

HMSS2: an advanced tool for the analysis of sulfur metabolism, including organosulfur compound transformation, in genome and metagenome assemblies

Tomohisa Tanabe¹ and Christiane Dahl¹

¹Rheinische Friedrich-Wilhelms-Universität Bonn

May 12, 2023

Abstract

The global sulfur cycle has implications for human health, climate change, biogeochemistry, and bioremediation. The organosulfur compounds that participate in this cycle not only represent a vast reservoir of sulfur, but are also used by prokaryotes as sources of energy and/or carbon. Closely linked to the inorganic sulfur cycle, it involves the interaction of prokaryotes, eukaryotes, and chemical processes. However, ecological and evolutionary studies of the conversion of organic sulfur compounds are hampered by the poor conservation of the relevant pathways and their variation even within strains of the same species. In addition, several proteins involved in the conversion of sulfonated compounds are related to proteins involved in sulfur dissimilation or turnover of other compounds. Therefore, the enzymes involved in the metabolism of organic sulfur compounds are usually not correctly annotated in public databases. To address this challenge, we have developed HMSS2, a profiled Hidden Markov Model-based tool for rapid annotation and synteny analysis of organic and inorganic sulfur cycle proteins in prokaryotic genomes. Compared to its previous version (HMS-S-S), HMSS2 includes several new features. HMM-based annotation is now supported by non-homology criteria and covers the metabolic pathways of important organosulfur compounds, including dimethylsulfopropionate, taurine, isethionate, and sulfoquinovose. In addition, the calculation speed has been increased by a factor of four and the available output formats have been extended to include iTol compatible datasets, and customised sequence FASTA files

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29

Title: HMSS2: an advanced tool for the analysis of sulfur metabolism, including organosulfur compound transformation, in genome and metagenome assemblies

Authors: Tomohisa Sebastian Tanabe* and Christiane Dahl*

Institut für Mikrobiologie & Biotechnologie, Rheinische Friedrich-Wilhelms-Universität Bonn,
Bonn, Germany

Running head: HMSS2, a sulfur metabolism analysis tool

*Correspondence: Tomohisa Sebastian Tanabe and Christiane Dahl
Institut für Mikrobiologie & Biotechnologie
Rheinische Friedrich Wilhelms-Universität Bonn
Meckenheimer Allee 168
53115 Bonn, Germany
Tel. +49-228-735591
Fax +49-228-737476
E-mails: s6totana@uni-bonn.de (T.S.T.), chdahl@uni-bonn.de
(C.D.)

Keywords: Hidden Markov model (HMM) database, organosulfur compounds, sulfur metabolism, dimethylsulphopropionate, sulfoquinovose

Abbreviations:

30 **Abstract**

31 The global sulfur cycle has implications for human health, climate change, biogeochemistry,
32 and bioremediation. The organosulfur compounds that participate in this cycle not only
33 represent a vast reservoir of sulfur, but are also used by prokaryotes as sources of energy
34 and/or carbon. Closely linked to the inorganic sulfur cycle, it involves the interaction of
35 prokaryotes, eukaryotes, and chemical processes. However, ecological and evolutionary
36 studies of the conversion of organic sulfur compounds are hampered by the poor conservation
37 of the relevant pathways and their variation even within strains of the same species. In
38 addition, several proteins involved in the conversion of sulfonated compounds are related to
39 proteins involved in sulfur dissimilation or turnover of other compounds. Therefore, the
40 enzymes involved in the metabolism of organic sulfur compounds are usually not correctly
41 annotated in public databases. To address this challenge, we have developed HMSS2, a
42 profiled Hidden Markov Model-based tool for rapid annotation and synteny analysis of organic
43 and inorganic sulfur cycle proteins in prokaryotic genomes. Compared to its previous version
44 (HMS-S-S), HMSS2 includes several new features. HMM-based annotation is now supported
45 by non-homology criteria and covers the metabolic pathways of important organosulfur
46 compounds, including dimethylsulfpopropionate, taurine, isethionate, and sulfoquinovose. In
47 addition, the calculation speed has been increased by a factor of four and the available output
48 formats have been extended to include iTol compatible datasets, and customised sequence
49 FASTA files.

50
51
52

53 1 INTRODUCTION

54 The global organic sulfur cycle occurs in both terrestrial and aquatic environments and
55 involves the interplay of prokaryotes, eukaryotes, and chemical processes. Millions of
56 megatonnes of sulfonated compounds are produced annually by biological and industrial
57 processes. These compounds not only represent a vast reservoir of sulfur but can also be used
58 by prokaryotes as a sources of energy and carbon (Moran & Durham, 2019). Understanding
59 the mechanisms and ecological interactions of prokaryotes in the organic sulfur cycle is of
60 great importance because the decomposition of organic sulfur compounds affects human
61 health, bacterial virulence in infection (Dhouib et al., 2021), global warming, bioremediation
62 processes such as wastewater treatment (Schäfer et al., 2010), and is linked to the
63 biogeochemical cycling of sulfur between habitats (Koch & Dahl, 2018).

64 Sulfonated compounds can range from small size with only a C₁ carbon skeleton up to
65 sulfonated lipids with long-chain alkanes, amino acids such as cysteine, or sulfur-containing
66 cofactors with complex structures such as lipoate (Boden & Hutt, 2019; Goddard-Borger &
67 Williams, 2017; Moran & Durham, 2019). While chemistry offers an infinite number of possible
68 sulfonated compounds and new ones are being discovered all the time, these compounds
69 often lack a described metabolic function or the pathways for their synthesis or degradation
70 have not been elucidated (Thume et al., 2018). Only the most abundant sulfonated
71 compounds, such as sulfoquinovose, dimethylsulfopropionate (DMSP), taurine, isethionate,
72 cysteine, and methionine, have been studied biochemically in terms of synthesis and
73 degradation pathways.

74 In aquatic environments, the anti-stress molecule DMSP is the most well-known
75 organosulfur compound (Kiene et al., 2000). Mainly produced by macroalgae and
76 phytoplankton, it is emitted by around 600 million tonnes per year. Bacterial DMSP
77 degradation in the oceans, salt marshes, and coastal regions is the major source of
78 dimethylsulfide (DMS), which is released at a rate of about 300 million tonnes per year (Moran
79 & Durham, 2019). As a volatile compound, DMS affects atmospheric chemistry and global
80 warming by forming cloud condensation nuclei that increase the reflection of solar radiation
81 (Schäfer et al., 2010). In the context of the global sulfur cycle, DMS acts as a link between the
82 terrestrial, atmospheric and aquatic environments (Lovelock et al., 1972). DMS-derived

83 carbon and sulfur are used as electron acceptors or donors during dissimilation, or are
84 assimilated via the intermediates dimethylsulfone and methanesulfinate (Fig. 1).

85 Sulfonated lipids are estimated to be the largest reservoir of sulfur in terrestrial
86 ecosystems (Goddard-Borger & Williams, 2017). Sulfoquinovose is a sulfonated glucose
87 derivate and the most common part of the head group of sulfolipids which are integral part of
88 thylakoid membranes of chloroplasts and photosynthetic systems. Mainly produced by plants,
89 algae, and cyanobacteria its turnover rate has been estimated at around 10 billion tonnes per
90 year (Goddard-Borger & Williams, 2017). The bacterial decomposition of sulfoquinovose
91 involves several different pathways similar to the degradation of glucose (Fig. 2a), with the
92 exception that smaller sulfonated compounds are often released, since complete utilisation
93 with release of free sulfur by a single organism is often not possible (Wei et al., 2022). Release
94 and scavenging of sulfonated intermediates is achieved by various transport systems (Fig. 2b).
95 Sulfoquinovose decomposition and release of inorganic sulfur is then completed by pathways
96 linked to taurine, isethionate and/or sulfoacetate (Fig. 2c). In summary, prokaryotic utilization
97 of these organic compounds as sources of sulfur, carbon, and energy is far from being a
98 uniform process and new metabolic pathways for the degradation of sulfonated compound
99 are constantly being discovered (Boden et al., 2010; Koch & Dahl, 2018; Sharma et al., 2022;
100 Wolf et al., 2022).

101 These processes are also closely linked to the availability of inorganic sulfur as the
102 released sulfur is either assimilated or excreted as sulfate (Ruff et al., 2003), sulfite (Koch &
103 Dahl, 2018; Li et al., 2022; Sharma et al., 2022), thiosulfate (De Zwart et al., 1997),
104 tetrathionate (Boden et al., 2010) or sulfide (Peck et al., 2019). Indeed, the complete
105 consumption of the volatile sulfonated C₁-compound DMS coupled with the oxidation of the
106 thiosulfate formed as an intermediate, has been reported for a single organism, providing a
107 new link between the organic and inorganic sulfur cycles (Koch & Dahl, 2018). However, the
108 fate of the sulfur released from sulfonated compounds is often not known or assumed to be
109 the same as in dissimilatory sulfur oxidation or reduction. The physiology and interactions of
110 bacterial communities that release sulfur from sulfonated carbon compounds have been
111 sparsely explored and the few existing studies are based on, or assume, sulfur cycling via
112 dissimilatory sulfite reductases (Burrichter et al., 2021; Hanson et al., 2021; Wolf et al., 2022).

113 Ecological studies of organic sulfur compounds are difficult because their metabolism is
114 poorly conserved across bacterial phylogeny and can even vary between strains of the same
115 species. Thus, even within a species, predictions based on taxonomic assignment are not
116 possible (Schäfer et al., 2010). As the functional annotation pipelines of public databases
117 mainly focus on the synthesis of methionine and cysteine, the enzymes involved in the
118 metabolism of organic sulfur compounds are usually not correctly annotated. Inaccurate
119 annotation in public databases is exacerbated by the fact that several proteins involved in the
120 conversion of sulfonated compounds are related to proteins involved in sulfur dissimilation or
121 the turnover of other compounds e.g. the DMSO reductase family (Leimkühler & Iobbi-Nivol,
122 2016) or quinone oxidoreductase complexes (Duarte et al., 2021). For these reasons, the
123 abundance of microbes utilising organic sulfur compounds is likely to be underestimated
124 (Carrion et al., 2019) and the role of sulfonated compounds is understudied (Wolf et al., 2022).
125 Thus, there is a knowledge gap of the link between inorganic and organic sulfur cycling in
126 ecological systems.

127 To fill this gap, we have extended HMS-S-S (Tanabe & Dahl, 2022). This tool was
128 originally developed for rapid detection and annotation of inorganic sulfur dissimilation in
129 prokaryotic genomes. With the substantial extension presented here, it now includes not only
130 inorganic sulfur metabolism enzymes, but also enzymes with characterized or at least strongly
131 indicated function in the metabolism of sulfonated sulfur compounds. These include
132 sulfoquinovose synthesis and degradation pathways, DMSP metabolism, taurine and
133 isethionate conversion, and transport systems for various sulfonated compounds. For all these
134 pathways, we developed individual profiled hidden Markov Models (HMM) and validated
135 score thresholds by cross-validation and with an independent test dataset. HMS-S-S itself has
136 been completely redesigned, improving usability and output formats, and extending the file
137 manipulation tool. By optimising the underlying algorithms, the overall computing speed has
138 been increased by a factor of four. Due to the complete overhaul, we have renamed the tool
139 “HMSS2”. HMSS2 now covers the known metabolism of inorganic and organic sulfur
140 compounds, facilitating the exploration of the microbe-driven natural sulfur cycle.

141 **2 METHODS**

142 **2.1 HMSS2 improvements and workflow**

143 Algorithmic improvements were made on the speed and user-friendliness by process
144 optimization and the implementation of additional features. HMSS2 algorithms are now
145 completely written in Python and precompiled versions are available. In this way, the number
146 of dependencies required to be installed by the user has been greatly reduced to just two
147 external programs. HMMER and Prodigal are still required but installing and configuring of
148 MySQL is no longer necessary. The installation was further simplified by preparation of a pre-
149 compiled executable, that will run directly on a Unix system.

150 HMSS2 includes the basic design of HMS-S-S with further automation. User-supplied
151 input requires a directory containing files in FASTA nucleotide format, consisting of scaffolds
152 or contigs. Alternatively, it is possible to provide amino acid sequences in FASTA files and the
153 corresponding features in GFF3 formatted files. All files in the directory will then be processed
154 in consecutive order. Nucleotide input files are first searched for open-reading frames and
155 translated into protein sequences by Prodigal. This step is omitted if protein sequences are
156 provided. Profile hidden Markov Models (HMM) are then queried against the protein
157 sequences of the current file with validated bit score cutoffs via hmmsearch. Hits are saved in
158 a local database together with corresponding genomic features and protein amino acid
159 sequences. The local database now uses the SQLite database engine and an improved
160 database table structure that allows to save multi-domain proteins with all domains. In the
161 next step, the detected proteins are searched for genetic co-localization. This is done via the
162 genomic features and a maximum nucleotide distance between two genes to be syntenic.
163 Syntenic gene clusters are then compared with a set of predefined and named gene patterns.
164 A new feature of HMSS2 is the detection of co-linear gene clusters. This is a special type of
165 synteny where the genes occur in exactly the same order as the gene pattern. Gene clusters
166 that are similar to the pattern(s) provided are then named by characteristic keywords. NCBI,
167 GTDB taxonomy files or custom files with a similar format can be used to assign taxonomic
168 information. As the taxonomy may change over time, it is recommended that the user updates
169 this information locally as required. Results can be retrieved from the local database filtered
170 by protein domains and/or keywords via HMSS2. The standard output now includes FASTA
171 formatted files and iTol datasets.

