

# 1 Recovering genomic clusters of secondary 2 metabolites from lakes: a Metagenomics 2.0 3 approach

4 Rafael R. C. Cuadrat<sup>1,2</sup>, Danny Ionescu<sup>1</sup>, Alberto M. R. Davila<sup>3</sup>, Hans-Peter Grossart<sup>1,4</sup> \*

5 <sup>1</sup> Leibniz-Institute of Freshwater Ecology and Inland Fisheries, Alte Fischerhuetten 2, OT Neuglobsow,  
6 16775, Stechlin, Germany

7 <sup>2</sup> Berlin Center for Genomics in Biodiversity Research, Königin-Luise-Straße 6-8, 14195, Berlin, Germany

8 <sup>3</sup> Computational and Systems Biology Laboratory, Oswaldo Cruz Institute, Fiocruz, Avenida Brasil 4365, Rio  
9 de Janeiro CEP 21040-360, Brazil

10 <sup>4</sup> Potsdam University, Institute for Biochemistry and Biology, Potsdam, Germany

11

12 \* Corresponding author

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

## 31 **Abstract**

### 32 **Background**

33 Metagenomic approaches became increasingly popular in the past decades due to  
34 decreasing costs of DNA sequencing and bioinformatics development. So far, however, the  
35 recovery of long genes coding for secondary metabolism still represents a big challenge.  
36 Often, the quality of metagenome assemblies is poor, especially in environments with a  
37 high microbial diversity where sequence coverage is low and complexity of natural  
38 communities high. Recently, new and improved algorithms for binning environmental  
39 reads and contigs have been developed to overcome such limitations. Some of these  
40 algorithms use a similarity detection approach to classify the obtained reads into  
41 taxonomical units and to assemble draft genomes. This approach, however, is quite limited  
42 since it can classify exclusively sequences similar to those available (and well classified) in  
43 the databases.

44 In this work, we used draft genomes from Lake Stechlin, north-eastern Germany, recovered  
45 by MetaBat, an efficient binning tool that integrates empirical probabilistic distances of  
46 genome abundance, and tetranucleotide frequency for accurate metagenome binning. These  
47 genomes were screened for secondary metabolism genes, such as polyketide synthases  
48 (PKS) and non-ribosomal peptide synthases (NRPS), using the Anti-SMASH and  
49 NAPDOS workflows.

### 50 **Results**

51 With this approach we were able to identify 243 secondary metabolite clusters from 121  
52 genomes recovered from the lake samples. A total of 18 NRPS, 19 PKS and 3 hybrid

53 PKS/NRPS clusters were found. In addition, it was possible to predict the partial structure  
54 of several secondary metabolite clusters allowing for taxonomical classifications and  
55 phylogenetic inferences.

## 56 **Conclusions**

57 Our approach revealed a great potential to recover and study secondary metabolites genes  
58 from any aquatic ecosystem.

59

60 Keywords: Metagenomics 2.0, secondary metabolites, nonribosomal peptide synthetase,  
61 polyketide synthase

62

63

64

65

66

67

68

69

70

71

## 72 **Background**

73 Metagenomics, also known as environmental genomics, describes the study of a microbial  
74 community without the need of *a priori* cultivation in the laboratory. It has the potential to  
75 explore uncultivable microorganisms by accessing and sequencing their nucleic acid [1]. In  
76 recent years, due to decreasing costs of DNA sequencing - metagenomic databases [2] (e.g.,  
77 MG-RAST) have rapidly grown and archive billions of short read sequences [3]. Many  
78 metagenomic tools and pipelines were proposed to better analyse these enormous datasets  
79 [4]. Additionally, these tools allow to (i) infer ecological patterns, alfa- and beta-diversity  
80 and richness [5]; (ii) assemble environmental contigs from the reads [6] and more recently,  
81 (iii) recover draft genomes from metagenomic bins [7][8][9]. By recovering a high number  
82 of draft genomes from these so far uncultivable organisms, it is now possible to screen for  
83 new genes and clusters, unlocking a previously underestimated metabolic potential such as  
84 secondary metabolite gene clusters by using a metagenomic approach called Metagenomics  
85 2.0 [10].

86 Polyketide synthases (PKS) and non-ribosomal peptide synthases (NRPS) are two families  
87 of modular mega-synthases, both are very important for the biotechnological and  
88 pharmaceutical industry due to their broad spectrum of products, spanning from antibiotics  
89 and antitumor drugs to food pigments. They act in an analogous way, producing  
90 polyketides (using acyl-coA monomers) and peptides (using aminoacyl monomers),  
91 respectively. Both families are broadly distributed in many taxonomical groups, ranging  
92 from bacteria (alphaproteobacteria, cyanobacteria, actinobacteria) to fungi [11][12].

93 PKS enzymes can be classified in types (I, II and III), where type I can be further classified  
94 into modular or iterative classes. The iterative PKS use the same domain many times,  
95 iteratively, to synthesize the polyketide. The modular PKS are large multi-domain enzymes  
96 in which each domain is used only once in the synthesis process [13][14]. The production  
97 of the polyketide follows the co-linearity rule, each module being responsible for the  
98 addition of one monomer to the growing chain [15].

99 Type I PKS are characterized by multiple domains in the same open reading frame (ORF)  
100 while in type II each domain is encoded in a separate ORF, acting interactively [16]. Type  
101 III is also known as Chalcone synthase and has different evolutionary origin from type I  
102 and II [17]. Type III PKSs are self-contained enzymes that form homodimers. Their single  
103 active site in each monomer catalyzes the priming, extension, and cyclization reactions  
104 iteratively to form polyketide products [17]. Hybrid PKS/NRPS and NRPS/PKS are also  
105 modular enzymes, encoding lipopeptides (hybrid between polyketides and peptides) and  
106 occur in bacterial as well as fungal genomes [18][19][20].

107 PKS and NRPS are very well explored in genomes from cultivable organisms, mainly  
108 *Actinomycetes* [21] and *Cyanobacteria* [22]. Recently, by using a metagenomic approach,  
109 studies have demonstrated the presence of these metabolite-genes in aquatic environments,  
110 as for example, Brazilian coast waters (in free living and particle-associated bacteria) and  
111 from the microbiomes of Australian marine sponges [23]. However, there are few  
112 metagenomic studies whose scope is to find these gene families in freshwater environments  
113 where most studies are based on isolation approaches [24][25].

114 In addition, due to the rather large size of genes involved in these pathways, yet, it is not  
115 possible to recover the full genes by using traditional read-based metagenomics or the

116 single sample assembly approach. Most of the studies aim to solely find specific domains,  
117 like Keto-synthase (KS) in PKS and Condensation domain (C) in NRPS, due to the high  
118 conservation of these domains [26].

119 We used a metagenomics 2.0 approach to overcome these limitations and improve the  
120 screening for secondary metabolism genes and clusters while evaluating the potential of  
121 microbial communities for future research on potential drugs. This study aims to (i)  
122 generate draft genomes from Lake Stechlin; (ii) to screen these genomes for new complete  
123 multi-modular enzymes from PKS and NRPS families, exploring their diversity and  
124 phylogeny.

125

126

127

## 128 **Methods**

### 129 **Sampling and sequencing**

130 A total of 26 metagenomic samples from Lake Stechilin, north-eastern Germany were used.  
131 Water was collected as metagenomic samples on several occasions (April, June 2013, July  
132 2014, Aug 2015) in sterile 2 L Schott bottles from Lake Stechlin (53°9'5.59N,  
133 13°1'34.22E). All samples, except those from Aug 2015, were filtered through 5 µm and  
134 subsequently 0.2 µm pore-size filters. The samples collected in Aug 2015 were not size-  
135 fractionated and directly filtered on a 0.2 µm pore size filter due to specific research

136 demands. Genomic DNA was extracted using a phenol/chloroform protocol as described in  
137 [27] and was sent for sequencing.

138 Sequencing was conducted at MrDNA (Shallowater, Texas) on an Illumina Hiseq 2500,  
139 using the V3 chemistry, following, fragmentation, adaptor ligation and amplification of 50  
140 ng genomic DNA from each sample, using the Nextera DNA Sample Preparation Kit.

141 Table S1 shows the general information about the 26 samples used in this study.

142

### 143 **Environmental draft genomes**

144 Briefly, all samples were pre-processed by Nesoni ([https://github.com/Victorian-](https://github.com/Victorian-Bioinformatics-Consortium/nesoni)  
145 [Bioinformatics-Consortium/nesoni](https://github.com/Victorian-Bioinformatics-Consortium/nesoni)) to remove low quality sequences and to trim adaptors,  
146 and afterwards assembled together using MegaHIT (default parameters) [6]. The reads from  
147 each sample were mapped back to these assembled contigs using BBMAP  
148 (<https://sourceforge.net/projects/bbmap/>) and then all data was binned using MetaBAT [7]  
149 to generate the draft genomes. The completeness and taxonomical classification were  
150 checked using CheckM [28].

