1    **ORIGINAL ARTICLE**

2    **A new multi-genomic approach for the study of biogeochemical cycles at global**
3    **scale: the molecular reconstruction of the sulfur cycle**

4

5    Valerie De Anda[1*], Icoquih Zapata-Peñasco[2], Bruno Contreras-Moreira[3, 4], Augusto Cesar

6    Poot-Hernandez[5] Luis E. Eguiarte[1] and Valeria Souza[1**]

7

8    *[1]Departamento de Ecología Evolutiva, Instituto de Ecología, Universidad Nacional Autónoma de*
9    *México, 70-275, México D.F. 04510*
10   *[2]Dirección de Investigación en Transformación de Hidrocarburos. Instituto Mexicano del Petróleo,*
11   *Eje Central Lázaro Cárdenas, Norte 152, Col. San Bartolo Atepehuacan, 07730, México*
12   *[3]Estación Experimental de Aula Dei, Consejo Superior de Investigaciones Científicas (EEAD-CSIC),*
13   *Avda. Montañana, 1005, Zaragoza 50059, Spain*
14   *[4]Fundación ARAID, calle María de Luna 11, 50018 Zaragoza, Spain.*
15   *[5]Departamento de Ingeniería Celular y Biocatálisis, Instituto de Biotecnología, Universidad*
16   *Nacional Autónoma de México, Cuernavaca, Morelos 62210, México*

17

18   *For correspondence. *E-mail valdeanda@ciencias.unam.mx **E-mail souza@unam.mx ; Tel. 925*
19   *423 2284; Fax 925 422 3160.*
20

21   **Summary**

22   Despite the great advances in microbial ecology and the explosion of high throughput

23   sequencing, our ability to understand and integrate the global biogeochemical cycles is still

24   limited. Here we propose a novel approach to summarize the complexity of the Sulfur cycle

25   based on the minimum ecosystem concept, the microbial mat model and the relative entropy

26   of protein domains involved in S-metabolism. This methodology produces a single value,

27   called the Sulfur Score (SS), which informs about the specific S-related molecular

28   machinery. After curating an inventory of microorganisms, pathways and genes taking part

29   in this cycle, we benchmark the performance of the SS on a collection of 2,107 non-

30   redundant RefSeq genomes, 900 metagenomes from MG-RAST and 35 metagenomes

31   analyzed for the first time. We find that the SS is able to correctly classify microorganisms

32   known to be involved in the S-cycle, yielding an Area Under the ROC Curve of 0.985.

33   Moreover, when sorting environments the top-scoring metagenomes were hydrothermal

34   vents, microbial mats and deep-sea sediments, among others. This methodology can be

35   generalized to the analysis of other biogeochemical cycles or processes. Provided that an

36   inventory of relevant pathways and microorganisms can be compiled, entropy-based scores

37    could be used to detect environmental patterns and informative samples in multi-genomic

38    scale.

39

**INTRODUCTION**

41    Despite their fundamental importance in sustaining life on Earth, understanding the fluxes of

42    fundamental elements (C, H, O, N, S, and P) through the Earth's surface has been challenging

43    for several reasons (Li *et al.*, 2012; Newman and Banfield, 2002). First, the global

44    biogeochemical cycles are enormously complex as they are an interconnected network of

45    biological, chemical and geophysical processes that have been coevolving in the biosphere

46    since the apparition of the first metabolic processes on Earth (~3.8 billion years ago)

47    (Falkowski *et al.*, 2008). Since then, the evolutionary history of life on Earth has been shaped

48    by a complex synergistic cooperation of multispecies a*s*semblages that differ in terms of

49    ecological and metabolic capabilities (Canfield 2005). Secondly, although these assemblages

50    are often spatially and temporally separated, they effectively couple electron transfer (i.e.,

51    redox) reactions that transform elements and energy derived from several abiotic processes

52    (Falkowski *et al.*, 2008). Thirdly, these abiotic processes involve the continuous supply and

53    removal of elements from various Earth surface reservoirs, such as geothermal processes

54    derived from mantle and crust, tectonics and rock weathering, and photochemical processes in

55    the atmosphere including the constant flux of solar energy (Hedges, 1992).

56        As a result of these challenges, the fluxes of fundamental atoms through Earth have

57    been studied and approached using different disciplines that address specific layers of

58    complexity. For instance, geochemistry and atmospheric sciences have been focused on

59    addressing the major abiotic process at global planetary scales, i.e., processes influencing the

60    flux of elements to and from the various Earth surface reservoirs and atmosphere (Canfield et

61    al., 2005, Falkowski *et al.*, 2008).

62        Moreover, microbial ecology has emphasized in understanding the links between

63    microbial catalyzed activities and ecosystem and biogeochemical processes (Morales and

64    Holben, 2011). Current approaches for establishing metabolic relationships *in situ* are based on

65    targeting coding-sequences involved in specific biochemical pathways related to carbon,

66    nitrogen, sulfur, and phosphorus cycling. In this way, DNA or RNA extracted directly from

67   natural environments is sequenced and quantified with conventional PCR and Sanger

68   sequencing or microarray analyses (Loy *et al.*, 2004; Khodadad and Foster, 2012) for example

69   the GeoChip (Tu *et al.*, 2014); or other 'omics' techniques such as metagenomics (Quaiser *et*

70   *al.*, 2011; Swingley *et al.*, 2012; Llorens-Marès *et al.*, 2015) or metatranscriptomics (Stewart *et*

71   *al.*, 2011; Chen *et al.*, 2015).

72        However, despite the great advances in microbial ecology and high throughput

73   sequencing, our ability to understand and integrate the global biogeochemical cycles is still

74   limited. Here, we propose a new comprehensive approach to evaluate, compare, and facilitate

75   the study of such cycles based on the metabolic machinery of the microorganisms responsible

76   for driving element fluxes throughout the Earth. The approach is based mainly on three aspects:

77   i) the minimum ecosystem concept: which considers the properties, forces (outside energy),

78   flow pathways (energy and matter), interactions, and feedback loops or circuits for the flow of

79   matter or energy (Odum,1993); ii)  the microbial mats, which are nearly closed and self-

80   sustaining ecosystems that comprise the major biogeochemical cycles, trophic levels and food

81   webs in a vertical laminate pattern (Bolhuis *et al.*, 2014); and iii) the mathematical

82   rationalization of Kullback-Leibler divergence, also known as relative entropy *H'* (Kullback

83   and Leibler, 1951). Relative entropy has been widely applied in physics, communication theory

84   and statistical inference, and is interpreted as a measure of disorder, information and

85   uncertainty, respectively (Commenges, 2015). Here we use the communication theory concept

86   of H' to summarize the information conveyed by the protein domains (metabolic machinery)

87   encoded by environmental DNA sequences. The application of this measure in biology was

88   originally developed by Stormo and colleagues to identify binding sites that regulate gene

89   transcription sites (Hertz and Stormo, 1999).

90        In this context, the compartmentalization of microbial mats provides clear, natural or

91   arbitrary boundaries that evoke the concept of "minimum ecosystem", which can be delimited

92   in a natural or arbitrary sense (Odum, 1993). Therefore, specific parts of the cycle may be seen

93   as parts of a whole. For instance, the redox level, reduced-oxidized compounds or even genes

94   and enzymes implicated in certain routes can be used as ecosystem boundaries. These

95   assemblies set up a unit that represents the minimum ecosystem with the minimum

96   requirements to be functional, therefore this can be applied to measure the information derived

97   from the complexity inside any biogeochemical cycle.

98    To test and evaluate the performance of our conceptualization, we focused on the

99    biogeochemical Sulfur cycle (from now on S-cycle), due to its critical role in the

100   biogeochemistry of the planet - *i.e.,* respiration of sedimentary organic matter, oxidation state

101   of the atmosphere and oceans, and the composition of seawater (Halevy *et al.,* 2012). Despite

102   the extensive biochemical knowledge of both oxidative and reductive microorganisms (Rabus

103   *et al.,* 2013; Canfield, 2015; Dahl, 2017), there are no studies aiming to integrate all the

104   microbiological and geochemical transformations and their corresponding metabolic pathways

105   of the sulfur cycle.  Our study proposes a general computational approach that can be easily

106   modified and used to compare and measure other biogeochemical cycles. This procedure

107   generates measurable scores to evaluate these cycles and their importance and ecological

108   weight on a global scale.

109

110   **MATERIALS AND METHODS**

111   The computational analysis addressing the different levels of complexity of the S-cycle was

112   divided into four stages illustrated in Figure 1. The corresponding scripts and documentation

113   are available for download from: https://github.com/eead-csic-compbio/metagenome_Pfam_score.

114

115   **STAGE 1: The biogeochemical complexity of S-cycle and 'omic'-datasets.**

116   *Taxonomic representatives of sulfur cycle: the microbial mat model.* According to the

117   minimum ecosystem concept, we consider microbial mats as models of a minimum ecological

118   system (Microbial Mat Model). Based on the metabolic guilds found in microbial mats and

119   other S-derived environments (i.e., hot springs, black smokers, sludge, etc.), we reviewed

120   primary literature and the MetaCyc database (Caspi *et al.*, 2012) to select S-based

121   microorganisms (at genus or species level) with experimental evidence of the physiology of

122   degradation, reduction, oxidation, or disproportionation of organic/inorganic S-compounds.

123   The list of S-based prokaryotes is found in Table S1. The non-redundant list of these S-based

124   microorganisms with fully sequenced genomes (December 2016) was called the Sulfur list

125   (Suli), which currently contains 161 genomes used as input of the pipeline.

126

127   *Random taxonomic representatives* (RList): In order to build negative control sets of organisms

128   that are not particularly enriched on metabolic preferences, 1,000 random lists of

4

129    microorganisms were drawn using the genomic dataset explained below as reference, with the

130    same number of microorganisms included in Suli.

131

132    *Metabolic pathways and genes:* We gathered and classified the metabolic pathways involved in

133    the S-cycle from the primary literature and expert-curated databases KEGG (Kanehisa and

134    Goto, 2000) and MetaCyc (Caspi *et al.*, 2012). This molecular information was integrated into

135    a single database named S̲ulfur cy̲cle (Sucy), which currently contains 152 genes and 48

136    enzyme classification numbers annotated in the Enzyme classification

137    (http://enzyme.expasy.org) (Table S2). The 152 FASTA sequences of the peptides encoded by

138    these genes were downloaded from UniProt (Magrane and Consortium, 2011) and used as the

139    input of the pipeline.