172 **2.1 Training dataset generation, annotation and HMM development**

173 Datasets were generated from genomic data downloaded from NCBI RefSeq (Haft et al.,
174 2018) or GenBank (Sayers et al., 2019) as of September 2022. The HMM training dataset
175 contained all assemblies from the NCBI RefSeq database with an assembly level of a complete
176 chromosome. The independent test data consisted of assemblies originating from GenBank,
177 again with an assembly level of the complete chromosome. GenBank covers a greater number
178 of phyla and a wider range of quality and is therefore not entirely similar to the training data
179 from RefSeq. Sequence annotation for Hidden-Markov-model generation was performed
180 using the training dataset and list of reference proteins for organic sulfur metabolism (Table
181 S1). Methods for annotating the training and independent test datasets and for HMM
182 generation were used as described previously (Tanabe & Dahl, 2022).

183 **2.4 Performance metric calculation**

184 Performance was determined using balanced accuracy (Brodersen et al., 2010), F1-score
185 (Forman & Scholz, 2010), and the Matthew-correlation-coefficient (MCC) (Chicco & Jurman,
186 2020). The metric values were additionally corrected for the dataset's skewness (Jeni et al.,
187 2013) (Table S2). Values for each Hidden Markov Model were calculated from a confusion
188 matrix obtained by comparing the annotation of the training/test dataset and annotation
189 assigned by the HMMs. Matching assignments were considered as true positives (TP), while
190 mismatching assignments were considered as false positives (FP), if the HMM recognised a
191 sequence unrelated to the HMM training sequences. All sequences that were not recognized
192 by the HMM but matched the annotation were counted as false negative (FN), and all other
193 sequences were recorded as true negatives (TN).

194 **2.5 Thresholding and cross-validation**

195 Thresholding and cross-validation were executed as previously described (Tanabe &
196 Dahl, 2022). For each HMM, bit scores for noise cutoff, trusted cutoff, and an optimized
197 threshold were determined prior to cross-validation. The noise cutoff corresponded to the
198 score of the lowest scoring TP hit. The trusted cutoff corresponded to the score of the highest
199 scoring FP hit. The optimized cutoff was computed during a nested cross-validation procedure
200 with a 10-fold outer loop and a 5-fold inner loop (Varma & Simon, 2006). The optimized cutoff

201 corresponded to the median of the thresholds with the highest F1 scores across all inner folds.
202 Outer folds were analyzed after all thresholds were set.

203 Each cross-validation fold was generated from the HMM training data. Sequences were
204 randomly sorted into the 10 outer folds of equal size, followed by the equal deviation of each
205 outer fold into 5 inner folds. A cross-validation procedure was then performed on all folds.
206 The inner folds were used to determine the optimized thresholds. The overall performance of
207 each HMM was then done with a confusion matrix created for the outer folds using the
208 optimized thresholds as a cutoff. Balanced accuracy was calculated as the average of all
209 accuracies from each fold. F1 score and MCC were calculated as the sum of the confusion
210 matrices from all folds (Forman & Scholz, 2010). The same procedure without fold generation
211 was performed for the independent test dataset (Chicco, 2017).

212 **2.6 Performance testing**

213 The performance of HMSS2 was compared with that of HMS-S-S version 1 (Tanabe &
214 Dahl, 2022). The HMM library included all 164 HMMs of the original library, detecting
215 dissimilatory sulfur metabolism. A quadratic increasing number of randomly selected
216 genomes ranging from 2 to 64 were chosen from the training dataset described for version 1
217 and used as input for the performance comparison. The input data were in FASTA nucleotide
218 format. Each run was repeated three times with newly randomised input data to reduce
219 performance bias caused by the input data. Both program versions were benchmarked for the
220 execution time required for the workflow from data entry to the final annotated hits with
221 appropriately named gene clusters, but without taxonomy assignment. Time was measured
222 as the required wall-clock runtime when running HMS-S-S or HMSS2 with four parallel threads
223 on an Intel Core i7-6700 CPU.

224 **3 RESULTS**

225 Here, we created a comprehensive database of reliable hidden Markov models (HMMs)
226 based on archaeal and bacterial proteins associated with organic sulfur metabolism. The same
227 approach has already been used for the enzymes of dissimilatory metabolism of inorganic sul-
228 fur compounds (Tanabe & Dahl, 2022). Not only sequence similarity, but also integrated
229 synteny was considered to assign a protein to a specific functional group. The HMMs created

230 here focus on the most abundant organic sulfur compounds in terrestrial and aquatic environ-
231 ments. The compounds covered here include dimethylsulfoniopropionate (DMSP), dimethyl
232 sulfide (DMS), dimethyl sulfoxide (DMSO), dimethyl sulfone (DMSO₂) (Fig. 1), 2,3-dihydroxy-
233 propane-1-sulfonate (DHPS), isethionate, taurine, and membrane sulfolipids (Fig. 2). The
234 HMMs for the enzymes of the metabolic pathways for degradation of individual compounds
235 are described in full below. Normally, prokaryotes do not code for the entire degradation
236 pathways, but only for parts of them.

237 **3.1 HMM Development: DMSP degradation**

238 DMSP is primarily produced by single-celled phytoplankton and algal seaweeds, where
239 it acts as an osmolyte and anti-stress molecule (Kiene et al., 2000). Degradation of DMSP
240 either requires a demethylation pathway or a DMSP lyase (Fig. 1). The demethylation pathway
241 is encoded by the *dmdABCD* gene cluster and starts with the demethylation of DMSP via DmdA
242 to form methylmercaptopropionate. This intermediate is further catabolized by DmdB, DmdC
243 and finally DmdD with the release of acetaldehyde and methanethiol (Bullock et al., 2014;
244 Reisch et al., 2011). For each of the enzymes, one HMM was generated, making four in total.
245 Several non-orthologous DMSP lyases, DddL, DddP, DddQ, DddW and DddY, have been
246 characterised which convert DMSP to acrylate with the release of DMS and acrylate. The latter
247 is then converted to 3-hydroxypropionate by AcuNK (Curson et al., 2011) or to propionyl-CoA
248 by AcuI (Todd et al., 2012). DMSP lyase DddD catalyzes formation of propionyl-CoA and DMS
249 from DMSP in a single reaction without the formation of an acrylate intermediate.
250 3-hydroxypropionate can be further converted to acetyl-CoA via DddA and DddC (Curson et
251 al., 2011). HMMs were generated for AcuI, AcuN, AcuK, DddA, and all DMSP lyases. As there
252 were less than ten sequences identified for DddQ, DddW and DddC, HMMs could not be
253 constructed for these three enzymes.

254 **3.2 HMM development: Assimilation of methanethiol and DMS**

255 DMS and methanethiol are C₁-organosulfur compounds derived mainly from the
256 degradation of DMSP. Both can be assimilated by bacteria as a source of sulfur and carbon,
257 where methanethiol is first converted to DMS, followed by oxidation and assimilation (Fig. 1).
258 The conversion of methanethiol to DMS is catalyzed by methanethiol S-methyltransferase,
259 MddA. This membrane-bound enzyme transfers a single sulfur atom from S-

260 adenosylmethionine to methanethiol (Carrion et al., 2015). The resulting DMS can be further
261 oxidized by either dimethylsulfide cytochrome *c* reductase, DdhABCD, also known as
262 dimethylsulfide dehydrogenase (McDevitt, Hanson, et al., 2002), or by multicomponent DMS
263 monooxygenase DsoABCDEF (Horinouchi et al., 1999). The periplasmic DdhABC
264 dimethylsulfide dehydrogenase couples the oxidation of DMS to the reduction of two *c*-type
265 cytochromes, producing DMSO as the final product. DdhD is a cytoplasmic protein that is not
266 part of the DMS dehydrogenase but has a proposed function in the assembly of the DdhAB
267 complex and its secretion via the Tat pathway (McDevitt, Hugenholtz, et al., 2002). For DdhA
268 and DdhB, it was possible to generate individual HMMs, while this was not the case for DdhC
269 and DdhD which had less than ten validly annotated sequences in the training dataset. The
270 multicomponent DMS monooxygenase DsoABCDEF oxidizes DMS in a two-step reaction to
271 DMSO₂ with DMSO as intermediate. As the sulfur moiety is specifically oxidised, this enzyme
272 is also referred to in the literature as assimilatory DMS S-monooxygenase (Boden & Hutt,
273 2019). A total of six HMMs were generated for this complex. After the oxidation of DMS to
274 DMSO₂, the next step in sulfur assimilation is the oxygen-dependent conversion of DMSO₂ to
275 methanesulfinate, catalyzed by FMN-dependent DMSO₂ monooxygenase SnfG (Wicht, 2016).
276 SnfG was represented by a single HMM. Methanesulfinate is chemically oxidized to
277 methanesulfonate, which is further oxidized to sulfite and formaldehyde by the assimilatory
278 methanesulfonate monooxygenase MsuDE in a NADH- and oxygen-dependent reaction. For
279 MsuDE, a HMM was trained for each subunit.

280 **3.3 HMM development: Dissimilation of DMSO₂**

281 Dimethylsulfone is mainly derived from oxidation of DMS. The degradation of dimethyl
282 sulfone (DMSO₂) begins with its reduction to dimethyl sulfoxide (DMSO) by a DMSO₂
283 reductase in an NADH-dependent reaction (Fig. 1). Although the activity has been measured
284 in crude extracts of some methylotrophic Actinobacteria and Alphaproteobacteria (Borodina et
285 al., 2000; Borodina et al., 2002), the enzyme has not been characterized. DMSO is then further
286 reduced to dimethylsulfide (DMS). Two types of DMSO reductases have so far been
287 characterized (Boden & Hutt, 2019). The first, membrane-bound enzyme is composed of the
288 three subunits, DmsABC, and uses electrons from the quinol pool for DMSO reduction (Bilous
289 & Weiner, 1985). For this enzyme one HMM for each subunit was trained. The second DMSO
290 reductase uses NADH for this purpose and probably consists of only one subunit with high

291 similarity to DmsA, indicated by its cross-reaction with DmsA antibodies. A separate HMM
292 could not be trained for this enzyme, because it is only known by its activity in crude extracts
293 (Borodina et al., 2002). In addition to the Dms-type DMSO reductases, a soluble periplasmic
294 DMSO reductase, DorCAD, has been characterized (A. G. McEwan et al., 1998). The
295 corresponding genes are regulated by DorS and DorR (Kappler & Schäfer, 2014). For each of
296 these five proteins/subunits, we constructed one HMM. The DMS, which is released by DMSO
297 reductase of both types, is oxidized to methanethiol (CH₃SH) and formaldehyde by a DMS
298 monooxygenase, DmoAB, in another NADH-consuming reaction (Boden et al., 2011). As only
299 *dmoA* has been validly identified so far, we trained a HMM specifically for DmoA, but not for
300 DmoB. Further oxidation of methanethiol by a methanethiol oxidase MtoX leads to the final
301 release of sulfide and another molecule of formaldehyde (Eyice et al., 2017). A single HMM
302 was trained for MtoX.

303 **3.4 HMM development: Dissimilation of methanesulfonate**

304 Methanesulfonate is formed by spontaneous chemical oxidation of DMS in the
305 atmosphere (Fig. 1). It is used by diverse aerobic bacteria as a sulfur source and by some
306 specialized methylotrophic prokaryotes as a source of carbon and energy (Kelly & Murrell,
307 1999). The dissimilatory methanesulfonate monooxygenase catalyzes the conversion of
308 methanesulfonate to formaldehyde and sulfite (Henriques & De Marco, 2015). This enzyme is
309 encoded by the *msmABCD* operon, which is often located adjacent to the *msmEFGH* operon,
310 usually in the opposite direction. The latter encodes a putative ABC-type transporter (Fig. 2b)
311 proposed to facilitate the import of methanesulfonate into to the cytoplasm (Henriques & De
312 Marco, 2015). Six HMMs were developed to represent each of these proteins. MsmC and
313 MsmD had to be excluded due to the small number of sequences in the training datasets.

