

1 **RESEARCH**

2 **MEBS, a software platform to evaluate large (meta)genomic collections**
3 **according to their metabolic machinery: unraveling the sulfur cycle**

4
5 **Authors**

6 Valerie De Anda^{1*}, Icoquih Zapata-Peñasco², Augusto Cesar Poot-Hernandez³ Luis E. Eguiarte¹,
7 Bruno Contreras-Moreira^{4,5*} and Valeria Souza^{1*}

8 **Affiliations:**

9 ¹*Departamento de Ecología Evolutiva, Instituto de Ecología, Universidad Nacional Autónoma de México,*
10 *70-275, Coyoacán 04510 México D.F.*

11 ²*Dirección de Investigación en Transformación de Hidrocarburos. Instituto Mexicano del Petróleo, Eje*
12 *Central Lázaro Cárdenas, Norte 152, Col. San Bartolo Atepehuacan, 07730, México*

13 ³*Departamento de Ingeniería de Sistemas Computacionales y Automatización. Sección de Ingeniería de*
14 *Sistemas Computacionales. Instituto de Investigaciones en Matemáticas Aplicadas y en Sistemas.*

15 ⁴*Estación Experimental de Aula Dei, Consejo Superior de Investigaciones Científicas (EEAD-CSIC), Avda.*
16 *Montañana, 1005, Zaragoza 50059, Spain*

17 ⁵*Fundación ARAID, calle María de Luna 11, 50018 Zaragoza, Spain.*

18
19 **Correspondence address:** Valerie de Anda. Instituto de Ecología, Universidad Nacional Autónoma de
20 México, 70-275, Coyoacán 04510; México D.F; Tel: (+52)5556229006; E-mail valdeanda@ciencias.unam.mx;
21 Bruno Contreras Moreira. Estación Experimental de Aula Dei, Avda. Montañana, 1005, Zaragoza 50059;
22 Spain; Tel:(+34) 976716089 ; E-mail bcontreras@eead.csic.es; Valeria Souza . Instituto de Ecología,
23 Universidad Nacional Autónoma de México, 70-275, Coyoacán 04510; México D.F; Tel: (+52)5556229006; E-
24 mail souza@unam.mx

25
26 **BACKGROUND:** The increasing number of metagenomic and genomic sequences has dramatically improved
27 our understanding of microbial diversity, yet our ability to infer metabolic capabilities in such datasets
28 remains challenging. **FINDINGS:** We describe the Multigenomic Entropy Based Score pipeline (MEBS), a
29 software platform designed to evaluate, compare and infer complex metabolic pathways in large ‘omic’
30 datasets, including entire biogeochemical cycles. MEBS is open source and available through
31 https://github.com/eead-csic-compbio/metagenome_Pfam_score. To demonstrate its use we modeled the
32 sulfur cycle by exhaustively curating the molecular and ecological elements involved (compounds, genes,
33 metabolic pathways and microbial taxa). This information was reduced to a collection of 112 characteristic
34 Pfam protein domains and a list of complete-sequenced sulfur genomes. Using the mathematical
35 framework of relative entropy (H'), we quantitatively measured the enrichment of these domains among
36 sulfur genomes. The entropy of each domain was used to both: build up a final score that indicates whether
37 a (meta)genomic sample contains the metabolic machinery of interest and to propose marker domains in
38 metagenomic sequences such as DsrC (PF04358). MEBS was benchmarked with a dataset of 2,107 non-
39 redundant microbial genomes from RefSeq and 935 metagenomes from MG-RAST. Its
40 performance, reproducibility, and robustness were evaluated using several approaches, including random
41 sampling, linear regression models, Receiver Operator Characteristic plots and the Area Under the Curve
42 metric (AUC). Our results support the broad applicability of this algorithm to accurately classify (AUC=0.985)
43 hard to culture genomes (e.g., *Candidatus Desulforudis audaxviator*), previously characterized ones and

44 metagenomic environments such as hydrothermal vents, or deep-sea sediment. **CONCLUSIONS:** Our
45 benchmark indicates that an entropy-based score can capture the metabolic machinery of interest and be
46 used to efficiently classify large genomic and metagenomic datasets, including uncultivated/unexplored
47 taxa.

48 **Keywords:**

49 Metabolic machinery, metagenomics, omic-datasets, Pfam domains, Relative entropy, sulfur cycle,
50 Multigenomic Entropy-based Score.

51

52 **Background**

53

54 Over the last 15 years, the enormous advances in high-throughput sequencing technologies have
55 revolutionized the field of microbial ecology, dramatically improving our understanding of life's
56 microbial diversity to an unprecedented level of detail [1–4].

57 Nowadays, accessing the total repertoire of genomes within complex communities by means of
58 metagenomics is becoming a standard and routine procedure in order to attain the full insight of
59 the diversity, ecology, evolution and functional makeup of the microbial world [5]. Furthermore,
60 the accurate reconstruction of microbial genomes and draft-populations from environmental
61 metagenomic studies has been shown to be a powerful approach [6–10], providing clues about the
62 potential metabolic strategies of hard-to-culture microbial lineages by linking the functional
63 mechanisms that support specific metabolisms with taxonomic, systematic, and ecological contexts
64 of that lineage [8].

65 Despite the accelerated accumulation of large collections of metagenomic and genomic sequences,
66 our ability to analyze, evaluate and compare complex metabolic capabilities in large-scale 'omic'
67 datasets remains biologically and computationally challenging [11]. Predicting the metabolic
68 potential is a key step in describing the relationship between a microbial community and its
69 ecosystem function. This is largely performed by mapping the protein coding genes of 'omic' data
70 onto reference pathway databases such as MetaCyc [12] or KEGG [13] based on their homology to
71 previously characterized genes [14]. The current available methods for metabolic pathway
72 prediction or reconstruction rely on the use of several metrics to infer the overall repertoire of
73 metabolic pathways present in a given metagenomic dataset (e.g., MinPath [14], HUMAnN[15],
74 PRMT [16], MetaPathways [17]).

75 However, due to the challenges involved in testing meaningful biological hypotheses with complex
76 data, only a small proportion of the metabolic information derived from these datasets is
77 eventually used to draw ecologically relevant conclusions. In this regard, most of the microbial
78 ecology-derived 'omic' studies have been mainly focused on either: i) developing broad description
79 of the metabolic pathways within a certain environment e.g., [18,19]; ii) analyzing the relative
80 abundance of marker genes involved in several metabolic processes and in certain ecosystems
81 (e.g., primary productivity, decomposition, biogeochemical cycling [20–24]; or iii) discovering
82 differentially abundant, shared or unique functional units (genes, proteins or metabolic pathways)
83 across several environmental metagenomic samples [25–27].

84 Therefore, in order to address some of the limitations of these methods, we propose a novel
85 approach to reduce the complexity of targeted metabolic pathways involved in several integral
86 ecosystem processes -- such as entire biogeochemical cycles -- into a single informative score,
87 called Multigenomic Entropy-Based Score (MEBS). This approach is based on the mathematical
88 rationalization of Kullback-Leibler divergence, also known as relative entropy H' [28]. Relative
89 entropy has been widely applied in physics, communication theory, and statistical inference, and it
90 is interpreted as a measure of disorder, information and uncertainty, respectively [29]. Here we
91 use the communication theory concept of H' to summarize the information derived from the
92 metabolic machinery encoded by the protein coding genes of 'omic' datasets. The application of
93 this metric in biology was originally developed by Stormo and colleagues identifying binding sites
94 that regulate gene transcription sites [30].

95 In order to evaluate the performance of our approach, we selected the sulfur cycle (from now on S-
96 cycle) because this is one of the most metabolically- and ecologically complex biogeochemical
97 cycles, but there are few studies analyzing the complete repertoire (genes, proteins, or metabolic
98 pathways) involved in the mobilization of inorganic-organic sulfur compounds through microbial-
99 catalyzed reactions at a planetary scale [20,31–35].

100

101

102

103

104 MEBS description

105

106 MEBS (Multigenomic Entropy-Based Score, RRID: 015708) runs in Linux systems and is available at
107 [36]. For practical purposes, the MEBS algorithm was divided into four stages summarized in Figure
108 1 and explained below.

109

110 **STAGE 1: Manual curation of Sulfur cycle and 'omic' datasets**

111 Sulfur taxonomic representatives. A data set comprehensively covering the currently known
112 representatives of the S-cycle was obtained from primary literature and the MetaCyc database
113 [12]. Each taxonomic representative (at genus or species level) was selected under the criteria of
114 having evidence suggesting their physiological and biochemical involvement in the degradation,
115 reduction, oxidation, or disproportionation of sulfur compounds. Then, each taxonomic
116 representative was scanned against our Genomic dataset (see further details below), in order to
117 obtain a list containing the completely sequenced and non-redundant genomes of the S-cycle. The
118 resulting Sulfur list (or 'Suli') currently contains 161-curated genomes, and was used as the first
119 input of the pipeline. Both the manually curated taxonomic representatives and Suli can be found
120 in Table S1.

121

122 Random taxonomic representatives (RList). As a negative control, we generated 1000 lists of
123 genomes that are not particularly enriched on sulfur metabolic preferences. Each list contains 161
124 random genomes, the same number of microorganisms included in Suli. These lists were obtained
125 by randomly subtracting from the Genomic dataset (see below) 161 Refseq accession numbers and
126 their corresponding names.

127

128 Metabolic pathways and genes. We gathered and classified the metabolic pathways involved in the
129 S-cycle from the primary literature and two experimentally validated curated databases: KEGG
130 (KEGG, RRID:SCR_012773) [13] and MetaCyc (MetaCyc, RRID:SCR_007778) [12]. All the molecular
131 information was then combined into a single database named Suci (for Sulfur cycle). Suci currently
132 contains 152 genes and 48 enzyme classification numbers annotated in the Enzyme classification

133 [37] (Table S2). The 152 FASTA sequences of the proteins encoded by these genes were
134 downloaded from UniProt [38] and used as the second input of the pipeline.

135

136 Genomic dataset (Gen). At the time of the analysis (December 21, 2016), a total of 4,158 genomes
137 were available from RefSeq database [39]. For comparative genomic purposes, we removed
138 redundancy in this large data set by using the Web interface [40] described in [41]. As
139 phylogenomic distance measure, we used a modified version of the Genomic Similarity Score
140 defined as GSSb in [41]; we selected the most tolerant threshold of 0.95 (so as not to drop many
141 sequenced genomes) and default parameters, resulting in 2,107 clusters containing similar
142 genomes, ordered by size (largest to smallest). Then, the largest genome representative for each
143 group was searched in the NCBI genome assembly summary file [42] and downloaded from the
144 NCBI FTP site [43].

145 Metagenomic dataset (Met). We used the Meta Genome Rapid Annotation using Sub- system
146 Technology server (MG-RAST, RRID:SCR_004814) [44] to download metagenomes that: i) were
147 publicly available; ii) contained associated metadata; and iii) had been isolated from well-defined
148 environments (i.e., rivers, soil, biofilms), discarding host associated microbiome sequences (i.e.,
149 human, cow, chicken). In addition we also included 35 unpublished metagenomes derived from
150 sediment, water and microbial mats from Cuatro Ciénegas, Coahuila (CCC), Mexico. The latter were
151 also submitted and annotated in the MG-RAST server, and will be described in depth elsewhere.
152 The resulting collection of 935 FASTA files (\approx 500 GB), containing gene-called protein sequences
153 (MG-RAST stage 350), were downloaded from the RESTful MG-RAST API
154 (<http://api.metagenomics.anl.gov/api.html>). While these metagenomes were evaluated and scored
155 in STAGE 4, they were also analyzed to estimate their mean sequence length, considering that the
156 fragmented nature of metagenomic sequences would have an impact on homology detection,
157 depending on the length of the reads [45,46]. Therefore, we measured the Mean Size Length (MSL)
158 of the peptide sequences of the 935 metagenomes in Met and the 152-curated proteins in Sucs,
159 which are summarized in Figure S1. It was observed that the MSL of Met varies broadly, with a
160 majority of metagenomic peptides with $MSL \leq 30$ aa, and that Sucs proteins range from 49 to 1,020

161 aa, with MSL=349 aa. According to this distribution, the metagenomes in Met were grouped into
162 seven well-defined categories: $MSL \leq 30$, ≤ 60 , ≤ 100 , ≤ 150 , ≤ 200 , ≤ 250 , ≤ 300 aa.