140

141    *Omic datasets.* In order to test, compare and evaluate the importance of the S-cycle in 'omic'

142    datasets, we generated the following databases:

143        i) *Genomic dataset* (Gen): Due to the redundancy of complete genomes deposited in

144    RefSeq (https://www.ncbi.nlm.nih.gov/genome/browse/reference; 4,158 genomes at the time of

145    the analysis, December 21, 2016), we decided to reduce the set of genomic data by using a

146    web-based tool, that uses "genomic similarity scores" (Moreno-Hagelsieb *et al.*, 2013).

147    Selecting values of genomic similarity of 0.95 and a DNA signature of 0.01, we obtained a total

148    of 2,107 non-redundant genomes in FASTA format.

149        ii) *Metagenomic dataset* (Met): We selected metagenomes from the MG-RAST server

150    version 3.6 that met the following conditions: i) publicly available; ii) metadata associated; and

151    iii) environmental samples (isolated from defined environments o features like rivers, soil,

152    biofilms), discarding microbiome host-associated metagenomes (*i.e.,* to human, cow, chicken).

153    In addition, we also included 35 unpublished metagenomes derived from sediment, water and

154    microbial mats from Cuatro Ciénegas (Coahuila, Mexico), which were also submitted to the

155    MG-RAST pipeline.

156        Using the above-mentioned conditions, a total of 935 metagenomes were downloaded in

157    FASTA format (http://api.metagenomics.anl.gov/api.html, coding regions within the reads,

158    December, 2016). Then, we measured the Mean Size Length (MSL) of the peptide coding

159    regions of the Met FASTA files (Figure S1A). Taking into account that the 152 sulfur proteins

160  (Sucy) have lengths ranging from 49 to 1,020 amino acid residues (aa), their detection in

161  metagenomic samples will likely be affected by MSL (Figure S1B). For example, the

162  identification of long proteins harboring several domains (i.e., catalytic, cofactor binding etc.)

163  might be impaired in metagenomes with short MSL.

164       *iii) Fragmented genomic dataset* (GenF). We used the genomic dataset as a reference

165  for benchmarking the detection limits of the protein families. The peptide FASTA-format files

166  of the 2,107 non-redundant genomes were *in silico* sheared into seven categories of increasing

167  fragment length, taking into account the variation in read sizes of the metagenomic dataset

168  (Figure S1A).

169

170  **STAGE 2: Domain composition of the sulfur proteins**

171  We used Interproscan 5.21-60.0 (Jones *et al.*, 2014) to annotate the protein domains encoded in

172  the 152 Sucy genes, according to the Pfam-A database v30 (Finn *et al.*, 2008). In total, 112

173  Pfam domains where identified and subsequently scanned against the 'omic' datasets with

174  HMMER 3.0 (Finn *et al.*, 2011).

175

176  **STAGE 3: Relative entropy and its use to detect informative sulfur-related protein**

177  **domains**

178  In order to obtain an estimate of how protein families are represented in S-based

179  microorganisms, we used a derivative of the Kullback-Leibler divergence (Kullback and

180  Leibler, 1951) — also known as relative entropy *H'(i)* — to measure the difference between

181  probabilities *P* and *Q* (see Eq. 1 below). In this context, *P(i)* represents the frequency of protein

182  domain *i* in Suli genomes (observed frequency), while *Q(i)* represents the frequency in the

183  complete genomic dataset (expected frequency). *H'*, in bits, captures the extent to which a

184  domain informs specifically about sulfur metabolism. *H'* values that are close to 1 correspond

185  to the most informative Pfam domains (enriched among S-based genomes), whereas low *H'*

186  values (close to zero) describe non-informative ones. Negative values correspond to proteins

187  observed less than expected.

188

189  $$H' = P(i) \log_2 \frac{P(i)}{Q(i)} \qquad \textbf{Eq. 1}$$

190 As a control, *H'* was recalculated both in Gen and GenF datasets replacing Suli with 1,000

191 equally sized lists of random-sampled genomes (Rlist). Using these procedures, we evaluated

192 the variation of relative entropy of each Pfam domain in order to i) short-list those that could be

193 used as markers in metagenomic dataset (regardless of their MSL) and ii) to generate a score

194 which could be used to compare the importance of S-metabolism in 'omic'-samples (either in

195 Gen or Met).

196

197 *Clustering of Pfam domains according to their entropy:* The following clustering algorithms

198 were tested: K-Means, Affinity propagation, Mean-shift Spectral, Ward hierarchical,

199 Agglomerative, DBSCAN and Birch. These methods are part of the scikit-learn Machine

200 Learning Python module (http://scikit-learn.org/stable/modules/clustering.html).

201

202 **STAGE 4: The Sulfur Score (*SS*): origin, interpretation, properties and benchmark**

203 We propose to evaluate the importance of biogeochemical S-cycle in 'omic'-datasets using a

204 single metric that we call "Sulfur Score" (*SS*) (Eq. 2). By this approach, sulfur informative

205 protein domains would contribute to higher *SS*, whereas non-informative ones would decrease

206 it. This is an extension of procedures originally developed for the alignment of DNA and

207 protein motifs, in which individual positions are independent and additive, and can be simply

208 summed up to obtain the total weight or information content (Hertz and Stormo, 1999). Instead

209 of aligning sequences, in our context we compare a presence/absence string of Pfam domains,

210 from which a total weight (*SS*) is computed.

211

212 $$SS = \sum_{i=1}^{112} H'$$ **Eq. 2**

213

214 If we compare total *SS* across several genomes or environments, those in which the majority of

215 metabolic pathways of S-cycle are represented will thus have a high *SS*; in contrast, low *SS*

216 values should be expected if proteins involved in the S-cycle are not particularly enriched.

217

218   *Calibration:* We took into account the MSL of each metagenome to compute the *SS*. Briefly;

219   we gathered the entropy values (*H'*) of the 112 Pfams in Gen and GenF (Figure S2A). *H'*

220   estimates vary among the different GenF categories, with the major differences observed for

221   fragments of 30 and 60 aa (Figure S2B) and estimates converging with increasing MSL.

222   Therefore, the *SS* algorithm considers the MSL of the 'omic'-sample and chooses the

223   appropriate baseline *H'* values pre-computed on the GenF dataset. The GenF fragment size

224   ranges (30, 60, 100, 150, 200, 250 and 300) match the observed MSL in real metagenomic sets

225   (see Figure S1).

226

227   *Properties and performance of SS:* Because the outcome of the *SS* depends on i) the list of S-

228   related Pfam domains and ii) the curated list of S-genomes, we measured its reproducibility

229   with several approaches. First, scores computed in 2014 (Pfam v27, 1528 non-redundant

230   genomes, 156 species in Suli) were compared to the current results (Pfam v30, 2017 non-

231   redundant genomes, 161 curated species in Suli). Second, we compared the outcomes of the *SS*

232   using a random sampling experiment. Briefly, we computed *SS* 1000 times both in the Gen and

233   the Met datasets sorted in terms of their GenF size category. In each test, ≈50% of the 112 S-

234   related Pfam domains were randomly selected to compute *SS*. Finally, we benchmarked the

235   predictive capacity of the *SS* in order to accurately classify the genomes of S-related organisms

236   (Suli, n=161, positive instances) in contrast with a larger set of non-redundant genomes (Gen -

237   Suli, n=1.946, negative instances). True Positive Rates (TPR), False Positive Rates (FPR) and

238   the resulting of Receiver Operating Characteristic (ROC) plots were produced with the scikit-

239   learn module (http://scikit-learn.org/stable/modules/model_evaluation.html), and finally the

240   Area Under the ROC Curve (AUC) was computed.

241

242   **RESULTS AND DISCUSSION**

243   **Defining the biogeochemical S-cycle**

244   What parameters define a biogeochemical cycle, and what are its limits and scope? Which

245   elements should be considered necessary for each cycle? The study of element fluxes between

246   rocks, atmosphere, oceans and biological activity can be extremely complex in terms of space

247   and time, ranging from single living cells to entire ecosystems, and can be completed in

248   microseconds or instead developed over geological time scales (Hedges, 1992; Falkowski *et*

8

249     *al.*, 2008; Zhao *et al.*, 2014; Olson *et al.*, 2016; Widder *et al.*, 2016). Currently, microbial

250     ecology techniques are insufficient to capture the integral complexity of the biogeochemical

251     cycles by selecting a few marker genes to evaluate the importance of any given element in the

252     environment. Here, we argue that a comprehensive description of the relationship between

253     complex biotic an abiotic process is crucial to describe and understand the importance of global

254     biogeochemical cycles and provide a method to do so by taking advantage of 'omic'-era data.

255          We propose a new approach to analyze, compare, evaluate and quantify the importance

256     of biogeochemical cycles in 'omic' datasets summarized in Figure 1, focusing, as a case-study,

257     on the S-cycle. The first step consists on the systematic acquisition of the molecular and

258     ecological information required to describe the cycle of interest. This information can be

259     considered an inventory (Odum, 1993).

260          With the manual curation effort, we obtained both: a list of microorganisms (Suli), and a

261     list of genes encoding enzymes (Sucy) related to the S-cycle. The Suli list includes the

262     microorganisms involved in the global S-cycle using the microbial mats as ecological model.

263     Microbial mats are compartmentalized organizations that were the first ecosystems to appear on

264     Earth and have evolved over more than three billion years into the complex ecosystems that we

265     know today (Herman EK and Kump LR, 2005). Functionally, microbial mats are self-sufficient

266     structures that support most of the major biogeochemical cycles within a vertical dimension of

267     only a few millimeters in a multilayer space (Pinckney and Paerl, 1997) (Figure 2A). These

268     assemblies represent ecosystems with the minimum requirements to be functional and therefore

269     can be used to explore the complexity of biogeochemical cycles. In contrast with the compact

270     nature of microbial mats, the distribution of the metabolic S-guilds is widely dispersed at

271     planetary scale, being found in rivers and estuaries, lagoons, oceans, sediments and deep

272     hydrothermal vents (Halevy *et al.*, 2012).