151

### 152 **Screening secondary metabolism genes and phylogenetic analysis of NRPS and PKS** 153 **domains**

154 DNA fasta files of the generated bins (288) were submitted to a locally installed version of  
155 Anti-SMASH (`--clusterblast -smcogs -limit 1500`) [29]. Using in-house ruby scripts, the  
156 domains from PKS and NRPS were parsed. The PKS KS domains and NRPS C domains

157 were submitted to NAPDOS for classification [30]. In addition, all the KS and C domains  
158 (trimmed by NAPDOS) were submitted to BLASTP against RefSeq database [31], using  
159 the default parameters. The 3 best hits of each domain were extracted and added to the  
160 original multi-fasta file with the environmental domains. The full set of KS and C domains  
161 (from bins and references obtained by the blast on RefSeq database) was submitted for  
162 NAPDOS for the phylogenetic analysis. The resulting alignment and tree were exported  
163 and the trees were manually checked and annotated.

164

### 165 **Relative abundance of bins in each sample**

166 The reads from each sample were mapped (using BBMAP) against each bin fasta file and  
167 an in-house ruby parser script was used to calculate the relative abundance of each bin in  
168 each sample, normalizing the read counts by the number of reads of each sample. The table  
169 with the results was loaded into STAMP [32] in order to analyse the significant differences  
170 of bin abundance over the samples.

171

172

## 173 **Results**

### 174 **Environmental draft genomes obtained (bins)**

175 Metagenomic binning resulted in 288 draft environmental genomes (called bins in this  
176 study). Of these, 45 had a predicted completion level of more than 75% according to  
177 CheckM.



178 Table S2 shows the general information about each bin, including completeness, genome  
179 size, number of open reading frames (ORFs) and taxonomical classifications (from  
180 CheckM).

181

182

183

#### 184 **Screening secondary metabolism genes and phylogenetic analysis**

185 By using Anti-SMASH, at least one secondary metabolite gene cluster was found in 121 of  
186 the bins, totalling 243 clusters and 2200 ORFs. From these 243 clusters, 125 (51.4%) were  
187 classified in the Terpene and 35 (14.40%) in the Bacteriocin pathway. In addition, a total of  
188 18 NRPS, 6 type I PKS and 3 hybrid PKS/NRPS clusters were found in 15 different bins  
189 (Figure 1a). The latest 3 obtained pathway clusters are the main focus of our study.

190 Figure 1b shows the taxonomical classification at phylum level for the bins showing NRPS,  
191 type I PKS and hybrid clusters. Supplementary table S3 shows the distribution of all  
192 clusters in all bins.

193

194 <FIGURE 1>

195

196 **Figure 1A: Abundance of secondary metabolite cluster types obtained with Anti-SMASH in the**  
197 **recovered 288 bins (environmental genomes). B: Taxonomical classification of bins (Phyla) in**

198 which NRPS, PKS and Hybrid PKS/NRPS clusters were found. Red bar and pie: NRPS; blue bar and  
199 pie: Type I PKS; green bars and pie: Hybrid clusters (NRPS-PKS and PKS-NRPS).

200

201 A total of 43 condensation (C) domains were obtained from NRPS clusters. All these  
202 sequences were submitted to NAPDOS analysis. Figure 2a shows the classification of C  
203 domains into classes. Most of the sequences were classified as LCL domains (58%). This  
204 kind of domain catalyzes the formation of a peptide bond between two L-amino acids.

205

206

207 <FIGURE 2>

208

209 **Figure 2A: NAPDOS classification of the NRPS KS domain.** Modular: possess a multidomain  
210 architecture consisting of multiple sets of modules; hybridKS: are biosynthetic assembly lines that  
211 include both PKS and NRPS components; PUFA: Polyunsaturated fatty acids (PUFAs) are long chain  
212 fatty acids containing more than one double bond, including omega-3-and omega-6- fatty acids;  
213 Eneidyne: a family of biologically active natural products. The Eneidyne core consists of two  
214 acetylenic groups conjugated to a double bond or an incipient double bond within a nine- or ten-  
215 membered ring. **2B: NAPDOS classification of NRPS C domain.** Cyc: cyclization domains catalyze  
216 both peptide bond formation and subsequent cyclization of cysteine, serine or threonine residues;  
217 DCL: link an L-amino acid to a growing peptide ending with a D-amino acid; Epim: epimerization  
218 domains change the chirality of the last amino acid in the chain from L- to D- amino acid; LCL:  
219 catalyze formation of a peptide bond between two L-amino acids; modAA: appear to be involved

220 in the modification of the incorporated amino acid; Start: first module of a Non-ribosomal peptide  
221 synthase (NRPS).

222

223

224

225 The screening for type I PKS resulted in 9 KS domain sequences. Most of them are  
226 classified as modular type I PKS (56%). All of them were submitted to NAPDOS and  
227 classified into 4 different classes (Figure 2b).

228

229 All the KS and C domains were also submitted to similarity analysis by using BLASTP  
230 against RefSeq database (table S4 and S5), and the best 3 hits of each sequence were  
231 extracted and used for phylogenetic analyses with NAPDOS. The trees for C and KS  
232 domains are shown in Figures 3 and 4, respectively.

233

234

235 <FIGURE 3>

236

237

238

239 **Figure 3: NAPDOS phylogenetic tree of C domains (environmental domains, the top 3 blast**  
240 **results on RefSeq and the NAPDOS reference sequences).** The shadow colours represent the  
241 domain classifications (LCL, CYC, Start domains, EPIM, ModAA, Dual, and DCL). The sidebars  
242 represent phyla (Proteobacteria, Cyanobacteria, Firmicutes, Actinobacteria, Verrucomicrobia). All  
243 sequences from environmental bins are in red.

244

245

246

247

248

249 <FIGURE 4>

250

251 **Figure 4: NAPDOS tree of KS domains (environmental domains, the top 3 blast results on RefSeq**  
252 **and the NAPDOS reference sequences).** The shadow colours represent the domain  
253 classifications (Modular, KS1, Iterative, Trans-AT, Hybrid, PUFA, Enediyenes, Type II  
254 and Fabs – Fatty acid synthase). The sidebars represent phyla (Cyanobacteria and  
255 Actinobacteria). All the sequences from environmental bins are in red.

256

257

258

259 **Relative abundance of bins on each sample**

260 The relative abundance of all bins in each sample was estimated by mapping the reads from  
261 each sample against the assembled bins. Table S6 shows the normalized bin abundance for  
262 every sample.

263 Due to differences between the filtration methods, we decided to classify the samples in 3  
264 groups: particle associated samples (PA) – filtered on 5  $\mu\text{m}$  membranes (and also samples  
265 from aggregates), free-living samples (FL) – pre-filtered through 5.0  $\mu\text{m}$  membranes and  
266 subsequently filtered on 0.22  $\mu\text{m}$  membranes, and non-size fractionated samples (NSF) –  
267 filtered direct on 0.22  $\mu\text{m}$  membranes (without previous filtering) retaining the whole  
268 bacterial community. In total, we obtained 7 samples in the PA group, 5 in the FL group,  
269 and 14 in the NSF group.

270 The table with the relative abundance of the bins in all samples was loaded on STAMP and  
271 an ANOVA test was conducted followed Games-Howell POST-HOC test and Benjamini-  
272 Hochberg FDR correction. Table S7 shows 158 bins for which the difference in relative  
273 abundance was statistically significant ( $p < 0.05$ ) between the 3 groups (FL, PA and NSF).

274 From the 15 bins containing NRPS and/or type I PKS clusters, only 4 showed significant  
275 difference between the 3 groups. Bins 1 and 2 are more abundant in PA samples and bins  
276 193 and 235 are more abundant in the FL samples.

277

278

279

280

281

282

283

### 284 **Exploring NRPS, type I PKS and hybrid clusters from draft genome bins**

285 We highlight 3 bins (with less than 35% contamination and more than 70% completeness)  
286 out of the 15 obtained type I PKS and/or NRPS and explore their clusters.

287 In bin 34 (*Pseudomonas*, 98.28% completeness) it was possible to retrieve 7 clusters,  
288 including 3 NRPS clusters (Figure 7) and 2 Bacteriocin clusters.

289 In cluster 2 (ctg181), multiple domains of NRPS (with the 3 minimal modules) and  
290 regulatory genes were identified, e.g., smCOG: SMCOG1057 (TetR family transcriptional  
291 regulator) (Figure 5, in green arrows).

292

293

294

295 <FIGURE 5>

296

297 **Figure 5: Bin 34, NRPS clusters detailed annotation and synteny.** The synteny of the clusters with  
298 a functional classification for each ORF is given. In addition, for the NRPS biosynthetic ORFS the

299 domain annotations are given. CAL: Co-enzyme A ligase domain, C: condensation, A: adenylation,  
300 E: epimerization, TE: Termination, KR: Ketoreductase domain, and ECH: Enoyl-CoA hydratase

301

302 All clusters show a high similarity with *Pseudomonas* proteins. Cluster 2 has a similarity of  
303 92% with *Pseudomonas synxantha* bg33r, conserving also the gene synteny.

304 The C domain sequences were submitted to NAPDOS analysis and 2 were classified as  
305 belonging to the Syringomycin pathway and the LCL class, and one was classified as  
306 belonging to the Microcystin pathway or/and the DCL class (links an L-amino acid to a  
307 growing peptide ending with a D-amino acid).