163

164 Fragmented genomic dataset (GenF). In order to simulate the observed variability of MSL across
165 metagenomes, protein sequences encoded in the genomic dataset (Gen, containing 2,107
166 genomes) were *in silico* sheared with Perl script *get_protein_fragments.pl* into the seven MSL
167 categories defined above (30 to 300). This produced the GenF dataset, which currently requires up
168 to 104GB of disk space.

169

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

171 The annotation of protein domains in Suci was conducted using Interproscan 5.21-60.0 [47]
172 against databases Pfam-A v30 (Pfam, RRID:SCR_004726) [48], TIGRFAM v13 (JCVI TIGRFAMS ,
173 RRID:SCR_005493) [49] and Superfamily v1.75 (SUPERFAMILY , RRID:SCR_007952) [50]. Then, the
174 Hidden Markov Models (HMMs) from matched Pfam domains (n=112) were extracted from Pfam-A
175 using script *extract_hmms.pl*. These selected HMMs were subsequently scanned against the
176 Genomic, Genomic Fragmented and Metagenomic datasets (from now on 'omic' datasets, see
177 subsequent stages) using HMMER 3.0 *hmmsearch --cut_ga* option [51].

178

179 **STAGE 3: Relative entropy and its use in detecting informative domains**

180 In order to detect protein domains enriched among sulfur-based microorganisms (Suli), we used a
181 derivative of the Kullback-Leibler divergence [28] — also known as relative entropy $H'(i)$ — to
182 measure the difference between probabilities P and Q (see Eq. 1 below). In this context, $P(i)$
183 represents the frequency of protein domain i in the 161 Suli genomes (observed frequency), while
184 $Q(i)$ represents its frequency in the 2,107 genomes in Gen (expected frequency). The script to
185 compute the entropy (*entropy.pl*) requires the list of the genomes of interest (Suli) and the tabular
186 output file obtained in from the scanning of Gen and GenF against Pfam-Suci database. The
187 obtained values of H' (in bits) capture to what extent a given Pfam domain informs about the
188 metabolism of interest. In this case, domains with H' values close or greater than one, correspond
189 to the most informative Pfam domains (enriched among S-based genomes), whereas low H' values

190 (close to zero) indicate non-informative ones. Negative values correspond to those observed less
191 than expected.

192

$$193 \quad H' = P(i) \log_2 \frac{P(i)}{Q(i)} \quad \text{Eq. 1}$$

194

195 As a negative control, the H' of the 112 Pfam domains were recalculated in both Gen and GenF
196 datasets, but replacing Suli with 1,000 equally sized lists of random-sampled genomes (Rlist).

197 We evaluated the impact of the MSL in the computed entropy values using Gen and GenF. First,
198 we focused on detecting informative Pfam domains that could be used as possible molecular
199 marker genes in variable length, metagenomic sequences. Specifically, we looked for domains
200 displaying stable H' values across both Gen and GenF by using the script
201 *plot_cluster_comparison.py*, which implements the following methods: K-Means, Affinity
202 propagation, Mean-shift Spectral, Ward hierarchical, Agglomerative, DBSCAN and Birch. All of
203 these are part of the scikit-learn Machine Learning Python module [52].

204

205 **STAGE 4: Final score, interpretation, properties and benchmark**

206 Peptide sequences from a given genome or metagenome of interest are evaluated by first scanning
207 their Pfam domains and then producing a final score, defined as the sum of the precomputed
208 entropies of matched S-related Pfam domains (see Equation 2). This score (Sulfur Score 'SS' in our
209 case) summarizes the information content of the metabolic machinery of interest. In this context,
210 informative sulfur protein domains would contribute to higher SS, whereas non-informative ones
211 would decrease it. This is an extension of procedures originally developed for the alignment of DNA
212 and protein motifs, in which individual positions are independent and additive, and can be simply
213 summed up to obtain the total weight or information content [30]. Instead of aligning sequences,
214 in our context we added up the entropy values of the Pfam domains matched in a given 'omic'
215 sample (resulting from scanning the sample of interest against Pfam-Sucy), from which a total
216 weight (SS) is computed by using script *pfam_score.pl*.

217

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

218

219 Datasets in which the majority of informative S-cycle protein domains are represented will yield a
220 high *SS*; in contrast, low *SS* values should be expected if proteins involved in the S-cycle are not
221 particularly enriched.

222

223 MSL. As the calculation of the *SS* depends on the MSL of the omic sample of interest, script
224 *pfam_score.pl* supports option `-size`, in amino acid residues (aa). In this way, appropriate
225 precomputed *H'* values for Pfam domains can be selected to produce the final score. Currently 30,
226 60, 100, 150, 200, 250, 300 and real sizes are supported.

227

228 Metabolic pathway completeness and KEGG visualization. The presence-absence patterns of Pfam
229 domains belonging to particular pathways can be exploited to compute metabolic completeness.
230 This optional task is invoked with parameter `-keggmap` and a TAB-separated file mapping Pfam
231 identifiers to KEGG Orthology entries (KO numbers) and the corresponding pathway in Sucey (see
232 Table S3). To compute completeness, the total number of domains involved in a given pathway
233 (i.e., sulfate reduction, sulfide oxidation) must be retrieved from the Sucey database (See Table S2).
234 Then, the protein domains currently present in any given sample are divided by the total number
235 of domains in the pre-defined pathway. The script produces: i) a detailed report of the metabolic
236 pathways of interest; and ii) a list of KO numbers with Hex color codes, corresponding to KO
237 matches in the omic sample, which can be exported to the KEGG Mapper – Search & Color
238 Pathway tool [53] (see Figure S2).

239

240 Properties and performance of *SS*. Since the outcome of the final score (*SS*) largely depends on the
241 list of microorganisms involved in the metabolism of interest (in our case Suli) and the Pfam
242 domains found in the input protein sequences ($n=112$), we evaluated its robustness and
243 reproducibility with several approaches. First, we compared our results with a benchmark
244 performed three years ago in which we used Pfam-A v27 (instead of version 30), a genomic dataset

245 containing 1,528 non-redundant genomes (579 less genomes than our current Genomic dataset),
246 and an input list of 156 genomes of interest (five less than our current Suli). Second, *SS* estimates
247 were compared with scores obtained by randomly selecting ≈50% of the 112 Pfam domains with
248 both Gen and Met. This analysis was performed a thousand times with *pfam_score.pl -random*.
249 Third, we benchmarked the predictive capacity of the *SS* in order to accurately classify genomes of
250 S-related organisms (Suli, n=161, positive instances), in contrast with a larger set of non-redundant
251 genomes (Gen - Suli, n=1.946, negative instances). Therefore, we computed the True Positive Rates
252 (TPR), False Positive Rates (FPR), Receiver Operating Characteristic (ROC) plots and the resulting
253 Area Under the Curve (AUC) using the scikit-learn module described in [52].

254

255 Results and discussion

256

257 We present MEBS a new open source software to evaluate, quantify, compare, and predict the
258 metabolic machinery of interest in large 'omic' datasets. The pipeline includes four stages. The first
259 one consists on the systematic and targeted acquisition of the molecular and ecological
260 information describing the metabolism of interest, represented by a list of curated microorganisms
261 and a FASTA file of proteins involved in that metabolic network. In the second stage, the domain
262 composition of the curated proteins is evaluated. Then, the domains enriched among the
263 microorganisms of interest are identified by using the mathematical framework of the relative
264 entropy (H' , third stage). Finally, the summation of the entropy of individual Pfam domains in a
265 given genome or metagenomic dataset yields the final score (see Figure 1).

266 To test the applicability of this approach, we evaluated the metabolic machinery of the S-cycle.
267 Due to its multiple redox states and its consequences on microbiological and geochemical
268 transformations, S-metabolism can be observed as a complex metabolic machinery, involving a
269 myriad of genes, enzymes, organic substrates and electron carriers, which largely depend on the
270 surrounding geochemical and ecological conditions. For these reasons, the complete repertory
271 involved in the metabolic machinery of S-cycle has remained underexplored despite the massive
272 data produced in 'omic' experiments. Here, we performed an integral curation effort to describe all
273 the elements involved in the S-cycle and then used, as explained in the following sections, to score
274 genomic and metagenomic datasets in terms of their Sulfur relevance.

275

276 **Manual curation: the complex metabolic machinery of the Sulfur cycle**

277 In order to integrate the complete biogeochemical S-cycle, we manually curated and modeled the
278 major processes involved in the mobilization and use of S-compounds through Earth biosphere.
279 This effort resulted in two comprehensive databases. The first one includes most of the known
280 microorganisms (with and without complete genomes) described in the literature to be closely
281 involved in the S-cycle (Table S1). In this database, we included representative taxa from the
282 following metabolic sulfur guilds: i) chemolithotrophic, colorless sulfur bacteria (CLSB: 24 genera);
283 ii) anaerobic phototrophs, purple sulfur bacteria (PSB:25 genera), and green sulfur bacteria (GSB:9
284 genera); iii) sulfate reducing bacteria (SRB: 40 genera); and iv) deep-branch sulfur
285 hyperthermophilic microorganisms, such as elemental sulfur reducing (SRM:19 genera) and
286 oxidizers (SO:4 genera). From all the microorganisms described to be involved in the S-cycle, at the
287 time of the analysis, a total of 161 were found to be completely sequenced and non-redundant
288 genomes, and were used as the first input of the pipeline (Suli).

289 The second database (Sucy) contains genes, proteins, and pathways with experimental evidence
290 linking them to the S-cycle. To compile this database, we first gathered the most important S-
291 compounds derived from biogeochemical processes and biological catalyzed reactions. Then we
292 classified each S-compound according to their chemical and thermodynamic nature (Gibbs free
293 energy of formation, GFEF). Finally, we classified whether each compound can be used as a source
294 of carbon, nitrogen, energy or electron donor, fermentative substrate, or terminal electron
295 acceptor in respiratory microbial processes. The schematic representation of the manual curated
296 effort summarizing the complexity of the sulfur biogeochemical cycle in a global scale is shown in
297 Figure 2.

298 Once we selected the microorganisms, genes, and biogeochemical processes involved, we
299 systematically divided the metabolic machinery of the S-cycle into 28 major metabolic pathways
300 described in Table 1. In general terms we included pathways involved in: i) the oxidation/reduction
301 of inorganic S-compounds, used as source of energy, electron donor or acceptor (P1-P7, P11 and
302 P20 and P21); ii) the degradation of organic S-compounds, such as aliphatic sulfonates, sulfur
303 amino acids, and organosulfonates (P8-P10, P12-P19, P22,P23,P27); iii) the methanogenesis from

304 methylated thiols, such as dimethyl sulfide DMS (P24), methylthio-propanoate (P25) and
305 methanethiol(P26), which are generated in nature by different biogeochemical processes [12]; and
306 finally, iv) the biosynthesis of sulfolipids (SQDG) (P28), because it has been observed that some
307 bacteria living in S-rich and P-lacking environments are able to synthesize sulfolipids, instead of
308 phospholipids, in the membrane as an adaptation to the selective pressures of these particular
309 environments [54].The synthetic pathway P29 is explained in further detail in the next sections
310 (Table 1).

311 After the comprehensive metabolic inventory was compiled, we linked all the elements in a single
312 network representation of the S-metabolic machinery (Figure 3). To the best of our knowledge, this
313 is the first molecular reconstruction of the cycle that considers all the sulfur compounds, genes,
314 proteins and the corresponding enzymatic steps resulting into higher order molecular pathways.
315 The latter representation also highlights the interconnection of pathways in terms of energy flow
316 and the interplay of the redox gradient (organic/inorganic) of the intermediate compounds that act
317 as key axes of organic and inorganic reactions (e.g., sulfite).