314 **3.5 HMM development: Alkanesulfonate oxidation and transporters**

315 The *ssuEADCB* gene cluster encodes the two-component alkanesulfonate
316 monooxygenase SsuDE and the alkanesulfonate ABC-transporter SsuABC (Fig. 2b).
317 Alkanesulfonate monooxygenase catalyzes the oxidation of various sulfonated alkanes as
318 substrates with variable affinity, including phenylated organic compounds like N-
319 phenyltaurine. After transport into the cell via SsuABC, the sulfonate is cleaved by SsuDE in a
320 reaction dependent on NADH and molecular oxygen (Eichhorn et al., 1999). Electrons are

321 provided by SsuE via an FMN cofactor. SsuD then cleaves the sulfonate group and oxidizes the
322 terminal carbon atom. For this pathway five HMMs, one for each encoded protein, were
323 created.

324 **3.6 HMM development: Sulfoquinovose synthesis**

325 Sulfoquinovose (SQ) is a sulfonated derivate of glucose where the 6-hydroxyl group is
326 substituted by a sulfonate group. SQ is a constituent of the unique head group of the
327 membrane-bound glycolipid sulfoquinovosyl diacylglycerol (SQDG) present in thylakoid
328 membranes and photosynthetic prokaryotes. On a genetic level, five genes *sqdA*, *sqdB*, *sqdC*,
329 *sqdD* and *sqdX* have been described to be involved in SQDG synthesis in bacteria so far
330 (Benning & Somerville, 1992a, 1992b; Guler et al., 2000; Rossak et al., 1995). The functions of
331 SqdA and SqdC have not been completely resolved (Benning & Somerville, 1992b; Rossak et
332 al., 1997). The synthesis begins with the exchange of the 6-hydroxyl group of uridine-
333 diphosphate (UDP)-glucose for a sulfonate group by UDP-sulfoquinovose synthase, SqdB. The
334 formation of SQDG is then catalyzed SQDG synthase, SqdD or SqdX (Rossak et al., 1995). A
335 total of five HMMs was trained to detect the enzymes of this pathway.

336 **3.7 HMM development: Sulfoquinovose degradation and transport**

337 As sulfoquinovose is a sulfonated derivate of glucose, it is catabolized in a similar manner
338 and can serve as a carbon and energy source (Hanson et al., 2021). Several pathways
339 resembling glucose degradation have been characterized, including the Sulfo-Embden-
340 Meyerhof-Parnas pathway (Denger et al., 2014), the Sulfo-Entner–Doudoroff pathway (Felux
341 et al., 2015), the transaldolase-based pathway related to the pentose phosphate pathway
342 (Frommeyer et al., 2020) and a complete degradation pathway based on a sulfoquinovose
343 monooxygenase (Sharma et al., 2022) (Fig. 2a).

344 The Sulfo-Embden-Meyerhof-Parnas pathway (Fig. 2a) begins with import of
345 sulfoquinovose by the transporter YihO. A sulfolipid α -glucosidase YihQ may also be involved
346 and other SQ derivatives may also be imported. Analogous to the EMP pathway, SQ is then
347 cleaved to dihydroxyacetonephosphate (DHAP) and 3-sulfolactaldehyde (SLA) via the
348 isomerase YihS, kinase YihV and aldolase YihT. In an NADH-dependent reaction, the reductase
349 YihU then reduces SLA to the final product 2,3-dihydroxypropane sulfonate (DHPS), which is

350 transported out of the cell again via YihP. A separate HMM was created for each of the Yih
351 proteins.

352 The Sulfo-Entner–Doudoroff is analogous to the ED pathway (Fig. 2a). As there was no
353 specific abbreviated name assigned to these enzymes by the original publication (Felux et al.,
354 2015), we assigned names to enhance HMSS2 output readability. SQ is cleaved by a
355 dehydrogenase SedA, a lactonase SedB, a dehydratase SedC and an aldolase SedD to pyruvate
356 and SLA. Another dehydrogenase, SedE, then oxidizes 3-sulfolactaldehyde (SLA) in an NAD-
357 dependent reaction to 3-sulfolactate (SL), which is then exported. A separate HMM was
358 generated for each of the proteins mentioned, for a total of 5 HMMs.

359 The third SQ degradation pathway contains a transaldolase as the key enzyme (Fig. 2a)
360 (Frommeyer et al., 2020). SQ is imported into this pathway via the transporter SftA and
361 converted to sulfofructose by the isomerase SftI. This product, together with glycerine-
362 aldehyde-3-phosphate, is then converted by the transaldolase SftT to SLA and fructose-6-
363 phosphate. SLA, in turn, is converted to SL in an NAD-dependent reaction by the
364 dehydrogenase SftD and exported via the transporter SftE or reduced to 2,3-
365 dihydroxypropane sulfonate (DHPS) in an NADH-dependent reaction by the reductase SftR. A
366 separate HMM was generated for each of the Sft proteins, for a total of 6 HMMs.

367 The fourth known degradation pathway for SQ (Fig. 2a) differs from the others described
368 so far, because it involves oxidation of the entire molecule, including cleavage of sulfur
369 (Sharma et al., 2022). The pathway described begins with the import of
370 sulfoquinovosylglycerol by an ABC transporter called SmoEFGH. In the cytoplasm,
371 sulfoquinovosyl glycerol is cleaved by the sulfoquinovosidase SmoI to SQ. In contrast to the
372 other pathways, SQ is now transformed to 6-oxo-glucose and sulfite by an alkanesulfonate
373 monooxygenase, SmoC. The electrons for this reaction come from NADPH via the flavin
374 reductase SmoA. 6-oxo-glucose is converted in another NADPH-dependent reaction by SmoB
375 into glucose, which is then available for glycolysis. Eight HMMs were generated for this
376 pathway, one for each protein. An additional HMM was trained for SmoD ,a putative regulator
377 encoded in the *smo* operon.

378 **3.8 HMM development: 2,3-dihydroxypropane sulfonate transporters and** 379 **degradation**

380 According to the postulated pathway for degradation of 2,3-dihydroxypropane sulfonate
381 (DHPS) (Fig. 2c), the compound is either taken up by the TRAP transporter HpsKLM or by HpsU
382 (Fig. 2b). The DHPS-3-dehydrogenase HpsN then converts (R)-DHPS to sulfolactate with
383 concomitant formation of two equivalents of NADH. For (S)-DHPS, it was postulated that this
384 compound is first converted to the (R)-DHPS enantiomer via (R)-DHPS-2-dehydrogenase HpsP
385 and (S)-DHPS-2-dehydrogenase HpsO (Mayer et al., 2010). The resulting (R)-sulfolactate can
386 be further converted in several ways: The (R)-sulfolactate sulfolyase SuyAB catalyzes a
387 desulfonation reaction, releasing sulfite and pyruvate. The (S)-enantiomer of sulfolactate is
388 first converted to sulfoxyacetate by SlcC and then to (R)-sulfolactate by ComC (Mayer et al.,
389 2010). Both enantiomers were postulated to be transported by the exporter SlcHFG (Mayer
390 et al., 2010) (Fig. 2b). On HMM was created for each protein/subunit of the DHPS degradation
391 pathway.

392 **3.9 HMM development: Isethionate and taurine degradation**

393 Isethionate and taurine are C₂-sulfonates which are produced by eukaryotes from
394 cysteine or methionine (Moran & Durham, 2019). Bacterial degradation of these compounds
395 includes sulfoacetaldehyde as an intermediate which is a point of convergence with
396 sulfoacetate degradation (Weinitschke, Hollemeyer, et al., 2010) (Fig. 2c). Two different
397 transporters are proposed for the import of isethionate (Fig. 2b). These are the TRAP
398 transporters IseKLM and IseU from the major facilitator superfamily. After import into the
399 cytoplasm, isethionate is oxidized to sulfoacetaldehyde by the isethionate dehydrogenase IseJ
400 (Weinitschke, Sharma, et al., 2010). In some organisms, isethionate is not converted, but the
401 sulfonate group is cleaved off by isethionate sulfite lyase IslAB, releasing sulfite and
402 acetaldehyde (Peck et al., 2019).

403 Taurine import is postulated to be facilitated by the ABC transporter TauAB1B2C or the
404 TRAP transporter TauKLM (Fig.2b). There are several possibilities for the further pathway.
405 Taurine can either be oxygenated by TauD to form 1-hydroxy-2-aminoethane sulfonic acid,
406 which decomposes to aminoacetaldehyde and sulfite (Eichhorn et al., 1999), or it is oxidized
407 in NADH-dependent reaction by the taurine dehydrogenase TauXY, which produces

408 sulfoacetaldehyde. The same product is also produced by the transfer of the amino group to
409 pyruvate by taurine:pyruvate aminotransferase Tpa (Bruggemann et al., 2004)) or to
410 2-oxoglutarate by taurine:2-oxoglutarate aminotransferase Toa (Krejci et al., 2010).

411 Sulfoacetaldehyde can be converted by the NADPH-dependent sulfoacetaldehyde
412 reductase IsfD to isethionate which is then exported by the IsfE transporter (Krejci et al.,
413 2010). Another possible fate of sulfoacetylaldehyde is desulfonation coupled to a
414 phosphorylation by sulfoacetaldehyde acetyltransferase Xsc to acetyl phosphate which is
415 further converted to acetyl-CoA by phosphate acetyltransferase Pta (Weinitschke, Sharma, et
416 al., 2010) Sulfite released in the each of these processes is exported via TauE (Weinitschke et
417 al., 2007). An individual HMM was developed for each individual protein/subunit mentioned
418 here. An exception was made for TauB1 and TauB2, which were combined into a single HMM
419 due to their similarity. Additionally, we trained an HMM for TauZ, a protein of unknown
420 function, and the regulator TauR. Both are commonly found genetically associated with other
421 *tau* genes.

422 **3.10 HMM development: Sulfoacetaldehyde formation**

423 Sulfoacetaldehyde is not only produced by taurine and isethionate degradation but also
424 by the dissimilation of sulfoacetate (Weinitschke, Hollemeyer, et al., 2010). The transporter
425 SauU is hypothesised to facilitate the entry of sulfoacetate into the cell (Fig. 2b). Subsequently,
426 sulfoacetate is activated by sulfoacetate-CoA ligase, SauT, and finally reduced to
427 sulfoacetaldehyde via sulfoacetaldehyde dehydrogenase, SauS, consuming NADPH. SauS,
428 SauT and SauU (Weinitschke, Hollemeyer, et al., 2010) were each represented by a HMM
429 respectively. Sulfoacetaldehyde can also be produced by decarboxylation of sulfopyruvate
430 (Fig. 2c) catalyzed by ComDE (Denger et al., 2009). These two subunits are each represented
431 by a HMM.

432 **3.11 HMM development: Cysteine synthesis**

433 Cysteine is an essential amino acid with a thiol side chain. Here, we started to cover the
434 relevant proteins with HMMs primarily based on knowledge collected with enterobacterial
435 model organisms. Biosynthesis begins with the import of sulfate or thiosulfate into the
436 bacterial cell via CysUWA (Aguilar-Barajas et al., 2011) or YeeE/YedE-like (Tanaka et al., 2020)
437 transporters. Sulfate is reduced to sulfide which is then incorporated into O-acetylserine to

438 synthesize cysteine (Kredich, 1996). In *E. coli*, sulfate is activated by ATP sulfurylase CysDN
439 (Leyh et al., 1988) to adenosine 5'-phosphosulfate (APS), which can be further activated by
440 APS kinase CysC to 3'-phosphoadenosine-5'-phosphosulfate (PAPS). PAPS reductase CysH then
441 reduces the activated compound to sulfite. In some bacteria, including most cyanobacteria,
442 APS can be reduced to sulfite directly, without phosphorylation to PAPS (Bick et al., 2000). The
443 assimilatory APS reductases catalyzing this reaction exhibit similarity to the assimilatory PAPS
444 reductases (Abola et al., 1999; Bick et al., 2000) and are covered by the same HMM (CysH) in
445 this work. In Enterobacteria, sulfite is reduced to sulfide via CysIJ. Finally, cysteine is
446 synthesized from sulfide and O-acetyl-L-serine by the cysteine synthase CysK. A total of 10
447 new HMMs was generated for the mentioned proteins/subunits. An HMM for YeeE/YedE-like
448 transporters was already available through HMS-S-S (Tanabe & Dahl, 2022)

449 **3.12 HMM validation: cross validation and independent test data set**

450 The HMMs developed were validated by cross-validation and with an independent test
451 data set. In cross-validation, sequences unrelated to the tested HMM training data were
452 added as true negative examples in addition to the omitted training sequences (Chicco, 2017;
453 Refaeilzadeh et al., 2009). The omitted sequences from each fold served as true positive
454 examples. Cross-validation was performed using the optimized thresholds calculated prior to
455 cross-validation. Thus, the threshold values should also be checked for their suitability.
456 Performance was measured using the Matthews Correlation Coefficient (MCC). This metric
457 ranges from -1 to 1, with 0 corresponding to random assignment, 1 corresponding to perfect
458 assignment with no misclassification, and -1 corresponding to complete misclassification.
459 Here, the individual occurrence of FP or FN lowers the score on the MCC, while the
460 combination of both misclassifications lowers the score more dramatically than the single
461 occurrence of either type of error (Chicco & Jurman, 2020).