273          The distribution of S-related bacterial taxa can be analyzed in terms of redox potential

274     and Gibbs Energy of Free Formation of S-compounds, resembling the compact layout of the

275     metabolic guilds in the microbial mat (Figure 2B). Therefore, Suli includes three main groups

276     of microorganisms belonging to the S-metabolic guilds in microbial mats i) chemolithotrophic

277     color less sulfur bacteria (CLSB), ii) anaerobic phototrophs: purple sulfur bacteria (PSB), green

278     sulfur bacteria, (GSB), and iii) sulfate reducing bacteria (SRB) as well as deep-branch sulfur

279  hyperthermophile microorganisms found in extreme conditions (hot springs, black smokers,

280  etc.), such as elemental sulfur reducing (SRM) and oxidizer (SO) microorganisms.

281      The other manually curated list of the inventory, Sucy, contains the metabolic pathways,

282  genes and enzyme activity numbers involved in the S-fluxes (Table S2). To the best of our

283  knowledge, this is the first attempt to integrate the biotic and abiotic processes involved in the

284  mobilization of inorganic/organic S-compounds through microbial-catalyzed reactions at global

285  scale; we gathered and classified all the metabolic pathways involved in the S-cycle according

286  to three key aspects described in Figure 2B, i) S-compound: either organic or inorganic, derived

287  from abiotic or biotic processes, ii) standard Gibbs free energy of formation (GFEF), and iii)

288  metabolic role of the S-compound. The metabolic pathways involved in these S-compounds

289  were systematically divided into 28 pathways (Table S3). We included the pathways involved

290  in a) the oxidation/reduction of inorganic S-compounds, used as source of energy, electron

291  donor or acceptor, b) the degradation of organic S-compounds  such as aliphatic sulfonates

292  sulfur amino acids, and organosulfonates, c) methanogenesis from methylated thiols such

293  dimethyl sulfide DMS, metylthio-propanoate and methanethiol, which are generated in nature

294  by different biogeochemical processes (Caspi *et al.*, 2012), and d) the biosynthesis of

295  sulfolipids (SQDG), because it has been observed that some bacteria living in S-rich and P-

296  lacking environments, are able to synthetized sulfolipids instead of phospholipids in the

297  membrane as an adaptation of the selective pressures of the environment (Alcaraz *et al.*, 2008).

298      Once we integrated the metabolic inventory (genes, enzymatic numbers and major

299  metabolic complexes), we linked all the enzymatic steps and S-compounds into a

300  comprehensive representation of the S-cycle in a single cell (Figure 3), with the following

301  features i) the comprehensive interconnection of pathways in terms of energy flow, ii) the

302  direction of the reactions of the important biogeochemical S-compounds, iii) the interplay of

303  the redox gradient (organic/inorganic) of the intermediate compounds that act as key axes of

304  organic and inorganic reactions (i.e., sulfite), and iv) the molecular reconstruction of S-cycle at

305  different levels (genes, abiotic sulfur-derived compounds and  enzymatic steps).

306      In order to benchmark the entropy-based approach described below, we used available

307  data in both genome and metagenome databases. We also included 35 unpublished

308  metagenomes derived from microbial mats, water and sediment from an ultra-oligotrophic

10

309    shallow lake in Cuatro Cienegas, Coahuila (CCC), Mexico. Altogether, these 935 metagenomes

310    set up the Met dataset. The Gen dataset was  sheared with different fragment sizes taking into

311    account the Mean Size Length (MSL) of Met (Figure S1), producing dataset GenF, which was

312    used for estimating the detection limits of sulfur protein families. We describe the computation

313    approach step-by-step below.

314

315    **Relative entropy as a way to identify protein domains that inform about the S-cycle**

316    After the first step of collecting the datasets, the second step involved the annotation of the

317    coding sequences of the 152 genes in Sucy, yielding a total of 112 Pfam domains (Pfam Sucy

318    in Figure 1). The third step consisted in measuring the relative entropy ($H'$) of each Pfam using

319    Equation 1. The presence of each Pfam domain in Suli (genome list) and in the genomic

320    datasets (Gen and GenF) was used as observed and expected frequencies respectively. The

321    obtained $H'$ values are shown in Figure S2.

322          As negative control, we tested to what extent those $H'$ values depend on the particular

323    genomes curated in Suli. To do so, we replaced Suli by 1,000 lists of random-sampled genomes

324    and used them to compute the observed frequencies. As expected, there was a clear difference

325    between the $H'$ computed using random genomes (Figure S3A), and those obtained with the

326    manually curated list of S- based microorganisms (Figure S3B). In particular, entropy values

327    derived from the random test were found to be approximately symmetric and consistently low

328    among the GenF size categories, yielding values of -0.09, and 0.1 as 5% and 95% percentiles,

329    respectively (Table S4).

330          We then evaluated the behavior of the $H'$ values in both Gen and GenF to test whether

331    informative Pfam domains can be used as molecular markers of S-cycle in metagenomic

332    sequences of variable MSL. In order to be considered as a marker gene, each Pfam domain had

333    to fulfill three requirements 1) produce consistently high mean $H'$ values in Gen and GenF, 2)

334    display low standard deviation (std), and 3) obtain $H'$ values clearly separated from the random

335    distribution. We tested several clustering methods, summarized in Figure S4; among them, the

336    Birch and Ward methods grouped together the informative protein domains with low std

337    (Figure S5). However, Ward clustering included a few protein domains relevant in the S-cycle,

338    which were otherwise discarded by the Birch method. Therefore, we selected the Ward method,

339    which produced three clusters of protein domains described in Figure 4.

340    Cluster 0 includes 94 domains with entropies in the range [-0.4, 0.4], thus overlapping

341    those obtained in the negative control explained earlier. Cluster 1 identifies 12 Pfam domains

342    listed in Table 1, with high entropy and low std, and can therefore be proposed as molecular

343    markers. Among the proposed molecular marker protein domains are APS-Reductase

344    (PF12139: $H'$=1.2), ATP-sulfurilase (PF01747: $H'$=1.03) and DsrC (PF04358: $H'$=0.52), key

345    protein families in metabolic pathways involved in both sulfur oxidation/reduction processes.

346    Finally, cluster 2 groups 6 domains (described in Table S5) with high entropy values and high

347    std, such as PUA-like domain (PF14306: $H'=1$). We presume that the protein domains listed in

348    Table S5 are key players in S-metabolism and their presence in almost all complete-sequenced

349    S-associated microorganisms suggests their important roles.

350    Despite their different properties, all the 112 Pfam domains mentioned in those clusters were

351    subsequently used in the next sections to detect peptides related to the S-cycle in '*omic*'

352    datasets.

353

Table 1 Informative Pfam domains with high H' and high std. Novel proposed molecular marker domains in metagenomic datasets of variable mean size length (MSL)

| Pfam ( Suli ocurrences) | $H'$ mean | $H'$ std | Description |
|---|---|---|---|
| PF12139 58/161 | 1.2 | 0.01 | **Adenosine-5'-phosphosulfate reductase beta subunit:** Key protein domain for both sulfur oxidation/reduction metabolic pathways. Has been widely studied in the dissimilatory sulfate reduction pathway. In all recognized sulfate-reducing prokaryotes, the dissimilatory process is mediated by three key enzymes: Sat, Apr and Dsr. Homologous proteins are also present in the anoxygenic photolithotrophic and chemolithotrophic sulfur-oxidizing bacteria (CLSB, PSB, GSB), in different cluster organization (Meyer and Kuever, 2007). |
| PF00374 135/161 | 1.1 | 0.09 | **Nickel-dependent hydrogenase:** Hydrogenases with S-cluster and selenium containing Cys-x-x-Cys motifs involved in the binding of nickel. Among the homologues of this hydrogenase domain is the alpha subunit of the sulfhydrogenase I complex of *Pyrococcus furiosus*, that catalyzes the reduction of polysulfide to hydrogen sulfide with NADPH as the electron donor (Pedroni *et al.*, 1995). |
| PF01747 103/161 | 1.03 | 0.06 | **ATP-sulfurylase:** Key protein domain for both sulfur oxidation and reduction processes. The enzyme catalyzes the transfer of the adenylyl group from ATP to inorganic sulfate, producing adenosine 5′-phosphosulfate (APS) and pyrophosphate, or the reverse reaction (Taguchi *et al.*, 2004). |
| PF02662 62/161 | 0.82 | 0.03 | **Methyl-viologen-reducing hydrogenase, delta subunit:** Is one of the enzymes involved in methanogenesis and encoded in the mth-flp-mvh-mrt cluster of methane genes in *Methanothermobacter thermautotrophicus*. No specific functions have been assigned to the delta subunit (Finn *et al.*, 2008). |
| PF10418 122/161 | 0.78 | 0.06 | **Iron-sulfur cluster binding domain of dihydroorotate dehydrogenase B:** Among the homologous genes in this family are *asrA* and *asrB* from *Salmonella enterica enterica serovar Typhimurium*, which encode 1) a dissimilatory sulfite reductase, 2) a gamma subunit of the sulfhydrogenase I complex of *Pyrococcus furiosus* and, 3) a gamma subunit of the sulfhydrogenase II complex of the same organism (Caspi *et al.*, 2012). |