Table 1. Metabolic pathways of global biogeochemical S-cycle

Pathway number	Metabolism ^a	Chemical process ^b	Sulfur compound	Type ^c	Chemical formula	Source ^d	Number of Pfam domains ^e
P1	DS	O	Sulfite		SO ₃ ²⁻	E	9
P2	DS	O	Thiosulfate		S ₂ O ₃ ²⁻	E	10
P3	DS	O	Tetrathionate		S ₄ O ₆ ²⁻	E	2
P4	DS	R	Tetrathionate		S ₄ O ₆ ²⁻	E	17
P5	DS	R	Sulfate		SO ₄ ²⁻	E	20
P6	DS	R	Elemental sulfur		S ⁰	E	20
P7	DS	D	Thiosulfate		S ₂ O ₃ ²⁻	E	9
P8	DS	O	Carbon disulfide	O	CS ₂	E	1
P9	A	DE	Alkanesulfonate	O	CH ₃ O ₃ SR	S	5
P10	A	R	Sulfate		SO ₄ ²⁻	S	20
P11	DS	O	Sulfide		H ₂ S	E/S	29
P12	A	DE	L-cysteate	O	C ₃ H ₆ NO ₅ S	C/E	1
P13	A	DE	Dimethyl sulfone	O	C ₂ H ₆ O ₂ S	C/E	3
P14	A	DE	Sulfoacetate	O	C ₂ H ₂ O ₅ S	C/E	2
P15	A	DE	Sulfolactate	O	C ₃ H ₄ O ₆ S	C/S	14
P16	A	DE	Dimethyl sulfide	O	C ₂ H ₆ S	C/S	16
P17	A	DE	Dimethylsulfoniopropionate	O	C ₅ H ₁₀ O ₂ S	C/S/E	12
P18	A	DE	Methylthiopropionate	O	C ₄ H ₇ O ₂ S	C/S	7
P19	A	DE	Sulfoacetaldehyde	O	C ₂ H ₃ O ₄ S	C/S	7
P20	DS	O	Elemental sulfur		S ⁰	C/S/E	7
P21	DS	D	Elemental sulfur		S ⁰	C/S/E	1
P22	A	DE	Methanesulfonate	O	CH ₃ O ₃ S	C/S/E	7
P23	A	DE	Taurine	O	C ₂ H ₇ NO ₃ S	C/S/E	11
P24	DS	M	Dimethyl sulfide	O	C ₂ H ₆ S	C	1
P25	DS	M	Methylthio-propionate	O	C ₄ H ₇ O ₂ S	C	1
P26	DS	M	Methanethiol	O	CH ₄ S	C	1
P27	A	DE	Homotaurine	O	C ₃ H ₉ NO ₃ S	N	1
P28	A	B	Sulfolipid	O	SQDG		4
P29			Markers		Markers		12

318 ^a Metabolism: Assimilative (A) inorganic compounds are reduced during biosynthesis; Dissimilative
 319 (DS) inorganic compounds used as electron acceptors in energy metabolism. A large amount of
 320 electron acceptor is reduced and the reduced product is secreted.

321 ^b Chemical Process: Oxidation (O): Reduction (R), Degradation (DE), Biosynthesis (B),
 322 Methanogenesis (M), Disproportionation (D).

323 ^c Compound Type: Organic (O): sulfur atoms with covalent bonds to carbon atoms. Inorganic (I):
 324 sulfur compounds with non-carbon atoms.

325
 326 ^d Source: sulfur compound used as source of energy (E), sulfur (S), carbon (C), nitrogen (N).

327
 328 ^e Number of Pfam domains belonging to each metabolic pathway described in Sucy (Table S2)

329
 330
 331
 332

333 **Annotation of Pfam domains within Sulfur proteins**

334 Our approach requires the detection of structural and evolutionary units, also known as domains,
335 in the curated list of protein sequences involved in the metabolism of interest (S-cycle in this case).
336 The annotation of protein domains against the Pfam-A database resulted in a total of 112 domains
337 identified in 147 proteins (out of 152). These 112 domains constitute the Pfam-Sucy database and
338 represent all the pathways listed in Table 1. Two other protein family databases were tested
339 (TGRFAM and Superfamily), but the number of proteins with positive matches was lower than with
340 Pfam (57 and 137, respectively) and thus were not further considered.

341

342 **Preparation of omic datasets: Gen, GenF and Met**

343 The genomic dataset required for computing domain entropies (Gen) was obtained from public
344 databases, as explained above in MEBS Description. A fragmented version of Gen, called GenF, was
345 generated by considering the Mean Size Length (MSL) distribution of metagenomic sequences
346 (Figure S1).

347 In order to benchmark MEBS with real environmental metagenomic samples, a collection of 900
348 public metagenomes was obtained from MG-RAST, to which we added 35 metagenomes sampled
349 from an ultra-oligotrophic shallow lake in México (CCC). Altogether, these 935 metagenomes set
350 up the Met dataset.

351

352 **Using the relative entropy to recognize S-cycle domains and candidate markers**

353 The next stage consists on the quantitative detection of informative domains (enriched among
354 organism in Suli), by computing its relative entropy (H') using Equation 1. The occurrences of each
355 of the 112 Pfam domains in Suli and the genomic datasets were taken as observed and expected
356 frequencies, respectively. Figure 4A summarizes the computed H' values in real (Gen) and
357 fragmented genomic sequences of increasing size (GenF). The results indicate that only a few Pfam
358 domains are equally informative regardless of the length of sequences. When H' values inferred
359 from real, full-length proteins are compared to those of fragmented sequences, it can be seen that
360 shorter sequences (MSL 30 & 60 aa) yield larger entropy differences than sequences of length >
361 100 aa (see in Figure 4B). Therefore, in order to shortlist candidate marker genes we selected those

362 Pfam domains displaying constant, high mean H' values in Gen and GenF, low H' standard deviation
363 (std) and a clear separation from the random distribution.

364 We tested several clustering methods, summarized in Figure S3, with Ward and Birch performing
365 best in grouping together informative protein domains with low std. However, the Ward
366 classification was eventually selected as Birch failed to include a few Pfam domains relevant in the
367 S-cycle (see Figure S4). By using Ward method, three well-defined clusters of Pfam domains were
368 generated, as observed in Figure 4C. Cluster 0 included 94 domains containing H' values ranging
369 from [-0.4, 0.4] and overlapping with the values obtained in the negative control explained in the
370 next section. Cluster 1 consistently grouped together 12 Pfam domains listed in Table 2 with high
371 entropy and low std, and can therefore be proposed as molecular markers in metagenomic
372 sequences of variable length. Among the proposed marker domains are APS-Reductase (PF12139:
373 H' =1.2), ATP-sulfurilase (PF01747: H' =1.03) and DsrC (PF04358: H' =0.52), key protein families in
374 metabolic pathways involved in both sulfur oxidation/reduction processes. Finally, cluster 2
375 includes Pfam domains displaying high entropy values and high std, such as the PUA-like domain
376 (PF14306: H' =1). We presume that domains within this cluster are also key players in S-
377 metabolism; however, their high std makes them unsuitable for markers, particularly with
378 metagenomic sequences of variable MSL. We suggest that further analyses will be required to test
379 the implication in S-energy conservation processes of proteins containing domains such as
380 PF03916, PF02665 or PF14697 (see complete list in Table S4).

Table 2 Informative Pfam domains with high H' and low std. Novel proposed molecular marker domains in metagenomic data of variable MSL

Pfam ID (Suli occurrences)	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 metabolism. 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 [35].
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 <i>Pyrococcus furiosus</i> , that catalyzes the reduction

			of polysulfide to hydrogen sulfide with NADPH as the electron donor [55].
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 [56].
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 <i>Methanothermobacter thermautotrophicus</i> . No specific functions have been assigned to the delta subunit [48].
PF10418 122/161	0.78	0.06	Iron-sulfur cluster binding domain of dihydroorotate dehydrogenase B: Among the homologous genes in this family are <i>asrA</i> and <i>asrB</i> from <i>Salmonella enterica enterica serovar Typhimurium</i> , which encode 1) a dissimilatory sulfite reductase, 2) a gamma subunit of the sulfhydrogenase I complex of <i>Pyrococcus furiosus</i> and, 3) a gamma subunit of the sulfhydrogenase II complex of the same organism [12].
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 <i>Rhodovulum sulfidophilum</i> and 3) sulfur reductase FeS subunit (sreB) of <i>Acidianus ambivalens</i> , involved in the sulfur reduction using H ₂ or organic substrates as electron donors [12].
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 [57]. 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 <i>Pyrococcus furiosus</i> [12].
PF01568 156/161	0.4	0.05	Molybdopterin dinucleotide binding domain: This domain corresponds to the C-terminal domain IV in dimethyl sulfoxide (DMSO) reductase [48].
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 <i>Chromatium vinosum</i> [48].
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., <i>Wolinella succinogenes</i> polysulfide reductase chain. <i>Salmonella typhimurium</i> 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 <i>Paracoccus parantropus</i> [12].

381

382

383 **Is the entropy affected by the input list of microorganisms? Negative control test**

384 In order to evaluate to what extent the H' values depend on the curated list of microorganisms, we
385 performed a negative control by replacing Suli in 1,000 lists of randomly-sampled genomes and
386 used them to compute the observed frequencies (see Equation 1). As expected, there was a clear
387 difference between both H' estimates (see Figure S5). In particular, entropy values derived from
388 the random test were found to be approximately symmetric and consistently low among the GenF
389 size categories (compared with the real values), yielding values of -0.09, and 0.1 as 5% and 95%
390 percentiles, respectively (Table S5).

391

392 **Sulfur Score and its predictive capacity to detect S-microbial players in a large genomic** 393 **dataset.**

394 To test whether Pfam entropies can be combined to capture the S-metabolic machinery in 'omic'-
395 samples, we calculated the final MEBS score, called in this case Sulfur Score (SS). We computed the
396 SS on each of the 2,107 non-redundant genomes in Gen with script `score_genomes.sh`. The
397 individual genomes along with their corresponding SS values and taxonomy according to NCBI are
398 found in Table S6.

399 For evaluation purposes, we classified and manually annotated all the genomes in Gen according
400 to their metabolic capabilities. First, we identified the 161-curated genomes belonging to Suli.
401 Then, we focused on the remaining genomes. A set of 192 genomes with $SS > 4$ were labeled as
402 Sulfur unconsidered or related microorganisms (Sur). Finally, the rest of genomes in Gen were
403 classified as NS (Non-Sulfur = Gen - (Suli + Sur)), including 1,754 genomes. The boxplots in Figure
404 5A summarize the scores obtained in these three subsets.

405 To double-check whether the Sur genomes -- selected due to their SS -- might be involved in the S-
406 cycle, we manually annotated all of them focusing on relevant genomic, biochemical, physiological
407 and environmental information that we might have missed since Suli was first curated (Table S7).
408 Out of 192 genomes, 68 are reported to metabolize S-compounds under culture conditions in the
409 literature. For instance, *Sideroxydans lithotrophicus ES-1*, a microaerophilic Fe-oxidizing bacterium,
410 has been observed to also grow in thiosulfate as an energy source [58]. Another 59 Sur organisms
411 have been isolated from Sulfur-rich environments, such as hot springs or solfataric muds.

412 Remarkably, some of this species include hard-to culture genomes reconstructed from
413 metagenomic sequences such as *Candidatus Desulfurudis audaxviator MP104C* isolated from
414 basalt-hosted fluids of the deep seafloor [6]; an unnamed endosymbiont of a scaly snail from a
415 black smoker chimney [59] and archaeon *Geoglobus ahangari*, sampled from a 2,000m depth
416 hydrothermal vent [60]. Furthermore, we also confirmed within Sur the implication of S-cycle of 20
417 species of the genus *Campylobacter*. These results are consistent with the ecological role of the
418 involved taxa, that along with SRB and methanogens inhabiting host-gastrointestinal and low
419 oxygen environments, where several inorganic (e.g., sulfates, sulfites) or organic (e.g., dietary
420 amino acids and host mucins) are highly metabolized by these metabolic guilds [61]. The
421 implication of *Campylobacter* species in the S-cycle is also supported by the fact that some of them
422 have been isolated from deep sea hydrothermal vents [62]. The remaining species in Sur were
423 classified in different categories, including bioremediation (7), Fe-environment (2), marine (2), peat
424 lands (2) and other environments (32, see Figure 5B).

