representing different regions of the arctic biome

194    (A), four from surface soils of Indian oil refineries (I), three from mangroves (M), seven from

195    surficial marine sediments (DWH), seven from oil sands cores (OSC), four from surface

196    waters of oil sands tailings ponds (OSTPu), three from oil sands tailings pond waters at

197    median depth (OSTPm) and four from deep oil sands tailings pond waters (OSTPd). We

198    deliberately kept the taiga and OSTP samples separate even though we expected high

199    amounts of similarity between them in certain aspects when compared to other samples,

200    due to evidence of ample distinctive characteristics in the said samples in their parent

201    studies [17,18]. All the 16S rRNA datasets used can be downloaded through the list of accession

202    numbers provided in Supplementary Table S1.  All datasets used in the study presented,

203    were sequenced in either Roche 454, Illumina or ABI Ion Torrent platforms. The 16S rRNA

204    datasets are described in greater detail in Table 1. The downloaded 16S rRNA datasets were

205    checked for quality using FastQC and filtered for high quality sequences in mothur using the

206    following criteria: minimum sequence length of 100 bp, sequences trimmed when average

207    quality drops below 20 in a sliding window of 15 bp, and a maximum of 2 mismatches in the

208    barcode-key-template region of the reads.

209

210    Analysis of microbial community structure and composition in 16S rRNA datasets

211

212     mothur was used to estimate abundances of bacterial taxa in the 16S rRNA datasets as

213     described above. Briefly, all the datasets containing high quality reads were aligned against

214     the mothur implementation of the SILVA 16S rRNA database, followed by removal of

215     chimeric sequences using the mothur implementation of Uchime. Classification of the

216     filtered sequences against the Greengenes database was carried out then, upon which

217     contaminating archaeal, chloroplast, mitochondrial, eukaryal or unknown sequences were

218     removed. Finally, OTUs were predicted from these high quality sequences. OTUs were again

219     mapped to the sequence taxonomy file generated previously in mothur to generate

220     comparative taxonomy data for the datasets. We also assessed the compositional similarity

221     between the soil samples from different sites. For doing this, we compared the pairwise

222     taxonomic abundances from each site against each other and within the datasets as well,

223     using Bray-Curtis measure for estimation of beta diversity [19]. The permutation-based

224     multivariate analysis of variance (PERMANOVA) was used to test the homogeneity of

225     taxonomic dispersion across samples along with concomitant estimation of 2D stress. The

226     resulting Bray-Curtis similarity distance matrix was used as input for ordination of the oil

227     contaminated samples through non-metric multidimensional scaling (NMDS) in PAST v3.11

228     [20].

229

230     Metagenome prediction and metabolic reconstruction of 16S rRNA datasets

231

232     Metagenomes were predicted from 16S rRNA data using PICRUSt [9]. OTU data generated in

233     mothur for all 16S rRNA datasets was used to prepare *.biom* files formatted as input for

234     PICRUSt v1.1.0 [9] with the *make.biom* script available in mothur. PICRUSt requires OTU

235    abundances mapped to Greengenes OTU IDs as input for prediction of corresponding

236    metagenomes. PICRUSt databases for 16S rRNA copy number normalization and KEGG

237    ortholog prediction were updated using publicly available information listed in Integrated

238    Microbial Genomes (IMG)[21] as on 4th April, 2016, according to the instructions (default

239    settings) provided in the Genome Prediction Tutorial for PICRUSt

240    (http://picrust.github.io/picrust/tutorials/genome_prediction.html#genome-prediction-

241    tutorial). The update involved the inclusion of 16S rRNA copy number information and KEGG

242    ortholog (KO) annotation data as per KEGG v77.1 [22] for ~34,000 bacterial and archaeal

243    genomes available in IMG. 16S rRNA copy numbers for 16S rRNA datasets were normalized

244    using the *normalize_by_copy_number.py* script. Metagenomes were predicted from the

245    copy number normalized 16S rRNA data in PICRUSt using the *predict_metagenomes.py*

246    script against the updated and PICRUSt-formatted, characterized protein functional

247    database of KEGG Orthology. Contributions of various taxa to different KOs were computed

248    with the script *metagenome_contributions.py* and visualized with the script

249    *plot_metagenome_contributions.R* (https://groups.google.com/forum/#!topic/picrust-

250    users/Hq9_G23J9W4) and ggplot [23] in R (http://www.R-project.org). Predicted

251    metagenomes were then used as inputs in HUMAnN2 [24] to estimate the relative abundances

252    of KEGG Pathways and/or KEGG modules. Based on the KO estimates, relative abundance

253    and coverage of KEGG Pathways was inferred by HUMAnN2. KO information was also used

254    by MinPath [25] to infer coverage and relative abundances of KEGG modules, which are

255    manually defined tight, functional units. KEGG Pathways and KEGG modules (KEGG v77.1)

256    data for HUMAnN2 were updated according to publicly available information in IMG [21] and

257    KEGG [22]. Relative abundances and coverages of KEGG modules were represented through

258    circular cladograms generated through GraPhlan.

259

260    Identification of metagenomic biomarkers

261

262    We furthered our study through detection of taxonomic clades, KEGG orthologs and

263    metabolic modules that are significantly over/under-represented in the individual oil

264    contaminated environments through statistical analyses carried out on the inferred relative

265    abundances. To this end, the procedure of linear discriminant analysis (LDA) effect size was

266    employed through LEfSe v1.0 [8] to identify differentially abundant features that can be used

267    as potential metagenomic biomarkers. For this analysis, the alpha parameter significance

268    threshold for the Krushkal-Wallis (KW) test implemented among classes in LEfSe was set to

269    0.01 and the logarithmic LDA score cut-off was set to 2.0, due to the relatively small sample

270    size under consideration. All analysis carried out through LEfSe was performed through the

271    Galaxy server [26]. Additionally, to estimate the associations between taxonomic and

272    functional enrichments in each oil polluted environment, we carried out tests of correlation

273    between abundances for KEGG orthologs and taxonomic clades using a non-parametric test

274    of Spearman's rank correlation. Detection of significant relationships, defined as a

275    correlation > 0.7 with a p-value < 0.001 and reaching a Benjamini-Hochberg false discovery

276    rate < 0.01 was carried out through the function *corr.test* implemented in the R package,

277    *psych* [27]. Correlations were only computed for oil polluted sites represented by greater than

278    6 samples. The resultant correlation network was visualized using the interactive platform,

279    Cytoscape v3.4.0 [28].

280

281    Detection of microbial interactions

282

283     Bacterial interactions in oil contaminated environments was investigated in the present

284     study through non-random bacterial co-occurrence and co-exclusion relationships within

285     individual soil sites. Only polluted sites consisting of more than 4 samples were subjected to

286     deductions of bacterial interactions. mothur implementation of the Sparse Correlations

287     for Compositional data algorithm (SparCC) [7], a tool capable of computing significant

288     correlations from compositional data while correcting for the effects of the same, was used

289     to detect significant co-occurrence and co-exclusion patterns. SparCC was run on absolute

290     count OTU tables generated by mothur for each sample, using the command *sparcc* with

291     default settings except a single non-default parameter of permutations=10,000. OTU

292     associations with an absolute SparCC correlation value above 0.6 with *p*-values < 0.01 were

293     considered statistically significant and incorporated into subsequent network construction.

294     The final network of significant SparCC correlations was built in Cytoscape 3.4.0 [28]. The

295     nodes in the reconstructed networks represent OTUs participating in robust, statistically

296     significant relationships (both positive and negative), which are in turn portrayed by edges

297     i.e. connections between the nodes.

298

299     Data Availability

300     16S rRNA amplicon sequencing data generated in this study were deposited in the NCBI

301     Sequence Read Archive (SRA) under accession numbers SRR3168574-SRR3168577. The

302     amplicon sequence data are bundled under NCBI BioProject number PRJNA306989.

303

304     **Results**

305

306     Bacterial community composition of oil polluted sediments in India

307

308      In the present study, we collected oil contaminated samples from sites subjected to regular

309      pollution events in state owned oil refineries at Guwahati and Barhola in the Indian state of

310      Assam to assess the *in situ* bacterial community composition of the same (Table 1). To our

311      knowledge, except for a metagenomic study on oil pipeline microbial populations by Joshi et

312      al. [29], this is the first study of its kind to be performed in India. To add diversity to our

313      samples, sampling was carried out from both surface and subsurface soils. 16S rRNA

314      amplicon sequencing data generated by pyrosequencing was further analyzed through

315      mothur, which classified the resulting OTUs into 465 phylotypes. This included 11 phyla, 29

316      orders and 28 families identified at ≥ 0.5% abundance in at least one of the samples (Fig. 1,

317      Supplementary Table S2). Proteobacteria was the most dominant phylum in all the samples

318      with an average relative abundance of ≥ 50% (Fig. 2). The relative abundance of

319      Proteobacteria (70%) was higher in the Barhola subsurface sample as compared to others

320      (Fig. S1). Acidobacteria was almost as abundant as Proteobacteria in the Noonmati samples

321      (~40%) while lower abundance was observed in the Barhola samples (~10%) (Supplementary

322      Fig. S1). Bacteroidetes and Chlorobi exhibited increase in abundance in the Barhola surface

323      sample when compared to others while Actinobacteria and Chloroflexi were seen to be

324      present in low but consistent abundance across all samples (Supplementary Fig. S1). At the

325      family level, bacterial community composition was much more divergent than at the phylum

326      level with some of them more abundant in particular samples. For instance,

327      Acidobacteraceae (33%), Xanthomonadaceae (10%) and Sphingomonadaceae (7.5%) were

328      detected at higher abundances in Noonmati surface than in others while Koribacteraceae

329      (25%) and Rhodocyclaceaea (~7%) were found to be more enriched in Noonmati subsurface

330      (Fig. 1). Additionally, Hydrogenophilaceae (~9%) and Hyphomicrobiaceae (~4%) exhibited

14

331    increased abundance in Barhola subsurface compared to other samples (Fig. 1).

332    Chitinophagaceaea, Ectothiorhodospiraceae and Ignavibacteraceae were more enriched in

333    the Barhola samples, while Acetobacteraceae was found in greater numbers in the

334    Noonmati samples (Fig. 1). All samples, except Noonmati surface, showed a high abundance

335    of Comamonadaceae whereas Weeksellaceae and Syntrophaceaea seemed to be specific to

336    Barhola surface (Fig. 1). Families as Sinobacteraceae, Microbacteriaceae and

337    Solibacteraceae were found in fairly consistent abundances across samples (Fig. 1). Overall,

338    at the phylum level the bacterial community composition of Noonmati samples was

339    observed to be more homogenous and less influenced by sampling depth than the Barhola

340    samples (Fig. 1, Fig. 2, Supplementary Fig. S1).

341

342    General characterization of bacterial community composition in petroleum hydrocarbon

343    polluted habitats

344

345    Comprehensive characterization of bacterial community composition in hydrocarbon

346    polluted environments was carried out using 61 publicly available and previously

347    validated/published 16S rRNA amplicon sequencing datasets distributed over 11 different

348    habitats (Table 1, Supplementary Table S1) along with 4 datasets generated in this study.

349    mothur analysis of all datasets led to the identification of 18 phyla, 38 orders and 39 families

350    at ≥ 2% abundance in at least one habitat (Fig. 2A, Fig. 2B). Proteobacteria dominated the

351    bacterial community composition at the phylum level with relative abundances ranging

352    from 20-77% across samples (Fig. 2A). Acidobacteria was detected in large numbers in all

353    samples with notably decreased abundances in OSC, OSTPu, OSTPm and OSTPd samples

354    (Fig. 2A, Supplementary Table S1). Actinobacteria and Chloroflexi were consistently

355    identified in all samples with significant increase in A samples, while Bacteroidetes showed

356    higher abundance in DWH and I samples (Fig. 2A, Supplementary Table S1). Similar to our

357    findings, an increase in abundance for the Actinobacteria was reported by Yergeau et al. in

358    diesel contaminated arctic soil biopiles [30]. Additionally, Chlorobi was detected in high

359    abundance only in M and I samples with increased Gemmatimonadetes abundance

360    identified in A, C and M samples (Fig. 2A, Supplementary Table S1). Verrucomicrobia

361    contribution in microbial community composition was higher in DWH, M and C, while

362    abundance of Firmicutes was higher in OSTP samples and Cyanobacteria in C samples as

363    compared to others (Fig. 2A, Supplementary Table S1). Order level clades with higher

364    abundances detected at ≥ 2% abundance in at least one habitat, tended to be more specific

365    to certain samples. For instance, Acidobacterales had a 21% abundance in I, while

366    Burkholderiales had an average abundance of 30% across OSTP samples and Caulobacerales

367    had an abundance of ~19% in taiga samples (Fig. 2B, Supplementary Table S1). Additionally,

368    Xanthomonadales showed high abundance (15-20%) in I and DWH samples and

369    Actinomycetales dominated A samples with an abundance of 24% (Fig. 2B, Supplementary

370    Table S1). In addition, Alteromonadales (15%) was found in increased abundance in DWH

371    samples, Ellin329 (15%) abundance was highly elevated in Taiga upper active layer (Tu), and

372    Burkholderiales (25%), Pseudomonadales (27%), Rhizobiales (22%) were enriched in OSC

373    (Fig. 2B, Supplementary Table S1). Bacterial families detected at ≥ 2% abundance in a

374    habitat also exhibited preferential sequestration to certain samples. While

375    Caulobacteraceae and Sphingomonadaceae were highly enriched in the taiga samples with

376    an average relative abundance of ~19% and ~29%, Comamonadaceae exhibited a highly

377    elevated mean abundance of 30% in the OSTP samples (Supplementary Table S1).

378    Additionally, Comamonadaceae dominated the I samples bacteriome with an abundance of

379     15% and contributed 10% of the bacteriome in A samples (Supplementary Table S1). Highly

380     specific increases in relative abundance as compared to other samples included

381     Microbacteriaceae (19%) for A samples, Alteromonadaceae (14%) and Xanthomonadaceae

382     (20%) for DWH samples, and Moraxellaceae (26%) for OSC samples (Supplementary Table

383     S1).

384

385     Similarity in bacterial community structure and detection of taxonomic biomarkers of oil

386     polluted environments

387

388     Bray-Curtis similarity scores were inferred from taxonomic data generated by mothur in

389     PAST v3.11 (Table 2) and consequently reduced to a two-dimensional space using NMDS

390     (Fig. 3) for estimation of structural similarity of bacteriomes from petroleum hydrocarbon

391     polluted environments. PERMANOVA tests carried out in PAST showed that taxonomic

392     composition of bacterial communities in the oil polluted environments were significantly

393     varied ($p$ = 0.05). However, there were three exceptions. The PERMANOVA results

394     demonstrated that the taiga samples and OSTP samples were not significantly different

395     among themselves ($p$ = 0.2-0.9) and that bacteriomes at these sites although separated by

396     depth shared substantial similarity. These observations indicated that unlike large distance

397     spatial separation i.e. geographical isolation, depth or local spatial separation is not a major

398     defining factor for effecting substantial dissimilarity. This is well supported by the Bray-

399     Curtis indices (Table 2) and NMDS plots of the same (Fig. 3) wherein all these samples

400     cluster fairly closely. Additionally, polluted mangrove sediments showed similarity with

401     OSTPm and Tp samples ($p$ = 0.057-0.09). Given the very low $p$ values these may be

402     aberrations and may have occurred due to preferences, assumptions, and thresholds set in

403    our analysis pipeline.  All habitats showed considerable conservation of taxonomic

404    composition within respective samples as described in Table 2. Among these intra-group

405    interactions, OSC samples were indeed clustered in very close proximity (Fig. 3) and

406    exhibited a Bray-Curtis similarity score of 0.85 ± 0.09, which was the highest among all inter

407    and intra-group comparisons (Table 2). Intra-group comparisons of taiga samples showed

408    lowest similarities (Bray-Curtis similarity score 0.45-0.57 ± 0.15) among all habitats, probably

409    due to sampling of source soil from 4 different regions of the China-Russia crude oil pipeline

410    (Table 1, Table 2). Among the inter group comparisons, lowest similarity was observed

411    among M and Tp samples (Bray-Curtis similarity score 0.31 ± 0.02). Apart from the taiga and

412    OSTP samples, which showed inter-group Bray-Curtis similarity score similar to intra-group

413    scores (Table 2), the highest inter-group similarity score of 0.54 ± 0.05 was seen between

414    the relatively related environments of M and DWH.