462 The majority of the HMMs developed showed high precision and recall in the cross-
463 validation and on the test dataset (Fig. 3). Of the 134 HMMs covering proteins of organic sulfur
464 compound metabolism, 127 stayed above an MCC of 0.80 during the cross-validation (Fig. 3,
465 Table S2). The evaluation of the 134 HMMs against the independent test dataset resulted in
466 120 HMMs with an MCC of 0.80 or higher. HMMs for the alkanesulfonate transporter subunits
467 SsuB and SsuC failed the cross-validation threshold of 0.8 slightly by 0.02 points but performed
468 better on the independent test dataset. These were the only cases where the cross-validation

469 performance was insufficient but the performance on the test dataset was above the
470 threshold. From the HMMs with an MCC > 0.8 during cross-validation, seven scored below 0.8
471 in the test dataset. These were MsmG with an MCC of 0.78, SmoI (0.76), MsmB (0.66), DddA
472 (0.62), DorA (0.46) and SftD (0.03). For SftD, MsmB, MsmG and DddA this was due to a high
473 number of sequences which were falsely classified as negative, probably due to a low training
474 sequence diversity. Thus, these HMMs had a high precision and did not generate high numbers
475 of false positive hits, but they performed low in recognition resulting in a high number of
476 unrecognized sequences. The opposite was the case for the DorA HMM, which generated too
477 many false positive hits but no false negative ones. Sulfoquinovosidase SmoI interfered in the
478 detection with sulfoquinovosidase named YihQ. The same holds true for transporters HpsU
479 and IseU. All sequences that were falsely classified by one of these two HMMs belonged to
480 the other HMM. Together these two HMMs performed well in detecting of isethionate and
481 DHPS transporters of the major facilitator superfamily. The situation was similar for YihO and
482 SftA which are both postulated sulfoquinovose importers that catalyse the same function in
483 the context of sulfoquinovose degradation. In summary, 112 of 134 HMMs were successfully
484 tested via cross-validation and with an independent dataset. Two other pairs of HMMs can be
485 used together, for the safe detection of sulfoquinovosidase and the transporters YihO and
486 SftA.

487 **3.13 HMM validation: Case study**

488 HMSS2 was also validated with 24 complete genomes from bacteria with organic sulfur
489 compound metabolism (Table S3), which were screened for the presence of enzymes for the
490 utilisation of taurine, isethionate, DHPS, sulfoquinovose and DMS (Fig. 4).

491 Proteins for taurine utilization were found mainly in the known taurine-utilizing genera
492 *Octadecabacter*, *Roseobacter*, *Roseovarius* and *Ruegeria* of the Rosebacterales, including the
493 taurine degraders *Roseovarius nubinhibens* (Denger et al., 2009) and *Ruegeria pomeroyi*
494 (Gorzynska et al., 2006). These strains encoded for the TauABC taurine importer, Tpa and Xsc
495 constituting the complete degradation pathway from free taurine via sulfoacetaldehyde to
496 acetyl phosphate with the release of sulfite. *Roseobacter denitrificans* additionally possessed
497 genes for the taurine dehydrogenase TauXY and the taurine:2-oxoglutarate aminotransferase
498 Toa, which can also convert taurine to sulfoacetaldehyde. The sulfoacetaldehyde
499 acetyltransferase Xsc was present in all genomes examined. This is probably due to the fact

500 that sulfoacetaldehyde is not exclusively an intermediate of taurine degradation but also of
501 isethionate, sulfoacetate and DHPS degradation, and possibly of other as yet unknown
502 pathways (Weinitschke, Hollemeyer, et al., 2010). In line with this possibility, genes encoding
503 isethionate dehydrogenase IseJ, which converts isethionate to sulfoacetaldehyde, were found
504 in almost all analyzed Rhodobacterales, Hyphomicrobiales and Gammaproteobacteria
505 genomes, consistent with earlier reports (Weinitschke, Sharma, et al., 2010). *Leminorella*
506 *grimontii*, *Hyphomicrobium denitrificans* and all *Methylophaga* species were exceptions,
507 consistent with the inability of *H. denitrificans* and *Methylophaga* to consume organosulfur
508 compounds with more than one carbon atom.

509 Isethionate desulfonation via isethionate sulfite-lyase IslAB has been found in
510 microcompartments of *Bilophila wadsworthia* (Burrichter et al., 2021). In accordance, HMSS2
511 detected the importer IseU and IslAB in this organism. A similar desulfonation pathway
512 without microcompartments was postulated for *Desulfovibrio alaskensis* and *D. desulfuricans*
513 (Burrichter et al., 2021). In *D. desulfuricans*, HMSS2 also found IseU and IslAB, suggesting that
514 this organisms, like *B. wadsworthia*, may scavenge free isethionate via IseU. In contrast, *D.*
515 *alaskensis* encodes IslAB but not IseU. Instead, it contains sulfocacetaldehyde reductase IsfD
516 (or SarD), which is also present in *Bilophila wadsworthia*. In both cases, this enzyme may
517 provide an endogenous source of isethionate (Burrichter et al., 2021).

518 Most analysed genomes possessed the potential for sulfopyruvate and (R)-sulfolactate
519 generation from DHPS and (L)-sulfolactate. The potential of (R)-DHPS oxidation via HpsN
520 generating 2 NADH equivalents was found in all analysed strains and most Iso encoded for
521 isomerization of (S)-DHPS to (R)-DHPS via HpsP (17/24 genomes). The predicted presence of
522 genes for desulfonation of sulfopyruvate by ComDE and sulfolactate by SuyAB as found here
523 is also in accordance with previous reports for the Roseobacterales clade (Chen et al., 2021;
524 Denger et al., 2009), the Hyphomicrobiales (Chen et al., 2021), *Desulfovibrio desulfuricans* and
525 *B. wadsworthia* (Hanson et al., 2021). Even without the ability to desulfonate sulfopyruvate
526 or sulfolactate, the conversion of DHPS to sulfopyruvate or sulfolactate and export of these as
527 end products provides 2-3 NADH equivalents and thus a growth advantage for the organism.

528 Sulfoquinovose degradation via the Sulfo-Entner-Doudoroff pathway is present in eight
529 bacteria, including *Pseudomonas putida* and other bacteria for which this pathway has been
530 described or postulated (Felux et al., 2015). The complete sulfoquinovose degradation

531 pathway based on a sulfoquinovose monooxygenase was found in seven proteobacteria in
532 accordance with previous reports (Sharma et al., 2022). The other known sulfoquinovose
533 degradation pathways were not detected, which is likely due to the presence of the Sulfo-
534 Embden-Meyerhof-Parnas pathway (Denger et al., 2014) primarily in Enterobacterales and the
535 transaldolase-dependent sulfoquinovose degradation in Firmicutes (Frommeyer et al., 2020).
536 Bacteria from these taxonomic groups were not included in the case study.

537 DMS degradation has been described for *Methylophaga thiooxydans*, *Methylophaga*
538 *sulfidovorans* (Kröber & Schäfer, 2019), *Hyphomicrobium denitrificans* (Koch & Dahl, 2018),
539 and *Hyphomicrobium sulfonivorans* (Boden et al., 2011). According to our HMSS2 analysis, *H.*
540 *sulfonivorans* encoded for DmoA, while all other three encoded only for methanethiol oxidase
541 MtoX. DmoA was missing and the organisms must contain a so far unknown DMS
542 monooxygenase. In accordance with previous reports, MtoX was also found in
543 *Methylacidiphilum fumariolicum* (Schmitz et al., 2022), and several Rosebacterales, including
544 *Ruegeria pomeroyi* (Eyice et al., 2017). The latter is a known degrader of DMSP to
545 methanethiol via DmdA, B, C and DmdD (Reisch et al., 2011) which were all detected by the
546 HMMs created here.

547 In summary, our case study on characterized organosulfur compound degraders has
548 shown that in all cases the detection by HMSS2 agrees with the published analyses of other
549 authors.

550 **3.14 HMSS2 improvements**

551 HMSS2 has a redesigned engine and additional features for protein annotation and
552 output format customisation (Fig. 5). Proteins with multiple domains are now stored with all
553 domains and not just the domain with the highest score. This was accomplished by improving
554 the local relational database structure. This requires that the recognised domain regions in
555 the primary sequence do not overlap, so that domains with high scores are not overwritten
556 by lower scores. On the other hand, high-scoring domains may still overwrite one or more
557 lower-scoring domains during annotation.

558 Gene arrangement can now be used by HMSS2 for annotation as a non-homologous
559 criterion. Hits below the threshold are also considered and annotated if they lie within a gene
560 cluster and the potentially assigned annotation would complete a known gene cluster

561 arrangement. Thus, a gene that highly likely occurs within a gene cluster must reach a lower
562 cutoff than normal to be detected if it is encoded within such a cluster.

563 The output formats have been greatly expanded, and new features were added to
564 improve usability and readability. It is still possible to retrieve sequences filtered by protein
565 type, the genomic proximity and the presence of proteins or gene clusters in the same
566 genome. HMSS2 automatically recovers a list of all hits with genomic features and a separate
567 protein sequence file in FASTA format. Additionally, two subsets of the latter file are created.
568 One subset includes all hits that are unique to their genome respectively, while another subset
569 includes all hits that occur at least twice in the same genome. Multi-domains proteins,
570 retrieved by the requested protein type, are listed separately if at least one other domain has
571 been annotated.

572 An output module for iTol compatible datasets was also included. This module
573 integrates the generation of iTol datasets for presence/absence of the keywords/domains for
574 each genome. Range datasets, which mark specific proteins in a phylogenetic tree, can now
575 also be generated by HMSS2, as well as iTol compatible datasets for displaying gene clusters.
576 HMSS2 also comes with several utilities to modify the output protein FASTA files. It is now
577 possible to assign the taxonomic names of the source organism to each sequence. Files can
578 now be filtered by length, merged without duplicating sequence identifiers and sequences
579 from multiple FASTA files originating from the same organism can be concatenated into a
580 single sequence. With a FASTA-formatted file as input, a list of neighboring genes is now
581 accessible to support searches for conserved but previously undiscovered gene constellations.

582 The execution time of the HMSS2 was compared to that of HMS-S-S to demonstrate the
583 scalability and efficiency of HMSS2. For this test, increasing numbers of genomes were
584 randomly selected from the assemblies of the training dataset and gene clusters were
585 annotated and determined with the 164 HMMs of the original library. Time measurements
586 were performed in triplicate with random selection of input assemblies for each replicate. The
587 execution time was then averaged over all replicates. Comparison between the two versions
588 showed a large difference in the required execution time (Fig. 6, Table S4). The observed
589 increase in execution speed for HMSS2 became more significant as the number of genomes
590 processed increased and scaled linearly with the number of input assemblies. While HMS-S-S
591 required around 26 minutes to process 64 assemblies, HMSS2 needed only 7 minutes for this

592 task. Thus, the introduced improvements led to a fourfold accelerated computation speed for
593 HMSS2.

594 **4 DISCUSSION**

595 Here, we present a substantial update that provides an HMM-based search tool for
596 proteins involved in the metabolism of inorganic and organic sulfur compounds. The high
597 accuracy of the advanced tool presented here provides a reliable basis for genome analysis
598 and is further supported by the genomic context detection. The HMSS2 algorithm now uses
599 homologous and non-homologous criteria already in the protein annotation step, not just for
600 the later identification of gene clusters. In addition, the overall execution time was accelerated
601 by fourfold compared to the previous version, further speeding up the detection of sulfur
602 metabolism pathways in genomes and metagenomes. With the increasing number of available
603 genomes, faster protein annotation is required to handle the immense amount of available
604 data.

605 We also significantly broadened the applicability of HMSS2 by adding the conversion of
606 sulfonated carbon compounds. HMSS2 now covers pathways from the entire sulfur cycle,
607 enabling studies on the link between the cycles of inorganic and organic sulfur compounds. In
608 addition to providing operon structure information to support equivalence prediction, the
609 accessibility and display of the annotated proteins has been greatly enhanced. Not only can
610 sequences now be filtered by annotation, but also the presence of genes and genomic context
611 can be displayed using other specialised applications, further extending the capabilities of
612 synteny analysis. Such analyses are not limited to studies of the ecological role of prokaryotes
613 but also include the evolution of metabolic pathways (Garcia et al., 2022), distribution of new
614 pathways (Sharma et al., 2022) and genomic context visualization (Garcia et al., 2019; Letunic
615 & Bork, 2021).