425 When the SS values of genomes in Sur are compared to the S-metabolic guilds represented in Suli
426 (e.g PSB, SRB, GSB), it can be seen that they are indeed similar and clearly separated from the rest
427 of NS genomes (Figure 5C). This strongly suggests that high scoring genomes are indeed
428 ecologically and metabolically implicated in the S-cycle.

429 Finally, in order to quantify the capacity of the SS to accurately classify S-related microorganisms,
430 we computed a Receiver Operator Characteristic (ROC) curve (for a detailed description of ROC
431 curves see [63]). We thus defined genomes annotated in Suli as positive instances, and the rest as
432 negative ones. The results are shown in Figure 5D, with an estimated Area Under the Curve (AUC)
433 of 0.985, and the corresponding cut-off values of SS for several False Positive Rates (FPR).
434 According to this test, a SS value of 8.705 is required to rule out all false positives in Gen, while
435 SS=5.231 is sufficient to achieve a FPR < 0.05.

436 Overall, these results indicate that MEBS is a powerful and broadly applicable approach to predict,
437 and classify microorganisms closely involved in the sulfur cycle even in hard-to culture microbial
438 lineages.

439

440 **Sulfur Score and its predictive capacity to detect S-related environments in a large**
441 **metagenomic collection.**

442 The SS was also computed for each metagenome in Met, using their corresponding MSL to choose
443 the appropriate entropies previously calculated in dataset GenF (Table S8). In order to test
444 whether SS values can be used to identify S-related environments, we performed the following
445 analyses. First, we use the geographical metadata associated with each metagenome to map the
446 global distribution of SS. In Figure 6A, SS values are colored from yellow to red. The most
447 informative S-environments (displaying SS values equal or greater than the 95th percentile of each
448 MSL category) are shown in blue.

449 Then, we sorted the metagenomes according to their environmental features as proposed by the
450 Genomic Standards Consortium [GSC] and implemented in MG-RAST. Each feature corresponds to
451 one of 13 environmental packages (EP) that standardize metadata describing particular habitats
452 that are applicable across all GSC checklists and beyond [64]. Therefore, each EP represents a
453 broad and general classification containing particular features. For example, the “water” EP
454 includes 330 metagenomes from our dataset, belonging to several features such as freshwater,
455 lakes, estuarine, marine, hydrothermal vents, etc. Since each of these features has different
456 ecological capabilities in terms of biogeochemical cycles, we can expect different behaviors among
457 SS values, as shown in Figure 6B. In general, all the metagenomes derived from hydrothermal vents
458 (2), marine benthic (6), intertidal (8) and our unpublished CCC microbial mats had SS values above
459 the 95th percentile, highlighting the importance of the S-cycle in these environments. In contrast,
460 the metagenomes belonging to features such as sub-terrestrial habitat (7), saline evaporation pond
461 (24) or organisms associated habitat (7) displayed consistently low or even negative SS values,
462 indicating a negligible presence of S-metabolic pathways in those environments. The remaining
463 features have intermediate median SS values and contain occasionally individual metagenomes
464 with SS values above the 95th percentile, such as freshwater, marine, ocean or biofilm
465 environments.

466 To validate the list of 50 high-scoring metagenomes (above the 95th percentiles), we double-
467 checked their annotations. According to the literature and associated metadata, all these
468 environments are closely involved in mineralization, uptake, and recycling processes of S-

469 compounds. For example, environmental sequences derived from costal Oligochaete worm *Olavius*
470 *algarvensis*, hydrothermal vents and marine deep-sea surface sediments around the Deep-Water
471 Horizon spill in the Gulf of Mexico. The complete list of annotated metagenomes, along with their
472 ecological capabilities, is found in Table S9.

473

474 **Evaluating the robustness of the Sulfur Score**

475 To test the reproducibility and robustness of MEBS final score (*SS*), we conducted two further
476 analyses. In the first one we compared *SS* estimates derived from Met dataset, computed with
477 Pfam entropies obtained in the first MEBS benchmark performed three years ago (2014) with the
478 current data described in this article (2017). Despite the changes of both databases (Pfam database
479 version and the Suli list), we found a strong correlation ($r^2=0.912$) between the *SS* outcomes (Figure
480 S6 A). A kernel density analysis of the latter comparison suggests a different behavior of low and
481 high *SS* scores, with the latter being more reproducible (see Figure S6B).

482 In the second analysis, we quantitatively tested to what extent the entropy estimates of the 112
483 Pfam domains directly affect the outcome of the *SS* in Gen and Met. We randomly subsampled
484 $\approx 50\%$ of those domains to compute the *SS* a thousand times for each genome and metagenome in
485 Gen and Met, respectively. The results, summarized in Table S10, confirm that *SS* values computed
486 with random subsets of Pfam domains are generally lower than *SS* derived from the full list ($n=112$)
487 of Suly-Pfam domains. To further inspect the distribution of *SS* values produced with random
488 subsets of domains (random *SS*), we focused on the particular case of the metagenomes belonging
489 to the category $MSL=60$. As expected, the distribution of random *SS* oscillates between negative
490 and positive values. Interestingly, metagenomes exhibiting only positive random *SS* are ranked
491 above the 95th percentile according to their real *SS* values (See Figure S7A). The latter indicates that
492 even a random subset of Pfam domains are used to compute the score, is more likely to high-rank
493 metagenomes containing the sulfur metabolic machinery (large number of high-entropy Pfam
494 domains), than those lacking the sulfur metabolism or displaying a large number of non-informative
495 Pfam domains. Furthermore, by comparing the median of random *SS* with the real scores, we
496 observe a clear separation between those distributions (see Figure S7B and Table S10).

497

498 **Completeness of S-metabolic pathways**

499 As we described above, the MEBS pipeline models a metabolic network as an array of S-related
500 protein domains (Sucy-Pfam), to ultimately use their entropies to produce the final score (*SS*). For a
501 closer look, we also dissected the total contribution of independent domains at the network level,
502 in order to assess whether *SS* depends on the partial or complete detection of S-pathways.
503 Consequently, we evaluated the pathway completeness in both genomic (Gen) and metagenomic
504 (Met) datasets (see Tables S11 and S12, respectively). Since the number of Pfam domains per
505 pathway goes from one to 29 (see Table 1 and Table S2), we suspect that pathways represented by
506 a single domain might not reflect their complete metabolic function. For example, the pathways
507 involved in the methanogenesis of compounds such as dimethylsulfide (DMS, P24), methyl-
508 thiolpropanoate (MTPA, P25), and methanethiol (MeSH, P26) are represented by the same protein
509 (MtsA, PF01208) in our Sucy database, as well as in Metacyc [12]. Therefore, we expect that
510 pathways P24-26 will have identical presence-absence patterns in Gen and Met.

511 The boxplots in Figure 7A and 7B summarize the distribution of completeness for each S-metabolic
512 pathway including the synthetic pathway (P29) composed by 12 candidate markers as described in
513 Table 2. As expected, the observed completeness per pathway was higher in Met than in Gen,
514 since microbial communities harbor a wider repertoire of metabolic functions than single genomes.
515 In the case of genomes, we noted that a few pathways were complete in most genomes, being the
516 majority involved in the usage of organic sulfur compounds such as alkanesulfonates (P9),
517 sulfoacetate (P14) and biosynthesis of sulfolipids (SQDG) and the single domain pathways P24-26.
518 Remarkably, we also detected a few organisms displaying the highest levels of metabolic
519 completeness in some S-energy based pathways. For example, we found that *Desulfosporosinus*
520 *acidiphilus SJ4* (*SS*=8,91) was the only genome harboring the complete repertoire of Pfam domains
521 described in Sucy for the sulfite oxidation (P1), strongly suggesting that it may oxidize sulfite.
522 However, this activity remains to be tested in culture [65]. In the case of thiosulfate oxidation (P3),
523 we detected three genomes displaying the highest levels of completeness, in agreement with their
524 ecological features: *Hydrogenobaculum sp. Y04AAS1* (*SS*=9,319) [66] and the CLSB: *Acidithiobacillus*
525 *calvus ATCC 51756* (*SS*=6,525) [67] and *Acidithiobacillus ferrivorans* (*SS*=7,436) [68]. For the sulfate
526 reduction dissimilative pathway (P5), out of 55 genomes displaying the higher completeness levels,

527 67% are actually SRB, 12% are Sur genomes, and the rest are sulfur oxidation microorganisms.
528 Furthermore, the PSB *Thioflavicoccus mobilis* 8321 (SS= 9,756), isolated from a microbial mat [69],
529 was the genome displaying the most complete sulfide oxidation pathway (P11). Elemental sulfur
530 disproportionation (P21) is represented by a single non-informative domain (PF07682, $H'=0.172$)
531 that remarkably is found in 14 sulfur respiring or related genomes such as *Sulfolobus tokodaii* str. 7
532 (SS= 5,341) and *Acidianus hospitalis* W1 (SS= 3,88). Finally, we identified six genomes encoding all
533 12 proposed markers. Among them, three were GSB (*Pelodictyon phaeoclathratiforme* BU-1,
534 SS=11,836, *Chlorobium chlorochromatii* CaD3, SS=11,625 and *Chlorobium tepidum* TLS, SS=
535 11,354), one CLSB (*Thiobacillus denitrificans* ATCC 25259 SS=11,61), another one PSB (*Thiocystis*
536 *violascens* DSM 198, SS=10,633) and finally one Sur (*Sedimenticola thiotaurini* SS=10,109). For a
537 complete description, see Table S13.

538 A global view of metabolic completeness was obtained by bulking the data from all pathways.
539 Linear regression models between mean completeness and SS were computed confirm the ,
540 yielding r^2 values of 0.003 and 0.627 for Gen and Met, respectively (See Figures 7C and 7D).
541 Moreover, we also assessed the relationship between the mean completeness of the synthetic
542 pathway of candidate markers (P29) and the SS. As expected, significant correlations were
543 obtained in both datasets ($r^2= 0.645$ and $r^2=0.881$ for Gen and Met, respectively; see Figure S8).

544 To get a more detailed insight of the completeness, we selected a few genomes and metagenomes
545 displaying high and low SS values. Specifically, from the Gen dataset we selected one
546 representative from the main S-guilds, one Sur genome and two genomes with low SS values (NS).
547 As observed in Figure 7, the low-scoring genomes *Enterococcus durans* (SS=-0,194), *Micrococcus*
548 *luteus* NCTC_2665 (SS=-3,588), and *Ruegeria pomeroyi* DSS-3 (SS=2,707) display unrelated patterns
549 of sulfur metabolic completeness, compared with the rest of genomes and therefore are
550 separated. In contrast, high-scoring S-respiring microorganisms *Desulfovibrio vulgaris* DP4 (SS=
551 11,442), *Sulfolobus acidocaldarius* DSM 639 (SS=5,457) and *Ammonifex degensii* KC4 (SS=12.508)
552 are clustered together. We also observed that mat-isolated cyanobacteria *Synechococcus* sp. JA-2-
553 3Ba 2-13, classified as NS with SS=3,704, was clustered together with other high-scoring genomes,
554 in agreement with the lack of correlation reported above.

555 In the case of metagenomes (see Figure 7E), we observed a clear correlation between SS and
556 completeness. For example, metagenomes 4440320.3 and 4489656.3, with the lowest scores
557 (SS=0.1 and SS=-2.649, respectively), also exhibit the largest number of incomplete pathways.
558 Similarly, high-scoring metagenomes derived from black smoker or marine sediment are grouped
559 together in terms of completeness.