415         To further investigate taxonomic apportionment and detect differentially abundant

416    clades in various oil polluted environments, we compared the abundances of clades

417    detected at an abundance of ≥ 0.5% in at least 5 samples, at each taxonomic level (Fig. 4).

418    The consequent taxonomic profile inferred for all samples (from domain to species level)

419    was then used by LEfSe to detect metagenomic biomarkers. In all, LEfSe detected 255

420    differentially abundant taxa including 66 families, 47 genera and 11 species level biomarkers

421    across all habitats (Fig. 4, Supplementary Table S2). The largest number of taxonomic

422    biomarkers were detected for the C samples (68) while the lowest were recorded for both

423    OSTPd and Tu (7). The very low number of detected taxonomic biomarkers for OSTPd and

424    Tu may be a fallout of the comparatively higher bacterial community structure similarity

425    between taiga and OSTP samples than others leading to smaller tally of unique and

426    significantly differential clades. Among the biomarkers detected at the family level, families

18

427    such as *Acetobacteraceae*, *Rhodospirillaceae*, *Ignavibacteriaceae* and *Chitinophagaceae*

428    were attributed to I samples, *Microbacteriaceae* to A, *Pirellulaceae* and *Planctomycetaceae*

429    to C, *Flavobacteriaceae*, *Rhodobacteraceae*, and *Xanthomonadaceae* to DWH,

430    *Erythrobacteraceae* and *Desulfuromonadaceae* to M, *Rhodocyclaceae* to OSTPd,

431    *Pseudomonadaceae*, *Anaerolinaceae* and *Syntrophaceae* to OSTPm, *Comamonadaceae* and

432    *Geobacteraceae* to OSTPu, *Caulobacteraceae*, *Bradyrhizobiaceae*, and *Hyphomicrobiaceae*

433    to Tb, *Methylobacteriaceae* to OSC, *Thermogemmatisporaceae*, *Alcaligenaceae* and

434    *Sphingomonadaceae* to Tp samples and *Burkholderiaceae*, *Nocardioidaceae*, and

435    *Micrococcaceae* to Tu samples (Fig. 4, Supplementary Table S2). At the genus level,

436    *Phenylobacterium* and *Novosphingobium* were detected as biomarkers for Tp samples,

437    while genera such as *Geobacter*, *Syntrophus*, *Microbacerium*, *Mycobacterium*, *HB2_32_21*,

438    *Candidatus_Koribacter*, *Methylobacterium*, *Caulobacter*, and *Rhodococcus* were attributed

439    as biomarkers for OSTPu, OSTPm, A, C, DWH, I, OSC, Tb, and Tu samples respectively (Fig. 4,

440    Supplementary Table S2). Interestingly, LEfSe detected 19 phylum level biomarkers which

441    indicate that preferential proliferation of bacterial lineages emanating from particular

442    higher level taxa, probably driven by hydrocarbon stress, is possible and may lead to

443    definitive compositional differences between oil polluted habitats. Moreover, candidate

444    phyla such as AC1, WS3 and WS6 were identified as biomarkers for OSTP samples which also

445    underline the uniqueness of these environments (Fig. 4, Supplementary Table S2).

446

447    Metabolic characterization and functional biomarkers of oil contaminated environments

448

449          For understanding the metabolic potential of oil polluted environments and

450    identifying differentially abundant functional features, metagenomes were predicted by

451    PICRUSt using the 16S rRNA gene amplicon data generated. Predicted proteins were

452    classified as KEGG orthologs (KOs) resulting in the identification of 7020 KOs across all

453    samples. Metabolic reconstruction of metagenomes predicted by PICRUSt was carried out in

454    HUMAnN2, which detected 585 KEGG modules across all samples. Among these functional

455    modules, 19 functional modules were present across all samples at a coverage of >90% and

456    were identified as core modules (Fig. 5, Supplementary Fig. S2, Table 3, Supplementary

457    Table S4). Most of the core modules identified are essential for sustenance of prokaryotic

458    life in the environments, such as translation (M00178), central carbon metabolism

459    (M00149), ATP synthesis (M00153, M00157) and nucleotide and amino acid metabolism

460    (M00005, M00020). Rest of the core modules identified were found to be involved in

461    various kinds of transport systems for cations, nutrients and peptides including iron,

462    phosphate, nickel, and amino acids (M00188, M00222, M00223, M00236, M00237,

463    M00239, M00240, M00250, M00254, M00255, M00256, M00258, M00320) (Fig. 5,

464    Supplementary Fig. S2, Table 3, Supplementary Table S4). This is important since these

465    resources are generally present in limiting quantities in nature and often determine the

466    survival and proliferation of microbes in the environment. Additionally, transport systems

467    for lipopolysaccharide (LPS), a principal component of the gram-negative bacterial cell wall,

468    were also understandably identified as core modules and included KEGG functional modules

469    for export of LPS across both cytoplasmic (M00250) and outer membranes (M00320) (Fig. 5,

470    Supplementary Fig. S2, Table 3, Supplementary Table S4). Furthermore, 56 differently

471    covered functional modules were detected across all oil contaminated samples (Fig. 5,

472    Supplementary Fig. S2, Supplementary Table S4). Among these, five modules were

473    completely covered in only one sample while being absent in all others (Fig. 5,

474    Supplementary Fig. S2, Supplementary Table S4). This included structural complexes for

475 Manganese/Iron transport (M00243), bacterial proteasomes (M00342) and putative

476 aldouronate transport (M00603), all of which were completely covered only in the C

477 samples (Fig. 5, Supplementary Fig. S2, Supplementary Table S4). This indicates that bacteria

478 in the C site are better equipped for transport of metallic cations, peptide utilization and

479 uptake of plant derived aldouronates than other sites. Furthermore, the presence of a

480 complete complement of D-Xylose transport system (M00215) in the C site also reaffirms

481 bacterial access to hemicellulosic plant material at this site (Fig. 5, Supplementary Fig. S2,

482 Supplementary Table S4). Additionally, glutamate transport system (M00233) was

483 completely covered at only the A site, and RstB-RstA stress response two component system

484 (M00446) at the OSC site (Fig. 5, Supplementary Fig. S2, Supplementary Table S4). The

485 bacteria at A site, thus are extremely capable of utilizing glutamate for growth, while

486 resident bacteria at OSC are better furnished with stress response mechanisms critical in

487 environmental adaptation and survival.

488 　　　　In addition to differently covered functional modules, 414 KEGG modules were

489 detected to be differentially abundant in at least one of the 12 contaminated environments

490 (Fig. 5, Supplementary Fig. S2, Supplementary Table S3). The largest number of differentially

491 abundant modules were attributed to the OSC samples (70) while the least (8) were

492 attributed to the OSTPm samples (Fig. 5, Supplementary Fig. S2, Supplementary Table S3).

493 The detection of a higher number of differentially abundant modules in the OSC samples is

494 possibly due to its highly extreme environment as compared to other samples, leading to

495 sequestration of a large number of convenient functions to optimize the use of available

496 resources and counteract distinct environmental stress conditions. On the contrary, similar

497 to the result for taxonomic biomarkers, the least number of differential functional modules

498 were detected in an OSTP sample (OSTPm), with the penultimate spot being taken by Tu

499  (13). As explained above, this is not surprising since both taiga and OSTP samples share

500  comparatively greater similarity between their habitats leading to an overlap of functional

501  capabilities and hence, fewer unique and over-represented functional modules (Fig. 5,

502  Supplementary Fig. S2, Supplementary Table S3). Most of the modules for metabolism of

503  aromatic hydrocarbons such as xylene degradation (M00537), toluene degradation

504  (M0053), benzoate degradation (M00540 and M00551), salicylate degradation (M00638)

505  and catechol ortho-cleavage (M00568) were also significantly associated with the OSC

506  samples (Fig. 5, Supplementary Fig. S2, Supplementary Table S3). A number of structural

507  complexes implicated in photosynthesis were found to be differentially abundant in C

508  samples which included Photosystems I and II (M00163, M00161), the cytochrome b6f

509  complex (M00162) and NADP(H): Quinone oxidoreductase for chloroplasts and

510  cyanobacteria (M00145) (Fig. 5, Supplementary Fig. S2, Supplementary Table S3).

511  Additionally, a plethora of amino acid biosynthesis modules were detected as functional

512  biomarkers in the taiga samples. For example, three different modules for lysine

513  biosynthesis (M00525-M00527), and one each for threonine, methionine and cysteine

514  biosynthesis (M00018, M00017, M00021) were significantly abundant in Tb samples while

515  modules for valine/isoleucine, phenylalanine, tyrosine, leucine and isoleucine biosynthesis

516  (M00019, M00024, M00026, M00040, M00432, M00535, M00570) were over-represented

517  in Tp samples (Fig. 5, Supplementary Fig. S2, Supplementary Table S3). The taiga samples

518  also exhibited an over-representation for modules involved in the biosynthesis of vitamins

519  and cofactors as heme, pantothenate, ubiquinone, tetrahydrofolate, thiamine and

520  ascorbate (M00127, M00129, M00121, M00119, M00128, M00126) (Fig. 5, Supplementary

521  Fig. S2, Supplementary Table S3).

522      Overall, all the sites were found to harbor a variety of differentially abundant

523      modules dedicated to the transport of saccharide, polyols, peptides, metallic cations,

524      vitamin, amino acid, mineral ions, organic ions, lipid and phosphate underlining the large

525      genetic investment of resident bacteria in the processing of environmental information

526      specific to the said site (Fig. 5, Supplementary Fig. S2, Supplementary Table S3). However,

527      while differentially over-represented transport systems for saccharides, polyols and lipids

528      were almost ubiquitously detected, significantly associated transport systems for other

529      substrates as phosphates, amino acids, peptides and organic ions were restrained to certain

530      sites. This may indicate differential availability of these nutrients resulting in preferential

531      dependence on certain substrates acquired from the environment and also reaffirms the

532      characteristically different nature of the environments under consideration. A large number

533      of differentially abundant biosynthetic pathways for sugars, amino acids and vitamins were

534      also detected along with a great diversity of two component systems catering to a range of

535      functions such as stress and redox response, quorum sensing, chemotaxis and heavy metal

536      tolerance across all sites (Fig. 5, Supplementary Fig. S2, Supplementary Table S3).

537      Additionally, some modules for atypical energy metabolism as denitrification, dissimilatory

538      nitrate reduction and dissimilatory sulfate reduction were also detected to be differentially

539      abundant and may be important biomarkers for the corresponding sites due to their

540      contribution in bacterial respiration (Fig. 5, Supplementary Fig. S2, Supplementary Table S3).

541      Finally, several modules describing microbial resistance to antibiotics and antimicrobial

542      peptides were detected to be over-represented at all sites (Fig. 5, Supplementary Fig. S2,

543      Supplementary Table S3). This is probably due to the method of ancestral state

544      reconstruction used by PICRUSt for genome prediction, that leads to these genes being

545      predicted for consequent metagenomes if input 16S rRNA data includes hits from bacteria

546 known to have antimicrobial resistance genes. The possession and even expression of these

547 genes probably will not have a significant selective advantage in environments already

548 undergoing natural selection due to oil pollution.

549

550 Associations between bacterial taxa and metagenomic gene families

551

552 Correlations between bacterial abundance and functions enriched at different sites

553 were evaluated following a statistical strategy similar to the approach described by Segata

554 et al. [31]. The results indicated strong and significant associations between a number

555 taxonomic clades and gene families predicted by PICRUSt (Supplementary Fig. S3). A subset

556 of these significant correlations included strong associations between previously detected

557 taxonomic biomarkers and over-represented KOs for each site, which further confirmed the

558 identified taxonomic biomarkers. For example, photosynthetic structural complex genes

559 *cpeA* (K05376) and *psb28-2* (K08904), found to be differentially abundant in C samples

560 exhibited strong positive association with an over-represented cyanobacterial order,

561 Oscillatoriophycideae (Spearman correlation > 0.7, *P*-value < 0.001) (Supplementary Fig. S3).

562 Additionally, an array of genes related to polycyclic aromatic hydrocarbon degradation such

563 as *nidABD*, *phdFGIEK,* and *phtAaBC* (K11943-48, K18251, K18255-57, K18275) were

564 differentially abundant in A samples and also significantly positively correlated to known

565 polyaromatic hydrocarbon degrader and taxonomic biomarker *Mycobacterium* (Spearman

566 correlation > 0.7, *P*-value < 0.001) [32] (Supplementary Fig. S3). In other observations,

567 *Methylobacterium* showed positive correlation with a number of genes associated with the

568 transport of sugars, saccharides and amino acids such as *ggtB-D* (K10232-34), *cebE-G*

569 (K10240-42), *chvE* (K10546), *gguA-B* (K10547-48), and *gluA-D* (K10005-08) in arctic samples

24

570     (Spearman correlation > 0.7, *P*-value < 0.001) (Supplementary Fig. S3). Hydrocarbon

571     degradation genes like *pcaG* (K00448), *bbsH* (K07546) and *pcaL* (K14727) were significantly

572     correlated to class Actinobacteria in a positive manner in the same samples (Spearman

573     correlation > 0.7, *P*-value < 0.001) (Supplementary Fig. S3). In the DWH sample,

574     Colwelliaceae exhibited positive correlations with both anaerobic C4-dicarboxylate

575     transporter (*dcuB*; K07792) and 2-oxopent-4-enoate/cis-2-oxohex-4-enoate hydratase

576     (*bphH, xylJ, tesE*; K18820), an enzyme implicated in oligosaccharide metabolism (Spearman

577     correlation > 0.7, *P*-value < 0.001) (Supplementary Fig. S3). Additionally, genus *HB2.32.21*,

578     associated positively with a number of genes involved in alginate production (*alg44*,

579     *algJXKFE*; K19291-3, K19295-6, K16081), flagellar synthesis/chemotaxis (*qseC*; K07645) and

580     aminobenzoate metabolism gene regulation (*feaR*; K14063) (Spearman correlation > 0.7, *P*-

581     value < 0.001) (Supplementary Fig. S3). Acidobacteria however, was found to be negatively

582     correlated with the *alkB1-2* gene (K00496) coding for alkane-1-monooxygenase (Spearman

583     correlation < - 0.7, *P*-value < 0.001) (Supplementary Fig. S3). In OSC samples, positive

584     correlations were detected between *Methylobacterium* and genes involved in furfural

585     degradation (*hmfABCDEF*; K16874-80) and benzoate degradation (*aliAB, badI*; K04116-17,

586     K07536) (Spearman correlation > 0.7, *P*-value < 0.001) (Supplementary Fig. S3).

587     *Methylobacterium*, although an aerobe [33], has been shown to possess anaerobic benzene

588     degradation genes in the genome annotation for *Methylobacterium extorquens* PA1 in KEGG

589     (http://www.genome.jp/kegg-bin/show_pathway?mex01220). Furthermore, a number of

590     two-component systems (TCS) showed strong positive association with *Acinetobacter* and

591     Enterobacteriaceae. *Acinetobacter* was positively correlated with the enrichment of

592     RstA/RstB stress response TCS (K07639, K07661), while Enterobacteriaceae showed

593    affirmative relationships with the aerobic stress response sensor kinase ArcB (K07648) and

594    nitrate/nitrite response regulator NarP (K07685) (Supplementary Fig. S3).

595         To further understand the association of bacterial clades with gene families

596    specifically with respect to hydrocarbon degradation, we categorized all taxa contributing to

597    the abundance of genes known to be involved in hydrocarbon degradation at the family and

598    genus level (Supplementary Fig. S4). The results showed that differences existed between

599    major contributors to the abundance of particular hydrocarbonoclastic genes at different

600    sites. For example, abundance for alkane-1-monooxygenase (K00496) was contributed

601    mainly by Alteromonadaceae in DWH samples, Comammonadaceae in I, Mycobacteriaceae

602    and Nocardiaceae in C, Propionibacteriaceae in OSC, and a mixture of Acetobacteraceae,

603    Mycobacteriaceae, Nocardiaceae and Rhodospirillaceae in the taiga samples

604    (Supplementary Fig. S4). Similarly, for protocatechuate-4,5-dioxygenase (K04100-01),

605    Alteromonadaceae were again the major contributors for DWH samples, Comamonadaceae

606    and Methylobacteriaceae for OSC, Rhodocyclaceae for I, Rhodocyclaceae and

607    Comamonadaceae in OSTP, and Comamonadaceae and Bradyrhizobiaceae for taiga samples

608    (Supplementary Fig. S4). These differences in patterns observed at the family level, were

609    even more stark at higher resolutions i.e. genus level, thus effectively differentiating such

610    metagenomic contributors from site to site. This was best demonstrated for the

611    hydrocarbonoclastic gene catechol-1,2-dioxygenase (K03381), for which Alteromonadaceae

612    was found to be the most dominant contributor in both DWH and M samples (Fig. S4).