308 In cluster 3 (ctg415) (Figure 5), in addition to the NRPS domains, the following transporter  
309 related genes are found: smCOG: SMCOG1288 (ABC transporter related protein) and  
310 SMCOG1051 (TonB-dependent siderophore receptor) (blue narrows). Nevertheless, this  
311 cluster is not complete and just one C domain was found (LCL class), which was also  
312 classified to the Syringomycin pathway.

313 Cluster 6 (ctg857 – Figure 5) shows many NRPS domains, regulatory factors and  
314 transporters genes, including drug resistance genes, e.g. SMCOG1005 (drug resistance  
315 transporter, EmrB/QacA), SMCOG1044 (ABC transporter, permease protein) and  
316 SMCOG1051 (TonB-dependent siderophore receptor) (blue arrows). Two C domains from  
317 this cluster were classified as belonging to the heterocyclization class. This class catalyzes  
318 both peptide bond formation and subsequent cyclization of cysteine, serine or threonine  
319 residues [33]. Both domains were classified in the Pyochelin pathway by NAPDOS. The

320 phylogenetic tree (Figure 3) confirms both the functional and taxonomical classification  
321 (confidence value 100).

322

323 <FIGURE 6>

324

325 **Figure 6: Bin 193 and 131 type I PKS clusters detailed annotation and synteny.** It is possible to see  
326 the synteny of the cluster with the functional classification for each ORF. In addition, for the PKS  
327 biosynthetic ORF the domain-specific annotations can be seen. KS: keto-synthase, AT:  
328 acyltransferase, KR: ketoreductase, E: epimerization, DH: dehydratase, and ER: enoylreductase.

329

330

331 In **bin 193** (*Mycobacterium*, 73.37 % completeness) a type I PKS cluster was identified  
332 (ctg514) (Figure 6). Five PKS domains were retrieved, including the minimal core from  
333 one of the ORFs on this contig. The KS domain BLASTP result shows 82% (and 99%  
334 coverage) similarity with *Mycobacterium kansasii*. The NAPDOS analysis from the KS  
335 domain suggests that it could be a modular (Epothilone pathway) or iterative type I PKS  
336 similar to the Calicheamicin pathway. However, by using the phylogenetic analysis it was  
337 clustered as the iterative clade (confidence value 98.4) together with the *Mycobacterium*  
338 *kansasii* sequence (confidence value 100) (Figure 4).

339 Two further clusters were recovered: one type III PKS and one unclassified one. All  
340 clusters show similarity with the *Mycobacterium* clusters.



341

342 Bin 131 (unclassified bacteria by CheckM) has 84.09% of completeness reported by  
343 CheckM. In this bin it was possible to find one cluster and 3 domains of type I PKS (KS,  
344 AT and KR) (Figure 6). The KS domain was classified by NAPDOS as belonging to the  
345 Maduropeptin and Neocarzinostatin pathways. Using clusterblast inside Anti-Smash it was  
346 not possible to find any similar cluster, but using BLASTP it was possible to find similarity  
347 with the cyanobacteria *Microcystis aeruginosa* (64% identity and 99% coverage on  
348 BLASTP search). In the phylogenetic tree, it was clustered within the Eneidiynes clade  
349 (Figure 4) and also with *Microcystis aeruginosa* (confidence value 99).

350

351 In addition, there are 12 more bins with NRPS or type I PKS clusters, but with less than  
352 70% of completeness or more than 35% contamination. The bins 1 (69.54% completeness)  
353 and 2 (16.52%) were classified as the genus *Anabaena*, showing NRPS and hybrid NRPS-  
354 type 1 PKS, respectively. Bins 6 (39.66% completeness), 7 and 8 were classified as the  
355 genus *Planktothrix* and show a high diversity of secondary metabolites: 3 NRPS clusters,  
356 one type I PKS and 2 NRPS-PKS hybrid (bin 6), 2 NRPS (bin 7) and 2 NRPS and one type  
357 I PKS (bin 8). The bin 8 also shows a Microviridin cluster.

358 Bins 73 and 217 are classified as *Acidobacteria* showing PKS and NRPS, respectively. Bin  
359 235 (*Burkholderiaceae* family) shows 3 NRPS clusters and bin 13 (*Comamonadaceae*  
360 family, also from *Burkholderiales* order) shows one NRPS cluster.

361 Bin 78 is classified as *Verrucomicrobiaceae* shows 1 NRPS cluster. Additionally, there is  
362 an unclassified *Archaea* (NRPS cluster) and one bin without any classification (bin 136,  
363 type I PKS).

364

365 The Anti-Smash results for all bins are available in the Supplemental Information (SI 1).

366

## 367 Discussion

368 The field of metagenomics has generated a vast amount of data in the last decades [34].  
369 Most of the data is poorly annotated and little quality controlled when loaded into the  
370 public databases, hence awaiting a more in-depth analysis [35]. There are many open  
371 challenges in this field, e.g., (i) the lack of representative genomic databases from  
372 uncultivable organisms to be used in a similarity-based annotation procedure; (ii) high  
373 confidence assembly of short reads from species-rich samples; (iii) obtaining high enough  
374 coverage for every organism in the sample, including those with a low abundance, etc. [36].  
375 Recently, some new algorithms have been proposed to overcome these limitations and to  
376 obtain partial or near complete genomes from environmental samples, e.g., MetaBat [7] and  
377 MetaWatt [8]. Most of them require many samples and high coverage sequencing per  
378 sample as an input. Recently, studies have been done to recover genomes even from rare  
379 bacteria [37]. The term Metagenomics 2.0 was introduced to describe this new generation  
380 of metagenomic analysis by Katherine McMahon [10] and most of the studies using this  
381 approach have been conducted to reveal ecological interactions and networks [38][39].

382 In this study, we recovered 288 environmental draft genomes using 26 samples from Lake  
383 Stechlin, a temperate oligo-mesotrophic lake. One of the advantages of this approach is to  
384 enable the recovery of large genomic clusters, especially the Megasyntases clusters of the  
385 secondary metabolism, e.g., involved in biosynthesis of antibiotics, including its regulatory  
386 and transporter genes. Here, we have used the Anti-SMASH and NAPDOS pipelines to  
387 identify, annotate, classify and to carry out the phylogenetic analysis of a total of 243  
388 clusters of known secondary metabolites. To our knowledge, this is the first study using the  
389 metagenomics 2.0 approach to recover Megasyntases clusters. A number of previous  
390 studies had been conducted using a traditional PCR based screening [40] and shotgun  
391 metagenomics approach [41][42] exploring the abundance and diversity of individual genes  
392 and domains, but these studies are missing the genomic context. By obtaining the entire  
393 genomic context it is possible, in future studies, to clone and to do heterologous expression  
394 for all the genes, including promoters and transporters.

395

396 Screening the bins for secondary metabolite clusters, we can see that the most abundant  
397 cluster belongs to the Terpene pathway (125 clusters) (Figure 1). This biosynthesis pathway  
398 is well known to be present in many plant and fungi genomes, but recently it was proposed  
399 to be also widely distributed in bacterial genomes. One study revealed 262 distinct terpene  
400 synthases in the bacterial domain of life [43]. Consequently, it can represent a fertile source  
401 of new natural products – yet, greatly underestimated. The second most abundant class of  
402 clusters belongs to the Bacteriocin pathway (35 clusters) (Figure 1). Bacteriocin is a group  
403 of ribosomal synthesized antimicrobial peptides which can kill or inhibit bacterial strains  
404 closely related or non-related to the Bacteriocin producing bacteria [44]. It has been

405 suggested as a viable alternative to traditional antibiotics and can be used as narrow-  
406 spectrum antibiotics [45]. Only a few studies have been conducted to screen for Bacteriocin  
407 genes by using a metatenuomic approach, and solely for the host-associated microbiome [46]  
408 or fermented food microbiome [47][48][49]. None of these studies was conducted for  
409 natural environments and none have used the metagenomics 2.0 approach.

410 In this study, we focused on 2 families of large modular secondary metabolite genes, type I  
411 PKS and NRPS. With our approach, it was possible to find a total of 18 NRPS, 6 type I  
412 PKS and 3 hybrid PKS/NRPS clusters. For NRPS clusters, it was possible to recover 43 C  
413 domains, most of them (58%) from the LCL class. An LCL domain catalyzes a peptide  
414 bond between two L-amino acids [50]. A previous study also found that the LCL class was  
415 the most abundant in another aquatic environment, dominated by gram-negative bacteria  
416 [42]. Many studies have shown that the LCL class in aquatic environments is limited to  
417 gram-negative bacteria [51]. Our results further support this as we also found the LCL class  
418 only in bins of gram-negative bacteria (Figure 3) with the only exception of the unclassified  
419 Archaea (bin 233), which should be further investigated in order to confirm the  
420 phylogenetic classification. It was also possible to recover 9 KS domains from the type I  
421 PKS clusters, 56% from the modular class and 22% from the hybrid PKS/NRPS class.  
422 Those classes are larger (with many copies of each domain) than the iterative ones,  
423 increasing the chances to be recovered by metagenomic approaches. Accordingly, the  
424 NRPS and PKS clusters were more in depth analyzed, including syntheny, domain  
425 phylogeny, and partial metabolite protein structure predictions.