560

561 Conclusions

562 Our study represents the first exploration of the Sulfur biogeochemical cycle in a large collection of
563 genomes and metagenomes. The manually curated effort resulted in an inventory of the
564 compounds, genes, proteins, molecular pathways, and microorganisms involved. This complex
565 universe of articulated data was reduced into a list of microorganisms and Pfam domains encoded
566 in the proteins that take part in that network. These domains were first ranked in terms of relative
567 entropy, and then summed to produce a single S-score representing the relevance of a given
568 genomic or metagenomic sample in terms of sulfur metabolic machinery. We took advantage of
569 the mathematical framework of information theory, which has been widely used in computational
570 biology.

571 The performance of the Multigenomic Entropy Based Score pipeline (MEBS) (designed for the
572 above mentioned tasks) was benchmarked on large genomic and metagenomic sets. Our results
573 support the broad applicability of this algorithm in order to classify annotated genomes as well as
574 newly sequenced environmental samples without prior culture. We also assessed to what extent
575 the final score depended on the partial or complete detection of pathways and observed a higher
576 completeness per pathway in metagenomic sequences than in individual genomes.

577 We demonstrated that a measurable score can be applied to evaluate any given metabolic
578 machinery or biogeochemical cycle in large (meta)genomic scale, holding the potential to
579 dramatically change the current view of inferring metabolic capabilities in the present 'omic'-era.

580

581 Availability and requirements

582 Project name: MEBS

583 Project home page: https://github.com/eead-csic-compbio/metagenome_Pfam_score

584 Operating system(s): Linux
585 Programming language: Python 3, Perl5, Bash,
586 Other requirements: HMMER
587 License: GNU General Public License (GPL)

588 Availability of supporting data

589 The datasets supporting the results of this article are available in the GigaDB repository [REF#]

590

591 Abbreviations

592 MEBS: Multigenomic Entropy Based Score ; S: Sulfur ; S-cycle: Sulfur cycle; SS: Sulfur Score; Suli:
593 Sulfur list ; Sucy: Sulfur cycle database; Rlist: Random list of taxonomic representatives ; MSL:
594 Mean Size Length, H': Relative Entropy ; Sur: Sulfur unconsidered ; NS: Non sulfur related
595 genomes; Gen: Genomic dataset; Met: Metagenomic dataset; GenF : Genomic Fragmented
596 dataset; CLSB: Color-less Sulfur Bacteria; SOM: Sulfur Oxidizing Microorganisms; GSB: Green Sulfur
597 Bacteria; PSB: Purple Sulfur Bacteria; SRB: Sulfate Reducing Bacteria; ESR: Elemental-Sulfur
598 Reducing microorganisms; CCC: Cuatro Ciénegas, Coahuila; HMM: Hidden Markov Models (HMMs);
599 ROC: Receiver-operating characteristic; AUC: Area Under the Curve; TPR True Positive Rates; FPR
600 False Positive Rates; GSC: Genomic Standards Consortium; EP: environmental packages.

601

602

603 Acknowledgments

604 The authors gratefully acknowledge Emilio Morella for their valuable support, feedback and comments
605 throughout the development of this project. We really appreciate the reviewers Daan Speth and Thulani
606 Makhwanyane whose valuable comments and suggestions critically improve the manuscript and software.
607 Our special thanks to acknowledge Carlos P Cantalapiedra, Seth Barribeau, Will Levitt and anonymous
608 reviewers for their comments on earlier versions of the manuscript who immensely improved the algorithm
609 and the final version of the article. The authors also thank the Laboratory of Computational and Structural
610 Biology (EEAD-CSIC) and Laboratorio de Evolución Molecular y Experimental, Instituto de Ecología, UNAM
611 for providing the computational resources described in the article. The paper was written during a
612 sabbatical leave of LEE and VSS in the University of Minnesota in Peter Tiffin and Michael Travisano
613 laboratories.

614 Funding

615 VDA is a doctoral student from Programa de Doctorado en Ciencias Biomédicas, Universidad Nacional
616 Autónoma de México (UNAM) and received fellowship 356832 from CONACYT. This research was also
617 supported by funding from WWF-Alianza Carlos Slim, Sep-Ciencia Básica Conacyt grant 238245 to both VS
618 and LEE and Spanish MINECO grant CSIC13-4E-2490. BCM was funded by Fundación ARAID. The sabbatical

619 leave of LEE and VSS at the University of Minnesota were supported by scholarships from PASPA, DGAPA,
620 UNAM.

621 Competing interest

622 The authors declare that they have no competing interest.

623

624 Author contribution

625

626 VDA, BCM and IZP wrote the paper. BCM, VDA and ACPH developed and wrote the software and performed
627 all the bioinformatics analyses. VDA produced all the figures and wrote the documentation of the software.
628 VDA and IZP conceived the manual curation of the Sulfur cycle inventory and the microbiological,
629 biogeochemical, and ecological interpretation. LE and VS provided the intellectual framework, expertise and
630 resources to develop and supervise the project. All the authors read and approved the final manuscript.
631

632 Endnotes

633 We are currently finishing the analyses to demonstrate the applicability of this approach to other
634 biogeochemical cycles (C, N, O, Fe, P). Thereby, we hope that the pipeline MEBS will facilitate analysis of
635 biogeochemical cycles or complex metabolic networks carried out by specific prokaryotic guilds, such as
636 bioremediation processes (i.e., degradation of hydrocarbons, toxic aromatic compounds, heavy metals etc.).
637 We look forward to collaborate and help other researchers by integrating comprehensive databases that
638 might be helpful to the scientific community. Furthermore, we are currently working to improve the
639 algorithm by using only a list of sequenced genomes involved in the metabolism of interest, in order to
640 reduce the manual curation effort. We are also considering taking *k-mers* instead of peptide Hidden Markov
641 Models to increase the speed of the pipeline. We anticipate that our platform will stimulate interest and
642 involvement among the scientific community to explore uncultured genomes derived from large
643 metagenomic sequences.

644

645 References

- 646 1. Thomas T, Gilbert J, Meyer F. Metagenomics - a guide from sampling to data analysis. *Microb.*
647 *Inform. Exp.* [Internet]. BioMed Central Ltd; 2012;2:3. Available from:
648 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3351745%7B&%7Dtool=pmcentrez%](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3351745%7B&%7Dtool=pmcentrez%7B&%7Drendertype=abstract)
649 [7B&%7Drendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3351745%7B&%7Dtool=pmcentrez%7B&%7Drendertype=abstract)
- 650 2. Oulas A, Pavloudi C, Polymenakou P, Pavlopoulos GA, Papanikolaou N, Kotoulas G, et al.
651 Metagenomics: Tools and insights for analyzing next-generation sequencing data derived from
652 biodiversity studies. *Bioinform. Biol. Insights.* 2015;9:75–88.
- 653 3. Morales SE, Holben WE. Linking bacterial identities and ecosystem processes: Can “omic”
654 analyses be more than the sum of their parts? *FEMS Microbiol. Ecol.* 2011;75:2–16.
- 655 4. Hug LA, Baker BJ, Anantharaman K, Brown CT, Probst AJ, Castelle CJ, et al. A new view of the tree
656 of life. *Nat. Microbiol.* [Internet]. 2016;1:16048. Available from:

- 657 <http://www.nature.com/articles/nmicrobiol201648>
- 658 5. Marco D. Metagenomics²: Current Innovations and Future Trends. Caister Academic Press; 2011.
- 659 6. Jungbluth SP, Glavina del Rio T, Tringe SG, Stepanauskas R, Rappé MS. Genomic comparisons of
660 a bacterial lineage that inhabits both marine and terrestrial deep subsurface systems. PeerJ
661 [Internet]. 2017;5:e3134. Available from: <https://peerj.com/articles/3134>
- 662 7. Iverson V, Morris RM, Frazar CD, Berthiaume CT, Morales RL, Armbrust E V. Untangling Genomes
663 from Metagenomes: Revealing an Uncultured Class of Marine Euryarchaeota. Science (80-.).
664 2012;335:587–90.
- 665 8. Sangwan N, Xia F, Gilbert JA. Recovering complete and draft population genomes from
666 metagenome datasets. Microbiome [Internet]. Microbiome; 2016;4:8. Available from:
667 <http://microbiomejournal.biomedcentral.com/articles/10.1186/s40168-016-0154-5>
- 668 9. Mehrshad M, Amoozegar MA, Ghai R, Shahzadeh Fazeli SA, Rodriguez-Valera F. Genome
669 reconstruction from metagenomic data sets reveals novel microbes in the brackish waters of the
670 Caspian Sea. Appl. Environ. Microbiol. 2016;82:1599–612.
- 671 10. Sharon I, Banfield JF. Genomes from Metagenomics. Science [Internet]. 2013;342:1057–8.
672 Available from: <http://www.ncbi.nlm.nih.gov/pubmed/24288324>
- 673 11. Hiraoka S, Yang C, Iwasaki W. Metagenomics and Bioinformatics in Microbial Ecology: Current
674 Status and Beyond. Microbes Environ. [Internet]. 2016;31:204–12. Available from:
675 https://www.jstage.jst.go.jp/article/jsme2/31/3/31_ME16024/_article
- 676 12. Caspi R, Altman T, Dreher K, Fulcher C a, Subhraveti P, Keseler IM, et al. The MetaCyc database
677 of metabolic pathways and enzymes and the BioCyc collection of pathway/genome databases.
678 Nucleic Acids Res. [Internet]. 2012 [cited 2013 May 27];40:D742-53. Available from:
679 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3245006&tool=pmcentrez&rendertype=abstract>
680
- 681 13. Kanehisa M, Goto S. KEGG: kyoto encyclopedia of genes and genomes. Nucleic Acids Res.
682 [Internet]. 2000;28:27–30. Available from:
683 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=102409%7B&%7Dtool=pmcentrez%7B&%7Drendertype=abstract>
684
- 685 14. Ye Y, Doak TG. A Parsimony Approach to Biological Pathway Reconstruction/Inference for
686 Metagenomes. Handb. Mol. Microb. Ecol. I Metagenomics Complement. Approaches. 2011;5:453–
687 60.
- 688 15. Abubucker S, Segata N, Goll J, Schubert AM, Izard J, Cantarel BL, et al. Metabolic reconstruction
689 for metagenomic data and its application to the human microbiome. PLoS Comput. Biol. [Internet].
690 2012 [cited 2014 Jan 23];8:e1002358. Available from:
691 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3374609&tool=pmcentrez&rendertype=abstract>
692
- 693 16. Larsen PE, Collart FR, Field D, Meyer F, Keegan KP, Henry CS, et al. Predicted Relative