613    However, at the genus level, it was seen that while *HB2-32-21* was the dominant effector

614    organism in DWH samples, *Marinobacter* was the largest metagenomic contributor for

615    K03381 in the M samples (Supplementary Fig. S4).

616

617      Bacterial interactions in oil polluted environments

618

619          To further understand complex ecological relationships in oil polluted environments,

620      bacterial association networks were deduced from estimated taxonomic profiles. For our

621      study, we concentrated on individual oil polluted habitats, i.e. Arctic, China oil refineries, oil

622      sands core and so on. The resulting bacterial correlation networks, inferred at or above the

623      species level, constituted 186 significant relationships among 115 phylotypes ($P < 0.01$) (Fig.

624      6). Among the associations deduced to be significant, 72.58% were detected to share

625      positive correlations while the rest shared antagonistic relationships. Almost half of the co-

626      occurrence patterns identified (46%) were observed between bacteria of the same phyla

627      while more than three-quarters of all negative correlations (78%) were detected between

628      bacteria belonging to distinct phyla (Fig. 6). Thus, our results from the inferred bacterial

629      correlation networks indicated that, co-occurrence of phylotypes was closely related to

630      sharing of evolutionary lineage. For example, in the OSC habitat, phylotypes belonging to

631      proteobacterial family Oxalobacteraceae shared positive pairwise correlations with

632      Moraxellaceae and Enterobacteriaceae phylotypes, both of which belong to phylum

633      Proteobacteria (Fig. 6). Additionally, similar co-occurrence patterns were observed between

634      phylotypes attributed to families belonging to the order Actinomycetales in the C samples.

635      Positive pairwise associations were observed in C samples between phylotypes from

636      families Micrococcaceae and Nocardioidaceae, Intrasporangiaceae and Mycobacteriaceae

637      with Solirubrobacteraceae, and Gaiellaceae and Geodermatophilaceae with

638      Microbacteriaceae, all of which belong to order Actinomycetales (Fig. 6). Furthermore,

639      genera *Arthrospira* and *Phenylobacterium*, both of which belong to family

640      Caulobacteraceae, co-occurred in the Tu samples (Fig. 6). Conversely, bacteria without

641 evolutionary commonalities tended to be negatively correlated. For example, in DWH

642 samples, antagonistic relationships were observed between phylotypes belonging to family

643 Flavobacteriaceae from phylum Bacteroidetes and proteobacterial families

644 Desulfuromonadaceae and Desulfobulbaceae (Fig. 6). Similarly, mutual exclusion was

645 observed between phylotypes belonging to family Weeksellaceae of phylum Bacteroidetes

646 and Xanthomonadaceae of phylum Proteobacteria in I samples. Additionally, negatively

647 correlated associations were observed between phylotypes belonging to genera

648 Pelotomaculum and Thiobacillus in OSTPu samples, the former of which belongs to phylum

649 Firmicutes and the latter to phylum Proteobacteria (Fig. 6).

650

651 **Discussion**

652

653 The advent of next-generation sequencing (NGS) technologies have revolutionized

654 investigative approaches into microbial processes. This has led to re-exploration of well-

655 known microbial processes as the nitrogen cycle [34], methane metabolism [35], sulfur cycle [36],

656 heavy metal remediation and petroleum bioremediation [37] along with examination of exotic

657 and extreme environments such as deep-sea hydrothermal vents [38], cold deserts as

658 Antarctica [39] and remote cave systems [40]. During this time, a large amount of work has been

659 done on the microbiology of hydrocarbon degradation using NGS technologies as well [41].

660 Most of these studies employed 16S rRNA based amplicon sequencing while some used

661 metagenomic shotgun sequencing for their enquiries. Although some of these studies have

662 concentrated on prediction of potential biomarkers for oil pollution in certain environments

663 [42,43], no investigative effort has been undertaken to use the large amounts of data

664 generated in oil pollution studies across the world to review, validate and further these

665   studies. In the present study, we report bacterial diversity in oil contaminated soil collected

666   from north-eastern India and describe taxonomic and functional characteristics of oil

667   polluted environments across the world in order to understand the differences and

668   similarities that exist between them. Additionally, we infer a large number of potential

669   biomarkers, both taxonomic and functional, along with co-occurrence networks, which

670   provide new insights into the process of oil bioremediation including taxa and metabolic

671   pathways critical to survival in different oil polluted ecosystems. To this end, we have used

672   65 16S rRNA datasets from different studies across the world (Table 1, Supplementary Table

673   S1), including 4 datasets generated in this study, and carried out robust *in-silico* analysis

674   with recently developed bioinformatics tools to compare and contrast the same. The

675   principal features and findings of our study are discussed below.

676

677   Taxonomic and functional features of oil contaminated soils in India.

678

679         To unravel the bacterial community structure, functional characteristics and complex

680   inter-relationships in oil contaminated environments in India, crude oil polluted samples

681   were collected from two oil refineries in north-east India. 16S rRNA amplicon sequencing

682   analysis of the said samples revealed an overall predominance of the metabolically versatile

683   phylum Proteobacteria (Fig. 2, Supplementary Fig. S1), as has been reported previously in

684   other oil polluted sites [44]. Additionally, phylum Acidobacteria was also detected to be

685   prevalent in these samples and was identified as a biomarker as well (Fig. 2, Supplementary

686   Fig. S1, Fig. 4, Supplementary Table S2). Further identification the acidobacterial genera

687   *Edaphobacter* and *Candidatus_Koribacer* as biomarkers underlines the importance of the

688   acidobacterial lineage in the oil contaminated India soils and also indicates a slightly acidic

689    environment (Fig. 4, Supplementary Table S2). *Methylibium petroleiphilum* was also

690    detected as a biomarker and contributed significantly to the hydrocarbon degradation

691    capabilities at these sites (Fig. 4, Supplementary Table S2, Supplementary Fig. S4).

692    *Methylibium petroleiphilum*, an aerobic bacterium, has previously been reported to degrade

693    hydrocarbons as methyl tert-butyl ether [45]. Moreover, anaerobic genera as the sulfate

694    dissimilating *Thiobacillus* and the photoautrophic *Ignavibacteriaceae*, were also identified as

695    biomarkers (Fig. 4, Supplementary Table S2). These observations, along with the detection

696    of differentially abundant KEGG modules as dissimilatory sulfate reduction, sulfate => H2S

697    (M00596) and NarX-NarL (nitrate respiration) two-component regulatory system (M00471;

698    also found to be present with complete coverage) (Supplementary Fig. S2, Supplementary

699    Table S3, Supplementary Table S4) indicate that anaerobic processes and taxa play a major

700    role in the oil contaminated India soils. Simultaneously, the identification of over-

701    represented aerobic pathways such as formaldehyde assimilation, serine pathway (M00346)

702    (Supplementary Fig. S2, Supplementary Table S3) and differentially abundant aerobic taxa as

703    *Methylibium* and Chitinophagaceae (Supplementary Fig. 4, Supplementary Table S2)

704    however indicate that these environments may be predominantly aerobic but partially

705    anaerobic or microaerophilic. Over-representation of Chitinophagaceae in these

706    environments also indicate possible availability of chitin as a carbon source [46]. NMDS

707    ordination of taxonomic profiles and estimation of Bray-Curtis similarity indices for I

708    samples show that they share much more similarity amongst themselves than with other

709    sites (Fig. 3, Table 2). I samples showed highest similarity with A samples, which may be due

710    to the enrichment of Acetobacteraceae and Comamonadaceae in both the sites possibly due

711    to oil contamination (Fig. 3, Table 2). This is well demonstrated in metagenomic

712    contributions of both families to hydrocarbonoclastic genes (Supplementary Fig. S4).

713    However, whereas both the families are the consistently major contributors for

714    hydrocarbon degrading genes across all I samples, other families are shown to also

715    contribute hydrocarbonoclastic capabilities to samples in A sites (Supplementary Fig. S4).

716    Furthermore, samples collected from India oil refineries show extensive adaptation to oil

717    pollution as revealed by the detection of KOs for known hydrocarbonoclastic genes such as

718    alkane-1-monooxygenase, protocatechuate-3,4-dioxygenase, catechol-1,2-dioxygenase and

719    protocatechuate-4,5-dioxygenase (Supplementary Table S3, Supplementary Fig. S4) along

720    with over-represented hydrocarbon degradation genes such as homogentisate-1,2-

721    dioxygenase and the toluene monooxygenase system (*tmoABCDEF / tbuA1A2BCUV/*

722    *touABCDEF*) (data not shown). Bacterial interaction networks inferred by SparCC show a

723    high proportion of co-exclusion relationships, most of which exist between taxa belonging

724    to different evolutionary lineages (Fig. 6). Competitive interactions were observed between

725    aerobic and anaerobic taxa, which indicates co-existence in relative proximity with

726    competition for resources and furthers our inference of a possibly microaerophilic or

727    partially anaerobic oil polluted environment (Fig. 6). One of only two positive correlations

728    was found to be shared between *Methylibium* and *Parvibaculum*, which may endow

729    *Methylibium* with a competitive edge and facilitate its enrichment in the I samples (Fig. 6).

730

731    Validation of bioinformatic pipeline.

732

733        To our knowledge this is the only study that has congregated existing 16S rRNA NGS

734    data generated during experiments on hydrocarbon pollution in different habitats around

735    the world to deduce possible biomarkers and associated bacterial characteristics and

736    interactions. The bioinformatics pipeline we designed to analyze this data employed

737    PICRUSt, which is a recently developed tool that uses 16S rRNA data to predict

738    metagenomes for corresponding samples along with LEfSe which predicts potential

739    biomarkers and HUMAnN2 for metabolic reconstruction of PICRUSt predicted

740    metagenomes. It is to be noted however, that KEGG orthologs and KEGG module databases

741    for PICRUSt and HUMAnN2 were meticulously updated (previously PICRUSt KEGG databases

742    included KOs only up to K15039 and HUMAnN had a KEGG module database represented

743    only up to M00378) to include currently available definitions of KEGG functional modules

744    and represent the metabolic terrain of petroleum hydrocarbon contaminated habitats in

745    totality, especially with respect to hydrocarbon degradation, functional modules for which

746    were absent in the original databases. To confidently interpret and infer our results, we

747    validated our findings in both taxonomic and functional aspects. For example, a complete

748    convergence of conclusion was observed when comparing our inferred taxonomic

749    compositions and biomarkers with the findings of Mason et al. [47] for the marine sediments

750    samples. Our analysis of the marine sediment samples identified a highly dominant

751    Gammaproteobacterial genus, HB2-32-21 (Greengenes OTU ID 248394) belonging to the

752    family Alteromonadaceae (Table S2), which was detected as a taxonomic biomarker (Fig. 4)

753    and also contributed significantly to the abundance of hydrocarbon degradation genes at

754    the site (Supplementary Fig. S4). Additionally, Colwelliaceae and Rhodobacteraceae were

755    also detected as over-represented taxonomic biomarkers at the Macondo oil contaminated

756    DWH sample sites (Fig. 4, Supplementary Table S2) with the latter contributing largely

757    abundance of the alkane degrading enzyme, alkane-1-monooxygenase (Supplementary Fig.

758    S4). All these observations, are extremely consistent with the findings of Mason et al. [47] in

759    the original article and furthers their study providing new insights. Moreover, we found

760    important similarities between conclusions inferred by An et al. [17] and our study, regarding

761    the oil sands core datasets. In the original study by An et al. [17], the oil sands core was

762    deduced as an aerobic environment with limited oxygen ingress in specific regions leading

763    to regional anaerobiasis. This theory of intermittent oxygen infusion in sections of the oil

764    sands core was strongly supported by the detection of both aerobic and anaerobic pathways

765    of hydrocarbon degradation in the oil sands core. For example, in the oil sands core samples

766    we detected differentially abundant KEGG modules for aerobic degradation of different

767    hydrocarbons as xylene, benzoate, toluene and cumate including metabolism of

768    corresponding intermediates as salicylate and catechol (M0537-40, M00568, M00638)

769    (Supplementary Fig. S2, Supplementary Table S3) [48] alongside a module implicated in

770    anaerobic degradation of benzoate (M00551) (Supplementary Fig. S2, Supplementary Table

771    S3) [49]. The larger number of aerobic hydrocarbonoclastic modules compared to anaerobic

772    modules therefore further validated our bioinformatic pipeline for consequent

773    interpretation of findings in the present study.

774

775    Metabolic reconstruction of oil polluted metagenomes reveals important functional

776    pathways in petroleum hydrocarbon contaminated habitats.

777

778    In order to understand the functional landscape of each oil polluted environment,

779    metagenomes were predicted by PICRUSt from 16S rRNA data and metabolic modules

780    detected using HUMAnN2. We identified 19 core modules which were present across all

781    habitats with a coverage of > 90%. Most of these are involved in processes central to

782    survival of bacteria in the environment. Furthermore, in order to identify preferential

783    genetic investments among resident bacteria at each habitat differentially abundant KOs

784    and KEGG modules were detected through LEfSe. Consequently, we analyzed over-

785    represented KOs and KEGG modules across all habitats to identify broad metabolic

786    signatures that may be indicative of important areas of genetic expenditure, especially

787    outside hydrocarbon degradation. As a result, we identified a number of differential

788    functional pathways dedicated to transport of certain sugars or lipids, biosynthesis of

789    particular biomolecules, stress response, quorum sensing, metabolism of polysaccharides,

790    assimilation and respiration of sulphur and/or nitrogen compounds besides hydrocarbon

791    degradation, across all sites. For example, a large number of putrescine transport complexes

792    (M00193, M00299, M00300) a transport system for arginine/ornithine (M00235) were

793    detected to be differentially abundant for the DWH samples (Fig. 5, Supplementary Fig. S2,

794    Supplementary Table S3). This sequestration of putrescine transporters along with

795    ornithine, which is readily converted by ornithine decarboxylase to putrescine indicates a

796    significant dependence of marine bacteria at an oil polluted site on putrescine. This can be

797    explained by the crucial role putrescine plays in bacteria as an osmoprotectant [50], and

798    therefore its prevalence in a marine oil polluted environment. Similarly, availability and

799    possible use of carbon sources besides hydrocarbons was apparent in the C samples. The

800    differential presence of a complete complement of D-xylose transport system (M00215) and

801    a putative aldouronate transport system (M00603) along with the over-representation of

802    KEGG module M00014 (Glucuronate pathway), strongly indicated that besides petroleum

803    hydrocarbons, plant wastes maybe available as possible sources of energy for resident soil

804    bacteria at the China oil refineries site (Fig. 5, Supplementary Fig. S2, Supplementary Table

805    S3, Supplementary Table S4).  Bacteria are known to extracellularly depolymerize

806    methylglucuronoxylan, a polysaccharide made of xylose that constitutes the hemicellulosic

807    component of terrestrial plants [51] leading to the production of aldouronates and

808    xylooligosaccharides. These compounds are taken up and normally converted intracellularly

809    to fermentable xylose, leading to generation of energy along with ethanol. Alternatively, D-

810    xylose can also be directly taken up from the environment. Also, two structural complexes

811    for transport of peptides/oligopeptides (M00239 & M00439) were detected to be

812    differentially abundant in the C samples along with bacterial proteasomes (M00342) (Fig. 5,

813    Supplementary Fig. S2, Supplementary Table S3). This indicates that acquisition of

814    environmental peptides and consequent proteasomal degradation of the same, may be the

815    dominant mechanism for obtaining amino acids for assimilatory purposes in the C samples.