426 From the bins showing secondary metabolite genes, the most complete was from bin 34  
427 (*Pseudomonas*). The 7 clusters on this genome vary from 8,675 base pairs (bp) to 52,516

428 bp in size, been only possible to be recovered due to the high completeness of the  
429 assembled genome (98.28%). The presence of a great diversity of clusters in *Pseudomonas*  
430 is expected, as many active secondary metabolites (encoded by NRPS) have been  
431 previously described in the *Pseudomonas* genus, ranging from antibiotics and antifungal to  
432 siderophores [52][53][54].

433 From those 7 clusters in bin 34, the NRPS clusters 2 and 3 showed a high similarity with  
434 Syringomycin (three domains) and Microcystin (one domain) pathways. The first one is  
435 found, for example, in the *Pseudomonas syringae* (a plant pathogen) genome, as a virulence  
436 factor (Syringomycin E) [55] which also has antifungal activity against *Saccharomyces*  
437 *cerevisiae* [56]. On the other hand, Microcystin is a class of toxins produced by freshwater  
438 Cyanobacteria species [57] and it can be produced in large quantities during massive bloom  
439 events [58]. Due to the taxonomical classification of the bin and the higher number of  
440 domains similar to Syringomycin, however, it is more likely that the product encoded by  
441 this cluster is functionally close to the latter pathway.

442 In bin 34 - cluster 6, both C domains were classified as belonging to the Pyochelin  
443 pathway. This peptide is a siderophore of *Pseudomonas aeruginosa* [59]. The presence of a  
444 TonB-dependent siderophore receptor in cluster 6 provides additional evidence about its  
445 functional classification. Additionally, two Bacteriocin clusters and one Aryl polyene  
446 cluster were found in bin 34. Aryl polyenes are structurally similar to the well-known  
447 carotenoids with respect to their polyene systems and it was recently demonstrated that it  
448 can protect bacteria from reactive oxygen species, similarly to what is known for  
449 carotenoids [60]. These results suggest that a wide range of metabolites is encoded in this

450 *Pseudomonas* genome, providing it an “arsenal” of secondary products, increasing the  
451 likelihood of the *Pseudomonas* species to succeed in aquatic systems.

452 Bin 131 (unclassified bacteria) shows a PKS cluster and 3 domains. It was classified as  
453 belonging to the Eneidyne pathway. These compounds are toxic to DNA and are under  
454 investigation as anti-tumor agents, with several compounds under clinical trials [61]. All  
455 are encoded by type I iterative PKS [62] and it was possible to recover the minimal core  
456 (KS, AT and KR) as well as transporter genes from the environmental genome. The most  
457 similar PKS I present in the public databases stems from *Microcystis aeruginosa*, but only  
458 with an identity of 64%, suggesting that it is encoding for a new compound, which has not  
459 previously been described.

460

461 In bin 193 (*Mycobacterium* - sister lineage *M. rhodesiae*), one of the 3 recovered clusters is a  
462 Type I PKS, similar to iterative PKS in the NAPDOS analysis (confirmed by the  
463 phylogenetic analysis). The most similar KS sequence belongs to *Mycobacterium kansasii*,  
464 with 82% similarity. *M. rhodesiae* and *M. kansasii* are both non-tuberculous mycobacteria  
465 (NTM) that can be found in different environments, but both can also be opportunistic  
466 pathogens and cause a chronic pulmonary infection in immunosuppressed patients [63].  
467 The species *M. kansasii* comprises various subtypes and some are often recovered from tap  
468 water and occasionally from river or lake water [64][65]. There is still controversy about  
469 how the transmission from environment to human host occurs and also about the  
470 implications on public health. The presence of PKS in *Mycobacterium* genus was  
471 discovered more than a decade ago and most of the polyketides encoded by different  
472 species of this genus play role in virulence and/or components of the extraordinarily

473 complex mycobacterial cell envelope [66]. Further studies must be done in order to  
474 investigate better the potential of this bin to cause infections on humans, i.e. by screening  
475 virulence factors on the full genome.

476

477 Five bins of the phyla Cyanobacteria contained PKS and NRPS clusters. In bins 1 and 2  
478 (*Anabaena*, now called “*Dolichospermum*”), it was possible to recover 2 NRPS and 1  
479 hybrid NRPS-PKS, respectively. The genus *Anabaena* is known to encode several toxins,  
480 including the dangerous Anatoxin-a, and to produce toxic blooms in lakes and reservoirs  
481 [67][68][69][70]. However, the Anatoxin-a is encoded by a type I PKS cluster [71], unlike  
482 the NRPS and Hybrid clusters found in the *Anabaena* bins from this study.

483 On the other hand, the hepatotoxic heptapeptide of the class Mycrocystin is present in many  
484 genera of Cyanobacteria, including *Anabaena*, and they are encoded by NRPS and also  
485 Hybrid NRPS-PKS clusters [72][73][74][75]. The results of NAPDOS reveal one C domain  
486 from bin 1 classified in the pathway of Mycrocystin with e-value 6e-83.

487 In bins 6, 7 and 8 (*Planktothrix*), it was possible to find several type I PKS, NRPS and  
488 hybrid clusters. In bin 6 there are 3 NRPS, one type 1 PKS and one hybrid cluster. The bin  
489 7 shows 2 NRPS and bin 8 shows one type I PKS and 2 NRPS clusters. The genus  
490 *Planktothrix* also can be producer of Anatoxin-a [75] and the presence of type I PKS cluster  
491 on these bins can be alarming. However, the 3 KS domains from type I PKS from bin 6  
492 reveal great similarity with the Epothilone pathway. The NapDOS analysis and the KS  
493 domain from bin 8 suggest a high similarity to the neurotoxin Jamaicamides pathway. A  
494 previous study showed in 2010 [76] the presence of an Anatoxin-a-producing

495 cyanobacterium in northeastern Germany Lake Stolpsee, rising concerns about the presence  
496 of these toxins in the waters of these lakes.

497

498 The absence of Anatoxin-a genes in the studied lake is in agreement with previous  
499 screening for a toxin screening in the lake [77].

500 In the bin 73 (*Acidobacteriales*) 1 PKS sequence was found. By the phylogenetic  
501 classification, its KS domain is clustered with *trans*-AT KS domains. The AT domains  
502 of *trans*-AT PKSs are not integrated into the assembly lines but expressed as free-standing  
503 polypeptides, unlike the more familiar *cis*-AT PKSs [78]. However, the NAPDOS result  
504 shows the AT domain of this bin in the same ORF with KS and KR domains, showing a  
505 syntheny that suggests a *cis*-AT PKS. In addition, the classification by similarity from  
506 NAPDOS suggests a polyunsaturated fatty acid (PUFA) but only with 31% of identity

507 To assess the life style of the bins (free-living or particle-associated), we calculated the  
508 relative abundance of the bins in every sample. A total of 158 bins with significant  
509 difference between the 3 groups were found (Table S7), however from the 15 bins on which  
510 this study focused (showing NRPS and/or type I PKS clusters), only 4 bins (26.6%) were  
511 significantly differently present in the life-styles. Bins 1 and 2 (*Anabaena* genus) are more  
512 abundant in the PA group, especially on samples B7 and B9 (and also the replicates Old\_b7  
513 and Old\_b9), accounting for 20-25% (bin 1) and 10-15% (bin 2) on these samples. The  
514 very high abundance of these bins on the samples can be explained, based on the long term  
515 monitoring program of IGB on Lake Stechlin, by the fact that these samples were collected  
516 during the occurrence of a massive cyanobacterial bloom.



517

518 From the other bins containing PKS/NRPS clusters, we can see that Bins 6, 7, and 8  
519 (*Planktothrix*), beside the lack of significant difference between FL and PA groups (p-value  
520  $> p 0.05$ ), they are clearly more abundant in NSF. The possible explanation for this notion  
521 is that the NSF samples were collected during a mesocosm experiment, whereas the other  
522 samples were directly derived from the lake.

523

524

525

## 526 Conclusions

527

528 Using the Metagenomics 2.0 approach, we were able to recover full megasynthases  
529 sequences and its genomic context from environmental draft genomes. However, there are  
530 limitations, e.g., the genomic coverage of less abundant organisms and the possibility of  
531 chimeras. Recently, it has been demonstrated that with an increasing number of samples, it  
532 will be possible to recover individual species genomes with a high confidence [7]. In the  
533 near future, with the advent of the 3<sup>rd</sup> generation sequencing, with longer reads, up to 100  
534 kilobases, it will be possible to further improve the quality of the assemblies [79]. These  
535 new approaches unlock the possibility of studying these newly recovered environmental  
536 pathways and their evolution in detail. Thus, allowing cloning and expressing these clusters  
537 will provide new insights on natural products of great interest for biotechnological and  
538 pharmaceutical industry. Moreover, studies have demonstrated the possibility to synthesise  
539 large functional DNA [80], and together with additional screening techniques, it will be

540 possible to obtain such sequences and to synthesise the full cluster for heterologous  
541 expression, skipping the cloning and functional screening process, saving considerable time  
542 and money. In addition, the current work highlights the great potential for the discovery of  
543 new metabolically active compounds in freshwaters such as oligo-mesotrophic Lake  
544 Stechlin. Further, the study of complete or near complete genomes from uncultivated  
545 bacteria in the natural environment will enable us to better understand the multiple forms of  
546 interactions between species and how they compete for the limiting natural resources.