- 694 Metabolomic Turnover (PRMT): determining metabolic turnover from a coastal marine
695 metagenomic dataset. *Microb. Inform. Exp.* [Internet]. BioMed Central Ltd; 2011;1:4. Available
696 from: <http://www.microbialinformatics.com/content/1/1/4>
- 697 17. Hanson NW, Konwar KM, Hawley AK, Altman T, Karp PD, Hallam SJ. Metabolic pathways for the
698 whole community. *BMC Genomics* [Internet]. 2014;15:619. Available from:
699 <http://www.biomedcentral.com/1471-2164/15/619>
- 700 18. Castañeda LE, Barbosa O. Metagenomic analysis exploring taxonomic and functional diversity
701 of soil microbial communities in Chilean vineyards and surrounding native forests. *PeerJ* [Internet].
702 2017;5:e3098. Available from: <https://peerj.com/articles/3098>
- 703 19. Fierer N, Leff JW, Adams BJ, Nielsen UN, Bates ST, Lauber CL, et al. Cross-biome metagenomic
704 analyses of soil microbial communities and their functional attributes. *Proc. Natl. Acad. Sci.*
705 [Internet]. 2012;109:21390–5. Available from:
706 <http://www.pnas.org/cgi/doi/10.1073/pnas.1215210110>
- 707 20. Llorens-Marès T, Yooseph S, Goll J, Hoffman J, Vila-Costa M, Borrego CM, et al. Connecting
708 biodiversity and potential functional role in modern euxinic environments by microbial
709 metagenomics. *ISME J.* [Internet]. 2015;1–14. Available from:
710 <http://www.nature.com/doi/10.1038/ismej.2014.254>
- 711 21. Quaiser A, Zivanovic Y, Moreira D, López-García P. Comparative metagenomics of bathypelagic
712 plankton and bottom sediment from the Sea of Marmara. *ISME J.* [Internet]. 2011 [cited 2014 Jul
713 16];5:285–304. Available from:
714 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3105693&tool=pmcentrez&rendertype=abstract>
715
- 716 22. Xie W, Wang F, Guo L, Chen Z, Sievert SM, Meng J, et al. Comparative metagenomics of
717 microbial communities inhabiting deep-sea hydrothermal vent chimneys with contrasting
718 chemistries. *ISME J.* [Internet]. Nature Publishing Group; 2011;5:414–26. Available from:
719 <http://dx.doi.org/10.1038/ismej.2010.144>
- 720 23. Delmont TO, Malandain C, Prestat E, Larose C, Monier J-M, Simonet P, et al. Metagenomic
721 mining for microbiologists. *ISME J.* [Internet]. 2011 [cited 2014 Apr 30];5:1837–43. Available from:
722 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3223302&tool=pmcentrez&rendertype=abstract>
723
- 724 24. Ganesh S, Parris DJ, DeLong EF, Stewart FJ. Metagenomic analysis of size-fractionated
725 picoplankton in a marine oxygen minimum zone. *ISME J.* [Internet]. Nature Publishing Group;
726 2014;8:187–211. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/24030599>
- 727 25. Segata N, Izard J, Waldron L, Gevers D, Miropolsky L, Garrett WS, et al. Metagenomic biomarker
728 discovery and explanation. *Genome Biol.* [Internet]. BioMed Central Ltd; 2011;12:R60. Available
729 from: <http://genomebiology.com/2011/11/6/R60>
- 730 26. Parks DH, Tyson GW, Hugenholtz P, Beiko RG. STAMP: Statistical analysis of taxonomic and
731 functional profiles. *Bioinformatics.* 2014;30:3123–4.

- 732 27. Sohn MB, Du R, An L. A robust approach for identifying differentially abundant features in
733 metagenomic samples. *Bioinformatics*. 2015;31:2269–75.
- 734 28. Kullback S, Leibler RA. On Information and Sufficiency. *Ann. Math. Stat.* 1951;22:79–86.
- 735 29. Commenges D. Information Theory and Statistics: an overview. 2015;1–22. Available from:
736 <http://arxiv.org/abs/1511.00860>
- 737 30. Hertz GZ, Stormo GD. Identifying DNA and protein patterns with statistically significant
738 alignments of multiple sequences. *Bioinformatics* [Internet]. 1999;15:563–77. Available from:
739 <http://www.ncbi.nlm.nih.gov/pubmed/10487864>
- 740 31. Dar S a, Yao L, van Dongen U, Kuenen JG, Muyzer G. Analysis of diversity and activity of sulfate-
741 reducing bacterial communities in sulfidogenic bioreactors using 16S rRNA and *dsrB* genes as
742 molecular markers. *Appl. Environ. Microbiol.* [Internet]. 2007;73:594–604. Available from:
743 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1796976%7B&%7Dtool=pmcentrez%7B&%7Drendertype=abstract>
- 744
- 745 32. Kerkhof LJ. Phylogeography of Sulfate-Reducing Bacteria among Disturbed Sediments ,
746 Disclosed by Analysis of the Dissimilatory Sulfite Reductase Genes (*dsrAB*) Phylogeography of
747 Sulfate-Reducing Bacteria among Disturbed Sediments , Disclosed by Analysis of the *Dissi*. 2005;
- 748 33. Loy A, Duller S, Baranyi C, Mussmann M, Ott J, Sharon I, et al. Reverse dissimilatory sulfite
749 reductase as phylogenetic marker for a subgroup of sulfur-oxidizing prokaryotes. *Environ.*
750 *Microbiol.* [Internet]. 2009;11:289–99. Available from:
751 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2702494%7B&%7Dtool=pmcentrez%7B&%7Drendertype=abstract>
- 752
- 753 34. Hügler M, Gärtner A, Imhoff JF. Functional genes as markers for sulfur cycling and CO₂ fixation
754 in microbial communities of hydrothermal vents of the Logatchev field. *FEMS Microbiol. Ecol.*
755 [Internet]. 2010 [cited 2013 Nov 3];73:526–37. Available from:
756 <http://www.ncbi.nlm.nih.gov/pubmed/20597983>
- 757 35. Meyer B, Kuever J. Molecular analysis of the diversity of sulfate-reducing and sulfur-oxidizing
758 prokaryotes in the environment, using *aprA* as functional marker gene. *Appl. Environ. Microbiol.*
759 [Internet]. 2007 [cited 2013 Jun 23];73:7664–79. Available from:
760 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2168068&tool=pmcentrez&rendertyp>
761 [e=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2168068&tool=pmcentrez&rendertyp)
- 762 36. MEBS [Internet]. Available from: [https://github.com/eead-csic-](https://github.com/eead-csic-compbio/metagenome_Pfam_score)
763 [compbio/metagenome_Pfam_score](https://github.com/eead-csic-compbio/metagenome_Pfam_score)
- 764 37. Enzyme Nomenclature [Internet]. Available from: <http://enzyme.expasy.org/>. Accesed 05 May
765 2016.
- 766 38. Magrane M, Consortium UP. UniProt Knowledgebase: A hub of integrated protein data.
767 *Database*. 2011;2011:1–13.
- 768 39. Reference and Representative Genomes [Internet]. Available from:

- 769 <https://www.ncbi.nlm.nih.gov/genome/browse/reference/>. Accessed 21 Dec. 2016
- 770 40. Genome clusters [Internet]. Available from:
771 <http://microbiome.wlu.ca/research/redundancy/redundancy.cgi>. Accessed 21 Dec. 2016
- 772 41. Moreno-Hagelsieb G, Wang Z, Walsh S, ElSherbiny A. Phylogenomic clustering for selecting
773 non-redundant genomes for comparative genomics. *Bioinformatics* [Internet]. 2013 [cited 2014
774 Sep 18];29:947–9. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/23396122>
- 775 42. NCBI genome assembly summary file [Internet]. Available from:
776 ftp://ftp.ncbi.nlm.nih.gov/genomes/refseq/assembly_summary_refseq.txt. Accessed 21 Dec. 2016.
- 777 43. NCBI. NCBI FTP site [Internet]. Available from: <ftp://ftp.ncbi.nlm.nih.gov/genomes/all/>. Accessed
778 21 Dec. 2016.
- 779 44. MG-RAST [Internet]. Available from: <http://metagenomics.anl.gov/>. Accessed 10 Nov. 2016.
- 780 45. Prakash T, Taylor TD. Functional assignment of metagenomic data: challenges and applications.
781 *Brief. Bioinform.* [Internet]. 2012;13:711–27. Available from:
782 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3504928%7B%7Dtool=pmcentrez%](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3504928%7B%7Dtool=pmcentrez%7B%7Drendertype=abstract)
783 [7B%7Drendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3504928%7B%7Dtool=pmcentrez%7B%7Drendertype=abstract)
- 784 46. Zhong C, Yang Y, Yooseph S. GRASP: Guided reference-based assembly of short peptides.
785 *Nucleic Acids Res.* 2015;43:e18.
- 786 47. Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, et al. InterProScan 5: Genome-scale
787 protein function classification. *Bioinformatics.* 2014;30:1236–40.
- 788 48. Finn RD, Tate J, Mistry J, Coghill PC, Sammut SJ, Hotz H-R, et al. The Pfam protein families
789 database. *Nucleic Acids Res.* [Internet]. 2008;36:D281–8. Available from:
790 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2238907%7B%7Dtool=pmcentrez%](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2238907%7B%7Dtool=pmcentrez%7B%7Drendertype=abstract)
791 [7B%7Drendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2238907%7B%7Dtool=pmcentrez%7B%7Drendertype=abstract)
- 792 49. Haft DH. The TIGRFAMs database of protein families. *Nucleic Acids Res.* [Internet]. 2003 [cited
793 2014 Aug 11];31:371–3. Available from:
794 <http://nar.oxfordjournals.org/lookup/doi/10.1093/nar/gkg128>
- 795 50. Gough J, Chothia C. SUPERFAMILY: HMMs representing all proteins of known structure. SCOP
796 sequence searches, alignments and genome assignments. *Nucleic Acids Res.* [Internet].
797 2002;30:268–72. Available from:
798 [http://www.ncbi.nlm.nih.gov/pmc/articles/PMC99153/%5Cnhttp://www.ncbi.nlm.nih.gov/pmc/ar](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC99153/%5Cnhttp://www.ncbi.nlm.nih.gov/pmc/articles/PMC99153/pdf/gkf097.pdf)
799 [ticles/PMC99153/pdf/gkf097.pdf](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC99153/pdf/gkf097.pdf)
- 800 51. Finn RD, Clements J, Eddy SR. HMMER web server: interactive sequence similarity searching.
801 *Nucleic Acids Res.* [Internet]. 2011;39:W29–37. Available from:
802 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3125773%7B%7Dtool=pmcentrez%](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3125773%7B%7Dtool=pmcentrez%7B%7Drendertype=abstract)
803 [7B%7Drendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3125773%7B%7Dtool=pmcentrez%7B%7Drendertype=abstract)
- 804 52. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine
805 Learning in {P}ython. *J. Mach. Learn. Res.* 2011;12:2825–30.

- 806 53. KEGG Mapper [Internet]. Available from: <http://www.genome.jp/kegg/mapper.html>. Accessed 4
807 Sept. 2017.
- 808 54. Alcaraz LD, Olmedo G, Bonilla G, Cerritos R, Hernández G, Cruz A, et al. The genome of *Bacillus*
809 *coahuilensis* reveals adaptations essential for survival in the relic of an ancient marine
810 environment. *Proc. Natl. Acad. Sci. U. S. A.* [Internet]. 2008;105:5803–8. Available from:
811 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2311347&tool=pmcentrez%](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2311347&tool=pmcentrez&rendertype=abstract)
812 [7B&%7Drendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2311347&tool=pmcentrez&rendertype=abstract)
- 813 55. Pedroni P, Volpe AD, Galli G, Mura GM, Pratesi C, Grandi G. Characterization of the locus
814 encoding the [Ni-Fe] sulfhydrogenase from the archaeon *Pyrococcus furiosus*: Evidence for a
815 relationship to bacterial sulfite reductases. *Microbiology*. 1995;141:449–58.
- 816 56. Taguchi Y, Sugishima M, Fukuyama K. Crystal Structure of a Novel Zinc-Binding ATP Sulfurylase
817 from *Thermus*. 2004;4111–8.
- 818 57. Santos AA, Venceslau SS, Grein F, Leavitt WD, Dahl C, Johnston DT, et al. A protein trisulfide
819 couples dissimilatory sulfate reduction to energy conservation. *Science (80-.)*. 2015;350:1541–5.
- 820 58. Emerson D, Field EK, Chertkov O, Davenport KW, Goodwin L, Munk C, et al. Comparative
821 genomics of freshwater Fe-oxidizing bacteria: implications for physiology, ecology, and
822 systematics. *Front. Microbiol.* [Internet]. 2013 [cited 2014 Jun 11];4:254. Available from:
823 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3770913&tool=pmcentrez&rendertyp](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3770913&tool=pmcentrez&rendertype=abstract)
824 [e=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3770913&tool=pmcentrez&rendertype=abstract)
- 825 59. Nakagawa S, Shimamura S, Takaki Y, Suzuki Y, Murakami S, Watanabe T, et al. Allying with
826 armored snails: the complete genome of gammaproteobacterial endosymbiont. *ISME J.* [Internet].
827 Nature Publishing Group; 2014;8:40–51. Available from:
828 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3869010&tool=pmcentrez&rendertyp](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3869010&tool=pmcentrez&rendertype=abstract)
829 [e=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3869010&tool=pmcentrez&rendertype=abstract)
- 830 60. Manzella MP, Holmes DE, Rocheleau JM, Chung A, Reguera G, Kashefi K. The complete genome
831 sequence and emendation of the hyperthermophilic, obligate iron-reducing archaeon “*Geoglobus*
832 *ahangari*” strain 234T. *Stand. Genomic Sci.* [Internet]. *Standards in Genomic Sciences*; 2015;10:77.
833 Available from:
834 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4600277&tool=pmcentrez%](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4600277&tool=pmcentrez&rendertype=abstract)
835 [7B&%7Drendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4600277&tool=pmcentrez&rendertype=abstract)
- 836 61. Carbonero F, Benefiel AC, Alizadeh-Ghamsari AH, Gaskins HR. Microbial pathways in colonic
837 sulfur metabolism and links with health and disease. *Front. Physiol.* 2012;3 NOV:1–11.
- 838 62. Nakagawa S, Takaki Y, Shimamura S, Reysenbach A-L, Takai K, Horikoshi K. Deep-sea vent
839 epsilon-proteobacterial genomes provide insights into emergence of pathogens. *Proc. Natl. Acad.*
840 *Sci. U. S. A.* 2007;104:12146–50.
- 841 63. Fawcett T. An introduction to ROC analysis. *Pattern Recognit. Lett.* 2006;27:861–74.
- 842 64. Field D, Sterk P, Kottmann R, Smet JW De, Amaral-zettler L, Cole JR, et al. Genomic Standards
843 Consortium Projects The Genomic Standards Consortium Initiating and Maintaining a Project