816    Interestingly, the DesK-DesR two-component system, implicated in regulation of the *des*

817    gene coding for a desaturase that helps control the saturation state of membrane lipids at

818    low temperatures [52] was detected to be differentially abundant in the arctic samples (Fig. 5,

819    Supplementary Fig. S2, Supplementary Table S3). Furthermore, the FitF-FitH two component

820    system, responsible for insecticidal toxin regulation [53], was over-represented in the urban

821    site of the Indian oil refinery samples (Fig. 5, Supplementary Fig. S2, Supplementary Table

822    S3). This makes sense, since it has previously been shown that relatively higher amount of

823    heat generation in cities compared to rural areas leads to sequestration of insects in urban

824    areas [54]. Sulfur assimilation in bacteria (M00616) was detected to be differentially abundant

825    in OSC samples along with a number of modules dedicated to transfer of sulfur compounds

826    (M00185, M00234, M00238, M00348, M00435-36) indicating a large genetic investment in

827    scavenging and metabolism of sulfur compounds in this site (Fig. 5, Supplementary Fig. S2,

828    Supplementary Table S3). Differential presence of a transport module for thiamine

829    (M00191), which is required for assimilation of sulfonate compounds, furthers affirms this

830    notion (Fig. 5, Supplementary Fig. S2, Supplementary Table S3). Additionally, differential

831    detection of assimilatory nitrate reduction module (M00531) also indicates the availability

832    and importance of nitrate ions in the sustenance of the OSC bacteriome (Fig. 5,

35

833 Supplementary Fig. S2, Supplementary Table S3). Sulfate and nitrate ions are also important

834 molecules in anaerobic respiration, and may therefore also play crucial roles in bacterial

835 survival in the anaerobic regions of the OSC. Interestingly, reduction of nitrate has been

836 reported to be closely linked to anaerobic degradation of benzene and concomitant growth

837 [55], functional modules for both of which have been differentially detected in OSC samples

838 (Fig. 5, Supplementary Fig. S2, Supplementary Table S3). Unlike other oil polluted sites, a

839 large number of hydrocarbonoclastic modules differentially detected in the OSC samples

840 (see previous section), whereas only two transport systems for small sugars (M00204,

841 M00215) and no major polysaccharide metabolism and/or transport pathways were

842 detected to be over-represented (Fig. 5, Supplementary Fig. S2, Supplementary Table S3).

843 This not only indicates at the extreme habitat with highly restricted availability of carbon

844 sources other than petroleum hydrocarbons in the OSC, it also explains the large clustering

845 of differentially abundant hydrocarbon degradation pathways in the OSC samples and

846 establishes petroleum hydrocarbons as the comprehensively dominant carbon

847 assimilation/energy production source for this site. A number of functional modules related

848 to methane metabolism, both methanogenic and methanotrophic, were detected to be

849 over-represented for the OSTP samples. For example, methanogenesis (M00356) was over-

850 represented in OSTPu, along with methane assimilation modules M00344-45 and M00608

851 detected to be differentially abundant in OSTPu and OSTPm respectively (Fig. 5,

852 Supplementary Fig. S2, Supplementary Table S3). Oil sands tailings ponds are known to be

853 important sources of methanogenesis and of methylotrophy [17], where deeper regions tend

854 to be highly anaerobic. Additionally, modules for copper processing (M00762) and copper

855 tolerance sensor (M00452) were detected in OSTPd (Fig. 5, Supplementary Fig. S2,

856 Supplementary Table S3). This is important, since copper is an essential component of

36

857     particulate methane monooxygenase (pMMO), and its availability can therefore determine

858     the survivability of methanotrophs [56] along with the ratio of soluble and particulate MMO in

859     the environment. A large number of modules dedicated to the biosynthesis of amino acids,

860     vitamins and co-factors were detected to be over-represented in the taiga samples. This

861     may be due to the poor availability of useful forms of amino acids and vitamins in the taiga

862     environment. Presence of differentially abundant modules for sulfur containing amino acid

863     biosynthesis (M00017, M00021) is supported by the detection of over-represented sulfur

864     assimilation modules (M00176, M00595) which involve biosynthesis of cysteine and

865     methionine as final/supplementary steps [57,58] (Fig. 5, Supplementary Fig. S2, Supplementary

866     Table S3). Additionally, presence of alternative carbon sources as pectin and component

867     sugars of other plant polysaccharides at the taiga sites can be inferred through the presence

868     of differentially abundant functional modules for pectin degradation (M00081), and uptake

869     and metabolism of other sugar and sugar derivatives as N-Acetylglucosamine, N, N'-

870     Diacetylchitobiose, D-glucuronate, aldouronates, D-galactouronate (M00606, M00205,

871     M00061, M00603, M00631) (Fig. 5, Supplementary Fig. S2, Supplementary Table S3). In the

872     M samples, two component systems for starvation of phosphate (M00434), a limiting

873     nutrient for mangroves [59] and metal tolerance (M00499) were detected as differentially

874     abundant (Fig. 5, Supplementary Fig. S2, Supplementary Table S3). Genetic investment in

875     metal tolerance should be important in M samples as mangroves in Brazil are routinely

876     subjected to pollution from factory effluents [42]. Furthermore, the clustering of differentially

877     abundant central carbohydrate metabolism pathways (M00001-2, M00004, M00009,

878     M00011) along with transport systems for fructose like sugars (M00273) in M samples,

879     indicate availability of simple sugars as carbon sources besides hydrocarbons (Fig. 5,

880     Supplementary Fig. S2, Supplementary Table S3). This is also supported by the differentially

881    abundant module for synthesis of trehalose (M00565), a known carbohydrate energy

882    storage compound and anti-desiccation agent [60], from glucose (Fig. 5, Supplementary Fig.

883    S2, Supplementary Table S3). All functional modules for degradation of aromatic

884    hydrocarbons were detected to be differentially abundant in OSC, DWH, taiga and OSTP (in

885    that order) samples (Fig. 5, Supplementary Fig. S2, Supplementary Table S3). This is probably

886    because these environments tend to be more extreme than other sites described in this

887    study and coupled with oil pollution, the bacterial metabolic pathways in these

888    environments have been further sculpted to rely greatly only on petroleum hydrocarbons

889    for growth. Additionally, sulfate and nitrate utilization modules have been identified in most

890    of these sites, which indicates an abundance of such ions in the environment and therefore

891    use of the same for anaerobic alkane degradation, as previously described [61]. Our results

892    thus indicate that for all habitats, genetic composition of the bacteriome is representative of

893    the immediate environment especially in terms of substrate usage, nutrient availability,

894    energy metabolism, biosynthesis of compounds, and survival strategies including quorum

895    sensing, chemotaxis, and stress response. Our findings reveal pathways differentially

896    important in these oil polluted environments, especially those not related to hydrocarbon

897    degradation and can therefore be used for differentiation between habitats of interest.

898    Further empirical studies will however be required strengthen these observations and

899    pinpoint functional biomarkers absolutely exclusive to oil polluted environments in specific

900    biomes.

901

902    Taxonomic biomarkers make important contributions to hydrocarbonoclastic and additional

903    functional capacities in oil polluted environments.

904

38

905    In order to identify taxonomic clades that may be differentially abundant in oil polluted sites

906    used in the present study, taxonomic profiles generated through analysis of 16S rRNA data

907    in mothur were examined using LEfSe. Additionally, to decipher functional associations of

908    taxonomic clades, direct correlations between KOs and taxa were determined along with

909    metagenomic contributions to hydrocarbonoclastic genes. Furthermore, bacterial co-

910    occurrence and co-exclusion networks were deduced to understand important bacterial

911    interactions in oil polluted sites. Our findings suggest that, taxonomic biomarkers inferred in

912    our study contribute significantly to important functions in the oil polluted metabolic

913    landscape and are often determined by their oil degradation capabilities. For example,

914    biomarkers for DWH samples *HB2.32.21* and Alteromonadaceae (Fig. 4, Supplementary

915    Table S2), were associated with over-represented KOs implicated in alginate biosynthesis

916    (Supplementary Fig. S3). Moreover, a two component pathway involved in the regulation of

917    alginate production (M00505) was also differentially abundant in DWH samples (Fig. 5,

918    Supplementary Fig. S2, Supplementary Table S3). Interestingly, previous studies have shown

919    that alginates provide increased mechanical stability to bacterial biofilms [62], and can

920    therefore be instrumental in aiding anchorage or adhesion of DWH Alteromonadaceae.

921    *HB.32.21* and Alteromonadaceae were found to be important contributors in

922    hydrocarbonoclastic properties of the DWH bacteriome (Supplementary Fig. S4) and the

923    former also exhibited strong associations with regulation of genes for aminobenzoate

924    metabolism through *feaR* (see Results). Furthermore, another taxonomic biomarker

925    identified for DWH samples, Colwelliaceae, was closely associated to the anaerobic C4-

926    dicarboxylate transporter DcuB, which is responsible for transport of molecules as fumarate,

927    succinate and malate [63]. This is important, as it may help the facultatively aerobic

928    Colwelliaceae to degrade alkanes anaerobically by addition of fumarates in marine

39

929    sediments [61]. Similarly, *Mycobacterium* was detected as a biomarker for C samples (Fig. 4,

930    Supplementary Table S2) and correlated strongly with KOs implicated in degradation of

931    hydrocarbons as naphthalene, benzoate and phthalate (Supplementary Fig. S3).

932    *Mycobacterium* have previously been shown to harbor the ability to degrade a variety of

933    aromatic hydrocarbons such as naphthalene, anthracene, phenanthrene, pyrene and so on

934    [32]. Phylum Cyanobacteria, a biomarker for C samples (Fig. 4, Supplementary Table S2),

935    strongly correlated with differentially abundant photosynthetic proteins *cpeA* (K05376) and

936    *psb28-2* (K08904) through an over-represented cyanobacterial order for C samples,

937    Oscillatoriophycideae (Fig. 4, Supplementary Table S2). Additionally, KEGG modules for

938    photosynthesis such as Photosystems I and II (M00163, M00161), cytochrome b6f complex

939    (M00162) and NADP(H):quinone oxidoreductase for chloroplasts and cyanobacteria

940    (M00145) were also found to be over-represented in C samples (Fig. 5, Supplementary Fig.

941    S2, Supplementary Table S3). These observations indicated important, differential and extra-

942    hydrocarbonoclastic contributions of Cyanobacteria in C samples. Furthermore,

943    *Microbacterium*, which is known to be a stringent chemoorganotroph [64], was identified to

944    be differentially abundant in A samples (Fig. 4, Supplementary Table S2). *Microbacterium*

945    was found to share close associations with KOs involved in over-represented transport

946    systems dedicated to acquisition of organic compounds such as cellobiose (M00206), alpha-

947    glucosides (M00201), glutamate (M00227, M00233) and multiple sugars (M00207, M00216,

948    M00221) (Supplementary Fig. S3), which can be used as possible sources of carbon and

949    energy and also indicates availability of the same in the environment. Phylum

950    Actinobacteria and class Actinobacteria, which were detected as biomarkers in the A

951    samples, exhibited significant correlations with almost all differentially abundant KOs for A

952    samples including hydrocarbon degrading genes as *pcaG*, phenol-2-monooxygenase, *bbsH*,

40

953    and *pcaL* (data not shown). In OSC samples, over-represented taxa *Methylobacterium* was

954    associated with genes involved in degradation of furfural and other hydrocarbons

955    (Supplementary Fig. S3). The presence of the strongly aerobic *Methylobacterium* [33] once

956    again reinforces the finding of ample availability of oxygen in the OSC. Interestingly,

957    differentially abundant taxa Enterobacteriaceae and *Acinetobacter* were detected to be

958    associated with a number of KOs implicated in stress response which included TCS KOs for

959    aerobic/anaerobic survival as ArcB and NarP and transcriptional regulation of the *mar-sox-*

960    *rob* regulon (Supplementary Fig. S3). The *mar-sox-rob* regulon has been reported in

961    coordinating survival against various environmental stresses activated by inducers as

962    paraquat, decanoate and intriguingly, salicylate [65], functional modules for which is

963    differentially abundant in OSC samples (Fig. 5. Supplementary Fig. S2, Supplementary Table

964    S3). Additionally, *Acinetobacter* was correlated with the stress response serine protease

965    DegS and iron starvation Fe/S biogenesis protein NfuA (Supplementary Fig. S3). Thus, these

966    biomarkers seem to contribute to important stress response pathways rather than

967    hydrocarbon degrading capabilities. Additionally, over-represented taxa such as

968    Oxalobacteraceae, *Cupriavidus*, Brucellaceae, and *Ochrobactrum* (Fig. 4, Supplementary

969    Table S2) were found to differentially contribute to the abundance of a number of

970    hydrocarbonoclastic genes (K00446, K00448-51, K03381) (Supplementary Fig. S4) in the OSC

971    samples. In taiga samples, a number of detected biomarkers such as *Phenylobacterium,*

972    Caulobacteraceae, Sphingomonadaceae, *Novosphingobium*, *Rhodococcus,* and

973    Burkholderiaceae (Fig. 4, Supplementary Table S2) were found to contribute heavily but

974    differently to the abundance of a plethora of hydrocarbonoclastic genes (Supplementary

975    Fig. S4).  Differentially abundant functional modules for the assimilation of sulphate,

976    transformation of thiosulphate to sulphate and regulation of the SOX complex responsible

41

977    for thiosulphate transformation (M00176, M00595, M00523) underline the preferential

978    sulphur usage in this site (Fig. 5, Supplementary Fig. S2, Supplementary Table S3). This is

979    well supported by the identification of *Bradyrhizobium*, *Caulobacter*, and *Burkholderia* as

980    biomarkers (Fig. 4, Supplementary Table S2), all which are known to be involved in sulfur

981    metabolism [66,67] and house homologous genes for the same. Interestingly, a number of

982    biomarkers identified here for the taiga samples such as *Phenylobacterium,*

983    Sphingomonadaceae, *Novosphingobium* and *Rhodococcus* were detected as "habitat

984    specialists" in oil contaminated taiga samples by Yang et al. [68]. Similar to the results of An et

985    al. [17], we encountered a significantly high proportion of anaerobic taxa in the OSTP samples,

986    among which Anaerolinaceae, Syntrophaceae, Desulfobulbaceae, Peptococcaceae,

987    Geobacteraceae, Syntrophorhabdaceae and the thermophilic Caldiserica [69] were detected

988    as biomarkers (Fig. 4, Supplementary Table S2). Detected taxonomic biomarkers such as

989    Anaerolinaceae and Comamonadaceae (Fig. 4, Supplementary Table S2) were found to

990    make significant contributions to the abundance of hydrocarbon degradation genes

991    (Supplementary Fig. S4) in OSTP samples. Other identified biomarkers such as

992    Geobacteraceae and *Thauera* (Fig. 4, Supplementary Table S2) are well known anaerobic

993    hydrocarbon degraders [70,71]. Additionally, another detected biomarker Nitrospirales (Fig. 4,

994    Supplementary Table S2), which is involved in nitrification [72] may contribute to nitrification,

995    for which over-represented module ammonia => nitrite transformation (M00528) was

996    identified in OSTP samples (Fig. 5, Supplementary Fig. S2, Supplementary Table S3).

997    Biomarkers of sulfate reducing bacteria such as Desulfuromonadales and Desulfobulbaceae

998    (Fig. 4, Supplementary Table S2), which is a known mesophilic/psychrophilic sulfate reducer

999    [73] may be involved in important sulfur metabolism pathways known to be important in

1000   OSTPs [74]. Interestingly, obligate anaerobes Anaerolinaceae have previously been associated

1001    with sulfate reducing conditions in the OSTPs [75]. Lastly, major contributions for

1002    hydrocarbonoclastic capabilities in OSTP samples was also observed from biomarkers

1003    *Pseudomonas* (K00446, K00448, K00449, K00496, K03381) and Rhodocyclaceae (K04100-01)

1004    (Fig. 4, Supplementary Table S2, Supplementary Fig. S4), furthering the hydrocarbon

1005    degradation capabilities of OSTPs.

1006         We also investigated significant bacterial associations in oil polluted sites to decipher

1007    important co-occurrence and co-exclusions. Our results showed that greater co-occurrence

1008    exists between phylotypes sharing an evolutionary lineage while more co-exclusions were

1009    observed between phylotypes from different ancestries. This observation has also been

1010    previously reported in microbial correlation studies in the environment [76,77]. Interestingly,

1011    not a large proportion of taxonomic biomarkers were observed to be represented in these

1012    significant correlations. This can possibly happen due to separation of niches due to various

1013    environmental and even temporal factors. For example, in the bacterial association network

1014    for DWH samples, biomarker Colwelliaceae was detected to participate in a significantly

1015    positive relationship with Desulfobulbaceae, a strictly anaerobic sulfate utilizing bacterial

1016    family (Fig. 6). The existence of this kind of a relationship, based on degradation of

1017    recalcitrant hydrocarbons, was inferred upon by the original authors too [47]. Strikingly

1018    however, the most abundant and robust hydrocarbon degrader i.e. *HB2.32.21*

1019    (Supplementary Fig. S4) was not detected to be involved in any significant associations. This

1020    observation can be explained by a possible individual capacity of survival for *HB2.32.21* due

1021    to its hydrocarbonoclastic capacities without extensive interactions with other resident

1022    bacteria, therefore occupying a separate niche in the oil polluted marine sediment site.