547

## 548 **Declarations**

### 549 **Available of data and materials**

550 The sequences generated for this study (metagenomic reads) were deposited in ENA  
551 (PRJEB22274 and PRJEB7963).

552

### 553 **Competing interests**

554 The authors declare no competing interests

### 555 **Funding**

556 This study was supported by the Science without Borders Program (Ciência Sem  
557 Fronteiras), CNPq. DI and HPG were funded by German science foundation (DFG)  
558 projects Aquameth (GR1540/21-1) and Aggregates (GR1540/28-1).

559

560 **Authors' contributions**

561 Conceived and designed the experiments: RRCC, DI, AMRD, and HPG. Performed the  
562 experiments: RRCC, DI. Analyzed the data: RRCC, DI, and HPG. Contributed  
563 reagents/materials/analysis tools: DI and HPG. All authors wrote the manuscript and  
564 revised it for significant intellectual content.

565

566 **Acknowledgements**

567 We thank Dr. Camila Mazzoni and all the team of Berlin Center for Genomics in  
568 Biodiversity Research (BeGenDiv) for allowing us to use the facilities and computational  
569 resources for the bioinformatics analyses. Elke Mach and the MIBI group are thanked for  
570 their technical support and fruitful discussions.

571

572

573

574

575

576

577

578

579

580

581

582

583

## 584 **References**

585

586 1. Rodríguez-Valera F: **Environmental genomics, the big picture?** *FEMS Microbiology*  
587 *Letters* 2004, **231**:153–158.

588

589 2. Vincent AT, Derome N, Boyle B, Culley AI, Charette SJ: **Next-generation sequencing**  
590 **(NGS) in the microbiological world: How to make the most of your money.** *Journal of*  
591 *Microbiological Methods* 2016.

592

593 3. Meyer F, Paarmann D, D'Souza M, Olson R, Glass E, Kubal M, Paczian T, Rodriguez  
594 A, Stevens R, Wilke A, Wilkening J, Edwards R: **The metagenomics RAST server – a**  
595 **public resource for the automatic phylogenetic and functional analysis of**  
596 **metagenomes.** *BMC Bioinformatics* 2008, **9**:386.

597

- 598 4. Huson DH, Auch AF, Qi J, Schuster SC: **MEGAN analysis of metagenomic data.**  
599 *Genome Research* 2007, **17**:377–386.  
600
- 601 5. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, Fierer  
602 N, Peña AG, Goodrich JK, Gordon JI, Huttley GA, Kelley ST, Knights D, Koenig JE, Ley  
603 RE, Lozupone CA, McDonald D, Muegge BD, Pirrung M, Reeder J, Sevinsky JR,  
604 Turnbaugh PJ, Walters WA, Widmann J, Yatsunenko T, Zaneveld J, Knight R: **QIIME**  
605 **allows analysis of high-throughput community sequencing data.** *Nature Methods* 2010,  
606 **7**:335–336.  
607
- 608 6. Li D, Liu C-M, Luo R, Sadakane K, Lam T-W: **MEGAHIT: an ultra-fast single-node**  
609 **solution for large and complex metagenomics assembly via succinct de Bruijn graph.**  
610 *Bioinformatics* 2015, **31**:1674–1676.  
611
- 612 7. Kang DD, Froula J, Egan R, Wang Z: **MetaBAT, an efficient tool for accurately**  
613 **reconstructing single genomes from complex microbial communities.** *PeerJ* 2015,  
614 **3**:e1165.  
615
- 616 8. Strous M, Kraft B, Bisdorf R, Tegetmeyer HE: **The Binning of Metagenomic Contigs**  
617 **for Microbial Physiology of Mixed Cultures.** *Frontiers in Microbiology* 2012, **3**.  
618
- 619 9. Wu Y-W, Tang Y-H, Tringe SG, Simmons BA, Singer SW: **MaxBin: an automated**  
620 **binning method to recover individual genomes from metagenomes using an**  
621 **expectation-maximization algorithm.** *Microbiome* 2014, **2**:1.

622

623

624 10. McMahon K: **“Metagenomics 2.0.”** *Environmental Microbiology Reports* 2015, **7**:38–

625 39.

626

627 11. Gokhale RS, Sankaranarayanan R, Mohanty D: **Versatility of polyketide synthases in**  
628 **generating metabolic diversity.** *Current Opinion in Structural Biology* 2007, **17**:736–743.

629

630 12. Koglin A, Walsh CT: **Structural insights into nonribosomal peptide enzymatic**  
631 **assembly lines.** *Natural Product Reports* 2009, **26**:987.

632

633 13. Lal R, Kumari R, Kaur H, Khanna R, Dhingra N, Tuteja D: **Regulation and**  
634 **manipulation of the gene clusters encoding type-I PKSs.** *Trends in biotechnology* 2000,  
635 **18**:264–274.

636

637

638 14. Cane DE, Walsh CT, Khosla C: **Harnessing the biosynthetic code: combinations,**  
639 **permutations, and mutations.** *Science* 1998, **282**:63–68.

640

641 15. Minowa Y, Araki M, Kanehisa M: **Comprehensive Analysis of Distinctive Polyke**  
642 **tide and Nonribosomal Peptide Structural Motifs Encoded in Microbial Genomes.**  
643 *Journal of Molecular Biology* 2007, **368**:1500–1517.

644

- 645 16. Sun W, Peng C, Zhao Y, Li Z: **Functional Gene-Guided Discovery of Type II**  
646 **Polyketides from Culturable Actinomycetes Associated with Soft Coral**  
647 ***Scleronephthya* sp.** *PLoS ONE* 2012, **7**:e42847.
- 648
- 649 17. Austin MB, Noel JP: **The chalcone synthase superfamily of type III polyketide**  
650 **synthases.** *Nat Prod Rep* 2003, **20**:79–110.
- 651
- 652
- 653 18. Fisch KM: **Biosynthesis of natural products by microbial iterative hybrid PKS–**  
654 **NRPS.** *RSC Advances* 2013, **3**:18228.
- 655
- 656 19. Mizuno CM, Kimes NE, López-Pérez M, Ausó E, Rodriguez-Valera F, Ghai R: **A**  
657 **Hybrid NRPS-PKS Gene Cluster Related to the Bleomycin Family of Antitumor**  
658 **Antibiotics in *Alteromonas macleodii* Strains.** *PLoS ONE* 2013, **8**:e76021.
- 659
- 660 20. Masschelein J, Mattheus W, Gao L-J, Moons P, Van Houdt R, Uytterhoeven B,  
661 Lamberigts C, Lescrinier E, Rozenski J, Herdewijn P, Aertsen A, Michiels C, Lavigne R: **A**  
662 **PKS/NRPS/FAS Hybrid Gene Cluster from *Serratia plymuthica* RVH1 Encoding the**  
663 **Biosynthesis of Three Broad Spectrum, Zeamine-Related Antibiotics.** *PLoS ONE* 2013,  
664 **8**:e54143.
- 665
- 666 21. Komaki H, Ichikawa N, Hosoyama A, Takahashi-Nakaguchi A, Matsuzawa T, Suzuki  
667 K, Fujita N, Gono T: **Genome based analysis of type-I polyketide synthase and**

668 **nonribosomal peptide synthetase gene clusters in seven strains of five representative**  
669 ***Nocardia* species. *BMC genomics* 2014, 15:1.**

670

671 22. Micallef ML, D'Agostino PM, Sharma D, Viswanathan R, Moffitt MC: **Genome**  
672 **mining for natural product biosynthetic gene clusters in the Subsection V**  
673 **cyanobacteria. *BMC Genomics* 2015, 16.**

674 23. Woodhouse JN, Fan L, Brown MV, Thomas T, Neilan BA: **Deep sequencing of non-**  
675 **ribosomal peptide synthetases and polyketide synthases from the microbiomes of**  
676 **Australian marine sponges. *The ISME journal* 2013, 7:1842–1851.**

677

678 24. Zothanpuia, Passari AK, Gupta VK, Singh BP: **Detection of antibiotic-resistant**  
679 **bacteria endowed with antimicrobial activity from a freshwater lake and their**  
680 **phylogenetic affiliation. *PeerJ* 2016, 4:e2103.**

681

682 25. Silva-Stenico ME, Silva CSP, Lorenzi AS, Shishido TK, Etchegaray A, Lira SP,  
683 Moraes LAB, Fiore MF: **Non-ribosomal peptides produced by Brazilian cyanobacterial**  
684 **isolates with antimicrobial activity. *Microbiological Research* 2011, 166:161–175.**

685

686

687 26. Selvin J, Sathiyarayanan G, Lipton AN, Al-Dhabi NA, Valan Arasu M, Kiran GS:  
688 **Ketide Synthase (KS) Domain Prediction and Analysis of Iterative Type II PKS Gene**  
689 **in Marine Sponge-Associated Actinobacteria Producing Biosurfactants and**  
690 **Antimicrobial Agents. *Front Microbiol* 2016, 7:2727–12.**



691

692

693 27. Ionescu D, Siebert C, Polerecky L, Munwes YY, Lott C, Häusler S, Bižić-Ionescu M,  
694 Quast C, Peplies J, Glöckner FO, Ramette A, Rödiger T, Dittmar T, Oren A, Geyer S, Stärk  
695 H-J, Sauter M, Licha T, Laronne JB, de Beer D: **Microbial and Chemical**  
696 **Characterization of Underwater Fresh Water Springs in the Dead Sea.** *PLoS ONE*  
697 2012, **7**:e38319–21.