- 844 within the GSC The GSC Project Description template provides a References☐: 2014;599–601.
- 845 65. Alazard D, Joseph M, Battaglia-Brunet F, Cayol JL, Ollivier B. Desulfosporosinus acidiphilus sp.
846 nov.: A moderately acidophilic sulfate-reducing bacterium isolated from acid mining drainage
847 sediments. *Extremophiles*. 2010;14:305–12.
- 848 66. Romano C, D’Imperio S, Woyke T, Mavromatis K, Lasken R, Shock EL, et al. Comparative
849 genomic analysis of phylogenetically closely related *Hydrogenobaculum* sp. isolates from
850 yellowstone national park. *Appl. Environ. Microbiol.* 2013;79:2932–43.
- 851 67. Chen L, Ren Y, Lin J, Liu X, Pang X, Lin J. Acidithiobacillus caldus sulfur oxidation model based on
852 transcriptome analysis between the wild type and sulfur oxygenase reductase defective mutant.
853 *PLoS One* [Internet]. 2012 [cited 2013 Apr 24];7:e39470. Available from:
854 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3440390&tool=pmcentrez&rendertype=abstract>
855
- 856 68. Liljeqvist M, Valdes J, Holmes DS, Dopson M. Draft genome of the psychrotolerant acidophile
857 *Acidithiobacillus ferrivorans* SS3. *J. Bacteriol.* 2011;193:4304–5.
- 858 69. Imhoff JF, Pfennig N. *Thioflavicoccus mobilis* gen. nov., sp. nov., a novel purple sulfur bacterium
859 with bacteriochlorophyll b. *Int. J. Syst. Evol. Microbiol.* 2001;51:105–10.

860

861 63. GigaDB DOI Citation

862

863

864 Figure Legends

865 **Figure 1.** Schematic representation of the four stages of the MEBS algorithm focusing on the S-cycle. The
866 first step consists on the systematic curation of a database containing the metabolic information of S-cycle,
867 which is reduced to a FASTA-file of proteins involved (Sucy) and a list of 161 related microorganisms (Suli). A
868 thousand lists of 161 random-sampled genomes were used as negative control (Rlist). The training dataset
869 comprises 2,107 genomes (Gen), which were fragmented in different sizes by considering the Mean Size
870 Length (MSL) of 935 metagenomes (Met). In the second stage the domain composition of Sucy proteins is
871 obtained by scanning Pfam-A, resulting in the Pfam-Sucy database. Then, the relative entropy (H') of each
872 Sucy-Pfam domain is obtained in the third stage. Finally, the precomputed entropies in Gen and GenF are
873 used to evaluate full-length genomic sequences (real) and metagenomic sequences of variable MSL (in this
874 example A, B and C).

875 **Figure 2.** Sulfur cycle at global scale. The most important organic and inorganic S-compounds derived from
876 biogeochemical processes are arranged according to the Standard Gibbs free energy of formation described
877 in Caspi et al., (2012). The left column indicates whether specific microorganisms are able to use those S-
878 compounds, as a source of Carbon (C), Nitrogen (N), Energy (E) or Electron donors (°). Double asterisks
879 indicate whether the S-compound is used as sole source, of C, N, or E. The corresponding electron acceptors
880 in redox-coupled reactions using the S-compound as electron donor are not shown. The right column

881 indicates whether the S-compound is used as fermentative substrate (F) or terminal electron acceptor in
882 respiratory processes (R). Colored boxes summarize the metabolic guilds involved in the metabolism of S-
883 compounds, in oxidation (i.e., CLSB, SOM, PSB, and GSB) or reduction (SR, SRB) processes. The complete list
884 of S-based microorganisms (Suli) is found in Table S1. Figure based on annotations from MetaCyc [12].

885 **Figure 3.** Comprehensive network representation of the machinery of the biogeochemical S-cycle in a single
886 cell. The 28 molecular pathways involved in the metabolism of sulfur compounds described in Table 1 are
887 included. The enzymatic steps are depicted as rectangles followed by arrows indicating the direction of the
888 reaction. Green hexagons represent metabolic links to other metabolisms. Bold dashed arrows indicate
889 bidirectional reactions. Inorganic S-compounds have been arranged according to their reduction potential,
890 from the most oxidized (yellow) to the most reduced (red) compounds. Grey rectangles indicate enzymes
891 acting in disproportionation processes in which a reactant is both oxidized and reduced in the same
892 chemical reaction, forming two separate compounds. Input biogeochemical S-compounds are shown
893 outside and connected with bold arrows. Dashed arrows indicate S-compounds excreted out of the cell. The
894 upper half of the modeled cell depicts the processes involved in the use of organic S-compounds (orange
895 circles) found in natural environments and used as source of carbon, sulfur and/or energy in several
896 aerobic/anaerobic strains described in Figure 2.

897 **Figure 4.** Entropy values of Sulfur-derived protein domains. A) Heatmap showing the entropy values (H') of
898 the 112 Pfam domains identified in proteins curated in SuCy. B) Difference between entropies estimated
899 from sizes categories of growing peptide size (GenF) and the real values measured within complete
900 genomes (Gen). Error bars show standard deviations. Both graphs were obtained with script
901 *plot_entropy.py*. Clustering of the Pfam relative entropies obtained in Gen and GenF produced with the
902 Ward method. Log frequency of the entropy values computed in the random test is colored in purple (see
903 scale bar). Cluster 0 (blue) groups protein domains with low relative entropy that overlap with the random
904 distribution. Cluster 1 (green) includes the Pfam domains that fulfill the requirements to be used as
905 molecular markers (high H' and low standard deviation, std). Red dots (cluster 2) correspond to Pfam
906 domains with high H' and std. The cluster was produced with script *F_meanVSstd.py*

907 **Figure 5.** Distribution of Sulfur Score (SS) in 2,107 non-redundant genomes (Gen). A) Subsets of genomes
908 annotated in Suli (n=161); ii) Sur, genomes not listed in Suli with SS > 4 and candidates to be S-related
909 microorganisms (n=192); iii) rest of the genomes in Gen (NS, n=1,754). According to the curated species,
910 True Positives can be defined as genomes with SS > max(SS_{NS}) distribution, whereas True Negatives are
911 those with SS < min(SS_{Suli}). B) Assignment of the 192 genomes in Sur to ecological categories based on
912 literature reports. C) Distribution of SS for different S-metabolic guilds, and the genomes in Sur. D) ROC
913 curve with Area Under the Curve (AUC) indicated together with thresholds for some False Positive Rates
914 (FPR).

915 **Figure 6.** Distribution of Sulfur Score (SS) in the metagenomic dataset Met. A) Geo-localized metagenomes
916 sampled around the globe are colored according to their SS values. The following cut-off values correspond
917 to the 95th percentiles of seven Mean Size Length classes (30, 60, 100, 150, 200, 250 and 300 aa): 7.66,
918 9.70, 8.81, 8.51, 8.18, 8.98 and 7.61, respectively. Circles with thick blue border indicate metagenomes with
919 SS ≥ the 95th percentile. B) Distribution of SS values observed in 935 metagenomes classified in terms of
920 features (X-axis) and colored according to their particular habitats Features are sorted according to their
921 median SS values. ccc: metagenomes from Cuatro Ciénegas, Coahuila, Mexico. Green lines indicate the
922 lowest and largest 95th percentiles observed across MSL classes.

923

924 **Figure 7.** Metabolic completeness of the metabolic pathways described in Table 1. A) Boxplot distribution of
925 the pathway completeness in genomic and B) metagenomic datasets. C) Linear regression models of the
926 Sulfur Score (*SS*) and the mean completeness in Gen and D) Met dataset. E) Heatmap showing the metabolic
927 completeness of the following genomes: *Desulfovibrio vulgaris DP4* (*SS*=11,442), *Ammonifex degensii KC4*
928 (*SS*=12.508); *Pelodictyon phaeoclathratiforme BU-1* (*SS*=11,836); *Thiobacillus denitrificans ATCC 25259*
929 (*SS*=11,61); PSB: *Allochromatium vinosum DSM 180* (*SS*=10.737); Sur *Methanosarcina barkeri MS* (*SS*=
930 5,93); *Sulfolobus acidocaldarius DSM 639* (*SS*=5,457); *Synechococcus sp. JA-2-3Ba 2-13* (*SS*=3,704);
931 *Hyphomicrobium denitrificans 1NES1* (*SS*= 3,236); *Ruegeria pomeroyi DSS-3* (*SS*=2,707); *Enterococcus*
932 *durans* (*SS*=-0,194); *Micrococcus luteus NCTC_2665* (*SS*=-3,588). F) Heatmap showing the metabolic
933 completeness of the metagenomes with the following MG-RAST ids and corresponding scores: 4489656.3
934 (*SS*=-2,649); 4440320.3(*SS*=0,1); 4441663.3(*SS*=9,986); 4510168.3 (*SS*=7,781) ; 4493725.3 (*SS*=9,547) ;
935 4461840.3 (*SS*=8,813); 4441599.3(*SS*=9,274); 4451035.3(*SS*=9,918); 4525341.3(*SS*=9,287);
936 4489328.3(*SS*=4,958); 4478222.3(*SS*=4,88). The color codes at the top of the heatmap correspond to
937 different environments. For a more detailed description of each metagenome see Table S8.

938

939 **Additional files -Supplementary Information**

940 The supplementary pdf file contains the following information:

941 **Supplementary figure S1.** Histogram distribution of the Mean Size Length of metagenomes in Met and the
942 input sulfur proteins.

943 **Supplementary figure S2.** Visualization of the Pfam domains mapped onto KEGG metabolic pathways

944 **Supplementary figure S3.** Comparison of clustering methods of the 112 Pfam entropies using script
945 *plot_cluster_comparison.py*

946 **Supplementary figure S4.** Clustering comparison between Birch and Ward clustering methods to stand out
947 the Pfam entropies with high *H'* and low std using the script

948 **Supplementary figure S5.** . Distribution of entropy values of 112 Pfam domains inferred from random-
949 sampled and Suli genomes.

950 **Supplementary figure S6.** Comparison of Sulfur Scores (*SS*) with data obtained three years ago (2014), with
951 the current data described in the article.

952 **Supplementary table S4.** Informative Pfam's with high *H'* and high std (not used as molecular marker
953 genes) in metagenomic fragmented data.

954 **Supplementary table S5.** Percentile distribution of the 112 Pfam entropies in the random test

955 **Supplementary table S10.** Statistics of *SS* computed on genomic (Gen, real sequences) and metagenomic
956 (Met, with increasing Mean Size Length, from 30 to 300aa) datasets

957 In separated excel files the following Supplementary tables are also provided:

958 **Supplementary table S1:** Table S1. Comprehensive list of the taxonomic representatives of sulfur cycle
959 including Sulfur list or 'Suli' containing 161-curated genomes used as input for the pipeline

960 **Supplementary table S2.** Sucey database containing the identifiers of the Sulfur proteins and their
961 corresponding annotations derived from Interproscan and manual curation.