1023    Thus, significant correlations (both positive and negative) may be driven by factors other

1024    than only oil pollution in oil contaminated sites with apparently benign taxa being involved

1025    in such interactions. This indicates that biomarkers and correlation networks must be

1026    studied in tandem to deduce meaningful conclusions. Even though empirical evidence in

1027    support of the natural presence of most microbial association networks is lacking, our study

1028    shows that they may be important to holistic interpretation of results as well as being a

1029    good starting point for further investigations.

1030        Our results therefore show that, detected biomarkers may contribute differently to

1031    strictly hydrocarbonoclastic properties when compared across sites, but their close

1032    association with a majority of differentially abundant KOs and as an extension a number of

1033    over-represented functional pathways for each site underlines their significance in these oil

1034    contaminated sites. We find that although many of the taxonomic biomarkers contribute to

1035    hydrocarbonoclastic capacities, some do not and can therefore contribute to other possibly

1036    important functions. These observations not only elucidate important taxa contributing

1037    functions more specific and essential to each site, but also shows that niches related to

1038    functions other than hydrocarbon degradation may significantly influence bacteriome

1039    structure in oil polluted sites, possibly more in sites with understandably lower degrees of

1040    contamination. This indicates clearly that while hydrocarbonoclastic capabilities may be a

1041    driving force for continued survival in these sites, other immediate factors including

1042    availability of different organic and inorganic compounds and environmental stress play

1043    heavily on the evolution of the bacteriome. Thus, we see that a combination of oil

1044    degradation capabilities and environmental factors shape the landscape for bacterial

1045    petroleum degradation. As an extension, it therefore becomes imperative to examine oil

1046    bioremediation processes, especially aimed at empirical identification of biomarkers, in

1047    totality with due comparison to similar studies and not in isolation as it may lead to

1048    misleading conclusions. This is well illustrated in some previous studies that have focused on

1049    predicting microbial markers or proxies for oil pollution in certain environments [42,43]. In the

1050    study on mangrove oil pollution and detection of microbial proxies by dos Santos et al. [43],

1051    *Marinobacter*, belonging to family Alteromonadaceae, was identified as a possible

1052    biomarker for oil pollution in mangroves. However, in our study when compared to other

1053    sites, Alteromonadaceae was detected to be differentially abundant in the DWH samples

1054    and *Marinobacter* was not identified as over-represented in any of the oil polluted sites (Fig.

1055    4, Supplementary Table S2).

1056

1057    **Conclusion**

1058

1059    In summary, our study showed that significant taxonomic and functional differences exist

1060    between geographically and/or spatially isolated oil polluted sites and that oil pollution is

1061    not the sole driving factor in determination of the bacteriome at these sites, even if maybe

1062    the most predominant one. In our study we have successfully detected functions that are

1063    important and contribute significantly to the sustenance of bacteriomes at these oil polluted

1064    sites. Additionally, we have detected taxa that are differentially abundant and also

1065    contribute to many of these functions.  Furthermore, we have also successfully shown that

1066    several of these important taxonomic clades and functional modules are often involved in

1067    extra-hydrocarbonoclastic activities, thus underlining the importance of these apparently

1068    peripheral niches related to endemic environmental responses such as variable resource

1069    utilization, alternate respiration and stress response in the survival of oil contaminated

1070    ecosystems. In the process, we therefore identified robust taxonomic and functional

1071    biomarkers, that are representative of an entire oil polluted environment and not only its

1072    hydrocarbonoclastic capabilities. These biomarkers can be implemented in monitoring of

1073    bacterial remediation processes as well as for distinguishing one of the 12 oil polluted

1074    habitats. Our results show that some parallels exist in the functional composition of oil

1075    polluted environments, mainly regarding adaptations to aerobic and anaerobic lifestyles

1076    depending on availability of limited resources compatible with sustenance of anaerobic

1077    growth (e.g. sulfate/nitrate related metabolism), but differences between them pertaining

1078    to transport of substrates, biosynthesis of biomolecules, response to various stress,

1079    degradation of diverse compounds and so on are significant and together with differences

1080    identified in taxonomic profiles and hydrocarbonoclastic capabilities enable us to truly

1081    differentiate between them. However, further studies are required which may involve

1082    simultaneous 16S rRNA based phylogenetic survey, metagenomic and metatranscriptomic

1083    investigations of oil polluted environments followed by empirical experimentations on the

1084    same to confirm robust and significantly differential taxonomic and functional biomarkers

1085    which may be used for monitoring. With the current sequencing technologies,

1086    bioinformatics tools and increasing data from other studies, it is highly possible that

1087    characteristic degradation profiles and subsequent remediation strategies for different

1088    environments across the world will be inferred fairly soon with unprecedented confidence.

1089    To our knowledge, this is the first population genomics study carried out on petroleum

1090    hydrocarbon polluted habitats. Our study presents novel understanding of oil contaminated

1091    habitats by showing how and in what manner petroleum hydrocarbons fashion the

1092    metagenomic fabric along with the evident effect of endemic factors characteristic of

1093    geographically separated diverse petroleum hydrocarbon contaminated environments,

1094    while interpreting the effects of both on the resident microbiome. This study therefore

1095    provides a foundation for future investigations into the microbial biodiversity and habitat

1096    function in oil polluted sites.

1097

**Acknowledgements**

1099

1113

**Author contributions**

1115

1116 D.C. and A.K.S. managed the project. A.M., D.C. and A.K.S. conceptualized and designed the

1117 experiments. B.C., J.L., A.K.S. and A.K.M. designed and conducted sampling for oil

1118 contaminated soil. A.M., B.C., J.L., P.B. and M.B. were involved in designing the sequencing

1119 strategy and conducting the same. A.M. designed the bioinformatic analysis strategy and

1120    conducted the same with assistance from A.P. A.M. and D.C. performed the data analyses.

1121    A.M, D.C. and A.K.S. prepared the manuscript. All authors reviewed the manuscript.

1122

1123    **Competing financial interests**

1124

1125    The authors declare no competing financial interests.

1126

1127    **References**

1128    1        Dean-Ross, D., Moody, J. & Cerniglia, C. E. Utilization of mixtures of polycyclic aromatic

1129             hydrocarbons by bacteria isolated from contaminated sediment. *FEMS microbiology ecology*

1130             **41**, 1-7, doi:10.1111/j.1574-6941.2002.tb00960.x (2002).

1131    2        Molina, M., Araujo, R. & Hodson, R. E. Cross-induction of pyrene and phenanthrene in a

1132             Mycobacterium sp. isolated from polycyclic aromatic hydrocarbon contaminated river

1133             sediments. *Canadian journal of microbiology* **45**, 520-529 (1999).

1134    3        Stringfellow, W. T. & Aitken, M. D. Competitive metabolism of naphthalene,

1135             methylnaphthalenes, and fluorene by phenanthrene-degrading pseudomonads. *Applied and*

1136             *environmental microbiology* **61**, 357-362 (1995).

1137    4        Bakken, L. R. Culturable and non-culturable bacteria in soil. *J.D. van Elsas, J.T. Trevor, E.M.H.*

1138             *Wellington (Eds.), Modern soil microbiology, Marcel Dekker, New York* 47-61 (1997).

1139    5        Gevers, D. *et al.* The Human Microbiome Project: a community resource for the healthy

1140             human microbiome. *PLoS biology* **10**, e1001377, doi:10.1371/journal.pbio.1001377 (2012).

1141    6        Segata, N. *et al.* Computational meta'omics for microbial community studies. *Molecular*

1142             *systems biology* **9**, 666, doi:10.1038/msb.2013.22 (2013).

1143    7        Friedman, J. & Alm, E. J. Inferring correlation networks from genomic survey data. *PLoS*

1144             *computational biology* **8**, e1002687, doi:10.1371/journal.pcbi.1002687 (2012).

1145  8    Segata, N. *et al.* Metagenomic biomarker discovery and explanation. *Genome biology* **12**,

1146       R60, doi:10.1186/gb-2011-12-6-r60 (2011).

1147  9    Langille, M. G. *et al.* Predictive functional profiling of microbial communities using 16S rRNA

1148       marker gene sequences. *Nature biotechnology* **31**, 814-821, doi:10.1038/nbt.2676 (2013).

1149  10   Weisburg, W. G., Barns, S. M., Pelletier, D. A. & Lane, D. J. 16S ribosomal DNA amplification

1150       for phylogenetic study. *Journal of bacteriology* **173**, 697-703 (1991).

1151  11   Lofgren, J. L. *et al.* Lack of commensal flora in Helicobacter pylori-infected INS-GAS mice

1152       reduces gastritis and delays intraepithelial neoplasia. *Gastroenterology* **140**, 210-220,

1153       doi:10.1053/j.gastro.2010.09.048 (2011).

1154  12   Andrews, S. FastQC: a quality control tool for high throughput sequence data. *Available*

1155       *online at: http://www.bioinformatics.babraham.ac.uk/projects/fastqc* (2010).

1156  13   Schloss, P. D. *et al.* Introducing mothur: open-source, platform-independent, community-

1157       supported software for describing and comparing microbial communities. *Applied and*

1158       *environmental microbiology* **75**, 7537-7541, doi:10.1128/AEM.01541-09 (2009).

1159  14   Edgar, R. C., Haas, B. J., Clemente, J. C., Quince, C. & Knight, R. UCHIME improves sensitivity

1160       and speed of chimera detection. *Bioinformatics* **27**, 2194-2200,

1161       doi:10.1093/bioinformatics/btr381 (2011).

1162  15   DeSantis, T. Z. *et al.* Greengenes, a chimera-checked 16S rRNA gene database and

1163       workbench compatible with ARB. *Applied and environmental microbiology* **72**, 5069-5072,

1164       doi:10.1128/AEM.03006-05 (2006).

1165  16   Asnicar, F., Weingart, G., Tickle, T. L., Huttenhower, C. & Segata, N. Compact graphical

1166       representation of phylogenetic data and metadata with GraPhlAn. *PeerJ* **3**, e1029,

1167       doi:10.7717/peerj.1029 (2015).

1168  17   An, D. *et al.* Metagenomics of hydrocarbon resource environments indicates aerobic taxa

1169       and genes to be unexpectedly common. *Environmental science & technology* **47**, 10708-

1170       10717, doi:10.1021/es4020184 (2013).

1171    18    Yang, S., Wen, X., Zhao, L., Shi, Y. & Jin, H. Crude oil treatment leads to shift of bacterial

1172        communities in soils from the deep active layer and upper permafrost along the China-

1173        Russia Crude Oil Pipeline route. *PloS one* **9**, e96552, doi:10.1371/journal.pone.0096552

1174        (2014).

1175    19    Bray, J. R. & Curtis, J. T. An ordination of the upland forest communities of southern

1176        Wisconsin. *Ecol Monographs* **27**, 325-349 (1957).

1177    20    Hammer, Ø., Harper, D. A. T. & Ryan, P. D. PAST: Paleontological Statistics Software Package

1178        for Education and Data Analysis. *Palaeontologia Electronica* **4**, 9pp (2001).

1179    21    Markowitz, V. M. *et al.* IMG: the Integrated Microbial Genomes database and comparative

1180        analysis system. *Nucleic acids research* **40**, D115-122, doi:10.1093/nar/gkr1044 (2012).

1181    22    Kanehisa, M., Goto, S., Furumichi, M., Tanabe, M. & Hirakawa, M. KEGG for representation

1182        and analysis of molecular networks involving diseases and drugs. *Nucleic acids research* **38**,

1183        D355-360, doi:10.1093/nar/gkp896 (2010).

1184    23    Wickham, H. ggplot2: Elegant Graphics for Data Analysis. *Springer-Verlag New York* (2009).

1185    24    Abubucker, S. *et al.* Metabolic reconstruction for metagenomic data and its application to

1186        the human microbiome. *PLoS computational biology* **8**, e1002358,

1187        doi:10.1371/journal.pcbi.1002358 (2012).

1188    25    Ye, Y. & Doak, T. G. A parsimony approach to biological pathway reconstruction/inference

1189        for genomes and metagenomes. *PLoS computational biology* **5**, e1000465,

1190        doi:10.1371/journal.pcbi.1000465 (2009).

1191    26    Goecks, J., Nekrutenko, A., Taylor, J. & Galaxy, T. Galaxy: a comprehensive approach for

1192        supporting accessible, reproducible, and transparent computational research in the life

1193        sciences. *Genome biology* **11**, R86, doi:10.1186/gb-2010-11-8-r86 (2010).

1194    27    Revelle, W. psych: Procedures for Personality and Psychological Research.  (2016).

1195    28    Shannon, P. *et al.* Cytoscape: a software environment for integrated models of biomolecular

1196        interaction networks. *Genome research* **13**, 2498-2504, doi:10.1101/gr.1239303 (2003).

1197  29  Joshi, M. N. *et al.* Metagenomic approach for understanding microbial population from

1198       petroleum muck. *Genome announcements* **2**, doi:10.1128/genomeA.00533-14 (2014).

1199  30  Yergeau, E., Sanschagrin, S., Beaumier, D. & Greer, C. W. Metagenomic analysis of the

1200       bioremediation of diesel-contaminated Canadian high arctic soils. *PloS one* **7**, e30058,

1201       doi:10.1371/journal.pone.0030058 (2012).

1202  31  Segata, N. *et al.* Composition of the adult digestive tract bacterial microbiome based on

1203       seven mouth surfaces, tonsils, throat and stool samples. *Genome biology* **13**, R42,

1204       doi:10.1186/gb-2012-13-6-r42 (2012).

1205  32  Kim, S. J., Kweon, O. & Cerniglia, C. E. Degradation of Polycyclic Aromatic Hydrocarbons by

1206       Mycobacterium Strains. *Handbook of Hydrocarbon and Lipid Microbiology*, 1865-1879.

1207  33  Green, P. N. Methylobacterium. *The Prokaryotes: Proteobacteria: Alpha and Beta Subclasses*

1208       **5**, 257-265, doi:10.1007/0-387-30745-1_14 (2006).

1209  34  Jurkowski, A., Reid, A. H. & Labov, J. B. Metagenomics: a call for bringing a new science into

1210       the classroom (while it's still new). *CBE life sciences education* **6**, 260-265 (2007).

1211  35  Guo, J. *et al.* Dissecting microbial community structure and methane-producing pathways of

1212       a full-scale anaerobic reactor digesting activated sludge from wastewater treatment by

1213       metagenomic sequencing. *Microbial cell factories* **14**, 33 (2015).

1214  36  He, Y., Feng, X., Fang, J., Zhang, Y. & Xiao, X. Metagenome and Metatranscriptome Revealed

1215       a Highly Active and Intensive Sulfur Cycle in an Oil-Immersed Hydrothermal Chimney in

1216       Guaymas Basin. *Frontiers in microbiology* **6**, 1236 (2015).

1217  37  Kuppusamy, S. *et al.* Pyrosequencing analysis of bacterial diversity in soils contaminated

1218       long-term with PAHs and heavy metals: Implications to bioremediation. *Journal of hazardous*

1219       *materials* **317**, 169-179 (2016).

1220  38  Cerqueira, T. *et al.* Microbial diversity in deep-sea sediments from the Menez Gwen

1221       hydrothermal vent system of the Mid-Atlantic Ridge. *Marine genomics* **24 Pt 3**, 343-355,

1222       doi:10.1016/j.margen.2015.09.001 (2015).

1223 39    Tytgat, B. *et al.* Bacterial diversity assessment in Antarctic terrestrial and aquatic microbial

1224        mats: a comparison between bidirectional pyrosequencing and cultivation. *PloS one* **9**,

1225        e97564, doi:10.1371/journal.pone.0097564 (2014).

1226 40    Barton, H. A. *et al.* Microbial diversity in a Venezuelan orthoquartzite cave is dominated by

1227        the Chloroflexi (Class Ktedonobacterales) and Thaumarchaeota Group I.1c. *Frontiers in*

1228        *microbiology* **5**, 615, doi:10.3389/fmicb.2014.00615 (2014).