616 The expansion to the metabolism of organic sulfur compounds resulted in the
617 generation of 134 additional HMMs in addition to the 164 HMMs previously included in
618 HMS-S-S, almost doubling the total number of proteins included. The accuracy of the newly
619 generated HMMs and the respective thresholds were demonstrated by cross-validation and a
620 test dataset. Observed deviations between both testing methods are likely due to an uneven
621 distribution and abundance of protein sequences influencing the number and diversity of

622 testable sequences. The quality of the 134 novel HMMs was ensured by selection of high-
623 quality genomes derived from the RefSeq and GenBank databases. The overall development
624 process had already been successfully applied for the proteins of inorganic sulfur metabolism
625 (Tanabe & Dahl, 2022). The test dataset was obtained from the full diversity of phyla accessible
626 from GenBank and should therefore reflect the widest possible range of sequence variation.
627 However, although the cutoff values have been validated, they are likely to need adjustment
628 for newly discovered phyla (Anantharaman et al., 2018; Jaffe et al., 2020).

629 The diversity of proteins involved in the metabolism of organic sulfur compounds
630 covered by HMSS2 also includes less prominent pathways for degradation and conversion of
631 compounds such as sulfoquinovose or DMS. Although a considerable proportion of sulfur in
632 the biosphere is bound in substrates or intermediates of these pathways, they are not
633 commonly included in annotation pipelines and often unrecognized or incorrectly annotated.
634 This is illustrated by fact that only 16 of the 124 proteins included here for the conversion of
635 sulfoquinovose, taurine, isethionate or DMSP have an exact counterpart in PFAM (El-Gebali et
636 al., 2019) or TIGRFAMs. In contrast, eight of ten HMMs covering sulfate assimilation for
637 cysteine biosynthesis have a TIGRFAM equivalent. A common problem in the functional
638 annotation of enzymes involved in metabolism of organic sulfur compounds are enzymes,
639 such as DmsA or DorA, that belong to the DMSO reductase superfamily. This family includes
640 tetrathionate reductase, polysulfide reductase and thiosulfate reductase, as well as several
641 other proteins unrelated to sulfur metabolism. Tertiary structure and complex composition is
642 conserved throughout all members of this family (Alastair G. McEwan et al., 2010) and
643 substrate specificity may only arise through a small number of conserved amino acids at the
644 active site (Struwe et al., 2021). The validation performed here showed that related complexes
645 in the DMSO reductase family did not negatively affect the HMMs for DmsA and DorA.
646 Furthermore, the reliability of prediction is raised when genomic context is paired with the
647 prediction made by the HMM detection as already discussed above.

648 **5 CONCLUSIONS**

649 In summary, HMSS2 is an advanced comprehensive HMM-based tool for annotation and
650 synteny analysis of prokaryotic sulfur metabolism. It has a higher speed and a much wider
651 coverage than its predecessor HMS-S-S and now includes proteins involved in the metabolism

652 of inorganic and organic sulfur compounds. The use of curated functionally equivalent
653 sequences for HMM training resulted in HMMs with high precision and recall. This also fills a
654 gap in the coverage of sulfur metabolism prediction by HMMs. The application possibilities
655 also include the combination with other HMMs from public databases or user-defined models
656 and can therefore be extended according to the user's needs. The improved output formats
657 are also applicable to ecology and evolutionary research.

658 **ACKNOWLEDGEMENTS**

659 This work was supported by the Deutsche Forschungsgemeinschaft (grant Da 351/13-1 to CD).
660 TST received a scholarship from the Studienstiftung des Deutschen Volkes.

661 **CONFLICT OF INTEREST**

662 The authors declare that they have no competing interests.

663 **AUTHOR CONTRIBUTIONS**

664 TST and CD conceived the study. TST developed and implemented the method and performed
665 the analyses. TST analysed and interpreted the data. Both authors wrote and approved the
666 final version of the manuscript.

667 **DATA AVAILABILITY STATEMENT**

668 HMSS2 program files are available at <https://github.com/TSTanabe/HMSS2>.

669 **References**

- 670 Abola, A. P., Willits, M. G., Wang, R. C., & Long, S. R. (1999). Reduction of adenosine-5'-phosphosulfate
671 instead of 3'-phosphoadenosine-5'-phosphosulfate in cysteine biosynthesis by *Rhizobium*
672 *meliloti* and other members of the family Rhizobiaceae. *Journal of Bacteriology*, *181*, 5280-
673 5287. <https://doi.org/10.1128/JB.181.17.5280-5287.1999>
- 674 Aguilar-Barajas, E., Diaz-Perez, C., Ramirez-Diaz, M. I., Riveros-Rosas, H., & Cervantes, C. (2011).
675 Bacterial transport of sulfate, molybdate, and related oxyanions. *Biometals*, *24*(4), 687-707.
676 <https://doi.org/10.1007/s10534-011-9421-x>
- 677 Anantharaman, K., Hausmann, B., Jungbluth, S. P., Kantor, R. S., Lavy, A., Warren, L. A., Rappé, M. S.,
678 Pester, M., Loy, A., Thomas, B. C., & Banfield, J. F. (2018). Expanded diversity of microbial
679 groups that shape the dissimilatory sulfur cycle. *ISME Journal*, *12*, 1715-1728.
680 <https://doi.org/10.1038/s41396-018-0078-0>
- 681 Benning, C., & Somerville, C. R. (1992a). Identification of an operon involved in sulfolipid biosynthesis
682 in *Rhodobacter sphaeroides*. *Journal of Bacteriology*, *174*(20), 6479-6487.
683 <https://doi.org/10.1128/jb.174.20.6479-6487.1992>
- 684 Benning, C., & Somerville, C. R. (1992b). Isolation and genetic complementation of a sulfolipid-deficient
685 mutant of *Rhodobacter sphaeroides*. *Journal of Bacteriology*, *174*(7), 2352-2360.
686 <https://doi.org/10.1128/jb.174.7.2352-2360.1992>

687 Bick, J. A., Dennis, J. J., Zylstra, G. J., Nowack, J., & Leustek, T. (2000). Identification of a new class of 5
688 'S-adenylylsulfate (APS) reductases from sulfate-assimilating bacteria. *Journal of Bacteriology*,
689 182, 135-142.

690 Bilous, P. T., & Weiner, J. H. (1985). Dimethyl sulfoxide reductase activity by anaerobically grown
691 *Escherichia coli* HB101. *Journal of Bacteriology*, 162(3), 1151-1155.
692 <https://doi.org/10.1128/jb.162.3.1151-1155.1985>

693 Boden, R., Borodina, E., Wood, A. P., Kelly, D. P., Murrell, J. C., & Schäfer, H. (2011). Purification and
694 characterization of dimethylsulfide monooxygenase from *Hyphomicrobium sulfonivorans*.
695 *Journal of Bacteriology*, 193(5), 1250-1258.

696 Boden, R., & Hutt, L. P. (2019). Bacterial metabolism of C₁ sulfur compounds. In F. Rojo (Ed.), *Aerobic*
697 *utilization of hydrocarbons, oils and lipids. Handbook of hydrocarbon and lipid microbiology*
698 (pp. 1-43). Cham: Springer Nature Switzerland AG.

699 Boden, R., Kelly, D. P., Murrell, J. C., & Schäfer, H. (2010). Oxidation of dimethylsulfide to tetrathionate
700 by *Methylophaga thiooxidans* sp. nov.: a new link in the sulfur cycle. *Environmental*
701 *Microbiology*, 12(10), 2688-2699. <https://doi.org/10.1111/j.1462-2920.2010.02238.x>

702 Borodina, E., Kelly, D. P., Rainey, F. A., Ward-Rainey, N. L., & Wood, A. P. (2000). Dimethylsulfone as a
703 growth substrate for novel methylotrophic species of *Hyphomicrobium* and *Arthrobacter*.
704 *Archives of Microbiology*, 173(5-6), 425-437.

705 Borodina, E., Kelly, D. P., Schumann, P., Rainey, F. A., Ward-Rainey, N. L., & Wood, A. P. (2002). Enzymes
706 of dimethylsulfone metabolism and the phylogenetic characterization of the facultative
707 methylotrophs *Arthrobacter sulfonivorans* sp. nov., *Arthrobacter methylotrophus* sp. nov., and
708 *Hyphomicrobium sulfonivorans* sp. nov. *Archives of Microbiology*, 177(2), 173-183.
709 <https://doi.org/10.1007/s00203-001-0373-3>

710 Brodersen, K. H., Ong, C. S., Stephan, K. E., & Buhmann, J. M. (2010). The balanced accuracy and its
711 posterior distribution. *2010 International Conference on Pattern Recognition*, 3121-3124.
712 <https://doi.org/10.1109/icpr.2010.764>

713 Bruggemann, C., Denger, K., Cook, A. M., & Ruff, J. (2004). Enzymes and genes of taurine and
714 isethionate dissimilation in *Paracoccus denitrificans*. *Microbiology (Reading)*, 150(Pt 4), 805-
715 816. <https://doi.org/10.1099/mic.0.26795-0>

716 Bullock, H. A., Reisch, C. R., Burns, A. S., Moran, M. A., & Whitman, W. B. (2014). Regulatory and
717 functional diversity of methylmercaptopropionate coenzyme A ligases from the
718 dimethylsulfoniopropionate demethylation pathway in *Ruegeria pomeroyi* DSS-3 and other
719 proteobacteria. *Journal of Bacteriology*, 196(6), 1275-1285.
720 <https://doi.org/10.1128/JB.00026-14>

721 Burrichter, A. G., Dorr, S., Bergmann, P., Haiss, S., Keller, A., Fournier, C., Franchini, P., Isono, E., &
722 Schleheck, D. (2021). Bacterial microcompartments for isethionate desulfonation in the
723 taurine-degrading human-gut bacterium *Bilophila wadsworthia*. *BMC Microbiology*, 21(1),
724 340. <https://doi.org/10.1186/s12866-021-02386-w>

725 Carrion, O., Curson, A. R. J., Kumaresan, D., Fu, Y., Lang, A. S., Mercade, E., & Todd, J. D. (2015). A novel
726 pathway producing dimethylsulphide in bacteria is widespread in soil environments. *Nature*
727 *Communications*, 6, 6579. <https://doi.org/10.1038/ncomms7579>

728 Carrion, O., Pratscher, J., Richa, K., Rostant, W. G., Farhan Ul Haque, M., Murrell, J. C., & Todd, J. D.
729 (2019). Methanethiol and dimethylsulfide cycling in Stiffkey saltmarsh. *Frontiers in*
730 *Microbiology*, 10, 1040. <https://doi.org/10.3389/fmicb.2019.01040>

731 Chen, X., Liu, L., Gao, X., Dai, X., Han, Y., Chen, Q., & Tang, K. (2021). Metabolism of chiral sulfonate
732 compound 2,3-dihydroxypropane-1-sulfonate (DHPS) by Roseobacter bacteria in marine
733 environment. *Environment International*, 157, 106829.
734 <https://doi.org/10.1016/j.envint.2021.106829>

735 Chicco, D. (2017). Ten quick tips for machine learning in computational biology. *BioData Mining*, 10,
736 35. <https://doi.org/10.1186/s13040-017-0155-3>

737 Chicco, D., & Jurman, G. (2020). The advantages of the Matthews correlation coefficient (MCC) over
738 F1 score and accuracy in binary classification evaluation. *BMC Genomics*, 21(1), 6.
739 <https://doi.org/10.1186/s12864-019-6413-7>

740 Curson, A. R., Todd, J. D., Sullivan, M. J., & Johnston, A. W. (2011). Catabolism of
741 dimethylsulphoniopropionate: microorganisms, enzymes and genes. *Nature Reviews*
742 *Microbiology*, 9(12), 849-859. <https://doi.org/10.1038/nrmicro2653>

743 De Zwart, J., Sluis, J., & Kuenen, J. G. (1997). Competition for dimethyl sulfide and hydrogen sulfide by
744 *Methylophaga sulfidovorans* and *Thiobacillus thioparus* T5 in continuous cultures. *Applied and*
745 *Environmental Microbiology*, 63(8), 3318-3322. [https://doi.org/10.1128/aem.63.8.3318-](https://doi.org/10.1128/aem.63.8.3318-3322.1997)
746 [3322.1997](https://doi.org/10.1128/aem.63.8.3318-3322.1997)

747 Denger, K., Mayer, J., Buhmann, M., Weinitschke, S., Smits, T. H., & Cook, A. M. (2009). Bifurcated
748 degradative pathway of 3-sulfolactate in *Roseovarius nubinhibens* ISM via sulfoacetaldehyde
749 acetyltransferase and (S)-cysteate sulfolyase. *Journal of Bacteriology*, 191(18), 5648-5656.
750 <https://doi.org/10.1128/JB.00569-09>

751 Denger, K., Weiss, M., Felux, A. K., Schneider, A., Mayer, C., Spitteller, D., Huhn, T., Cook, A. M., &
752 Schleheck, D. (2014). Sulphoglycolysis in *Escherichia coli* K-12 closes a gap in the
753 biogeochemical sulphur cycle. *Nature*, 507(7490), 114-117.
754 <https://doi.org/10.1038/nature12947>