962 **Supplementary table S3.** Sulfur Pfam domains (Pfam-Sucey), and their corresponding mapping into KEGG (KO
963 number), and the manual assignment into sulfur metabolic pathways

964 **Supplementary table S6.** Gen dataset containing their corresponding SS and taxonomy assignment.

965 **Supplementary table S7.** Manual annotation of Sulfur unconsidered or related microorganisms (Sur) with
966 SS>4 Supplementary table S8

967 **Supplementary table S8.** Met dataset with their corresponding SS values and metadata.

968 **Supplementary table S9.** Manually annotated high scoring metagenomes along with their ecological
969 capabilities in terms of sulfur cycle
970

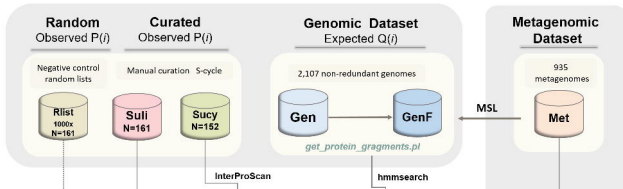
971 **Supplementary table S11.** Metabolic completeness in Gen dataset for each of the 28 metabolic pathways of the S-
972 cycle described in Table 1. (Pathway 29 contains the proposed marker genes)
973

974 **Supplementary table S12.** Metabolic completeness in Men dataset for each of the 28 metabolic pathways of the S-
975 cycle described in Table 1. (Pathway 29 contains the proposed marker genes)
976

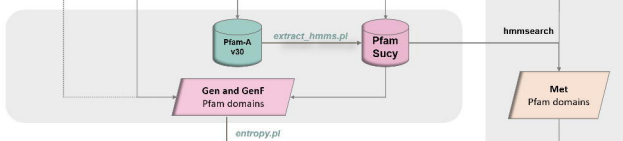
977 **Supplementary table S13.** Frequency and description of the most complete genomes in terms of S-cycle
978 metabolic pathways
979

980
981

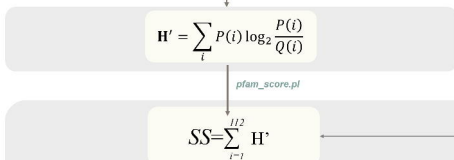
Stage 1 Curation and Omic datasets



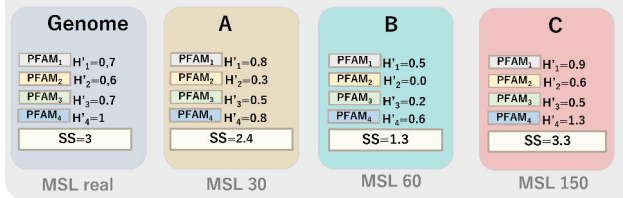
Stage 2 Pfam domain composition



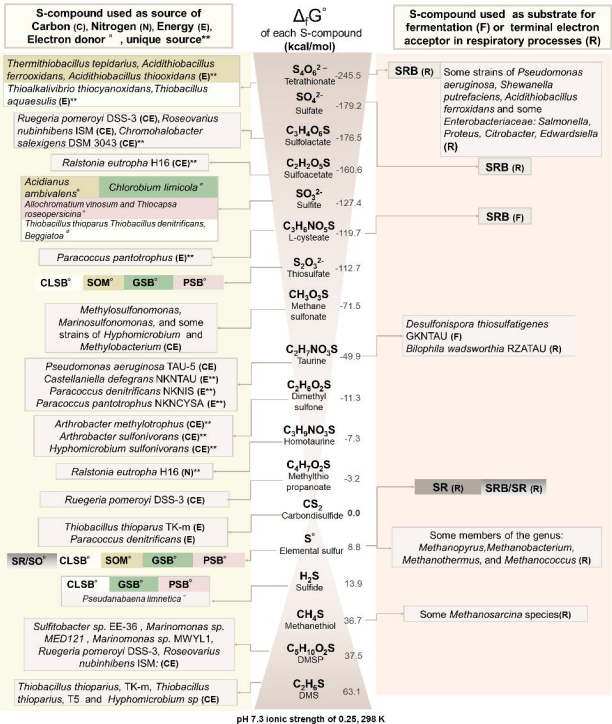
Stage 3 Relative entropy



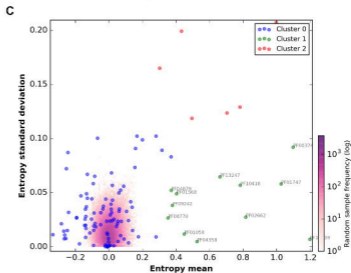
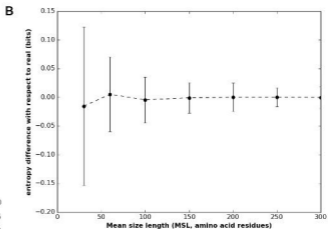
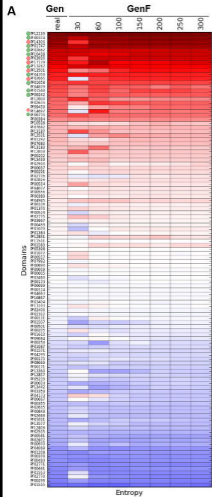
Stage 4 MEBS Sulfur Score (SS)

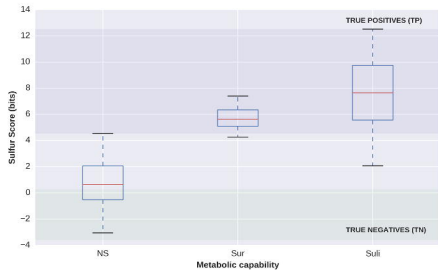
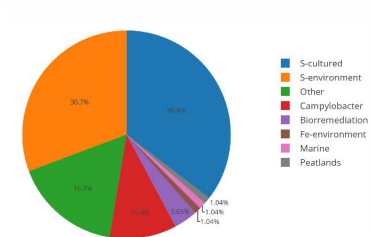
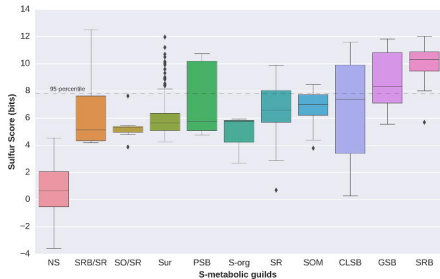
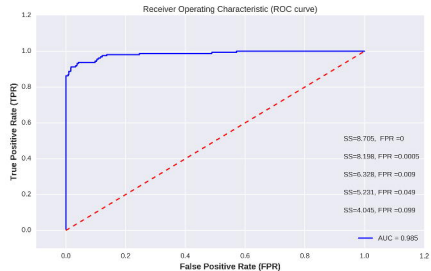


Biogeochemical Sulfur cycle at global scale

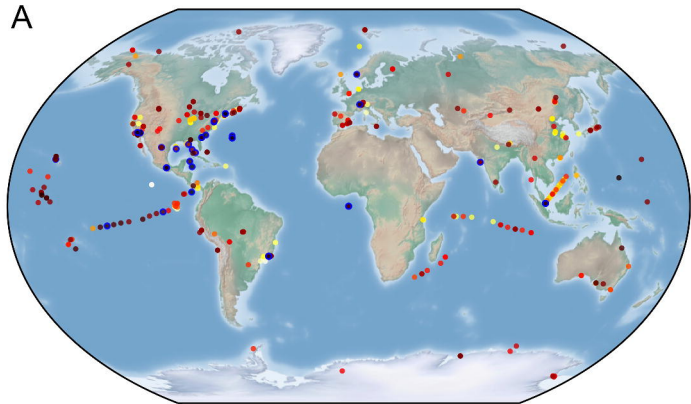


CLSB	Color-less Sulfur Bacteria: 24 genera (<i>i.e</i> <i>Beggiatoa</i> , <i>Thiomargarita</i> , <i>Thiobacillus</i>)
SOM	Sulfur Oxidizing Microorganisms: 12 genera (<i>i.e</i> <i>Thermithiobacillus</i> , <i>Acidithiobacillus</i>)
GSB	Green Sulfur Bacteria: 9 genera (<i>i.e</i> <i>Clorobaculum</i> , <i>Chloroflexus</i> , <i>Chlorobium</i>)
PSB	Purple Sulfur Bacteria: 25 genera (<i>i.e</i> <i>Ectothiorhodospira</i> , <i>Cromatium</i>)
SRB	Sulfate Reducing Bacteria: 40 genera (<i>i.e</i> <i>Desulfivibrio</i> , <i>Desulfotomaculum</i> , <i>Desulfotignum</i>)
SR	Elemental-Sulfur Reducing microorganisms: 19 genera (<i>i.e</i> <i>Sulfospirillum</i> , <i>Desulfurella</i> , <i>Thermoproteus</i>)
ESR/SO	Elemental-Sulfur Reducing and Sulfur Oxidizing microorganisms: includes 4 genera: <i>Sulfolobus</i> , <i>Acidianus</i> , <i>Aquifex</i> , <i>Thermoplasma</i>
SRB/ESR	Sulfate Reducing Bacteria and Elemental Sulfur Reducing microorganism : 42 genera (<i>i. e</i> <i>Thermovirga</i> , <i>Pyrolobus</i>)

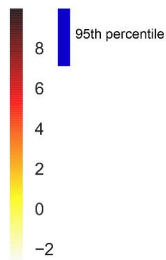


A**B****C****D**

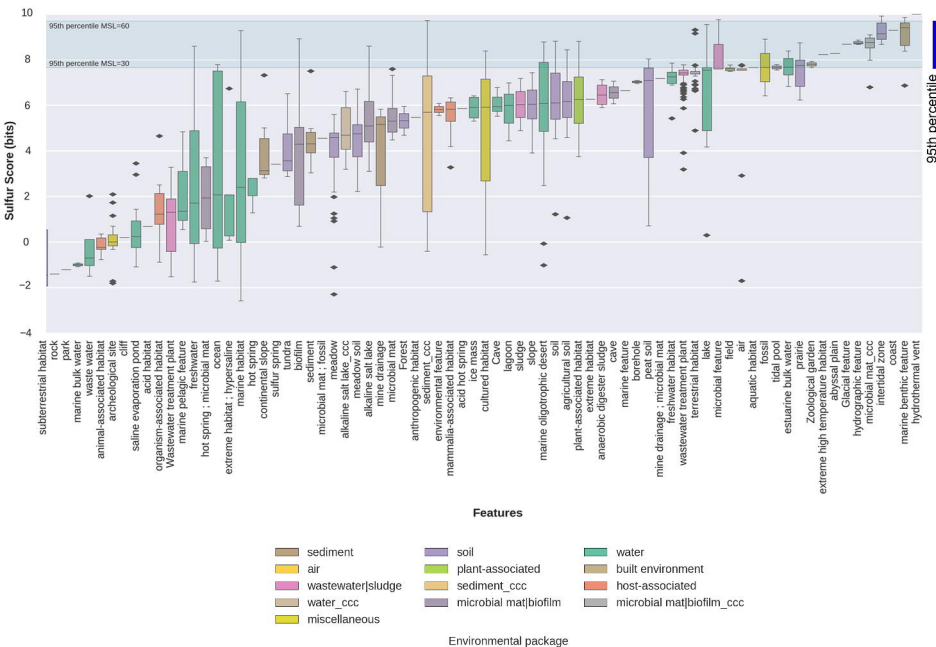
A



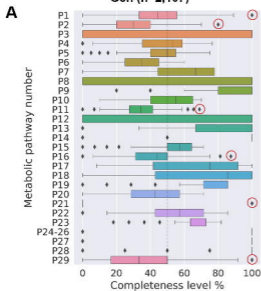
Sulfur Score (SS)



B



Gen (n=2,107)



Met (n=935)