698

699

700 28. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW: **CheckM: assessing**  
701 **the quality of microbial genomes recovered from isolates, single cells, and**  
702 **metagenomes.** *Genome research* 2015, **25**:1043–1055.

703

704

705 29. Weber T, Blin K, Duddela S, Krug D, Kim HU, Bruccoleri R, Lee SY, Fischbach MA,  
706 Müller R, Wohlleben W, Breitling R, Takano E, Medema MH: **antiSMASH 3.0—a**  
707 **comprehensive resource for the genome mining of biosynthetic gene clusters.** *Nucleic*  
708 *Acids Research* 2015, **43**:W237–W243.

709 30. Ziemert N, Podell S, Penn K, Badger JH, Allen E, Jensen PR: **The Natural Product**  
710 **Domain Seeker NaPDoS: A Phylogeny Based Bioinformatic Tool to Classify**  
711 **Secondary Metabolite Gene Diversity.** *PLoS ONE* 2012, **7**:e34064.

712

- 713 31. O'Leary NA, Wright MW, Brister JR, Ciufu S, Haddad D, McVeigh R, Rajput B,  
714 Robbertse B, Smith-White B, Ako-Adjei D, Astashyn A, Badretdin A, Bao Y, Blinkova O,  
715 Brover V, Chetvernin V, Choi J, Cox E, Ermolaeva O, Farrell CM, Goldfarb T, Gupta T,  
716 Haft D, Hatcher E, Hlavina W, Joardar VS, Kodali VK, Li W, Maglott D, Masterson P, et  
717 al.: **Reference sequence (RefSeq) database at NCBI: current status, taxonomic**  
718 **expansion, and functional annotation.** *Nucleic Acids Res* 2016, **44**:D733–D745.
- 719 32. Parks DH, Tyson GW, Hugenholtz P, Beiko RG. **STAMP: statistical analysis of**  
720 **taxonomic and functional profiles.** *Bioinformatics.* 2014;30(21):3123-3124.  
721 doi:10.1093/bioinformatics/btu494.
- 722
- 723
- 724 33. Di Lorenzo M, Stork M, Naka H, Tolmasky ME, Crosa JH: **Tandem heterocyclization**  
725 **domains in a nonribosomal peptide synthetase essential for siderophore biosynthesis**  
726 **in *Vibrio anguillarum*.** *BioMetals* 2008, **21**:635–648.
- 727
- 728
- 729 34. **National Research Council (US) Committee on Metagenomics: Challenges and**  
730 **Functional Applications. The New Science of Metagenomics: Revealing the Secrets of**  
731 **Our Microbial Planet.** Washington (DC): National Academies Press (US); 2007. 5, Data  
732 Management and Bioinformatics Challenges of Metagenomics.
- 733
- 734 35. Gilbert JA, Meyer F, Bailey MJ: **The future of microbial metagenomics (or is**  
735 **ignorance bliss?).** *ISME Journal-International Society for Microbial Ecology* 2011, **5**:777.

736

737 36. Teeling H, Glockner FO: **Current opportunities and challenges in microbial**  
738 **metagenome analysis--a bioinformatic perspective.** *Briefings in Bioinformatics* 2012,  
739 **13:728–742.**

740

741 37. Albertsen M, Hugenholtz P, Skarshewski A, Nielsen KARL, Tyson GW, Nielsen PH:  
742 **Genome sequences of rare, uncultured bacteria obtained by differential coverage**  
743 **binning of multiple metagenomes.** *Nature Biotechnology* 2013, **31:533–538.**

744

745 38. Sangwan N, Xia F, Gilbert JA: **Recovering complete and draft population genomes**  
746 **from metagenome datasets.** *Microbiome* 2016, **4.**

747

748 39. Vanwonterghem I, Jensen PD, Rabaey K, Tyson GW: **Genome-centric resolution of**  
749 **microbial diversity, metabolism and interactions in anaerobic digestion: Genome-**  
750 **centric resolution through deep metagenomics.** *Environmental Microbiology* 2016,  
751 **18:3144–3158.**

752

753 41. Amos GCA, Borsetto C, Laskaris P, Krsek M, Berry AE, Newsham KK, Calvo-Bado L,  
754 Pearce DA, Vallin C, Wellington EMH: **Designing and Implementing an Assay for the**  
755 **Detection of Rare and Divergent NRPS and PKS Clones in European, Antarctic and**  
756 **Cuban Soils.** *PLOS ONE* 2015, **10:e0138327.**

757

758 41. Foerstner KU, Doerks T, Creevey CJ, Doerks A, Bork P: **A Computational Screen for**  
759 **Type I Polyketide Synthases in Metagenomics Shotgun Data.** *PLoS ONE* 2008, **3:e3515.**

760

761 42. Cuadrat R, Cury J, Dávila A: **Metagenomic Analysis of Upwelling-Affected**  
762 **Brazilian Coastal Seawater Reveals Sequence Domains of Type I PKS and Modular**  
763 **NRPS**. *International Journal of Molecular Sciences* 2015, **16**:28285–28295.

764

765 43. Yamada Y, Kuzuyama T, Komatsu M, Shin-ya K, Omura S, Cane DE, Ikeda H:  
766 **Terpene synthases are widely distributed in bacteria**. *Proceedings of the National*  
767 *Academy of Sciences* 2015, **112**:857–862.

768

769 44. Yang S-C, Lin C-H, Sung CT, Fang J-Y: **Antibacterial activities of bacteriocins:**  
770 **application in foods and pharmaceuticals**. *Frontiers in Microbiology* 2014, **5**.

771

772 45. Cotter PD, Ross RP, Hill C: **Bacteriocins — a viable alternative to antibiotics?**  
773 *Nature Reviews Microbiology* 2012, **11**:95–105.

774

775 46. Zheng J, Gänzle MG, Lin XB, Ruan L, Sun M: **Diversity and dynamics of**  
776 **bacteriocins from human microbiome: Bacteriocins of human microbiome**.  
777 *Environmental Microbiology* 2015, **17**:2133–2143.

778

779 47. Illeghems K, Weckx S, De Vuyst L: **Applying meta-pathway analyses through**  
780 **metagenomics to identify the functional properties of the major bacterial communities**  
781 **of a single spontaneous cocoa bean fermentation process sample**. *Food Microbiology*  
782 2015, **50**:54–63.

783

784 48. Więckowicz M, Schmidt M, Sip A, Grajek W: **Development of a PCR-based assay**  
785 **for rapid detection of class IIa bacteriocin genes: Detection of class IIa bacteriocins.**  
786 *Letters in Applied Microbiology* 2011, **52**:281–289.

787

788 49. Escobar-Zepeda A, Sanchez-Flores A, Quirasco Baruch M: **Metagenomic analysis of a**  
789 **Mexican ripened cheese reveals a unique complex microbiota.** *Food Microbiology*  
790 2016, **57**:116–127.

791

792 50. Rausch C, Hoof I, Weber T, Wohlleben W, Huson DH: **Phylogenetic analysis of**  
793 **condensation domains in NRPS sheds light on their functional evolution.** *BMC*  
794 *Evolutionary Biology* 2007, **7**:78.

795

796 51. Woodhouse JN, Fan L, Brown MV, Thomas T, Neilan BA: **Deep sequencing of non-**  
797 **ribosomal peptide synthetases and polyketide synthases from the microbiomes of**  
798 **Australian marine sponges.** *The ISME Journal* 2013, **7**:1842–1851.

799

800 52. Esmael Q, Pupin M, Kieu NP, Chataigné G, Béchet M, Deravel J, Krier F, Höfte M,  
801 Jacques P, Leclère V: ***Burkholderia* genome mining for nonribosomal peptide**  
802 **synthetases reveals a great potential for novel siderophores and lipopeptides synthesis.**  
803 *MicrobiologyOpen* 2016, **5**:512–526.

804

805 53. Van Der Voort M, Meijer HJG, Schmidt Y, Watrous J, Dekkers E, Mendes R,  
806 Dorrestein PC, Gross H, Raaijmakers JM: **Genome mining and metabolic profiling of the**

807 **rhizosphere bacterium *Pseudomonas* sp. SH-C52 for antimicrobial compounds.**

808 *Frontiers in Microbiology* 2015, **6**.