| PF13247 149/161 | 0.66 | 0.06 | **4Fe-4S dicluster domain:** Homologues of this family include: 1) DsrO, a ferredoxin-like protein, related to the electron transfer subunits of respiratory enzymes, 2) dimethylsulfide dehydrogenase β subunit (ddhB ), involved in dimethyl sulfide degradation in *Rhodovulum sulfidophilum* and 3) sulfur reductase FeS subunit (sreB) of *Acidianus ambivalens*, involved in the sulfur reduction using $H_2$ or organic substrates as electron donors (Caspi *et al.*, 2012). |
|---|---|---|---|
| PF04358 73/161 | 0.52 | 0 | **DsrC like protein:** DsrC is present in all organisms encoding a dsrAB sulfite reductase (sulfate/sulfite reducers or sulfur oxidizers). The physiological studies suggest that sulfate reduction rates are determined by cellular levels of this protein. The dissimilatory sulfate reduction couples the four-electron reduction of the DsrC trisulfide to energy conservation (Santos *et al.*, 2015). DsrC was initially described as a subunit of DsrAB, forming a tight complex; however, it is not a subunit, but rather a protein with which DsrAB interacts. DsrC is involved in sulfur-transfer reactions; there is a disulfide bond between the two DsrC cysteines as a redox-active center in the sulfite reduction pathway. Moreover, DsrC is among the most highly expressed sulfur energy metabolism genes in isolated organisms and meta-transcriptomes (Santos et al., 2015). |
| PF01058 158/161 | 0.45 | 0.01 | **NADH ubiquinone oxidoreductase, 20 Kd subunit:** Homologous genes are found in the delta subunits of both sulfhydrogenase complexes of *Pyrococcus furiosus* (Caspi *et al.*, 2012). |
| PF01568 156/161 | 0.4 | 0.05 | **Molydopterin dinucleotide binding domain:** This domain corresponds to the C-terminal domain IV in dimethyl sulfoxide (DMSO) reductase (Finn *et al.*, 2008). |
| PF09242 39/161 | 0.38 | 0.04 | **Flavocytochrome c sulphide dehydrogenase, flavin-binding:** Enzymes found in S-oxidizing bacteria such as the purple phototrophic bacteria *Chromatium vinosum* (Finn *et al.*, 2008). |
| PF04879 151/161 | 0.37 | 0.05 | **Molybdopterin oxidoreductase Fe4S4 domain:** Is found in a number of reductase/dehydrogenase families, which include the periplasmic nitrate reductase precursor and the formate dehydrogenase alpha chain, *i.e., Wolinella succinogenes* polysulfide reductase chain. *Salmonella typhimurium* thiosulfate reductase (gene phsA). |
| PF08770 45/161 | 0.35 | 0.03 | **Sulphur oxidation protein SoxZ:** SoxZ sulfur compound chelating protein, part of the complex known as the Sox enzyme system (for sulfur oxidation) that is able to oxidize thiosulfate to sulfate with no intermediates in *Paracoccus parantropus* (Caspi *et al.*, 2012) . |

354

## Identification of S-based genomes using the Sulfur Score

To test whether Pfam entropies can be combined to capture the S-related metabolism, we calculated the Sulfur Score (*SS*) with Equation 2 for all non-redundant genomes in dataset Gen. The obtained *SS* and the corresponding taxonomy according in NCBI for each genome can be found in Table S6. Then we classified all the genomes according to their metabolic capabilities in three subsets: Suli) containing manually 161 curated genomes; Sur) Sulfur unconsidered or related microorganisms not included in Suli with *SS* > 4 (in total 192), which were curated afterwards, and NS) including 1,754 Non-Sulfur species, comprehends the subset Gen – (Suli + Sur). Boxplots summarizing the scores for these subsets are plotted in Figure 5A.

In order to measure the predictive value of *SS*, we computed a Receiver Operator Characteristic (ROC) curve by calling positive instances those annotated in Suli and negative the rest of the genomes, while iterating along increasing values of *SS*. The results are shown in Figure 5B, with an estimated Area Under the Curve (AUC) of 0.985, and the cut-off values of *SS* for several False Positive Rates (FPR). With this training Gen dataset, *SS*=8.705 is required

369  to rule out all false positives. However, *SS*=5.231 is sufficient to achieve a FPR < 0.05. Figure

370  5C breaks down the species in Suli according to the metabolic guilds of the microbial mat

371  model observing the clear difference between the distribution of *SS* in NS and Sulfur-related

372  genomes (Suli and Sur). Finally, the *SS* values of candidate genomes in Sur are also plotted to

373  show that they are comparable to the metabolic guilds of the S-cycle. Figure 5D shows the

374  result of manually assigning candidate genomes in Sur to classes in terms of their ecological

375  capabilities (see Table S7).

376       Out of 192 Sur genomes, 68 are known to metabolize S-compounds under culture

377  conditions in literature reports. For instance, *Sideroxydans lithotrophicus* ES-1, a

378  microaerophilic Fe-oxidizing bacterium has been observed to grow using thiosulfate as an

379  energy source (Emerson *et al.*, 2013). Another 59 Sur organisms were isolated from Sulfur-rich

380  environments such as hot springs or solfataric muds, including uncultured species with

381  genomes assembled from metagenomic sequences. For instance, *Candidatus Desulforudis*

382  *audaxviator* MP104C was isolated from basalt-hosted fluids of the deep subseafloor (Jungbluth

383  *et al.,* 2016). Moreover, an unnamed endosymbiont of a scaly snail was sampled from a black

384  smoker chimney (Nakagawa *et al.*, 2014). Finally, the archaon *Geoglobus ahangari* was

385  isolated from a 2,000m depth hydrothermal vent (Manzella *et al.*, 2015).

386  Combining these two subsets, 68% of species in Sur were confirmed by curation to be S-

387  based.

388       We additionally confirmed the importance of S-cycle in gastrointestinal microbes of the

389  genus *Campylobacter* by detecting 20 genomes with high SS values. This result is consistent

390  with the implication of S-metabolism in the low oxygen environment of the host guts, where

391  several inorganic (e.g., sulfates, sulfites) or organic (e.g., dietary amino acids and host mucins)

392  S-compounds originate and are metabolized by several microorganisms. Among the microbes

393  involved in colonic S-metabolism are SRB, many methanogens and *Campylobacter* genus

394  (Carbonero *et al.*, 2012). Furthermore, some species of *Campylobacter* have been isolated from

395  deep sea hydrothermal vents (Nakagawa *et al.*, 2007), suggesting that this genus plays an

396  important role in the S cycle. The remaining species in Sur were classified in these categories:

397  biorremediation (7), Fe-environment (2), marine (2), peatlands (2) and other environments

398  (32).

399       Overall, these results highlight the broad applicability of our proposed entropy-based

400    score to accurately predict, classify and measure the importance of the S-cycle in both in

401    culture-derived and novel sequenced genomes without prior culture and biochemical

402    knowledge.

403

**404    Detecting key sulfur metagenomic environments with the Sulfur Score**

405    Encouraged by the genomic benchmark results described above, we set out to evaluate the

406    importance of the S-cycle across 935 metagenomes in dataset Met. We calculated the $SS$ for each

407    metagenome taking into account its Mean Size Length (MSL) and the corresponding entropies

408    calculated in dataset GenF (Table S8). The global distribution of Met is mapped in Figure 6A,

409    with $SS$ scores colored from yellow to red. To discriminate the most important S-related

410    environments, those with $SS$ values equal or larger than the 95[th] percentile of the corresponding

411    MSL category are shown with blue borders.

412       In order to analyze some ecological patterns of the metagenomes they were further sorted

413    by their environmental features as proposed by the Genomic Standards Consortium [GSC] and

414    implemented in MG-RAST. Each feature corresponds to one of 13 environmental packages (EP)

415    that standardize metadata describing particular habitats that are applicable across all GSC

416    checklists and beyond (Field *et al.*, 2014). The EPs represent a broad and general classification

417    containing particular features. For example, the "water" EP includes 330 metagenomes from our

418    dataset belonging to several features such as freshwater, lakes, estuarine, marine, hydrothermal

419    vents, etc. Each of these features has different ecological capabilities in terms of biogeochemical

420    cycles; therefore, will likely yield different $SS$ values. The results are shown in Figure 6B. In

421    general, all the metagenomes derived from hydrothermal vents (2), marine benthic (6), intertidal

422    (8), and our CCC microbial mats had $SS$ values above the 95[th] percentile, highlighting the

423    importance of the S-cycle in these environments. In contrast, the metagenomes belonging to

424    features such as sub-terrestrial habitat (7), saline evaporation pond (24) or organisms associated

425    habitat (7) displayed consistently low or even negative $SS$ values, indicating an insignificant

426    presence of S-metabolic pathways in those environments. The remaining features have

427    intermediate median $SS$ values and contain occasionally individual metagenomes with $SS$ values

428    above the 95[th] percentile, such as freshwater, marine, ocean or biofilm environments.

429     Using our approach, we identified and annotated a total of 50 high-scoring metagenomes

430     (Table S9). According to their corresponding literature and associated metadata, all these

431     environments can be described as Sulfur-related as they are reported to be involved in

432     mineralization, uptake, and recycling processes of S-compounds, for example:

433          i) Metagenome 4525341.3 (MSL=172aa, *SS*=9.287) sequenced from costal Oligochaete

434     worm *Olavius algarvensis*, from which metabolic pathway reconstruction revealed the presence of

435     gamma proteobacteria symbionts that are S-oxidizing and SRB. The chemoautotrophic symbionts

436     provide their hosts with multiple sources of nutrition such as organic carbon from autotrophic $CO_2$

437     fixation driven by oxidation of reduced inorganic S-compounds (Woyke *et al.*, 2006).

438          ii) Metagenome 4441663.3 (MSL=158aa, *SS*=9.986) sampled from an hydrothermal vent

439     in the Mariana Trough in 2003 (depth: 2,850 m, fluid temperature:106°C)  (Nakai *et al.*, 2011).

440     The hydrothermal vents are highly productive ecosystems fueled by a number of reduced

441     inorganic substances (*e.g.,* reduced S-compounds, hydrogen or methane) contained in the deep

442     hydrothermal fluids. Through the oxidation of such compounds, chemolithoautotrophic

443     microorganisms gain energy, which can be used for the fixation of inorganic carbon, mediating the

444     transfer of energy from the geothermal source to higher trophic levels and thus form the basis of

445     the unique food chains existing in these environments (Hügler *et al.*, 2010) .

446          iii) Metagenomes 4510162.3, 4510168.3 and 4510170.3, with MSL=32, and *SS* 7.676,

447     7.781 and 7.772, respectively, were sampled from the marine deep-sea surface sediments around

448     the Deep-water horizon spill in the Gulf of Mexico. They belong to ocean feature of EP sediment.

449     The activity of key hydrocarbon degradation pathways was confirmed in these metagenomic deep-

450     sea sediments and the presence of metabolic pathways involved in C, N and S cycles were also

451     confirmed in the metagenomic analysis described in (Mason *et al.*, 2014).

452          iv) Microbial mats from CCC were also detected above the 95th percentile (MSL: 100 and

453     *SS*: 8.945, 9.093, 9.0, and 8.978). Samples were obtained from an ultra-oligotrophic shallow lake

454     recovered from a desiccation event due to water over-exploitation. The sequenced microbial mats

455     showed a clear layered visual pattern structure following the S-metabolic guilds described in

456     Figure 2A. These were assigned to EP and feature "microbial_mat_ccc".

457          v) Metagenomes 4516637.3 and 4516803.3 (MSL=30, SS=7.762 and 7.753 respectively),

458     belong to EP and feature "air". Their high *SS* are consistent with the biogeochemical S-cycle,

459     since the importance of gas-phase reactions of S-compounds and the formation and subsequent

16

460    involvement of sulfate aerosols as cloud-forming nuclei is well established. While these reactions

461    can be carried out without microbial intervention, it has been suggested that microbial

462    communities might contribute to the degradation of some of these S-compounds (Cao *et al.*,

463    2014).