755 Dhoub, R., Nasreen, M., Othman, D., Ellis, D., Lee, S., Essilfie, A. T., Hansbro, P. M., McEwan, A. G., &
756 Kappler, U. (2021). The DmsABC sulfoxide reductase supports virulence in non-typeable
757 *Haemophilus influenzae*. *Frontiers in Microbiology*, 12, 686833.
758 <https://doi.org/10.3389/fmicb.2021.686833>

759 Duarte, A. G., Barbosa, A. C. C., Ferreira, D., Manteigas, G., Domingos, R. M., & Pereira, I. A. C. (2021).
760 Redox loops in anaerobic respiration - The role of the widespread NrfD protein family and
761 associated dimeric redox module. *Biochimica et Biophysica Acta Bioenergetics*, 1862(7),
762 148416. <https://doi.org/10.1016/j.bbabi.2021.148416>

763 Eichhorn, E., van der Ploeg, J. R., & Leisinger, T. (1999). Characterization of a two-component
764 alkanesulfonate monooxygenase from *Escherichia coli*. *Journal of Biological Chemistry*,
765 274(38), 26639-26646. <https://doi.org/10.1074/jbc.274.38.26639>

766 El-Gebali, S., Mistry, J., Bateman, A., Eddy, S. R., Luciani, A., Potter, S. C., Qureshi, M., Richardson, L. J.,
767 Salazar, G. A., Smart, A., Sonnhammer, E. L. L., Hirsh, L., Paladin, L., Piovesan, D., Tosatto, S. C.
768 E., & Finn, R. D. (2019). The Pfam protein families database in 2019. *Nucleic Acids Research*,
769 47(D1), D427-D432. <https://doi.org/10.1093/nar/gky995>

770 Eyice, O., Myronova, N., Pol, A., Carrión, O., Todd, J. D., Smith, T. J., Gurman, S. J., Cuthbertson, A.,
771 Mazard, S., Mennink-Kersten, M. A., Bugg, T. D., Andersson, K. K., Johnston, A. W., Op den
772 Camp, H. J., & Schäfer, H. (2017). Bacterial SBP56 identified as a Cu-dependent methanethiol
773 oxidase widely distributed in the biosphere. *ISME Journal*, 12(1), 145-160.
774 <https://doi.org/10.1038/ismej.2017.148>

775 Felux, A. K., Spitteller, D., Klebensberger, J., & Schleheck, D. (2015). Entner-Doudoroff pathway for
776 sulfoquinovose degradation in *Pseudomonas putida* SQ1. *Proceedings of the National*
777 *Academy of Sciences of the United States of America*, 112(31), E4298-4305.
778 <https://doi.org/10.1073/pnas.1507049112>

779 Forman, G., & Scholz, M. (2010). Apples-to-apples in cross-validation studies. *ACM SIGKDD*
780 *Explorations Newsletter*, 12(1), 49-57. <https://doi.org/10.1145/1882471.1882479>

781 Frommeyer, B., Fiedler, A. W., Oehler, S. R., Hanson, B. T., Loy, A., Franchini, P., Spitteller, D., &
782 Schleheck, D. (2020). Environmental and intestinal phylum Firmicutes bacteria metabolize the
783 plant sugar sulfoquinovose via a 6-deoxy-6-sulfofructose transaldolase pathway. *iScience*,
784 23(9), 101510. <https://doi.org/10.1016/j.isci.2020.101510>

785 Garcia, P. S., D'Angelo, F., Ollagnier de Choudens, S., Dussouchaud, M., Bouveret, E., Gribaldo, S., &
786 Barras, F. (2022). An early origin of iron-sulfur cluster biosynthesis machineries before Earth
787 oxygenation. *Nature Ecology & Evolution*, 6(10), 1564-1572. [https://doi.org/10.1038/s41559-](https://doi.org/10.1038/s41559-022-01857-1)
788 [022-01857-1](https://doi.org/10.1038/s41559-022-01857-1)

789 Garcia, P. S., Jauffrit, F., Grangeasse, C., & Brochier-Armanet, C. (2019). GeneSpy, a user-friendly and
790 flexible genomic context visualizer. *Bioinformatics*, 35(2), 329-331.
791 <https://doi.org/10.1093/bioinformatics/bty459>

792 Goddard-Borger, E. D., & Williams, S. J. (2017). Sulfoquinovose in the biosphere: occurrence,
793 metabolism and functions. *Biochemical Journal*, 474(5), 827-849.
794 <https://doi.org/10.1042/BCJ20160508>

795 Gorzyska, A. K., Denger, K., Cook, A. M., & Smits, T. H. M. (2006). Inducible transcription of genes
796 involved in taurine uptake and dissimilation by *Silicibacter pomeroyi* DSS-3^T. *Archives of*
797 *Microbiology*, 185(5), 402-406. <https://doi.org/10.1007/s00203-006-0106-8>

798 Guler, S., Essigmann, B., & Benning, C. (2000). A cyanobacterial gene, *sqdX*, required for biosynthesis
799 of the sulfolipid sulfoquinovosyldiacylglycerol. *Journal of Bacteriology*, 182(2), 543-545.
800 <https://doi.org/10.1128/JB.182.2.543-545.2000>

801 Haft, D. H., DiCuccio, M., Badretdin, A., Brover, V., Chetvernin, V., O'Neill, K., Li, W., Chitsaz, F.,
802 Derbyshire, M. K., Gonzales, N. R., Gwadz, M., Lu, F., Marchler, G. H., Song, J. S., Thanki, N.,
803 Yamashita, R. A., Zheng, C., Thibaud-Nissen, F., Geer, L. Y., Marchler-Bauer, A., & Pruitt, K. D.
804 (2018). RefSeq: an update on prokaryotic genome annotation and curation. *Nucleic Acids*
805 *Research*, 46(D1), D851-D860. <https://doi.org/10.1093/nar/gkx1068>

806 Hanson, B. T., Dimitri Kits, K., Loffler, J., Burrichter, A. G., Fiedler, A., Denger, K., Frommeyer, B.,
807 Herbold, C. W., Rattei, T., Karcher, N., Segata, N., Schleheck, D., & Loy, A. (2021).
808 Sulfoquinovose is a select nutrient of prominent bacteria and a source of hydrogen sulfide in
809 the human gut. *ISME J.* <https://doi.org/10.1038/s41396-021-00968-0>

810 Henriques, A. C., & De Marco, P. (2015). Methanesulfonate (MSA) catabolic genes from marine and
811 estuarine bacteria. *PLoS One*, 10(5), e0125735.
812 <https://doi.org/10.1371/journal.pone.0125735>

813 Horinouchi, M., Yoshida, T., Nojiri, H., Yamane, H., & Omori, T. (1999). Polypeptide requirement of
814 multicomponent monooxygenase DsoABCDEF for dimethyl sulfide oxidizing activity.
815 *Bioscience, biotechnology and biochemistry*, 63(10), 1765-1771.
816 <https://doi.org/10.1271/bbb.63.1765>

817 Jaffe, A. L., Castelle, C. J., Matheus Carnevali, P. B., Gribaldo, S., & Banfield, J. F. (2020). The rise of
818 diversity in metabolic platforms across the Candidate Phyla Radiation. *BMC Biology*, 18(1), 69.
819 <https://doi.org/10.1186/s12915-020-00804-5>

820 Jeni, L. A., Cohn, J. F., & De La Torre, F. (2013). Facing imbalanced data recommendations for the use
821 of performance metrics. *Int Conf Affect Comput Intell Interact Workshops, 2013*, 245-251.
822 <https://doi.org/10.1109/ACII.2013.47>

823 Kappler, U., & Schäfer, H. (2014). Transformations of dimethylsulfide. *Metal Ions in Life Sciences*, 14,
824 279-313.

825 Kelly, D. P., & Murrell, J. C. (1999). Microbial metabolism of methanesulfonic acid. *Archives of*
826 *Microbiology*, 172(6), 341-348. <https://doi.org/10.1007/s002030050770>

827 Kiene, R. P., Linn, L. J., & Bruton, J. A. (2000). New and important roles for DMSP in marine microbial
828 communities. *Journal of Sea Research*, 43(3-4), 209-224. [https://doi.org/10.1016/S1385-1101\(00\)00023-X](https://doi.org/10.1016/S1385-1101(00)00023-X)

830 Koch, T., & Dahl, C. (2018). A novel bacterial sulfur oxidation pathway provides a new link between the
831 cycles of organic and inorganic sulfur compounds. *ISME Journal*, 12(10), 2479-2491.
832 <https://doi.org/10.1038/s41396-018-0209-7>

833 Kredich, N. M. (1996). Biosynthesis of cysteine. In F. C. Neidhardt (Ed.), *Escherichia coli and Salmonella*
834 *typhimurium. Cellular and molecular biology* (pp. 514-527). Washington D.C.: American Society
835 for Microbiology.

836 Krejci, Z., Hollemeyer, K., Smits, T. H. M., & Cook, A. M. (2010). Isethionate formation from taurine in
837 *Chromohalobacter salexigens*: purification of sulfoacetaldehyde reductase. *Microbiology*
838 *(Reading)*, 156(Pt 5), 1547-1555. <https://doi.org/10.1099/mic.0.036699-0>

839 Kröber, E., & Schäfer, H. (2019). Identification of proteins and genes expressed by *Methylophaga*
840 *thiooxydans* during growth on dimethylsulfide and their presence in other members of the
841 genus. *Frontiers in Microbiology*, *10*, 1132. <https://doi.org/10.3389/fmicb.2019.01132>
842 Leimkühler, S., & Iobbi-Nivol, C. (2016). Bacterial molybdoenzymes: old enzymes for new purposes.
843 *FEMS Microbiology Reviews*, *40*(1), 1-18. <https://doi.org/10.1093/femsre/fuv043>
844 Letunic, I., & Bork, P. (2021). Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree
845 display and annotation. *Nucleic Acids Research*, *49*(W1), W293-W296.
846 <https://doi.org/10.1093/nar/gkab301>
847 Leyh, T. S., Taylor, J. C., & Markham, G. D. (1988). The sulfate activation locus of *Escherichia coli* K12:
848 cloning, genetic, and enzymatic characterization. *Journal of Biological Chemistry*, *263*, 2409-
849 2416.
850 Li, J., Koch, J., Flegler, W., Garcia Ruiz, L., Hager, N., Ballas, A., Tanabe, T. S., & Dahl, C. (2022). A
851 metabolic puzzle: consumption of C₁ compounds and thiosulfate in *Hyphomicrobium*
852 *denitrificans* X^T. *Biochimica et Biophysica Acta (BBA) - Bioenergetics*, *1864*, 148932.
853 <https://doi.org/10.1016/j.bbabi.2022.148932>
854 Lovelock, J. E., Maggs, R. J., & Rasmussen, R. A. (1972). Atmospheric dimethyl sulphide and the natural
855 sulphur cycle. *Nature*, *237*(5356), 452-453. <https://doi.org/10.1038/237452a0>
856 Mayer, J., Huhn, T., Habeck, M., Denger, K., Hollemeyer, K., & Cook, A. M. (2010). 2,3-
857 Dihydroxypropane-1-sulfonate degraded by *Cupriavidus pinatubonensis* JMP134: purification
858 of dihydroxypropanesulfonate 3-dehydrogenase. *Microbiology (Reading)*, *156*(Pt 5), 1556-
859 1564. <https://doi.org/10.1099/mic.0.037580-0>
860 McDevitt, C. A., Hanson, G. R., Noble, C. J., Cheesman, M. R., & McEwan, A. G. (2002). Characterization
861 of the redox centers in dimethyl sulfide dehydrogenase from *Rhodovulum sulfidophilum*.
862 *Biochemistry*, *41*(51), 15234-15244.
863 McDevitt, C. A., Hugenholtz, P., Hanson, G. R., & McEwan, A. G. (2002). Molecular analysis of dimethyl
864 sulphide dehydrogenase from *Rhodovulum sulfidophilum*: its place in the dimethyl sulphoxide
865 reductase family of microbial molybdopterin-containing enzymes. *Molecular Microbiology*,
866 *44*(6), 1575-1587.
867 McEwan, A. G., Hanson, G. R., & Bailey, S. (1998). Dimethylsulphoxide reductase from purple
868 phototrophic bacteria: structures and mechanism(s). *Biochemical Society Transactions*, *26*,
869 390-396.
870 McEwan, A. G., Ridge, J. P., McDevitt, C. A., & Hugenholtz, P. (2010). The DMSO reductase family of
871 microbial molybdenum enzymes; molecular properties and role in the dissimilatory reduction
872 of toxic elements. *Geomicrobiology Journal*, *19*(1), 3-21.
873 <https://doi.org/10.1080/014904502317246138>
874 Moran, M. A., & Durham, B. P. (2019). Sulfur metabolites in the pelagic ocean. *Nature Reviews*
875 *Microbiology*. <https://doi.org/10.1038/s41579-019-0250-1>
876 Peck, S. C., Denger, K., Burrichter, A., Irwin, S. M., Balskus, E. P., & Schleheck, D. (2019). A glycy radical
877 enzyme enables hydrogen sulfide production by the human intestinal bacterium *Bilophila*
878 *wadsworthia*. *Proceedings of the National Academy of Sciences of the United States of*
879 *America*, *116*(8), 3171-3176. <https://doi.org/10.1073/pnas.1815661116>
880 Refaeilzadeh, P., Tang, L., & Liu, H. (2009). Cross-validation. In L. Liu & T. Özsu (Eds.), *Encyclopedia of*
881 *database systems* (pp. 532-538). Boston, MA: Springer.
882 Reisch, C. R., Stoudemayer, M. J., Varaljay, V. A., Amster, I. J., Moran, M. A., & Whitman, W. B. (2011).
883 Novel pathway for assimilation of dimethylsulphoniopropionate widespread in marine
884 bacteria. *Nature*, *473*(7346), 208-211. <https://doi.org/10.1038/nature10078>
885 Rossak, M., Schafer, A., Xu, N., Gage, D. A., & Benning, C. (1997). Accumulation of sulfoquinovosyl-1-
886 O-dihydroxyacetone in a sulfolipid-deficient mutant of *Rhodobacter sphaeroides* inactivated in
887 *sqdC*. *Archives of Biochemistry and Biophysics*, *340*(2), 219-230.
888 <https://doi.org/10.1006/abbi.1997.9931>