809

810 54. Pan H-Q, Hu J-C: **Draft genome sequence of the novel strain *Pseudomonas* sp.**

811 **10B238 with potential ability to produce antibiotics from deep-sea sediment.** *Mar*

812 *Genomics* 2015, **23**:55–57.

813

814 55. Scholz-Schroeder BK, Soule JD, Gross DC: **The *sypA*, *sypB*, and *sypC* synthetase**

815 **genes encode twenty-two modules involved in the nonribosomal peptide synthesis of**

816 **syringopeptin by *Pseudomonas syringae* pv. *syringae* B301D.** *Molecular Plant-Microbe*

817 *Interactions* 2003, **16**:271–280.

818

819 56. Stock SD, Hama H, Radding JA, Young DA, Takemoto JY: **Syringomycin E**

820 **inhibition of *Saccharomyces cerevisiae*: requirement for biosynthesis of sphingolipids**

821 **with very-long-chain fatty acids and mannose-and phosphoinositol-containing head**

822 **groups.** *Antimicrobial agents and chemotherapy* 2000, **44**:1174–1180.

823

824 57. Dawson RM: **The toxicology of microcystins.** *Toxicon* 1998, **36**:953–962.

825

826 58. Bouhaddada R, Nelieu S, Nasri H, Delarue G, Bouaicha N: **High diversity of**

827 **microcystins in a *Microcystis* bloom from an Algerian lake.** *Environ Pollut* 2016,

828 **216**:836–844.

829

830 59. Brandel J, Humbert N, Elhabiri M, Schalk IJ, Mislin GLA, Albrecht-Gary A-M:  
831 **Pyochelin, a siderophore of *Pseudomonas aeruginosa*: Physicochemical**  
832 **characterization of the iron(iii), copper(ii) and zinc(ii) complexes.** *Dalton Transactions*  
833 2012, **41**:2820.

834

835 60. Schoner TA, Gassel S, Osawa A, Tobias NJ, Okuno Y, Sakakibara Y, Shindo K,  
836 Sandmann G, Bode HB: **Aryl Polyenes, a Highly Abundant Class of Bacterial Natural**  
837 **Products, Are Functionally Related to Antioxidative Carotenoids.** *ChemBiochem* 2016,  
838 **17**:247–253.

839

840 61. Jones GB, Fouad FS: **Designed enediyne antitumor agents.** *Curr Pharm Des* 2002,  
841 **8**:2415–2440.

842

843 62. Ahlert J, Shepard E, Lomovskaya N, Zazopoulos E, Staffa A, Bachmann BO, Huang K,  
844 Fonstein L, Czisny A, Whitwam RE, Farnet CM, Thorson JS: **The calicheamicin gene**  
845 **cluster and its iterative type I enediyne PKS.** *Science* 2002, **297**:1173–1176.

846

847

848 63 - Fedrizzi T, Meehan CJ, Grottola A, Giacobazzi E, Serpini GF, Tagliazucchi S, Fabio  
849 A, Bettua C, Bertorelli R, De Sanctis V, Rumpianesi F, Pecorari M, Jousson O, Tortoli E,  
850 Segata N: **Genomic characterization of Nontuberculous Mycobacteria.** *Nature*  
851 *Publishing Group* 2017:1–14.

852

- 853 64. van der Wielen PWJJ, Heijnen L, van der Kooij D: **Pyrosequence Analysis of the**  
854 **hsp65 Genes of Nontuberculous Mycobacterium Communities in Unchlorinated**  
855 **Drinking Water in the Netherlands.** *Applied and Environmental Microbiology* 2013,  
856 **79**:6160–6166.
- 857
- 858 65. Bakula Z, Safianowska A, Nowacka-Mazurek M, Bielecki J, Jagielski T: **Short**  
859 **Communication: Subtyping of Mycobacterium kansasii by PCR-Restriction Enzyme**  
860 **Analysis of the hsp65 Gene.** *BioMed Research International* 2013, **2013**:1–4.
- 861
- 862 66. Quadri LEN: **Biosynthesis of mycobacterial lipids by polyketide synthases and**  
863 **beyond.** *Critical Reviews in Biochemistry and Molecular Biology* 2014, **49**:179–211.
- 864
- 865 67. Brown NM, Mueller RS, Shepardson JW, Landry ZC, Morré JT, Maier CS, Hardy FJ,  
866 Dreher TW: **Structural and functional analysis of the finished genome of the recently**  
867 **isolated toxic Anabaena sp. WA102.** *BMC Genomics* 2016:1–18.
- 868
- 869 68. Li X, Dreher TW, Li R: **An overview of diversity, occurrence, genetics and toxin**  
870 **production of bloom-forming Dolichospermum (Anabaena) species.** *Harmful Algae*  
871 2016, **54**:54–68.
- 872



873 69. Calteau A, Fewer DP, Latifi A, Coursin T, Laurent T, Jokela J, Kerfeld CA, Sivonen K,  
874 Piel J, Gugger M: **Phylum-wide comparative genomics unravel the diversity of**  
875 **secondary metabolism in Cyanobacteria.** *BMC Genomics* 2014, **15**:977–14.

876

877

878 70. Carmichael WW, Biggs DF, Gorham PR. **Toxicology and pharmacological action**  
879 **of *Anabaena flos-aquae* toxin.** *Science* 1975, ;187:542-544

880

881 71. Méjean A, Paci G, Gautier V, Ploux O: **Biosynthesis of anatoxin-a and analogues**  
882 **(anatoxins) in cyanobacteria.** *Toxicon* 2014, **91**(C):15–22.

883

884 72. **Structural organization of microcystin biosynthesis in *Microcystis aeruginosa***  
885 **PCC7806: an integrated peptide<sup>^</sup>polyketide synthetase system.** 2000:1–12.

886

887 73. Rouhiainen L, Vakkilainen T, Siemer BL, Buikema W, Haselkorn R, Sivonen K: **Genes**  
888 **Coding for Hepatotoxic Heptapeptides (Microcystins) in the Cyanobacterium**  
889 ***Anabaena* Strain 90.** *Applied and Environmental Microbiology* 2004, **70**:686–692.

890

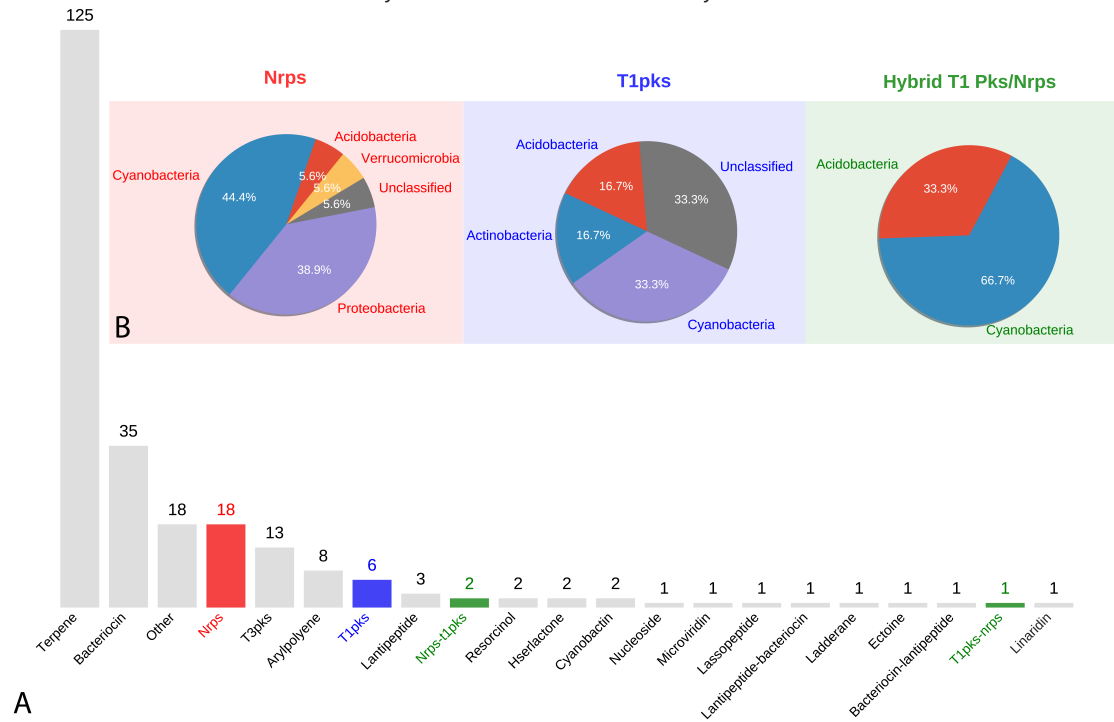
891 74. Rastogi RP, Madamwar D, Incharoensakdi A: **Bloom Dynamics of Cyanobacteria**  
892 **and Their Toxins: Environmental Health Impacts and Mitigation Strategies.** *Front*  
893 *Microbiol* 2015, **6**:223–22.