464          To test the reproducibility and robustness of the Sulfur Score, we conducted two further

465    analyses. In the first one, summarized in Figure S6, we compared *SS* estimates of the Met dataset

466    which combined Pfam entropies computed in 2014 and 2017. Despite the changes of both, the

467    Pfam database and the Suli list, overall we find a strong correlation, yielding an $R^2$=0.912 (Figure

468    S6 A). A kernel density analysis of the comparison suggests a different behavior of low and high

469    *SS* scores, with the latter being more reproducible (see Figure S6B). In the second analysis, we

470    quantitatively tested to what extent the entropy estimates of the set of 112 Pfam domains affect the

471    outcome of the *SS* in Gen and Met. To do so, we subsampled randomly ≈50% of those domains to

472    compute the *SS* 1,000 times for each of the 2,017 nr-genomes and 935 metagenomes. The results,

473    summarized in Table S10, confirm that *SS* values computed with random subsets of Pfam domains

474    are generally lower than *SS* derived from the full list (n=112) of S-related Pfam domains in Sucy,

475    regardless of MSL. However, as shown in Figure S7 for MSL=60, the overall ranking of

476    metagenomes produced by computing *SS* values only with a subset of Pfam domains is broadly

477    equivalent. Altogether, these analyses suggest that the choice of protein domains does affect the

478    absolute scores obtained. However, high scoring metagenomes are ranked highly, even when only

479    a reduced set of Pfam domains are employed.

480    The latter result confirms that a comprehensive database of protein-coding genes derived from

481    a systematic reconstruction of global biogeochemical cycles, is necessary to apply the

482    algorithm to quantify the importance of any given elements biogeochemical cycle in a given

483    environment.

484

485    **CONCLUSIONS**

486    In this study, we proposed a computational approach to address the complexity of global

487    biogeochemical cycles on a multi-genomic scale. This methodology requires the curation of an

488    inventory of biotic players (including genes, molecular pathways and microorganisms) and

489    abiotic compounds involved in the cycle of interest. We focused on the S-cycle due to is

490  importance on the biogeochemistry of the planet, but in principle this approach could be

491  applied to any other cycle with microbial participation.

492  For the first time, we systematically build two databases that describe the complexity of

493  the S-cycle based on the minimum ecosystem concept and the microbial mat model. We

494  compiled available non-redundant fully sequenced S-based microorganisms (Suli), and all the

495  known metabolic pathways involved in the S-cycle, taking into account also relevant abiotic S-

496  compounds (Sucy). We used these databases as the input of a computational entropy-based

497  approach that works in two stages. First, we used the individual entropies of protein domains

498  annotated in Sucy to propose a list of 12 molecular markers of the S-cycle and then combined

499  these in order to produce the Sulfur Score, a measure that informs about the presence of the

500  molecular machinery involved in S metabolism. We benchmarked the predictive value of the

501  Sulfur Score by producing a ROC curve, and tested its robustness with simulations against

502  randomly subsampled subsets of the curated protein domains.

503  Altogether, we propose that this method can be used to evaluate other global

504  biogeochemical cycles or complex molecular pathways across genomic and metagenomic

505  sequence datasets, therefore allowing the detection of environmental patterns and informative

506  samples using a single score.

507

513

514  **Conflicts of interest**

515  The authors declare no conflict of interest

516

517  **Figure Legends**

518  **Figure 1.** Computational pipeline for the analysis of the different levels of complexity of the

519  sulfur cycle (S-cycle). The first step is to obtain the datasets. The biogeochemical processes of

18

520    the S-cycle were compiled in two lists. First, the most important S-compounds involved in

521    abiotic or abiotic reactions were curated to produce a database of metabolic pathways, redox

522    reactions, enzymatic numbers and their corresponding coding genes (Sucy). Second, a list of S-

523    based microorganisms (Suli) was also curated and compiled using the information of the

524    Microbial Mat Model described in Figure 2. Then a total of 1,000 list of random genomes were

525    used as negative control (Rlist). The information gathered was then used to evaluate the 'omic'

526    datasets (Gen, GenF, Met). Stages 2 to 4 summarize the determination of the Sulfur Score ($SS$)

527    with sets of peptide sequences of increasing mean size length (MSL, in this example A, B and

528    C). $SS$ is calculated by summation of the entropy of 112 scanned protein domains, each one

529    with a certain entropy value ($H'$).

530    **Figure 2.** Sulfur cycle in small and planetary scale. A) Microbial Mat Model: Simplified

531    scheme of the relevant reactions carried out in microbial mats according to the redox potential.

532    The redox couples (at pH 7) are approximate; the exact values depend upon how the reactions

533    are coupled. 1) oxygenic photosynthesis by Cyanobacteria 2) chemosynthesis by

534    chemolithotrophic color-less sulfur bacteria (CLSB), anoxygenic photosynthesis by 3) purple

535    sulfur bacteria (PSB) and 4) green sulfur bacteria (GSB). sulfate reduction performed by 5)

536    sulfate reducing bacteria (SRB). B) Sulfur cycle at planetary scale. Most important organic and

537    inorganic S-compounds derived from biogeochemical processes arranged according to the

538    Standard Gibbs free energy of formation described in Caspi *et al.*, (2012). The left column

539    indicates whether specific microorganisms are able to use those S-compounds, as a source of

540    Carbon (C), Nitrogen (N), Energy (E) or Electron donors (°). Double asterisks indicate if the S-

541    compound is used as sole source, of C, N, or E. The corresponding electron acceptors in redox-

542    coupled reactions using the S-compound as electron donor are not show. The right column

543    indicates whether the S-compounds are used as fermentative substrate (F) or terminal electron

544    acceptor in respiratory processes (R). Colored boxes summarized the metabolic guilds involved

545    in the metabolism of S-compounds, in oxidation (*i.e.*, CLSB, SOM, PSB, and GSB) or

546    reduction (SR, SRB) process. Some redox coupling reactions carried out by the latter metabolic

547    guilds are showed in panel A. The complete list of S-based microorganisms (Suli) is found in

548    Table S1. Figure based on annotations from MetaCyc (Caspi *et al.*, 2012).

549    **Figure 3**. Comprehensive representation of the global biogeochemical S-cycle assembled from

550    many metabolic pathways found in a variety of organisms combined in a single cell. To the

551    best of our knowledge, all the molecular pathways involved in the metabolism of sulfur

552    compounds, described in Figure 2B, are included. The enzymatic steps are depicted as

553    rectangles followed by arrows indicating the direction of the reaction. Green hexagons

554    represent metabolic links to other metabolisms. Bold dashed arrows indicate bidirectional

555    reactions. Inorganic S-compounds have been arranged according to their reduction potential,

556    from the most oxidized (yellow) to the most reduced (red) compounds. Grey rectangles indicate

557    enzymes acting in disproportionation processes in which a reactant is both oxidized and

558    reduced in the same chemical reaction, forming two separate compounds. Input biogeochemical

559    S-compounds are shown outside and connected with bold arrows. Dashed arrows indicate S-

560    compounds excreted out of the cell. The upper half of the modelled cell depicts the processes

561    involved in the use of organic S-compounds (orange circles) found in natural environments and

562    used as source of carbon, sulfur and/or energy in several aerobic/anaerobic strains described in

563    Figure 2.

564    **Figure 4**. Clustering of the Pfam relative entropies obtained in Gen and GenF produced with

565    the Ward method. Log frequency of the entropy values computed in the random test is colored

566    in purple (see scale bar). Cluster 0 (blue) groups protein domains with low relative entropy that

567    overlap with the random distribution. Cluster 1 (green) includes the Pfam domains that fulfill

568    the requirements to be used as molecular markers (high H' and low standard deviation, std).

569    Red dots (cluster 2) correspond to Pfam domains with high H' and std.

570

571    **Figure 5.**  Distribution of Sulfur Score (*SS*) in genomes of dataset Gen. A) Subsets of non-

572    redundant genomes: i) genomes annotated in Suli (n=161); ii) Sur, genomes not listed in Suli

573    with $SS > 4$ and candidates to be S-related microorganisms (n=192); iii) rest of the genomes in

574    Gen (NS, n=1,754). According to the curated species, True Positives can be defined as

575    genomes with $SS > \max (SS_{NS})$ distribution, whereas True Negatives are those with $SS <$

576    $\min(SS_{Suli})$.  B) ROC curve with Area Under the Curve (AUC) indicated together with

577    thresholds for some False Positive Rates (FPR). C) Distribution of SS for different S-metabolic

578    guilds according to the microbial mat model (Figure 2A) and also the genomes in Sur. D)

579    Assignment of the 192 genomes in Sur to ecological categories based on literature reports.

580

581    **Figure 6.** Distribution of Sulfur Score (*SS*) in the metagenomic dataset Met. A) Geo-localized

582    metagenomes sampled around the globe are colored according to their *SS* values. The following

583    cut-off values correspond to the 95[th] percentiles of seven Mean Size Length classes (30, 60,

584    100, 150, 200, 250 and 300 aa): 7.66, 9.70, 8.81, 8.51, 8.18, 8.98 and 7.61, respectively. Circles

585    with thick blue border indicate metagenomes with SS $\geq$ the 95[th] percentile. B) Distribution of

586    SS values observed in 935 metagenomes classified in terms of features (X-axis) and colored

587    according to their environmental packages. Features are sorted according to their median *SS*

588    values. ccc: metagenomes from Cuatro Cienegas, Coahuila, Mexico. Green lines indicate the

589    lowest and largest 95[th] percentiles observed across MSL classes.

590

591

592    **References**

593    Alcaraz, L.D., Olmedo, G., Bonilla, G., Cerritos, R., Hernández, G., Cruz, A., et al. (2008) The

594        genome of Bacillus coahuilensis reveals adaptations essential for survival in the relic of an

595        ancient marine environment. *Proc. Natl. Acad. Sci. U. S. A.* **105**: 5803–8.