1229 41    Mukherjee, A. & Chattopadhyay, D. Exploring environmental systems and processes through

1230        next-generation sequencing technologies: insights into microbial response to petroleum

1231        contamination in key environments. *The Nucleus*, doi:10.1007/s13237-016-0190-3 (2016).

1232 42    Andreote, F. D. *et al.* The microbiome of Brazilian mangrove sediments as revealed by

1233        metagenomics. *PloS one* **7**, e38600, doi:10.1371/journal.pone.0038600 (2012).

1234 43    dos Santos, H. F. *et al.* Mangrove bacterial diversity and the impact of oil contamination

1235        revealed by pyrosequencing: bacterial proxies for oil pollution. *PloS one* **6**, e16943,

1236        doi:10.1371/journal.pone.0016943 (2011).

1237 44    Hernandez-Raquet, G. *et al.* Molecular diversity studies of bacterial communities of oil

1238        polluted microbial mats from the Etang de Berre (France). *FEMS microbiology ecology* **58**,

1239        550-562, doi:10.1111/j.1574-6941.2006.00187.x (2006).

1240 45    Kane, S. R. *et al.* Whole-genome analysis of the methyl tert-butyl ether-degrading beta-

1241        proteobacterium Methylibium petroleiphilum PM1. *Journal of bacteriology* **189**, 1931-1945,

1242        doi:10.1128/JB.01259-06 (2007).

1243 46    Rosenberg, E. The Family Chitinophagaceae. *The Prokaryotes: Other Major Lineages of*

1244        *Bacteria and The Archaea*, 493-495, doi:10.1007/978-3-642-38954-2_137 (2014).

1245 47    Mason, O. U. *et al.* Metagenomics reveals sediment microbial community response to

1246        Deepwater Horizon oil spill. *The ISME journal* **8**, 1464-1475, doi:10.1038/ismej.2013.254

1247        (2014).

1248   48   Jindrova, E., Chocova, M., Demnerova, K. & Brenner, V. Bacterial aerobic degradation of
1249         benzene, toluene, ethylbenzene and xylene. *Folia microbiologica* **47**, 83-93 (2002).

1250   49   Pelletier, D. A. & Harwood, C. S. 2-Hydroxycyclohexanecarboxyl coenzyme A dehydrogenase,
1251         an enzyme characteristic of the anaerobic benzoate degradation pathway used by
1252         Rhodopseudomonas palustris. *Journal of bacteriology* **182**, 2753-2760 (2000).

1253   50   Wood, J. M. Osmosensing by bacteria: signals and membrane-based sensors. *Microbiology*
1254         *and molecular biology reviews : MMBR* **63**, 230-262 (1999).

1255   51   Chow, V., Nong, G. & Preston, J. F. Structure, function, and regulation of the aldouronate
1256         utilization gene cluster from Paenibacillus sp. strain JDR-2. *Journal of bacteriology* **189**, 8863-
1257         8870, doi:10.1128/JB.01141-07 (2007).

1258   52   Mansilla, M. C., Cybulski, L. E., Albanesi, D. & de Mendoza, D. Control of membrane lipid
1259         fluidity by molecular thermosensors. *Journal of bacteriology* **186**, 6681-6688,
1260         doi:10.1128/JB.186.20.6681-6688.2004 (2004).

1261   53   Kupferschmied, P., Pechy-Tarr, M., Imperiali, N., Maurhofer, M. & Keel, C. Domain shuffling
1262         in a sensor protein contributed to the evolution of insect pathogenicity in plant-beneficial
1263         Pseudomonas protegens. *PLoS pathogens* **10**, e1003964, doi:10.1371/journal.ppat.1003964
1264         (2014).

1265   54   Meineke, E. K., Dunn, R. R., Sexton, J. O. & Frank, S. D. Urban warming drives insect pest
1266         abundance on street trees. *PloS one* **8**, e59687, doi:10.1371/journal.pone.0059687 (2013).

1267   55   Burland, S. M. & Edwards, E. A. Anaerobic benzene biodegradation linked to nitrate
1268         reduction. *Applied and environmental microbiology* **65**, 529-533 (1999).

1269   56   Lieberman, R. L. & Rosenzweig, A. C. Biological methane oxidation: regulation, biochemistry,
1270         and active site structure of particulate methane monooxygenase. *Critical reviews in*
1271         *biochemistry and molecular biology* **39**, 147-164, doi:10.1080/10409230490475507 (2004).

1272   57   Wei, J. *et al.* Cysteine biosynthetic enzymes are the pieces of a metabolic energy pump.
1273         *Biochemistry* **41**, 8493-8498 (2002).

1274    58    Sekowska, A., Kung, H. F. & Danchin, A. Sulfur metabolism in Escherichia coli and related

1275          bacteria: facts and fiction. *Journal of molecular microbiology and biotechnology* **2**, 145-177

1276          (2000).

1277    59    Chakraborty, A. *et al.* Changing bacterial profile of Sundarbans, the world heritage

1278          mangrove: impact of anthropogenic interventions. *World journal of microbiology &*

1279          *biotechnology* **31**, 593-610, doi:10.1007/s11274-015-1814-5 (2015).

1280    60    Elbein, A. D., Pan, Y. T., Pastuszak, I. & Carroll, D. New insights on trehalose: a

1281          multifunctional molecule. *Glycobiology* **13**, 17R-27R, doi:10.1093/glycob/cwg047 (2003).

1282    61    Rojo, F. Degradation of alkanes by bacteria. *Environmental microbiology* **11**, 2477-2490,

1283          doi:10.1111/j.1462-2920.2009.01948.x (2009).

1284    62    Garrett, T. R., Bhakoo, M. & Zhang, Z. B. Bacterial adhesion and biofilms on surfaces. *Prog*

1285          *Nat Sci* **18**, 1049-1056, doi:10.1016/j.pnsc.2008.04.001 (2008).

1286    63    Ullmann, R., Gross, R., Simon, J., Unden, G. & Kroger, A. Transport of C(4)-dicarboxylates in

1287          Wolinella succinogenes. *Journal of bacteriology* **182**, 5757-5764 (2000).

1288    64    Suzuki, K. I. & Hamada, M. Microbacterium. *Bergey's Manual of Systematics of Archaea and*

1289          *Bacteria*, 1-52, doi:10.1002/9781118960608.gbm00104 (2015).

1290    65    Chubiz, L. M., Glekas, G. D. & Rao, C. V. Transcriptional cross talk within the mar-sox-rob

1291          regulon in Escherichia coli is limited to the rob and marRAB operons. *Journal of bacteriology*

1292          **194**, 4867-4875, doi:10.1128/JB.00680-12 (2012).

1293    66    Elsen, S., Swem, L. R., Swem, D. L. & Bauer, C. E. RegB/RegA, a highly conserved redox-

1294          responding global two-component regulatory system. *Microbiology and molecular biology*

1295          *reviews : MMBR* **68**, 263-279, doi:10.1128/MMBR.68.2.263-279.2004 (2004).

1296    67    Lochowska, A. *et al.* Regulation of sulfur assimilation pathways in Burkholderia cenocepacia

1297          through control of genes by the SsuR transcription factor. *Journal of bacteriology* **193**, 1843-

1298          1853, doi:10.1128/JB.00483-10 (2011).

1299    68    Yang, S. *et al.* Hydrocarbon degraders establish at the costs of microbial richness, abundance

1300          and keystone taxa after crude oil contamination in permafrost environments. *Scientific*

1301          *Reports* **6**, 37473, doi:10.1038/srep37473 (2016).

1302    69    Mori, K., Yamaguchi, K., Sakiyama, Y., Urabe, T. & Suzuki, K. Caldisericum exile gen. nov., sp.

1303          nov., an anaerobic, thermophilic, filamentous bacterium of a novel bacterial phylum,

1304          Caldiserica phyl. nov., originally called the candidate phylum OP5, and description of

1305          Caldisericaceae fam. nov., Caldisericales ord. nov. and Caldisericia classis nov. *International*

1306          *journal of systematic and evolutionary microbiology* **59**, 2894-2898,

1307          doi:10.1099/ijs.0.010033-0 (2009).

1308    70    Childers, S. E., Ciufo, S. & Lovley, D. R. Geobacter metallireducens accesses insoluble Fe(III)

1309          oxide by chemotaxis. *Nature* **416**, 767-769, doi:10.1038/416767a (2002).

1310    71    Macy, J. M. *et al.* Thauera selenatis gen. nov., sp. nov., a member of the beta subclass of

1311          Proteobacteria with a novel type of anaerobic respiration. *International journal of systematic*

1312          *bacteriology* **43**, 135-142, doi:10.1099/00207713-43-1-135 (1993).

1313    72    Daims, H. *et al.* Complete nitrification by Nitrospira bacteria. *Nature* **528**, 504-509,

1314          doi:10.1038/nature16461 (2015).

1315    73    Kuever, J. The Family Desulfobulbaceae. *The Prokaryotes: Deltaproteobacteria and*

1316          *Epsilonproteobacteria*, 75-86, doi:10.1007/978-3-642-39044-9_267 (2014).

1317    74    Warren, L. A., Kendra, K. E., Brady, A. L. & Slater, G. F. Sulfur Biogeochemistry of an Oil Sands

1318          Composite Tailings Deposit. *Frontiers in microbiology* **6**, 1533,

1319          doi:10.3389/fmicb.2015.01533 (2015).

1320    75    Penner, T. J. & Foght, J. M. Mature fine tailings from oil sands processing harbour diverse

1321          methanogenic communities. *Canadian journal of microbiology* **56**, 459-470,

1322          doi:10.1139/w10-029 (2010).

1323    76    Xu, Z., Hansen, M. A., Hansen, L. H., Jacquiod, S. & Sorensen, S. J. Bioinformatic approaches

1324          reveal metagenomic characterization of soil microbial community. *PloS one* **9**, e93445,

1325          doi:10.1371/journal.pone.0093445 (2014).

1326    77    Barret, M. *et al.* Emergence shapes the structure of the seed microbiota. *Applied and*

1327          *environmental microbiology* **81**, 1257-1266, doi:10.1128/AEM.03722-14 (2015).

1328    78    Bell, T. H. *et al.* Predictable bacterial composition and hydrocarbon degradation in Arctic

1329          soils following diesel and nutrient disturbance. *The ISME journal* **7**, 1200-1210,

1330          doi:10.1038/ismej.2013.1 (2013).

1331    79    Sun, W. *et al.* Microbial communities inhabiting oil-contaminated soils from two major

1332          oilfields in Northern China: Implications for active petroleum-degrading capacity. *Journal of*

1333          *microbiology* **53**, 371-378, doi:10.1007/s12275-015-5023-6 (2015).

1334    **Legends**

1335    **Figure 1. Taxonomic composition of bacterial communities in Noonmati and Barhola oil**

1336    **contaminated soil.** Taxonomic cladogram showing all taxa detected at a relative abundance

1337    ≥ 0.5% in at least one of the samples. The four rings of the cladogram represent phyla

1338    (innermost), class, order and family (outermost) respectively. Circles in the cladogram depict

1339    detected taxonomic clades and are colored according to corresponding phyla. Outermost

1340    circular rings (external to the cladogram), show rectangular heatmaps depicting abundance

1341    of corresponding families in the cladogram.  Increasing abundance is represented by

1342    increasing opacity.

1343    **Figure 2. Taxonomic distribution of bacterial communities in oil contaminated**

1344    **environments.** Taxonomic clades detected at an average relative abundance ≥ 2% in at least

1345    one of 12 oil contaminated habitats, at the phylum level (A), and at the order level (B).

1346 **Figure 3. Non-metric multidimensional scaling (NMDS) plot of taxonomic composition of**

1347 **all oil contaminated samples of all habitats.** NMDS ordination of 65 oil contaminated

1348 samples across 12 habitats was carried out based on Bray-Curtis similarity distances

1349 calculated from pairwise taxonomic profile comparisons between all samples. Taxonomic

1350 clades present in at least one sample at a relative abundance ≥ 0.5% were used as input. A

1351 shorter linear distance between two samples denote greater similarity between the

1352 corresponding samples. Samples from 12 environments are depicted by different colors.

1353 **Figure 4. Taxonomic biomarkers of bacterial communities from oil polluted habitats.**

1354 Cladogram showing all taxonomic clades detected at a relative abundance ≥ 0.5% in at least

1355 five samples across all habitats. These were used as inputs for LEfSe. Seven rings of the

1356 cladogram represent phylum (innermost), class, order, family, genus and species

1357 (outermost), respectively. Enlarged circles represent differentially abundant taxa detected

1358 as taxonomic biomarkers and are colored corresponding to the individual soil habitat

1359 wherein they are over-represented among 12 oil polluted ecosystems (see legend).

1360 **Figure 5. Metabolic reconstruction of metagenomes from oil polluted habitats.** Cladogram

1361 showing KEGG BRITE hierarchical structures denoted by innermost four rings as inferred

1362 against detected KEGG metabolic modules for all oil contaminated samples. Outermost ring

1363 represents KEGG functional modules that have been detected in at least one of the 65

1364 PICRUSt predicted metagenomes as reconstructed by HUMAnN2. Over-represented

1365 metabolic modules inferred by LEfSe are depicted by enlarged circles and are colored

1366 corresponding to the oil contaminated habitat they have been identified to be differentially

1367 abundant in. Outermost rings (external to the cladogram) depict heat-maps representing

1368 significantly or non-significantly present/absent modules across all oil polluted

1369    environments. Presence is represented by ≥ 90% coverage and absence by ≤ 10% coverage

1370    of KEGG module as estimated by HUMAnN2. Varied modules are labeled and included in

1371    legend.


1372    **Figure 6. SparCC network plot of global microbial interactions in individual oil polluted**

1373    **habitats.** Significant bacterial associations captured by SparCC ($p$-value < 0.01) with an

1374    absolute correlation magnitude of ≥ 0.6 are presented. Nodes represent detected

1375    phylotypes (OTU clustered at 97% similarity) involved in either significant co-occurrence

1376    (green edges) or co-exclusion (red edges) relationships. Border coloration depicts taxonomic

1377    affiliation of nodes at the phylum level. Node size is proportional to the connectivity of the

1378    node (both positive and negative relationships).