889 Rossak, M., Tietje, C., Heinz, E., & Benning, C. (1995). Accumulation of UDP-sulfoquinovose in a
890 sulfolipid-deficient mutant of *Rhodobacter sphaeroides*. *Journal of Biological Chemistry*,
891 270(43), 25792-25797. <https://doi.org/10.1074/jbc.270.43.25792>

892 Ruff, J., Denger, K., & Cook, A. M. (2003). Sulphoacetaldehyde acetyltransferase yields acetyl
893 phosphate: purification from *Alcaligenes defragrans* and gene clusters in taurine degradation.
894 *Biochemical Journal*, 369(Pt 2), 275-285. <https://doi.org/10.1042/BJ20021455>

895 Sayers, E. W., Cavanaugh, M., Clark, K., Ostell, J., Pruitt, K. D., & Karsch-Mizrachi, I. (2019). GenBank.
896 *Nucleic Acids Research*, 47(D1), D94-D99. <https://doi.org/10.1093/nar/gky989>

897 Schäfer, H., Myronova, N., & Boden, R. (2010). Microbial degradation of dimethylsulphide and related
898 C₁-sulphur compounds: organisms and pathways controlling fluxes of sulphur in the biosphere.
899 *Journal of Experimental Botany*, 61(2), 315-334. <https://doi.org/10.1093/jxb/erp355>

900 Schmitz, R. A., Mohammadi, S. S., van Erven, T., Berben, T., Jetten, M. S. M., Pol, A., & Op den Camp,
901 H. J. M. (2022). Methanethiol consumption and hydrogen sulfide production by the
902 thermoacidophilic methanotroph *Methylacidiphilum fumariolicum* SolV. *Frontiers in*
903 *Microbiology*, 13, 857442. <https://doi.org/10.3389/fmicb.2022.857442>

904 Sharma, M., Lingford, J. P., Petricevic, M., Snow, A. J. D., Zhang, Y., Jarva, M. A., Mui, J. W., Scott, N. E.,
905 Saunders, E. C., Mao, R., Epa, R., da Silva, B. M., Pires, D. E. V., Ascher, D. B., McConville, M. J.,
906 Davies, G. J., Williams, S. J., & Goddard-Borger, E. D. (2022). Oxidative desulfurization pathway
907 for complete catabolism of sulfoquinovose by bacteria. *Proceedings of the National Academy*
908 *of Sciences of the United States of America*, 119(4). <https://doi.org/10.1073/pnas.2116022119>

909 Struwe, M. A., Kalimuthu, P., Luo, Z., Zhong, Q., Ellis, D., Yang, J., Khadanand, K. C., Harmer, J. R., Kirk,
910 M. L., McEwan, A. G., Clement, B., Bernhardt, P. V., Kobe, B., & Kappler, U. (2021). Active site
911 architecture reveals coordination sphere flexibility and specificity determinants in a group of
912 closely related molybdoenzymes. *Journal of Biological Chemistry*, 296, 100672.
913 <https://doi.org/10.1016/j.jbc.2021.100672>

914 Tanabe, T. S., & Dahl, C. (2022). HMS-S-S: a tool for the identification of sulphur metabolism-related
915 genes and analysis of operon structures in genome and metagenome assemblies. *Molecular*
916 *Ecology Resources*, 22(7), 2758-2774. <https://doi.org/10.1111/1755-0998.13642>

917 Tanaka, Y., Yoshikaie, K., Takeuchi, A., Ichikawa, M., Mori, T., Uchino, S., Sugano, Y., Hakoshima, T.,
918 Takagi, H., Nonaka, G., & Tsukazaki, T. (2020). Crystal structure of a YeeE/YedE family protein
919 engaged in thiosulfate uptake. *Science Advances*, 6(35), eaba7637.
920 <https://doi.org/10.1126/sciadv.aba7637>

921 Thume, K., Gebser, B., Chen, L., Meyer, N., Kieber, D. J., & Pohnert, G. (2018). The metabolite
922 dimethylsulfoxonium propionate extends the marine organosulfur cycle. *Nature*, 563(7731),
923 412-415. <https://doi.org/10.1038/s41586-018-0675-0>

924 Todd, J. D., Curson, A. R., Sullivan, M. J., Kirkwood, M., & Johnston, A. W. (2012). The *Ruegeria pomeroyi*
925 *acul* gene has a role in DMSP catabolism and resembles *yhdH* of *E. coli* and other bacteria in
926 conferring resistance to acrylate. *PLoS One*, 7(4), e35947.
927 <https://doi.org/10.1371/journal.pone.0035947>

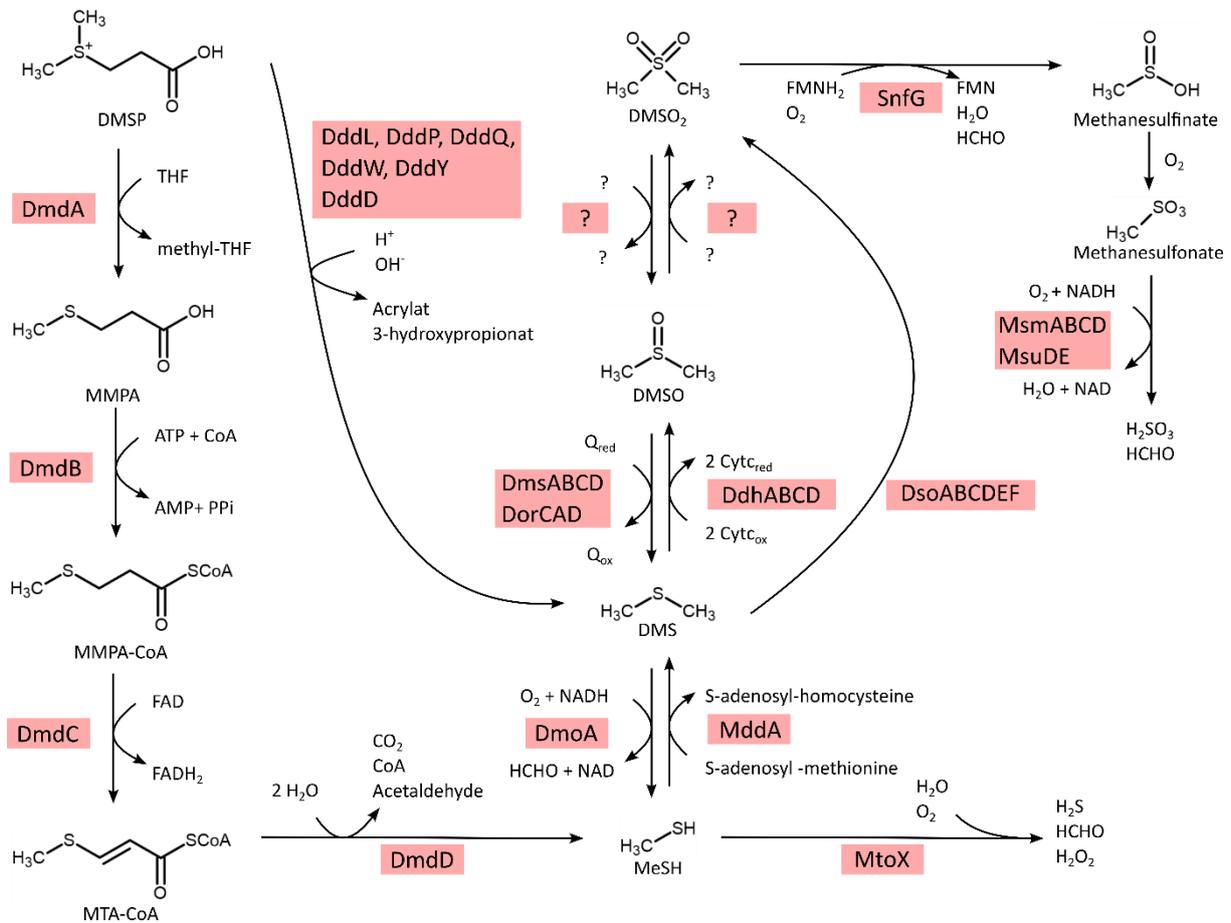
928 Varma, S., & Simon, R. (2006). Bias in error estimation when using cross-validation for model selection.
929 *BMC Bioinformatics*, 7, 91. <https://doi.org/10.1186/1471-2105-7-91>

930 Wei, Y., Tong, Y., & Zhang, Y. (2022). New mechanisms for bacterial degradation of sulfoquinovose.
931 *Bioscience Reports*, 42(10). <https://doi.org/10.1042/BSR20220314>

932 Weinitschke, S., Denger, K., Cook, A. M., & Smits, T. H. M. (2007). The DUF81 protein TauE in
933 *Cupriavidus necator* H16, a sulfite exporter in the metabolism of C₂ sulfonates. *Microbiology*,
934 153, 3055-3060. <https://doi.org/10.1099/mic.0.2007/009845-0>

935 Weinitschke, S., Hollemeyer, K., Kusian, B., Bowien, B., Smits, T. H., & Cook, A. M. (2010). Sulfoacetate
936 is degraded via a novel pathway involving sulfoacetyl-CoA and sulfoacetaldehyde in
937 *Cupriavidus necator* H16. *Journal of Biological Chemistry*, 285(46), 35249-35254.
938 <https://doi.org/10.1074/jbc.M110.127043>

939 Weinitschke, S., Sharma, P. I., Stingl, U., Cook, A. M., & Smits, T. H. (2010). Gene clusters involved in
940 isethionate degradation by terrestrial and marine bacteria. *Applied and Environmental*
941 *Microbiology*, 76(2), 618-621. <https://doi.org/10.1128/AEM.01818-09>
942 Wicht, D. K. (2016). The reduced flavin-dependent monooxygenase SfnG converts dimethylsulfone to
943 methanesulfinate. *Archives of Biochemistry and Biophysics*, 604, 159-166.
944 <https://doi.org/10.1016/j.abb.2016.07.001>
945 Wolf, P. G., Cowley, E. S., Breister, A., Matatov, S., Lucio, L., Polak, P., Ridlon, J. M., Gaskins, H. R., &
946 Anantharaman, K. (2022). Diversity and distribution of sulfur metabolic genes in the human
947 gut microbiome and their association with colorectal cancer. *Microbiome*, 10(1), 64.
948 <https://doi.org/10.1186/s40168-022-01242-x>
949
950