596    Bolhuis, H., Cretoiu, M.S., and Stal, L.J. (2014) Molecular Ecology of Microbial Mats. *FEMS

597        Microbiol. Ecol.* **3**: 1–16.

598    Canfiel. DE, Thamdrup B, Kristensen E. (2005) Aquatic Geomicrobiology, Volume 48

599        (Advances in Marine Biology). *Elsevier Academic Press*.

600    Cao, C., Jiang, W., Wang, B., Fang, J., Lang, J., Tian, G., et al. (2014) Inhalable

601        Microorganisms in Beijing ' s PM 2.5 and PM 10 Pollutants during a Severe Smog Event.

602        *Enviromental Sci. Technol.* **48**: 1499–1507.

603    Carbonero, F., Benefiel, A.C., Alizadeh-Ghamsari, A.H., and Gaskins, H.R. (2012) Microbial

604        pathways in colonic sulfur metabolism and links with health and disease. *Front. Physiol.*

605        **3**: 1–11.

606    Caspi, R., Altman, T., Dreher, K., Fulcher, C. a, Subhraveti, P., Keseler, I.M., et al. (2012) The

607      MetaCyc database of metabolic pathways and enzymes and the BioCyc collection of
608      pathway/genome databases. *Nucleic Acids Res.* **40**: D742-53.

609  Chen, L., Hu, M., Huang, L., Hua, Z., Kuang, J., Li, S., and Shu, W. (2015) Comparative
610      metagenomic and metatranscriptomic analyses of microbial communities in acid mine
611      drainage. *ISME J.* **9(7)**: 1579–1592.

612  Commenges, D. (2015) Information Theory and Statistics: an overview. *ARXIV* preprint
613      arXiv:1511.00860.

614  Dahl (2017) Modern Topics in the Phototrophic Prokaryotes. Metabolism, Bioenergetics and
615      Omics. *Springer International Publishing* pp 27-66, 10.1007/978-3-319-51365-2_2.

616  Emerson, D., Field, E.K., Chertkov, O., Davenport, K.W., Goodwin, L., Munk, C., et al. (2013)
617      Comparative genomics of freshwater Fe-oxidizing bacteria: implications for physiology,
618      ecology, and systematics. *Front. Microbiol.* **4**: 254.

619  Falkowski, P.G., Fenchel, T., and Delong, E.F. (2008) The microbial engines that drive Earth's
620      biogeochemical cycles. *Science* **320**: 1034–9.

621  Field, D., Sterk, P., Kottmann, R., Smet, J.W. De, Amaral-zettler, L., Cole, J.R., et al. (2014)
622      Genomic Standards Consortium Projects The Genomic Standards Consortium Initiating
623      and Maintaining a Project within the GSC The GSC Project Description template provides
624      a References : *Stand. Genomic Sci.* 599–601.

625  Finn, R.D., Clements, J., and Eddy, S.R. (2011) HMMER web server: interactive sequence
626      similarity searching. *Nucleic Acids Res.* **39**: W29-37.

627  Finn, R.D., Tate, J., Mistry, J., Coggill, P.C., Sammut, S.J., Hotz, H.-R., et al. (2008) The Pfam
628      protein families database. *Nucleic Acids Res.* **36**: D281-8.

629  Halevy, I., Peters, S.E., and Fischer, W.W. (2012) Sulfate burial constraints on the Phanerozoic
630      sulfur cycle. *Science* **337**: 331–4.

631  Hedges, J.I. (1992) Global biogeochemical cycles: progress and problems. *Mar. Chem.* **39**: 67–
632      93.

633  Herman EK and Kump LR (2005) Biogeochemistry of microbial mats under Precambrian

634      environmental conditions : a modelling study. *Geob* **3**: 77–92.

635      Hertz, G.Z. and Stormo, G.D. (1999) Identifying DNA and protein patterns with statistically

636      significant alignments of multiple sequences. *Bioinformatics* **15**: 563–77.

637      Hügler, M., Gärtner, A., and Imhoff, J.F. (2010) Functional genes as markers for sulfur cycling

638      and CO2 fixation in microbial communities of hydrothermal vents of the Logatchev field.

639      *FEMS Microbiol. Ecol.* **73**: 526–537.

640      Jones, P., Binns, D., Chang, H.Y., Fraser, M., Li, W., McAnulla, C., et al. (2014) InterProScan

641      5: Genome-scale protein function classification. *Bioinformatics* **30**: 1236–1240.

642      Jungbluth SP, Glavina del Rio T, Tringe SG, Stepanauskas R, Rappé MS. (2016) Genomic

643      comparisons of a bacterial lineage that inhabits both marine and terrestrial deep subsurface

644      systems. PeerJ Preprints 4:e2592v1https://doi.org/10.7287/peerj.preprints.2592v1

645      Kanehisa, M. and Goto, S. (2000) KEGG: kyoto encyclopedia of genes and genomes. *Nucleic*

646      *Acids Res.* **28**: 27–30.

647      Khodadad, C.L.M. and Foster, J.S. (2012) Metagenomic and metabolic profiling of

648      nonlithifying and lithifying stromatolitic mats of Highborne Cay, The Bahamas. *PLoS One*

649      **7**: e38229.

650      Kullback, S. and Leibler, R.A. (1951) On Information and Sufficiency. *Ann. Math. Stat.* **22**:

651      79–86.

652      Li, Y., Yu, S., Strong, J., and Wang, H. (2012) Are the biogeochemical cycles of carbon,

653      nitrogen, sulfur, and phosphorus driven by the "FeIII–FeII redox wheel" in dynamic redox

654      environments? *J. Soils Sediments* **12**: 683–693.

655      Llorens-Marès, T., Yooseph, S., Goll, J., Hoffman, J., Vila-Costa, M., Borrego, C.M., et al.

656      (2015) Connecting biodiversity and potential functional role in modern euxinic

657      environments by microbial metagenomics. *ISME J.* **9(7)**: 1579–92.

658      Loy, A., Ku, K., Lehner, A., Drake, H.L., and Wagner, M. (2004) Microarray and Functional

659      Gene Analyses of Sulfate-Reducing Prokaryotes in Low-Sulfate , Acidic Fens Reveal

660      Cooccurrence of Recognized Genera and Novel Lineages. *Appl. Environ. Microbiol.* **70**:

23

661      6998–7009.

662  Magrane, M. and Consortium, U.P. (2011) UniProt Knowledgebase: A hub of integrated
663      protein data. *Database* **2011**: 1–13.

664  Manzella, M.P., Holmes, D.E., Rocheleau, J.M., Chung, A., Reguera, G., and Kashefi, K.
665      (2015) The complete genome sequence and emendation of the hyperthermophilic, obligate
666      iron-reducing archaeon "Geoglobus ahangari" strain 234T. *Stand. Genomic Sci.* **10**: 77.

667  Mason, O.U., Scott, N.M., Gonzalez, A., Robbins-Pianka, A., Bælum, J., Kimbrel, J., et al.
668      (2014) Metagenomics reveals sediment microbial community response to Deepwater
669      Horizon oil spill. *ISME J.* **8**: 1464–75.

670  Meyer, B. and Kuever, J. (2007) Molecular analysis of the diversity of sulfate-reducing and
671      sulfur-oxidizing prokaryotes in the environment, using aprA as functional marker gene.
672      *Appl. Environ. Microbiol.* **73**: 7664–79.

673  Morales, S.E. and Holben, W.E. (2011) Linking bacterial identities and ecosystem processes:
674      Can "omic" analyses be more than the sum of their parts? *FEMS Microbiol. Ecol.* **75**: 2–
675      16.

676  Moreno-Hagelsieb, G., Wang, Z., Walsh, S., and ElSherbiny, A. (2013) Phylogenomic
677      clustering for selecting non-redundant genomes for comparative genomics. *Bioinformatics*
678      **29**: 947–9.

679  Nakagawa, S., Shimamura, S., Takaki, Y., Suzuki, Y., Murakami, S., Watanabe, T., et al.
680      (2014) Allying with armored snails: the complete genome of gammaproteobacterial
681      endosymbiont. *ISME J.* **8**: 40–51.

682  Nakagawa, S., Takaki, Y., Shimamura, S., Reysenbach, A.-L., Takai, K., and Horikoshi, K.
683      (2007) Deep-sea vent epsilon-proteobacterial genomes provide insights into emergence of
684      pathogens. *Proc. Natl. Acad. Sci. U. S. A.* **104**: 12146–12150.

685  Nakai, R., Abe, T., Takeyama, H., and Naganuma, T. (2011) Metagenomic Analysis of 0.2-µm-
686      Passable Microorganisms in Deep-Sea Hydrothermal Fluid. *Mar. Biotechnol.* **13**: 900–
687      908.

688  Newman, D.K. and Banfield, J.F. (2002) Geomicrobiology: how molecular-scale interactions
689       underpin biogeochemical systems. *Science* **296**: 1071–7.

690  Odum EP (1993) Ecology and our endangered life-support systems, 2nd edn. *Sinauer*
691  *Associates Inc., Sunderland, Massachusetts.* 301 pp.

692

693  Olson, K.R., Straub, K.D., and Straub, K.D. (2016) The Role of Hydrogen Sulfide in Evolution
694       and the Evolution of Hydrogen Sulfide in Metabolism and Signaling The Role of
695       Hydrogen Sulfide in Evolution and the Evolution of Hydrogen Sulfide in Metabolism and
696       Signaling. *Physiology* **31**: 60–72.

697  Pedroni, P., Volpe, A.D., Galli, G., Mura, G.M., Pratesi, C., and Grandi, G. (1995)
698       Characterization of the locus encoding the [Ni-Fe] sulfhydrogenase from the archaeon
699       Pyrococcus furiosus: Evidence for a relationship to bacterial sulfite reductases.
700       *Microbiology* **141**: 449–458.

701  Pinckney, J.L. and Paerl, H.W. (1997) Anoxygenic photosynthesis and nitrogen fixation by a
702       microbial mat community in a bahamian hypersaline lagoon. *Appl. Environ. Microbiol.*
703       **63**: 420–6.

704  Quaiser, A., Zivanovic, Y., Moreira, D., and López-García, P. (2011) Comparative
705       metagenomics of bathypelagic plankton and bottom sediment from the Sea of Marmara.
706       *ISME J.* **5**: 285–304.