**Table 1. Summary of datasets used in the study. (For additional details, refer to Supplementary data Table S1).**

| Biome Type | ID | Sequencing Platform | Original ID | Location | Depth of sample collection (cm below surface) | Source material for sequencing | Predominant contaminant/hydrocarbon | Reference |
|---|---|---|---|---|---|---|---|---|
| Urban | I1 | 454 GS Junior | Noonmati_Surface | Guwahati, Assam, India | 0-10 | *In situ* soil | Crude oil | This study |
| Urban | I2 | 455 GS Junior | Noonmati_Deep | Guwahati, Assam, India | 20-30 | *In situ* soil | Crude oil | This study |
| Urban | I3 | 456 GS Junior | Barhola_Surface | Barhola, Assam, India | 0-10 | *In situ* soil | Crude oil | This study |
| Urban | I4 | 457 GS Junior | Barhola_Deep | Barhola, Assam, India | 20-30 | *In situ* soil | Crude oil | This study |
| Arctic | A1 | Ion Torrent PGM | AH1d1, AH1d2 | Axel Heiburg, Nunavut, Canada | 0-15 | Treated microcosm sediment | Diesel oil | 78 |
| Arctic | A2 | Ion Torrent PGM | AK1d1, AK1d2, AK1d3 | Toolik Lake, Alaska, USA | 0-15 | Treated microcosm sediment | Diesel oil | 78 |

| Arctic | A3 | Ion Torrent PGM | AL1d1, AL1d2, AL1d3 | Alert, Nunavut, Canada | 0-15 | Treated microcosm sediment | Diesel oil | [78] |
|--------|-----|-----------------|---------------------|-------------------------|------|----------------------------|-----------|------|
| Arctic | A4 | Ion Torrent PGM | AVKd1, AVKd2, AVKd3 | Akulivik, Quebec, Canada | 0-15 | Treated microcosm sediment | Diesel oil | [78] |
| Arctic | A5 | Ion Torrent PGM | BDEd1, BDEd2, BDEd3 | Baie Deception, Quebec, Canada | 0-15 | Treated microcosm sediment | Diesel oil | [78] |
| Arctic | A6 | Ion Torrent PGM | BY1d1, BY1d2, BY1d3 | Bylot Island, Nunavut, Canada | 0-15 | Treated microcosm sediment | Diesel oil | [78] |
| Arctic | A7 | Ion Torrent PGM | EBAd1, EBAd2, EBAd3 | East Bay, Nunavut, Canada | 0-15 | Treated microcosm sediment | Diesel oil | [78] |
| Arctic | A8 | Ion Torrent PGM | IQAd1, IQAd2, IQAd3 | Iqaluit, Nunavut, Canada | 0-15 | Treated microcosm sediment | Diesel oil | [78] |
| Arctic | A9 | Ion Torrent PGM | NORd1, NORd2, NORd3 | Tromso, Norway | 0-15 | Treated microcosm sediment | Diesel oil | [78] |
| Arctic | A10 | Ion Torrent PGM | RANd1, RANd2, RANd3 | Rankin Inlet, Nunavut, Canada | 0-15 | Treated microcosm sediment | Diesel oil | [78] |
| Arctic | A11 | Ion Torrent PGM | RUSd1, RUSd2, RUSd3 | Yamal, Russia | 0-15 | Treated microcosm sediment | Diesel oil | [78] |

| Arctic | A12 | Ion Torrent PGM | THUd1, THUd2, THUd3 | Thule, Greenland | 0-15 | Treated microcosm sediment | Diesel oil | [78] |
|--------|-----|-----------------|---------------------|------------------|------|----------------------------|------------|------|
| Urban | C1 | Illumina Miseq | CQM1 | Changqing, China | 2-10 | *In situ* soil | Crude oil | [79] |
| Urban | C2 | Illumina Miseq | CQM2 | Changqing, China | 2-10 | *In situ* soil | Crude oil | [79] |
| Urban | C3 | Illumina Miseq | CQM3 | Changqing, China | 2-10 | *In situ* soil | Crude oil | [79] |
| Urban | C4 | Illumina Miseq | DQM3 | Daqing, China | 2-10 | *In situ* soil | Crude oil | [79] |
| Urban | C5 | Illumina Miseq | DQM4 | Daqing, China | 2-10 | *In situ* soil | Crude oil | [79] |
| Urban | C6 | Illumina Miseq | DQM5 | Daqing, China | 2-10 | *In situ* soil | Crude oil | [79] |
| Urban | C7 | Illumina Miseq | DQM12 | Daqing, China | 2-10 | *In situ* soil | Crude oil | [79] |
| Urban | C8 | Illumina Miseq | DQM50 | Daqing, China | 2-10 | *In situ* soil | Crude oil | [79] |
| Urban | C9 | Illumina Miseq | DQM200 | Daqing, China | 2-10 | *In situ* soil | Crude oil | [79] |
| Mangrove | M1 | 454 GS FLX | T23 2% I, T23 2% II | Restinga da Marambaia, Rio de Janeiro, Brazil | 0-20 | Treated microcosm sediment | Crude oil | [43] |
| Mangrove | M2 | 455 GS FLX | T66 2% I, T66 2% II | Restinga da Marambaia, Rio de | 0-20 | Treated microcosm sediment | Crude oil | [43] |

| | | | | Janeiro, Brazil | | | | |
|---|---|---|---|---|---|---|---|---|
| Mangrove | M3 | 456 GS FLX | T23 5% I, T23 5% II | Restinga da Marambaia, Rio de Janeiro, Brazil | 0-20 | Treated microcosm sediment | Crude oil | 43 |
| Marine sediment | DWH1 | Illumina | SE-20101001-GY-LBNL1-BC-120 | Gulf of Mexico | 0-1# | *In situ* soil | Crude oil | 47 |
| Marine sediment | DWH2 | Illumina | SE-20101001-GY-ALTNF001-BC-139 | Gulf of Mexico | 0-1# | *In situ* soil | Crude oil | 47 |
| Marine sediment | DWH3 | Illumina | SE-20101001-GY-NF006MOD-BC-143 | Gulf of Mexico | 0-1# | *In situ* soil | Crude oil | 47 |
| Marine sediment | DWH4 | Illumina | SE-20101017-GY-D031S-BC-278 | Gulf of Mexico | 0-1# | *In situ* soil | Crude oil | 47 |
| Marine sediment | DWH5 | Illumina | SE-20101017-GY-D040S-BC-315 | Gulf of Mexico | 0-1# | *In situ* soil | Crude oil | 47 |
| Marine sediment | DWH6 | Illumina | SE-20101017-GY-D038SW-BC-331 | Gulf of Mexico | 0-1# | *In situ* soil | Crude oil | 47 |
| Marine sediment | DWH7 | Illumina | SE-20101017-GY-D042S-BC-350 | Gulf of Mexico | 0-1# | *In situ* soil | Crude oil | 47 |
| Taiga | Tu1 | 456 GS FLX | 1330 | Walagan, China | 20-30 | Treated microcosm sediment | Crude oil | 18 |
| Taiga | Tu2 | 456 GS FLX | 2330 | Walagan North, China | 20-30 | Treated microcosm sediment | Crude oil | 18 |

| Taiga | Tu3 | 456 GS FLX | 3330 | Taiyuan, China | 20-30 | Treated microcosm sediment | Crude oil | [18] |
|-------|-----|------------|------|----------------|-------|---------------------------|-----------|------|
| Taiga | Tu4 | 456 GS FLX | 4330 | Jiagedaqi, China | 20-30 | Treated microcosm sediment | Crude oil | [18] |
| Taiga | Tb1 | 456 GS FLX | 1530 | Walagan, China | 70-80 | Treated microcosm sediment | Crude oil | [18] |
| Taiga | Tb2 | 456 GS FLX | 2530 | Walagan North, China | 70-80 | Treated microcosm sediment | Crude oil | [18] |
| Taiga | Tb3 | 456 GS FLX | 3530 | Taiyuan, China | 70-80 | Treated microcosm sediment | Crude oil | [18] |
| Taiga | Tb4 | 456 GS FLX | 4530 | Jiagedaqi, China | 70-80 | Treated microcosm sediment | Crude oil | [18] |
| Taiga | Tp1 | 456 GS FLX | 1630 | Walagan, China | 140-150 | Treated microcosm sediment | Crude oil | [18] |
| Taiga | Tp2 | 456 GS FLX | 2630 | Walagan North, China | 140-150 | Treated microcosm sediment | Crude oil | [18] |
| Taiga | Tp3 | 456 GS FLX | 3630 | Taiyuan, China | 140-150 | Treated microcosm sediment | Crude oil | [18] |
| Taiga | Tp4 | 456 GS FLX | 4630 | Jiagedaqi, China | 140-150 | Treated microcosm sediment | Crude oil | [18] |

| Arctic | OSC1 | 456 GS FLX | 2010Suncor Oil Sands Core Run 11 Subsample 1 | Alberta, Canada | 2,985-2,990 | *In situ* soil | Oil sands bitumen | [17] |
|---|---|---|---|---|---|---|---|---|
| Arctic | OSC3 | 456 GS FLX | 2010Suncor Oil Sands Core Run 11 Subsample 3 | Alberta, Canada | 2,985-2,990 | *In situ* soil | Oil sands bitumen | [17] |
| Arctic | OSC4 | 456 GS FLX | 2010Suncor Oil Sands Core Run 11 Subsample 4 | Alberta, Canada | 2,985-2,990 | *In situ* soil | Oil sands bitumen | [17] |
| Arctic | OSC5 | 456 GS FLX | 2010Suncor Oil Sands Core Run 11 Subsample 5 | Alberta, Canada | 2,985-2,990 | *In situ* soil | Oil sands bitumen | [17] |
| Arctic | OSC7 | 456 GS FLX | 2010Suncor Oil Sands Core Run 11 Subsample 7 | Alberta, Canada | 2,985-2,990 | *In situ* soil | Oil sands bitumen | [17] |
| Arctic | OSC9 | 456 GS FLX | 2010Suncor Oil Sands Core Run 11 Subsample 9 | Alberta, Canada | 2,985-2,990 | *In situ* soil | Oil sands bitumen | [17] |
| Arctic | OSC12 | 456 GS FLX | 2010Suncor Oil Sands Core Run 11 Subsample 13 | Alberta, Canada | 2,985-2,990 | *In situ* soil | Oil sands bitumen | [17] |
| Arctic | OSTPu1 | 456 GS FLX | 2009Suncor Tailings Pond 5 (6.5 ft) | Alberta, Canada | 200 | *In situ* soil | Oil sands bitumen | [17] |
| Arctic | OSTPu2 | 456 GS FLX | 2009Suncor Tailings Pond 5 (8.0 ft) | Alberta, Canada | 240 | *In situ* soil | Bitumen and various other hydrocarbons | [17] |
| Arctic | OSTPu3 | 456 GS FLX | 2010Syncrude Tailings Pond MLSB M MFT - 1m | Alberta, Canada | 100 | *In situ* soil | Bitumen and various other hydrocarbons | [17] |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Arctic | OSTPu4 | 456 GS FLX | 2010Syncrude Tailings Pond MLSB N MFT - 1.1m | Alberta, Canada | 110 | *In situ* soil | Bitumen and various other hydrocarbons | 17 |
| Arctic | OSTPm2 | 456 GS FLX | 2009Suncor Tailings Pond 5 (25 ft) | Alberta, Canada | 750 | *In situ* soil | Bitumen and various other hydrocarbons | 17 |
| Arctic | OSTPm4 | 456 GS FLX | 011Suncor Tailings Pond 6 (7 m) | Alberta, Canada | 700 | *In situ* soil | Bitumen and various other hydrocarbons | 17 |
| Arctic | OSTPm6 | 456 GS FLX | 2010Syncrude Tailings Pond MLSB N MFT - 6.1m | Alberta, Canada | 610 | *In situ* soil | Bitumen and various other hydrocarbons | 17 |
| Arctic | OSTPd1 | 456 GS FLX | 2009Suncor Tailings Pond 5 (40 ft) | Alberta, Canada | 1220 | *In situ* soil | Bitumen and various other hydrocarbons | 17 |
| Arctic | OSTPd2 | 456 GS FLX | 2009SuncorTailings Pond 5 (45 ft) | Alberta, Canada | 1370 | *In situ* soil | Bitumen and various other hydrocarbons | 17 |
| Arctic | OSTPd3 | 456 GS FLX | 2010Suncor Tailings Pond 6 (12 m) | Alberta, Canada | 1200 | *In situ* soil | Bitumen and various other hydrocarbons | 17 |
| Arctic | OSTPd4 | 456 GS FLX | 2011Suncor Tailings Pond 6 (13 m) | Alberta, Canada | 1300 | *In situ* soil | Bitumen and various other hydrocarbons | 17 |

#All samples collected at an average of ~1500 metres below sea level, depth given is from surface of ocean floor marine sediment.

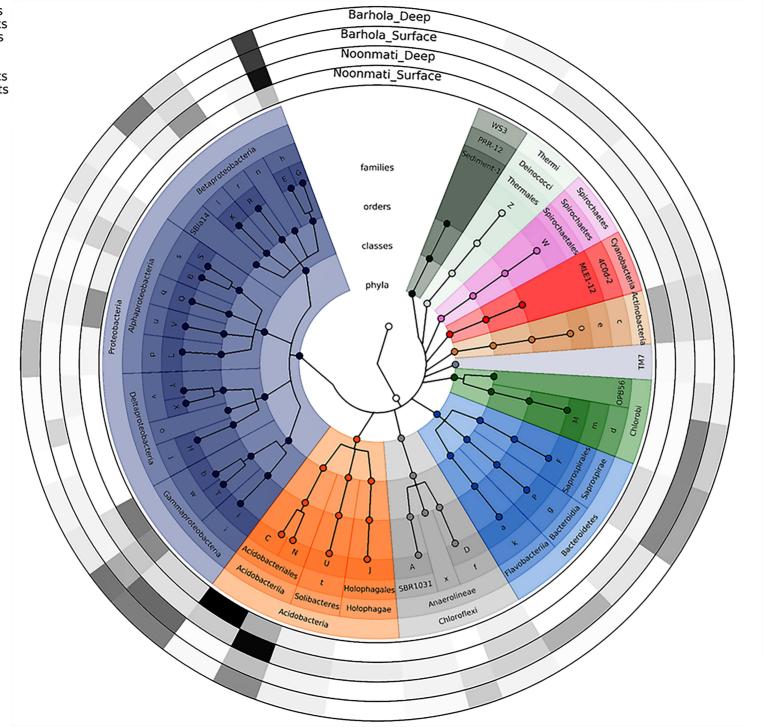**Table 2. Similarities of bacterial community structure within a habitat and between pairs of habitats.**

| Habitat | India oil refineries | Arctic | China oil refineries | Mangrove | Marine sediments | Taiga upper active layer | Taiga bottom active layer | Taiga permafrost layer | Oil sands core | Oil sands tailings pond upper | Oil sands tailings pond median | Oil sands tailings pond deep |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| India oil refineries | **0.63 ± 0.06** | 0.51 ± 0.07 | 0.39 ± 0.03 | 0.44 ± 0.04 | 0.47 ± 0.04 | 0.41 ± 0.06 | 0.41 ± 0.04 | 0.36 ± 0.05 | 0.48 ± 0.04 | 0.46 ± 0.08 | 0.46 ± 0.08 | 0.48 ± 0.08 |
| Arctic | 0.51 ± 0.07 | **0.72 ± 0.07** | 0.48 ± 0.05 | 0.46 ± 0.03 | 0.42 ± 0.06 | 0.45 ± 0.07 | 0.41 ± 0.05 | 0.39 ± 0.05 | 0.49 ± 0.05 | 0.39 ± 0.05 | 0.39 ± 0.04 | 0.41 ± 0.06 |
| China oil refineries | 0.39 ± 0.03 | 0.48 ± 0.05 | **0.69 ± 0.08** | 0.48 ± 0.02 | 0.40 ± 0.05 | 0.38 ± 0.06 | 0.35 ± 0.05 | 0.34 ± 0.06 | 0.37 ± 0.06 | 0.32 ± 0.02 | 0.35 ± 0.04 | 0.34 ± 0.05 |
| Mangrove | 0.44 ± 0.04 | 0.46 ± 0.03 | 0.48 ± 0.02 | **0.83 ± 0.02** | 0.54 ± 0.05 | 0.34 ± 0.03 | 0.34 ± 0.02 | 0.31 ± 0.02 | 0.36 ± 0.02 | 0.41 ± 0.02 | 0.43 ± 0.04 | 0.41 ± 0.05 |
| Marine sediments | 0.47 ± 0.04 | 0.42 ± 0.06 | 0.40 ± 0.05 | 0.54 ± 0.05 | **0.77 ± 0.09** | 0.35 ± 0.05 | 0.35 ± 0.02 | 0.33 ± 0.05 | 0.43 ± 0.02 | 0.38 ± 0.03 | 0.39 ± 0.03 | 0.42 ± 0.04 |
| Taiga upper active layer | 0.41 ± 0.06 | 0.45 ± 0.07 | 0.38 ± 0.06 | 0.34 ± 0.03 | 0.35 ± 0.05 | **0.52 ± 0.18** | 0.59 ± 0.17 | 0.52 ± 0.18 | 0.45 ± 0.09 | 0.34 ± 0.06 | 0.35 ± 0.05 | 0.37 ± 0.07 |
| Taiga bottom active layer | 0.41 ± 0.04 | 0.41 ± 0.05 | 0.35 ± 0.05 | 0.34 ± 0.02 | 0.35 ± 0.02 | 0.59 ± 0.17 | **0.57 ± 0.12** | 0.55 ± 0.20 | 0.45 ± 0.05 | 0.33 ± 0.04 | 0.34 ± 0.04 | 0.37 ± 0.06 |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Taiga permafrost layer | 0.36 ± 0.05 | 0.39 ± 0.05 | 0.34 ± 0.06 | 0.31 ± 0.02 | 0.33 ± 0.05 | 0.52 ± 0.18 | 0.55 ± 0.20 | **0.45 ± 0.22** | 0.42 ± 0.10 | 0.32 ± 0.06 | 0.33 ± 0.05 | 0.36 ± 0.08 |
| Oil sands core | 0.48 ± 0.04 | 0.49 ± 0.05 | 0.37 ± 0.06 | 0.36 ± 0.02 | 0.43 ± 0.02 | 0.45 ± 0.09 | 0.45 ± 0.05 | 0.42 ± 0.10 | **0.85 ± 0.09** | 0.45 ± 0.05 | 0.45 ± 0.04 | 0.56 ± 0.10 |
| Oil sands tailings pond upper | 0.46 ± 0.08 | 0.39 ± 0.05 | 0.32 ± 0.02 | 0.41 ± 0.02 | 0.38 ± 0.03 | 0.34 ± 0.06 | 0.33 ± 0.04 | 0.32 ± 0.06 | 0.45 ± 0.05 | **0.67 ± 0.12** | 0.64 ± 0.09 | 0.63 ± 0.13 |
| Oil sands tailings pond median | 0.46 ± 0.08 | 0.39 ± 0.04 | 0.35 ± 0.04 | 0.43 ± 0.04 | 0.39 ± 0.03 | 0.35 ± 0.05 | 0.34 ± 0.04 | 0.33 ± 0.05 | 0.45 ± 0.04 | 0.64 ± 0.09 | **0.56 ± 0.05** | 0.61 ± 0.12 |
| Oil sands tailings pond deep | 0.48 ± 0.08 | 0.41 ± 0.06 | 0.34 ± 0.05 | 0.41 ± 0.05 | 0.42 ± 0.04 | 0.37 ± 0.07 | 0.37 ± 0.06 | 0.36 ± 0.08 | 0.56 ± 0.10 | 0.63 ± 0.13 | 0.61 ± 0.12 | **0.61 ± 0.15** |