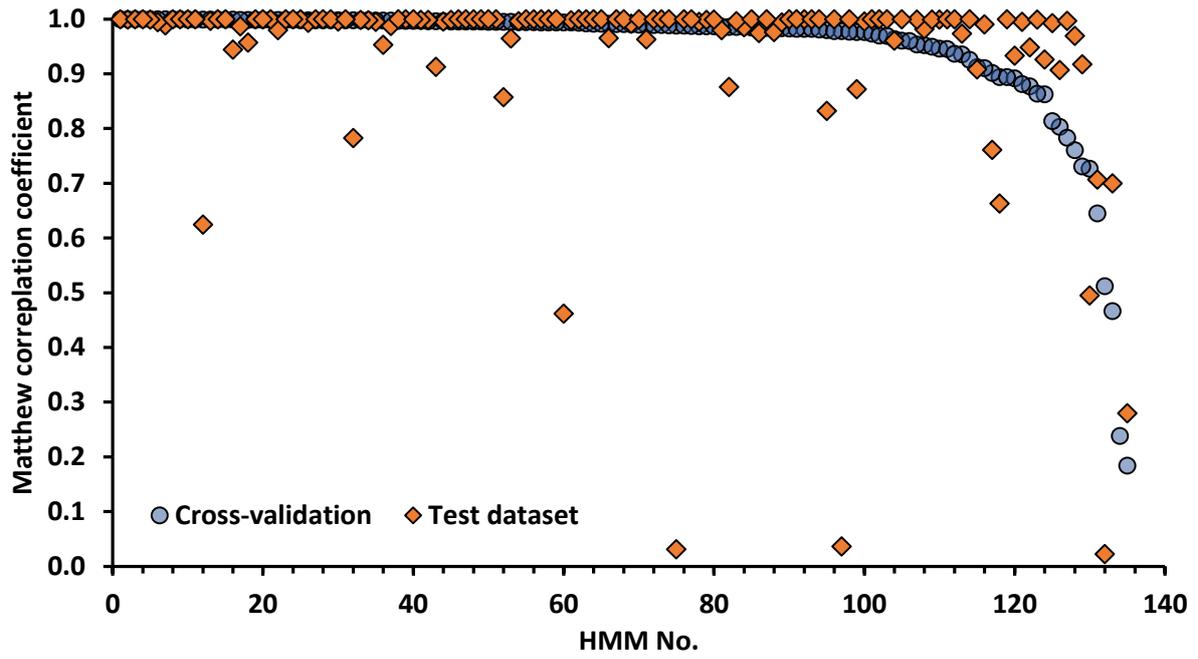


951

952 **Figure 1 Prokaryotic metabolism of C₁ organosulfur compounds.** All proteins shown have a corre-
 953 sponding HMM in HMSS2. Cyt_c, Cytochrome c; DMSP, dimethylsulfoniopropionate; DHPS, 2,3-dihydrox-
 954 ypropane-1-sulfonate; DMS, dimethylsulfide; DMSO, dimethylsulfone, DMSO₂ dimethylsulfoxide; FMN,
 955 flavin mononucleotide; FMNH₂, reduced flavin mononucleotide; MeSH, methanethiol; MMPA, methyl-
 956 mercaptopropionate; MMPA-CoA, 3-methylmercaptopyrionyl-CoA; MTA-CoA, methylthioacryloyl-
 957 CoA; THF, tetrahydrofolate.

958

966

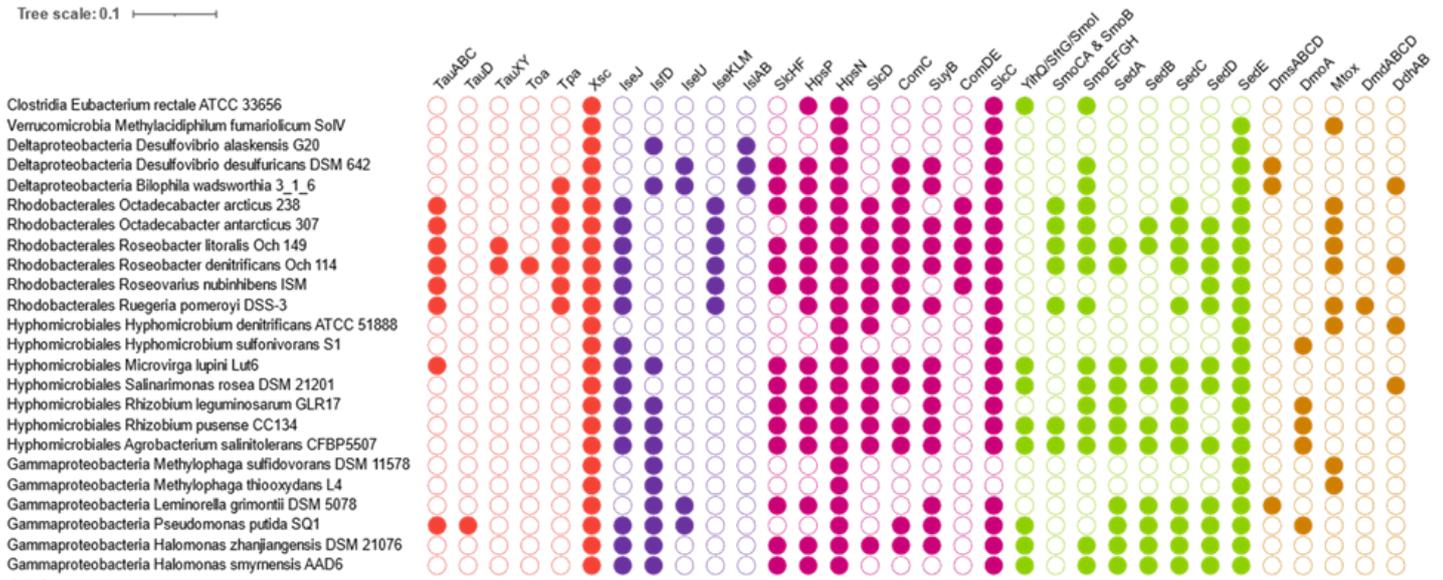


967

968 **Figure 3 Validation of the 134 HMMs generated in this work.** Performance was assessed by cross-
969 validation (blue dots) and on an independent test dataset (red diamonds). For each HMM Matthew
970 correlation coefficient was calculated. HMMs were ranked by their performance in cross-validation.

971

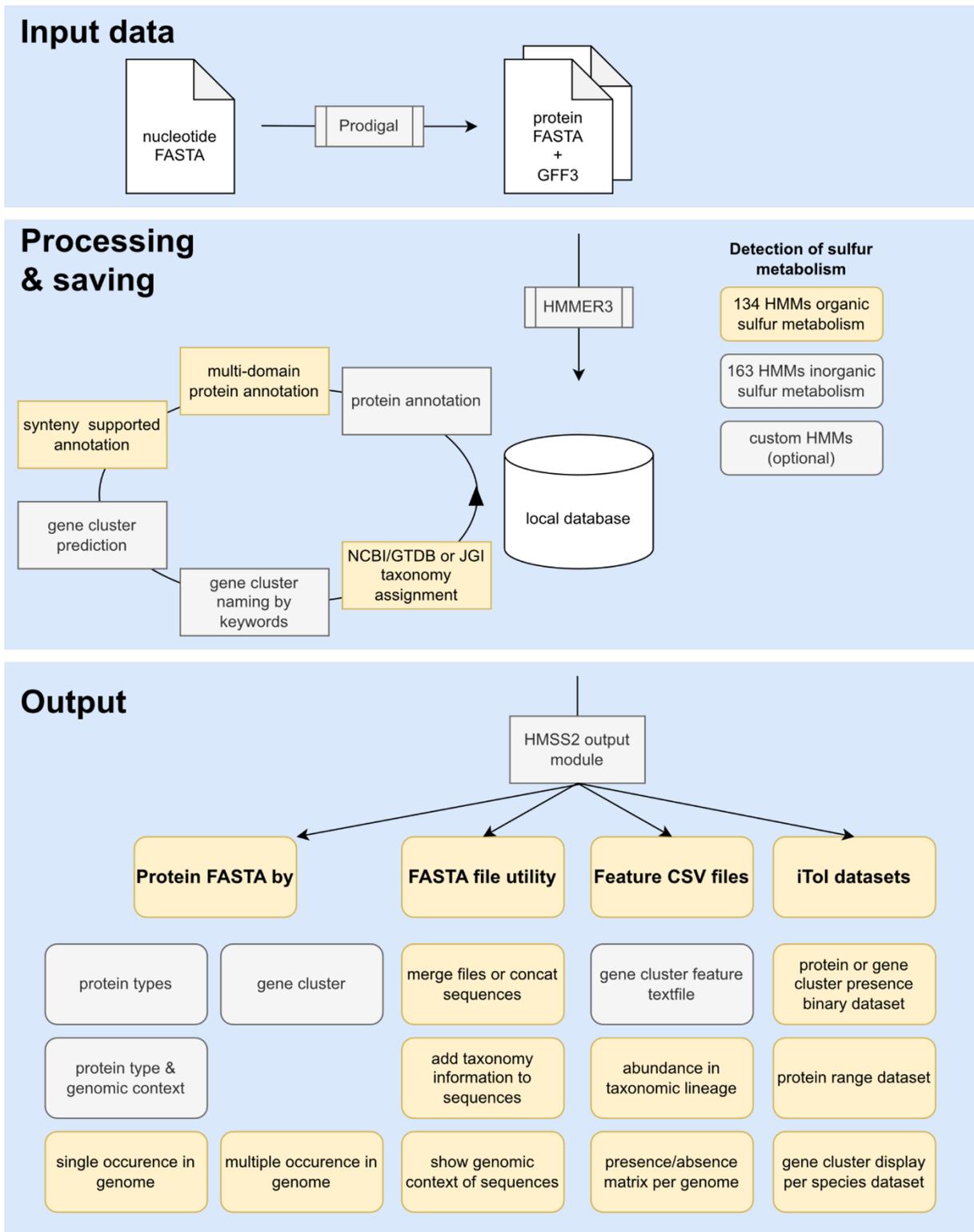
Tree scale: 0.1



972

973 **Figure 4. Presence/absence of proteins involved in the metabolism of organic sulfur compounds.** Oc-
 974 currence of genes for proteins involved in taurine degradation, isethionate degradation, 2,3-dihydrox-
 975 ypropane-1-sulfonate, sulfoquinovose and DMS metabolism, is indicated by filled orange, violet, pur-
 976 ple, green and light brown circles, respectively. The function of the individual proteins can be deduced
 977 from Figures 1 and 2.

978

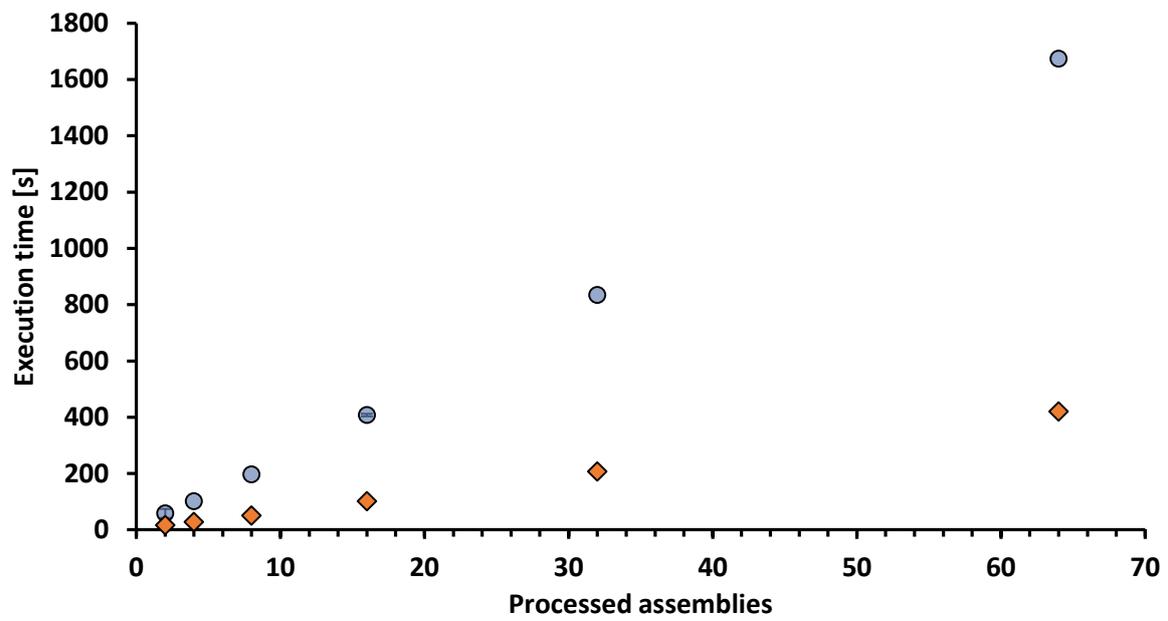


979

980 **Figure 5. Algorithm overview of HMSS2.** New features added in HMSS2 are highlighted in yellow. The
 981 only external programs required are HMMER3 and Prodigal.

982

983



984

985 **Figure 6. Computing time required by HMS-S-S compared to HMSS2.** Test were performed in triplicate
 986 with defined numbers of randomly selected sulfur-oxidizing or sulfur-reducing prokaryotes and 164
 987 HMMs. Blue circles: HMS-S-S, orange diamonds: HMSS2

988

989

- 990 **Supplementary Tables and Figures**
- 991 **Table S1. Reference proteins for dataset annotation**
- 992 **Table S2. HMM performance evaluation**
- 993 **Table S3. HMS-S-S vs. HMSS2 Benchmark**
- 994 **Table S4. Organisms for case study**