707  Rabus, R., Hansen, T., and Widdel, F. (2013) "Dissimilatory sulfate- and sulfur-reducing
708       prokaryotes," in The Prokaryotes, eds E. Rosenberg, E. Delong, S. Lory, E. Stackebrandt,
709       and F. Thompson. *Heidelberg: Springer*. 309–404.

710  Santos, A.A., Venceslau, S.S., Grein, F., Leavitt, W.D., Dahl, C., Johnston, D.T., and Pereira,
711       I.A.C. (2015) A protein trisulfide couples dissimilatory sulfate reduction to energy
712       conservation. *Science (80-. ).* **350**: 1541–1545.

713  Stewart, F.J., Dmytrenko, O., Delong, E.F., and Cavanaugh, C.M. (2011) Metatranscriptomic
714       analysis of sulfur oxidation genes in the endosymbiont of solemya velum. *Front.*
715       *Microbiol.* **2**: 134.

716  Swingley, W.D., Meyer-Dombard, D.R., Shock, E.L., Alsop, E.B., Falenski, H.D., Havig, J.R.,
717      and Raymond, J. (2012) Coordinating environmental genomics and geochemistry reveals
718      metabolic transitions in a hot spring ecosystem. *PLoS One* **7**: e38108.

719  Taguchi, Y., Sugishima, M., and Fukuyama, K. (2004) Crystal Structure of a Novel Zinc-
720      Binding ATP Sulfurylase from Thermus. *Biochemistry* **43**: 4111–4118.

721  Tu, Q., Yu, H., He, Z., Deng, Y., Wu, L., Van Nostrand, J.D., et al. (2014) GeoChip 4: A
722      functional gene-array-based high-throughput environmental technology for microbial
723      community analysis. *Mol. Ecol. Resour.* **14**: 914–928.

724  Widder, S., Allen, R.J., Pfeiffer, T., Curtis, T.P., Wiuf, C., Sloan, W.T., et al. (2016)
725      Challenges in microbial ecology: building predictive understanding of community
726      function and dynamics. *ISME J.* **10**: 2557–2568.

727  Woyke, T., Teeling, H., Ivanova, N.N., Huntemann, M., Richter, M., Gloeckner, F.O., et al.
728      (2006) Symbiosis insights through metagenomic analysis of a microbial consortium.
729      *Nature* **443**: 950–5.

730  Zhao, M., Xue, K., Wang, F., Liu, S., Bai, S., Sun, B., et al. (2014) Microbial mediation of
731      biogeochemical cycles revealed by simulation of global changes with soil transplant and
732      cropping. *ISME J.* **8(10)**: 2045–55.

733

# S-cycle in microbial mats

**A**

Oxygen

$E_0$ (v)

+0.82

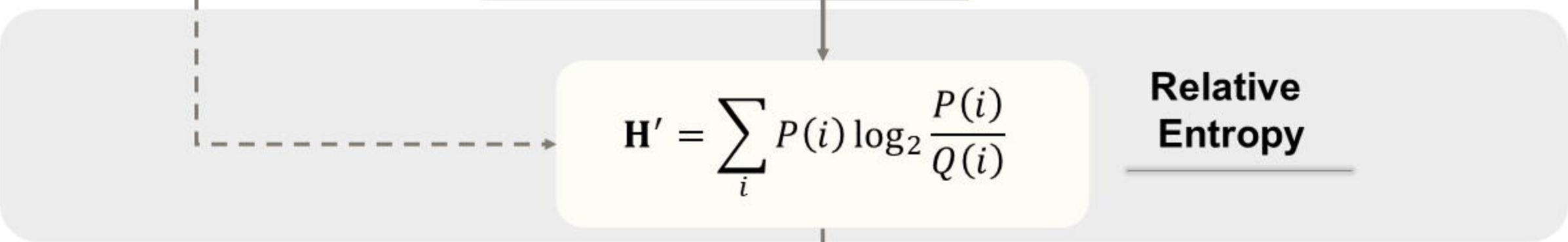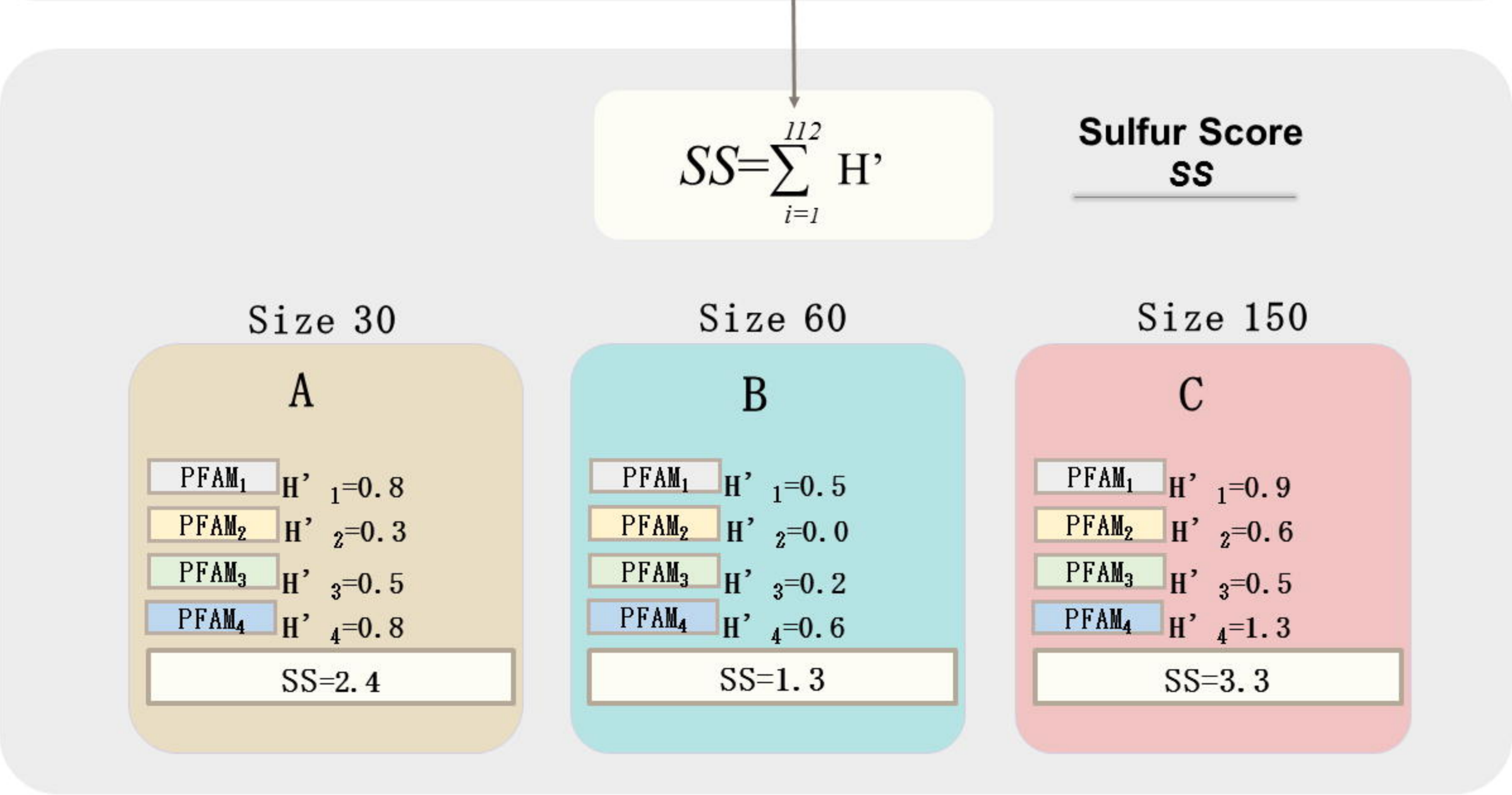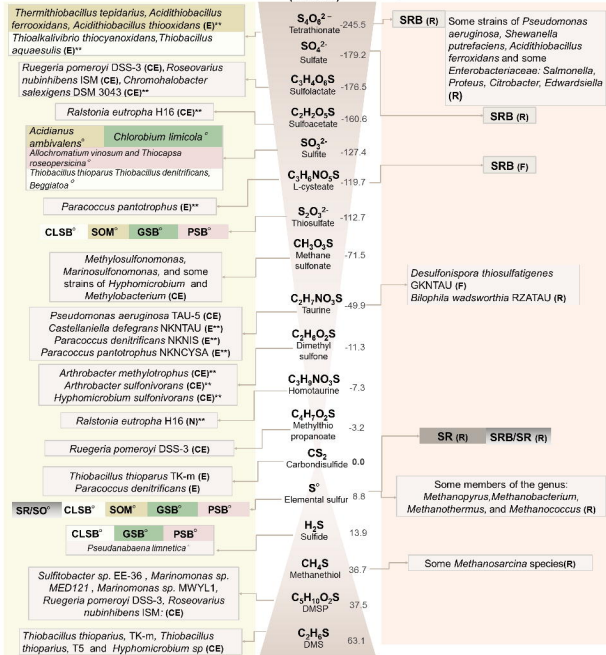| Cyanobacteria Oxygenic photosynthesis | $CO_2+H_2O \Rightarrow (CH_2O)+O_2$ |
| CLSB Chemosynthesis | $2CO_2+4H_2S+6O_2+2H_2O \Rightarrow 2(CH_2O)+4H_2SO_4$ |
| PSB Anoxygenic Photosynthesis | $2CO_2+H_2S+2H_2O \Rightarrow 2(CH_2O)+4H_2SO_4$ |
| GSB Anoxygenic photosynthesis | $2CO_2+H_2S+2H_2O \Rightarrow 2(CH_2O)+4H_2SO_4$ |
| SRB Dissimilatory sulfate reduction | $2CH_3CHOHCOOH+3H_2SO_4 \Rightarrow 6CO^2+3H_2S+6H_2O$ |