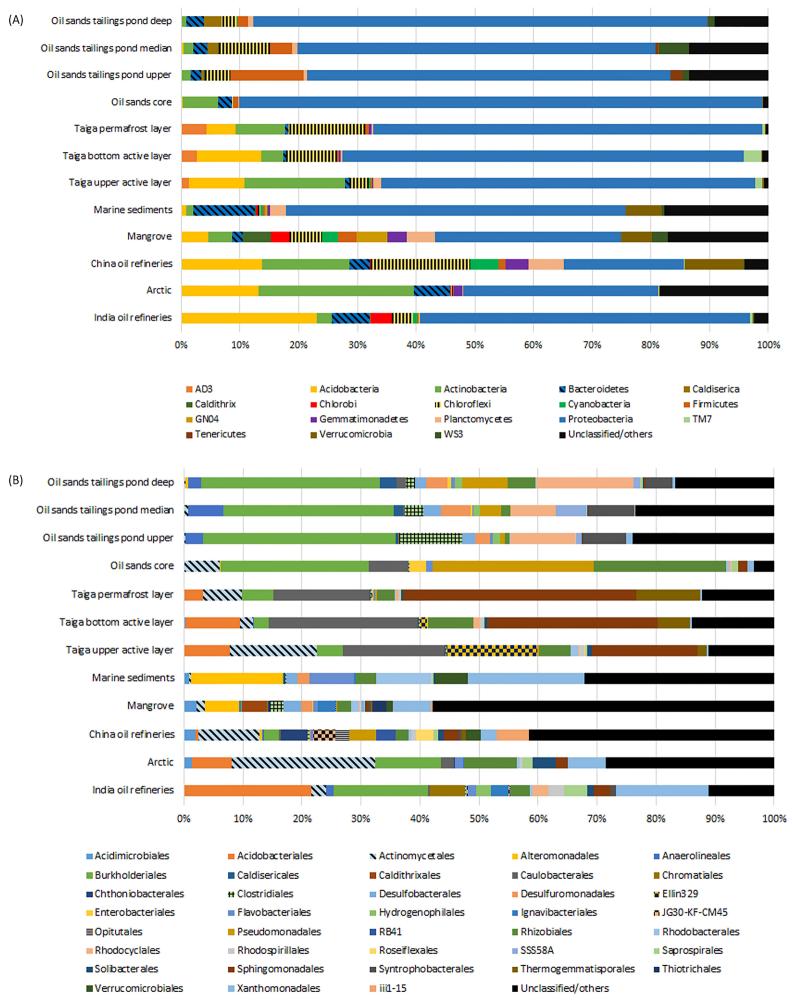**Table 3: Core modules shared between habitats as detected by HUMAnN2.**

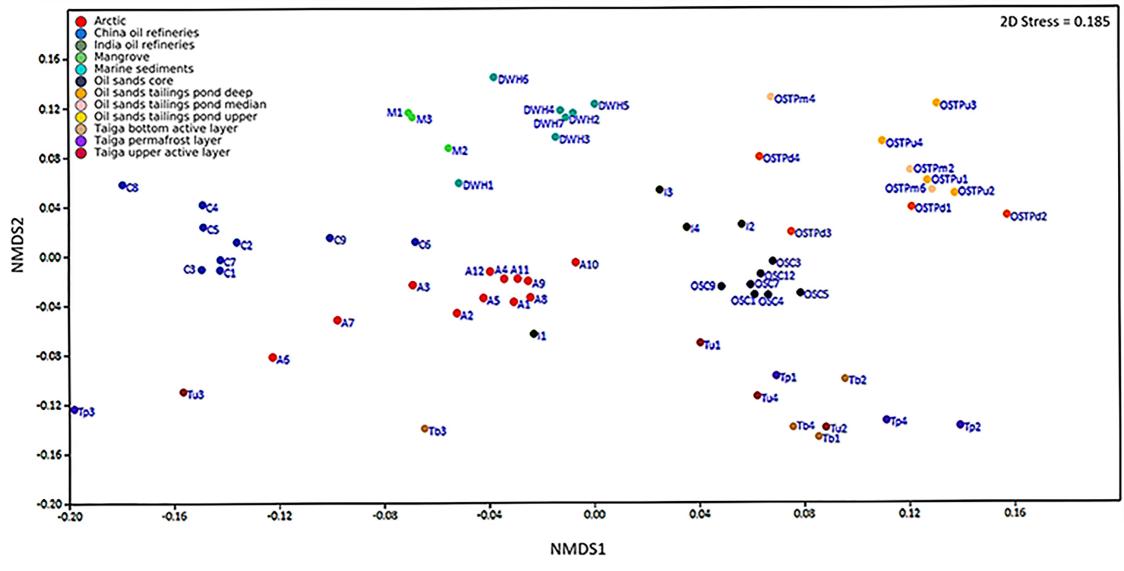| Module ID | Definition of modules in KEGG |
|---|---|
| M00005 | PRPP biosynthesis, ribose 5P => PRPP |
| M00020 | Serine biosynthesis, glycerate-3P => serine |
| M00149 | Succinate dehydrogenase, prokaryotes |
| M00153 | Cytochrome d ubiquinol oxidase |
| M00157 | F-type ATPase, prokaryotes and chloroplasts |
| M00178 | Ribosome, bacteria |
| M00188 | NitT/TauT family transport system |
| M00222 | Phosphate transport system |
| M00223 | Phosphonate transport system |
| M00236 | Putative polar amino acid transport system |
| M00237 | Branched-chain amino acid transport system |
| M00239 | Peptides/nickel transport system |
| M00240 | Iron complex transport system |
| M00250 | Lipopolysaccharide transport system |
| M00254 | ABC-2 type transport system |
| M00255 | Lipoprotein-releasing system |
| M00256 | Cell division transport system |
| M00258 | Putative ABC transport system |
| M00320 | Lipopolysaccharide export system |

(A)

| Legend (A) | | |
|---|---|---|
| AD3 | Acidobacteria | Actinobacteria |
| Bacteroidetes | Caldiserica | Caldithrix |
| Chlorobi | Chloroflexi | Cyanobacteria |
| Firmicutes | GN04 | Gemmatimonadetes |
| Planctomycetes | Proteobacteria | TM7 |
| Tenericutes | Verrucomicrobia | WS3 |
| Unclassified/others | | |

(B)

| Legend (B) | | |
|---|---|---|
| Acidimicrobiales | Acidobacteriales | Actinomycetales |
| Alteromonadales | Anaerolineales | Burkholderiales |
| Caldisericales | Caldithrixales | Caulobacterales |
| Chromatiales | Chthoniobacterales | Clostridiales |
| Desulfobacterales | Desulfuromonadales | Ellin329 |
| Enterobacteriales | Flavobacteriales | Hydrogenophilales |
| Ignavibacteriales | JG30-KF-CM45 | Opitutales |
| Pseudomonadales | RB41 | Rhizobiales |
| Rhodobacterales | Rhodocyclales | Rhodospirillales |
| Roseiflexales | SSS58A | Saprospirales |
| Solibacterales | Sphingomonadales | Syntrophobacterales |
| Thermogemmatisporales | Thiotrichales | Verrucomicrobiales |
| Xanthomonadales | iii1-15 | Unclassified/others |

**Legend (top left):**
- Arctic
- China oil refineries
- India oil refineries
- Mangrove
- Marine sediments
- Oil sands core
- Oil sands tailings pond deep
- Oil sands tailings pond median
- Oil sands tailings pond upper
- Taiga bottom active layer
- Taiga permafrost layer
- Taiga upper active layer

a:Sphingomonadaceae
b:Acetobacteraceae
c:Rhodospirillaceae
d:Rhodobacteraceae
e:Caulobacteraceae
f:Geobacteraceae
g:Desulfuromonadaceae
h:Syntrophorhabdaceae
i:Syntrophaceae
j:Myxococcaceae
k:Desulfobulbaceae
l:.Opitutaceae
m:Chthoniobacteraceae
n:Verrucomicrobiaceae
o:Peptococcaceae
p:Chitinophagaceae
q:Weeksellaceae
r:Cytophagaceae
s:Dolo_23
t:Kouleothrixaceae
u:Thermogemmatisporaceae
v:Anaerolinaceae
w:Phycisphaeraceae
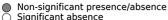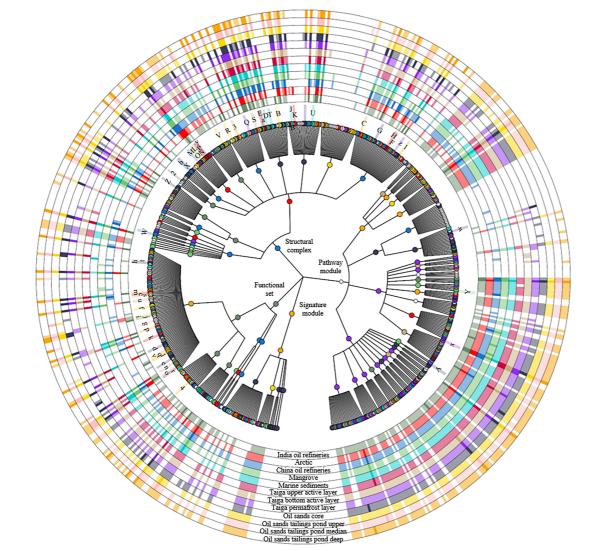x:Gemmataceae
y:Planctomycetaceae
z:Pirellulaceae

A:Spirochaetes
AA:Solirubrobacteraceae
AB:Gaiellaceae
AC:Iamiaceae
AD:Microbacteriaceae
AE:Sporichthyaceae
AF:Nocardiaceae
AG:Micromonosporaceae
AH:Mycobacteriaceae
AI:Nocardioidaceae
AJ:Intrasporangiaceae
AK:Propionibacteriaceae
AL:.Micrococcaceae
AM:Geodermatophilaceae
AN:Dietziaceae
AO:Chlorobi
AP:Ignavibacteriaceae
AQ:Caldiserica
AR:GN04
AS:Nitrospirae
AT:WS3
AU:Gemmatimonadetes
AV:WS6
AW:mb2424
AX:Ellin6075
AY:Koribacteraceae
AZ:Acidobacteriaceae
B:Spirochaetaceae
C:Cyanobacteria
D:Oscillatoriophycideae
E:Synechococcophycideae
F:OP8
G:TM7
H:Colwelliaceae
I:Alteromonadaceae
J:Marinicellaceae
K:Enterobacteriaceae
L:Xanthomonadaceae
M:Sinobacteraceae
N:Pseudomonadaceae
O:Moraxellaceae
P:Alcaligenaceae
Q:Burkholderiaceae
R:Comamonadaceae
S:Oxalobacteraceae
T:Rhodocyclaceae
U:Hydrogenophilaceae
V:Brucellaceae
W:Bradyrhizobiaceae
X:Methylobacteriaceae
Y:Hyphomicrobiaceae
Z:Erythrobacteraceae

**Legend (top left):**
- Arctic
- China oil refineries
- India oil refineries
- Mangrove
- Marine sediments
- Oil sands core
- Oil sands tailings pond deep
- Oil sands tailings pond median
- Oil sands tailings pond upper
- Taiga bottom active layer
- Taiga permafrost layer
- Taiga upper active layer

**Node legend:**
- Non-significant presence/absence
- Significant absence
- Significant presence

- Differential abundance

**Tree labels:**
Structural complex; Pathway module; Functional set; Signature module

**Outer ring labels:**
India oil refineries; Arctic; China oil refineries; Mangrove; Marine sediments; Taiga upper active layer; Taiga bottom active layer; Taiga permafrost layer; Oil sands core; Oil sands tailings pond upper; Oil sands tailings pond median; Oil sands tailings pond deep

**Right-hand legend:**
1: Glycerol transport system
2: Multidrug resistance, efflux pump VexEF-TolC
3: Copper-processing system
4: Multidrug resistance, efflux pump GesABC
A: Histidine degradation, histidine => N-formiminoglutamate => glutamate
B: Tungstate transport system
C: Putative sn-glycerol-phosphate transport system
D: Glycine betaine/proline transport system
E: Osmoprotectant transport system
F: Putative ABC transport system
G: D-Xylose transport system
H: Multiple sugar transport system
I: Glutamate/aspartate transport system
J: General L-amino acid transport system
K: Glutamate transport system
L: Zinc transport system
M: Manganese/iron transport system
N: Putative zinc/manganese transport system
O: Cobalt/nickel transport system
P: Putative ABC transport system
Q: Capsular polysaccharide transport system
R: Lipooligosaccharide transport system
S: Heme transport system
T: Spermidine/putrescine transport system
U: Urea transport system
V: Adhesin protein transport system
W: Bacterial proteasome
X: Microcin C transport system
Y: C10-C20 isoprenoid biosynthesis, archaea
Z: Cytochrome o ubiquinol oxidase
a: Sulfonate transport system
b: Oligopeptide transport system
c: SenX3-RegX3 (phosphate starvation response) two-component regulatory system
d: EnvZ-OmpR (osmotic stress response) two-component regulatory system
e: RstB-RstA two-component regulatory system
f: CreC-CreB (phosphate regulation) two-component regulatory system
g: BaeS-BaeR (envelope stress response) two-component regulatory system
h: QseC-QseB (quorum sensing) two-component regulatory system
i: TctE-TctD (tricarboxylic acid transport) two-component regulatory system
j: NarX-NarL (nitrate respiration) two-component regulatory system
k: DesK-DesR (membrane lipid fluidity regulation) two-component regulatory system
l: AlgZ-AlgR (alginate production) two-component regulatory system
m: GlnL-GlnG (nitrogen regulation) two-component regulatory system
n: NtrY-NtrX (nitrogen regulation) two-component regulatory system
o: HydH-HydG (metal tolerance) two-component regulatory system
p: PilS-PilR (type 4 fimbriae synthesis) two-component regulatory system
q: GlrK-GlrR (amino sugar metabolism) two-component regulatory system
r: DctB-DctD (C4-dicarboxylate transport) two-component regulatory system
s: CckA-CtrA/CpdR (cell cycle control) two-component regulatory system
t: ChvG-ChvI (acidity sensing) two-component regulatory system
u: RegB-RegA (redox response) two-component regulatory system
v: FixL-FixJ (nitrogen fixation) two-component regulatory system
w: Benzoate degradation, benzoate => catechol / methylbenzoate => methylcatechol
x: Nucleotide sugar biosynthesis, galactose => UDP-galactose
y: Biotin transport system
z: Putative aldouronate transport system

A: Arctic
C: China oil refineries
I: India oil refineries
DWH: Marine sediments
OSC: Oil sands core
OSTPd: Oil sands tailings pond deep
OSTPu: Oil sands tailings pond upper
Tb: Taiga bottom active layer
Tp: Taiga permafrost layer
Tu: Taiga upper active layer

Planctomycetes
Proteobacteria
Verrucomicrobia
Cyanobacteria
Gemmatimonadetes
Chloroflexi
AD3
Spirochaetes
Chlorobi
Thermotogae
AC1
Firmicutes
Acidobacteria
Actinobacteria
Bacteroidetes