-0.22

---

# S-cycle at global scale

**B**

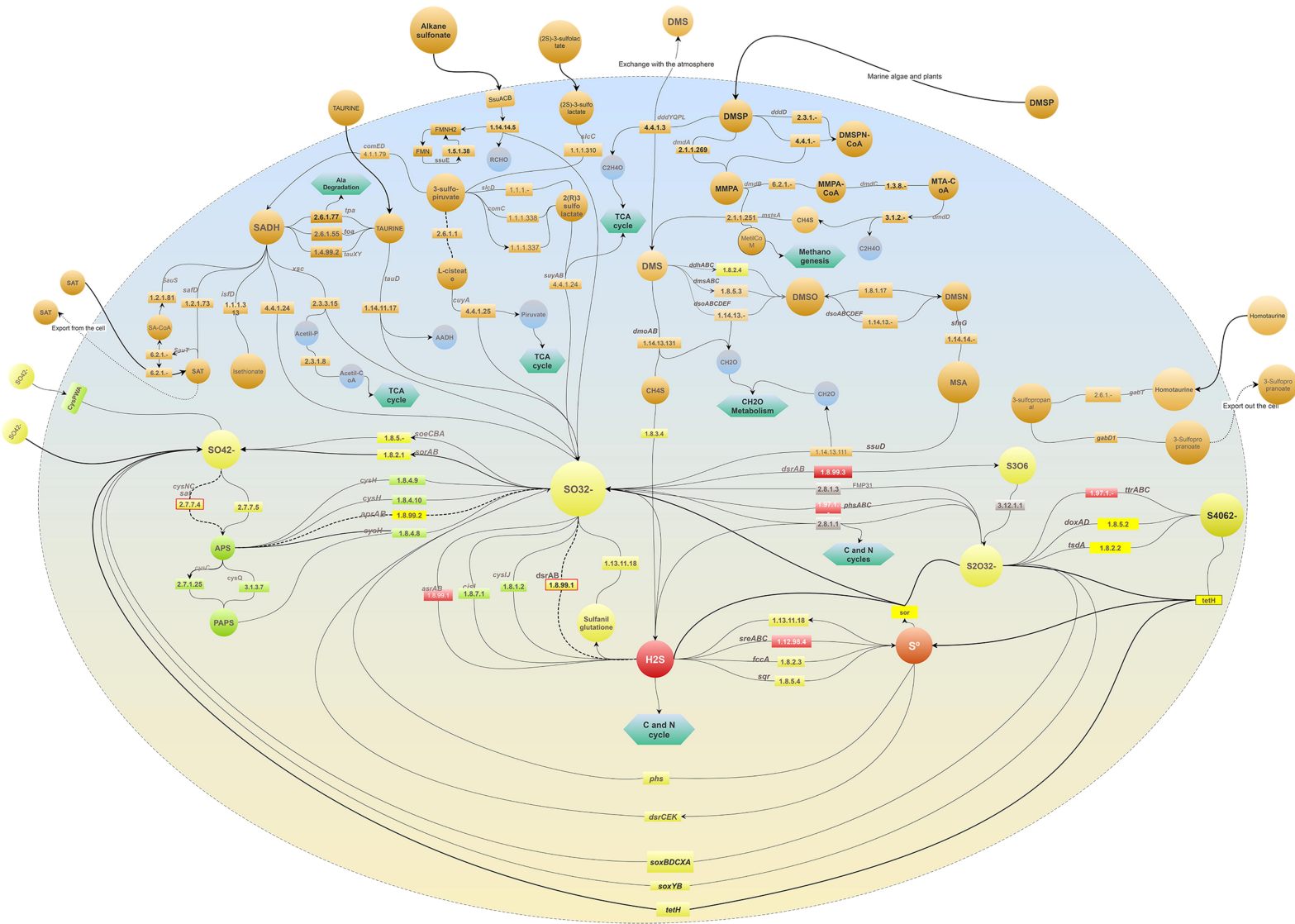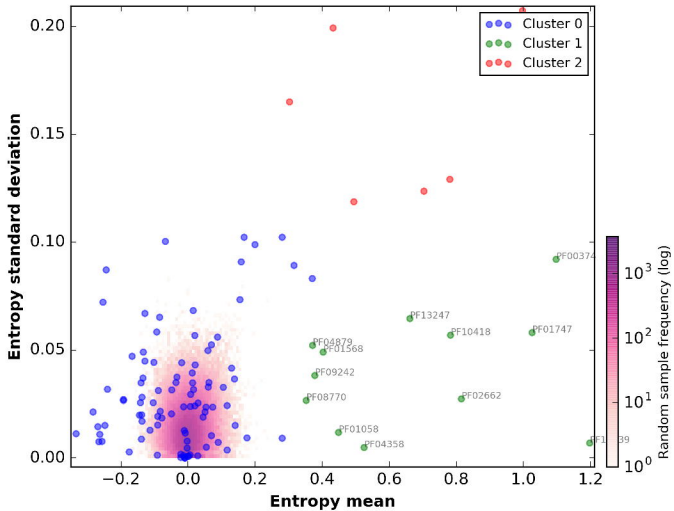**S-compound used as source of Carbon (C), Nitrogen (N), Energy (E), Electron donor °, unique source\*\***

$\Delta_f G°$ of each S-compound (kcal/mol)

**S-compound used as substrate for fermentation (F) or terminal electron acceptor in respiratory process (R)**

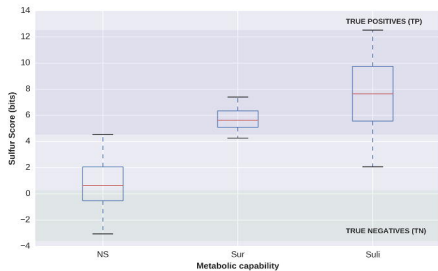| Left box | S-compound / $\Delta_f G°$ | Right box |
|---|---|---|
| *Thermithiobacillus tepidarius, Acidithiobacillus ferrooxidans, Acidithiobacillus thiooxidans* (E)\*\* *Thioalkalivibrio thiocyanodixans, Thiobacillus aquaesulis* (E)\*\* | $S_4O_6^{2-}$ Tetrathionate  -245.5 | **SRB (R)** Some strains of *Pseudomonas aeruginosa, Shewanella putrefaciens, Acidithiobacillus ferrooxidans* and some Enterobacteriaceae: *Salmonella, Proteus, Citrobacter, Edwardsiella* **(R)** |
| *Ruegeria pomeroyi* DSS-3 (CE), *Roseovarius nubinhibens* ISM, *Chromohalobacter salexigens* DSM 3043 (CE)\* | $SO_4^{2-}$ Sulfate  -179.2 | |
| | $C_3H_4O_6S$ Sulfolactate  -176.5 | |
| *Ralstonia eutropha* H16 (CE)\*\* | $C_2H_2O_5S$ Sulfoacetate  -160.6 | **SRB (R)** |
| *Acidianus ambivalens*ᵃ *Chlorobium limicola*ᶜ *Allochromatium vinosum* and *Thiocapsa roseopersicina*ᵃ *Thiobacillus thioparus Thiobacillus denitrificans, Beggiatoa*° | $SO_3^{2-}$ Sulfite  -127.4 | **SRB (F)** |
| | $C_3H_6NO_5S$ L-cysteate  -119.7 | |
| *Paracoccus pantotrophus* (E)\*\* | $S_2O_3^{2-}$ Thiosulfate  -112.7 | |
| **CLSB**° **SOM**° **GSB**° **PSB**° | | |
| *Methylosulfonomonas, Marinosulfonomonas,* and some strains of *Hyphomicrobium* and *Methylobacterium* (CE) | $CH_3O_3S$ Methane sulfonate  -71.5 | *Desulfonispora thiosulfatigenes* GKNTAU **(F)** *Bilophila wadsworthia* RZATAU **(R)** |
| *Pseudomonas aeruginosa* TAU-5 (CE) *Castellaniella defegrans* NKNTAU (E\*\*) *Paracoccus denitrificans* NKNIS (E\*\*) *Paracoccus pantotrophus* NKNCYSA (E\*\*) | $C_2H_7NO_3S$ Taurine  -49.9 | |
| | $C_2H_6O_2S$ Dimethyl sulfone  -11.3 | |
| *Arthrobacter methylotrophus* (CE)\*\* *Arthrobacter sulfonivorans* (CE)\* *Hyphomicrobium sulfonivorans* (CE)\* | $C_3H_9NO_3S$ Homotaurine  -7.3 | |
| *Ralstonia eutropha* H16 (N)\*\* | $C_4H_7O_2S$ Methylthio propanoate  -3.2 | **SR** **SRB/SR** |
| *Ruegeria pomeroyi* DSS-3 (CE) | $CS_2$ Carbondisulfide  0.0 | |
| *Thiobacillus thioparus* TK-m (E) *Paracoccus denitrificans* (E) | $S°$ Elemental sulfur  8.8 | Some members of the genus: *Methanopyrus, Methanobacterium, Methanothermus,* and *Methanococcus* **(R)** |
| **SR/SO**° **CLSB**° **SOM**° **GSB**° **PSB**° | | |
| **CLSB**° **GSB**° **PSB**° *Pseudanabaena limnetica*ᶜ | $H_2S$ Sulfide  13.9 | |
| *Sulfitobacter* sp. EE-36 , *Marinomonas* sp. MED121, *Marinomonas* sp. MWYL1, *Ruegeria pomeroyi* DSS-3, *Roseovarius nubinhibens* ISM: (CE) | $CH_4S$ Methanethiol  36.7 | Some *Methanosarcina* species**(R)** |
| | $C_3H_{10}O_2S$ DMSP  37.5 | |
| *Thiobacillus thioparus,* TK-m, *Thiobacillus thioparius,* T5 and *Hyphomicrobium* sp (CE) | $C_2H_6S$ DMS  63.1 | |

pH 7.3 ionic strength of 0.25, 298 K

A

Metabolic capability: NS, Sur, Suli

TRUE POSITIVES (TP)

TRUE NEGATIVES (TN)

Sulfur Score (bits)

B

Receiver Operating Characteristic (ROC curve)

True Positive Rate (TPR) vs False Positive Rate (FPR)

SS=8.705, FPR =0
SS=8.198, FPR =0.0005
SS=6.328, FPR =0.009
SS=5.231, FPR =0.049
SS=4.045, FPR =0.099

AUC = 0.985

C

95-percentile

Sulfur Score (bits)

S-metabolic guilds: NS, SRB/SR, SO/SR, Sur, PSB, S-org, SR, SOM, CLSB, GSB, SRB

D

- S-cultured — 35.4%
- S-environment — 39.7%
- Other — 16.2%
- Campylobacter — 5.6%
- Biorremediation — 1.04%
- Fe-environment — 1.04%
- Marine — 1.04%
- Peatlands — 1.04%

A

B

Sulfur Score (SS)

Features

Environmental package