894

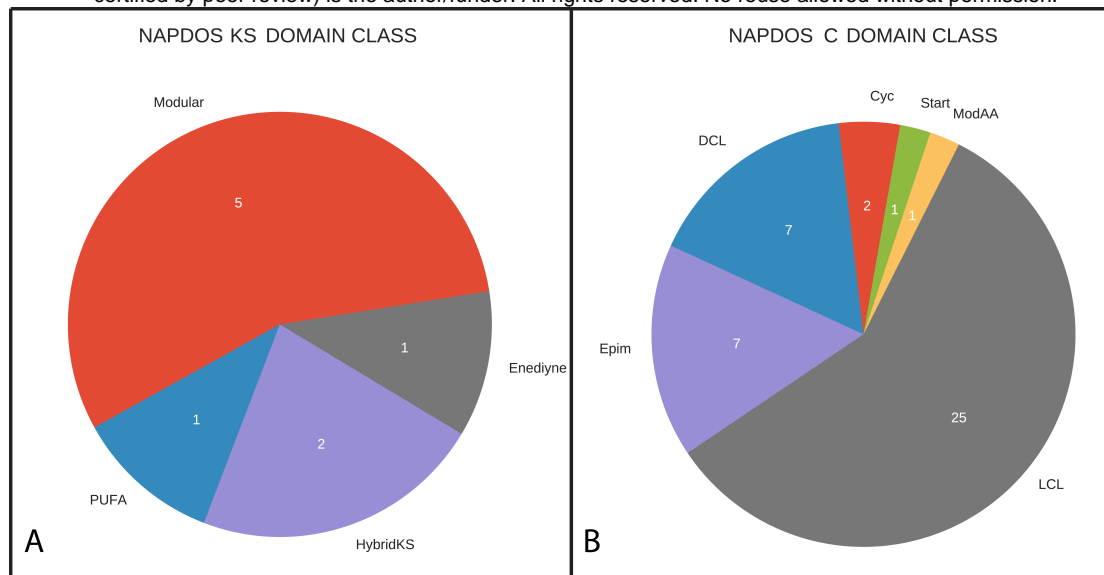
- 895 75. Viaggiu E, Melchiorre S, Volpi F, Di Corcia A, Mancini R, Garibaldi L, Crichigno G,  
896 Bruno M: **Anatoxin-a toxin in the cyanobacterium *Planktothrix rubescens* from a**  
897 **fishing pond in northern Italy.** *Environ Toxicol* 2004, **19**:191–197.
- 898
- 899 76. Ballot A, Fastner J, Lentz M, Wiedner C: **First report of anatoxin-a-producing**  
900 **cyanobacterium *Aphanizomenon issatschenkoi* in northeastern Germany.** *Toxicon*  
901 2010, **56**:964–971.
- 902
- 903 77. Dadheech P, Selmeczy G, Vasas G, Padišák J, Arp W, Tapolczai K, Casper P, Krienitz  
904 L: **Presence of Potential Toxin-Producing Cyanobacteria in an Oligo-Mesotrophic**  
905 **Lake in Baltic Lake District, Germany: An Ecological, Genetic and Toxicological**  
906 **Survey.** *Toxins* 2014, **6**:2912–2931.
- 907
- 908 78 . Weissman KJ: **The structural biology of biosynthetic megaenzymes.** *Nat Chem Biol*  
909 2015, **11**:660–670.
- 910 79. Frank JA, Pan Y, Tooming-Klunderud A, Eijsink VGH, McHardy AC, Nederbragt AJ,  
911 Pope PB: **Improved metagenome assemblies and taxonomic binning using long-read**  
912 **circular consensus sequence data.** *Scientific Reports* 2016, **6**:25373.
- 913
- 914 80. Hutchison CA, Chuang R-Y, Noskov VN, Assad-Garcia N, Deerinck TJ, Ellisman MH,  
915 Gill J, Kannan K, Karas BJ, Ma L, others: **Design and synthesis of a minimal bacterial**  
916 **genome.** *Science* 2016, **351**:aad6253.



Secondary metabolites clusters obtained by Anti-SMASH



**Figure 1A: Abundance of secondary metabolite cluster types obtained with Anti-SMASH in the recovered 288 bins (environmental genomes). B: Taxonomical classification of bins (Phyla) in which NRPS, PKS and Hybrid PKS/NRPS clusters were found. Red bar and pie: NRPS; blue bar and pie: Type I PKS; green bars and pie: Hybrid clusters (NRPS-PKS and PKS-NRPS).**



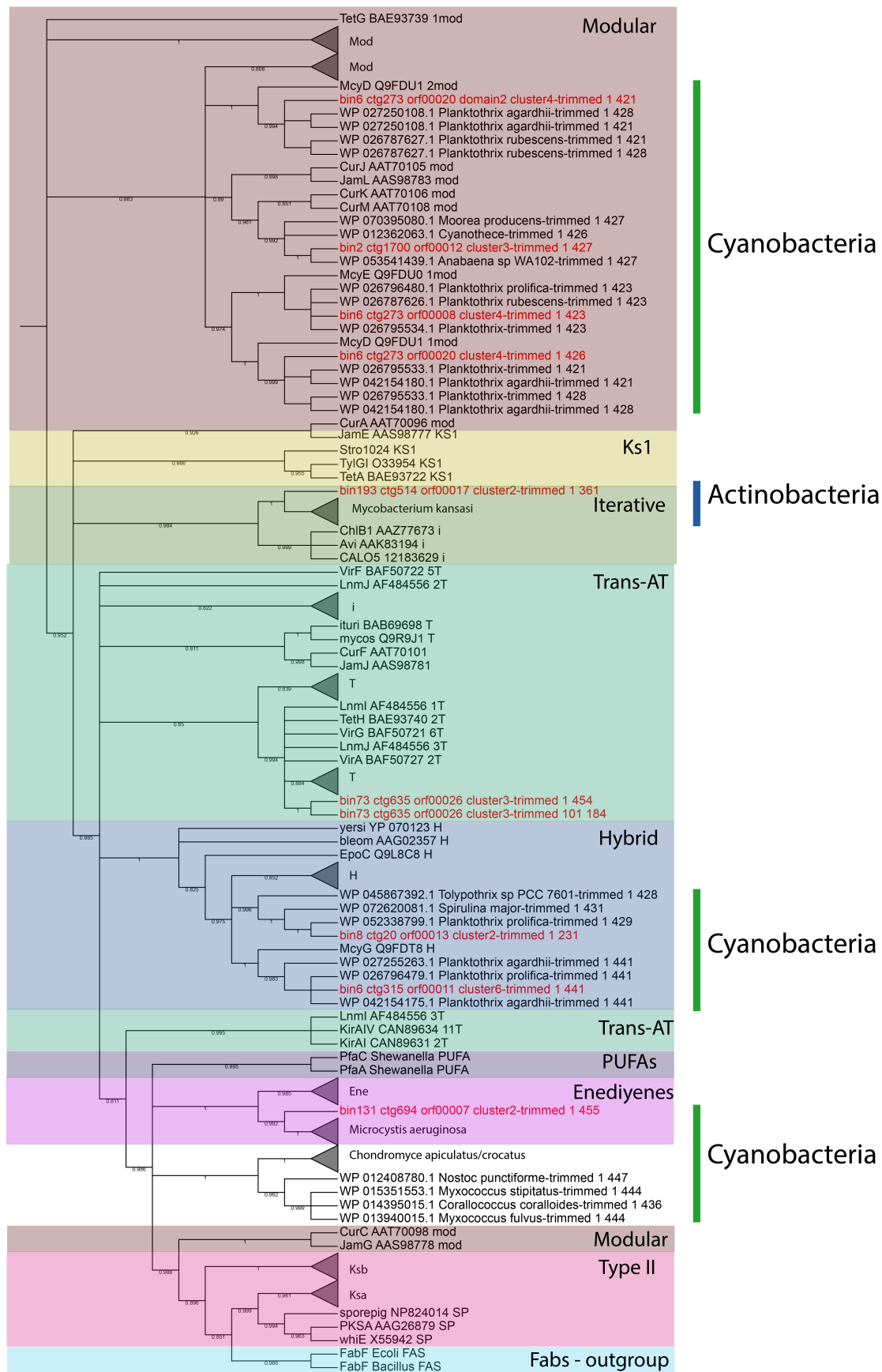
**Figure 2A: NAPDOS classification of the NRPS KS domain.** Modular: possess a multidomain architecture consisting of multiple sets of modules; hybridKS: are biosynthetic assembly lines that include both PKS and NRPS components; PUFA: Polyunsaturated fatty acids (PUFAs) are long chain fatty acids containing more than one double bond, including omega-3-and omega-6- fatty acids; Ene diyne: a family of biologically active natural products. The Ene diyne core consists of two acetylenic groups conjugated to a double bond or an incipient double bond within a nine- or ten-membered ring. **2B: NAPDOS classification of NRPS C domain.** Cyc: cyclization domains catalyze both peptide bond formation and subsequent cyclization of cysteine, serine or threonine residues; DCL: link an L-amino acid to a growing peptide ending with a D-amino acid; Epim: epimerization domains change the chirality of the last amino acid in the chain from L- to D- amino acid; LCL: catalyze formation of a peptide bond between two L-amino acids; modAA: appear to be involved in the modification of the incorporated amino acid; Start: first module of a Non-ribosomal peptide synthase (NRPS).



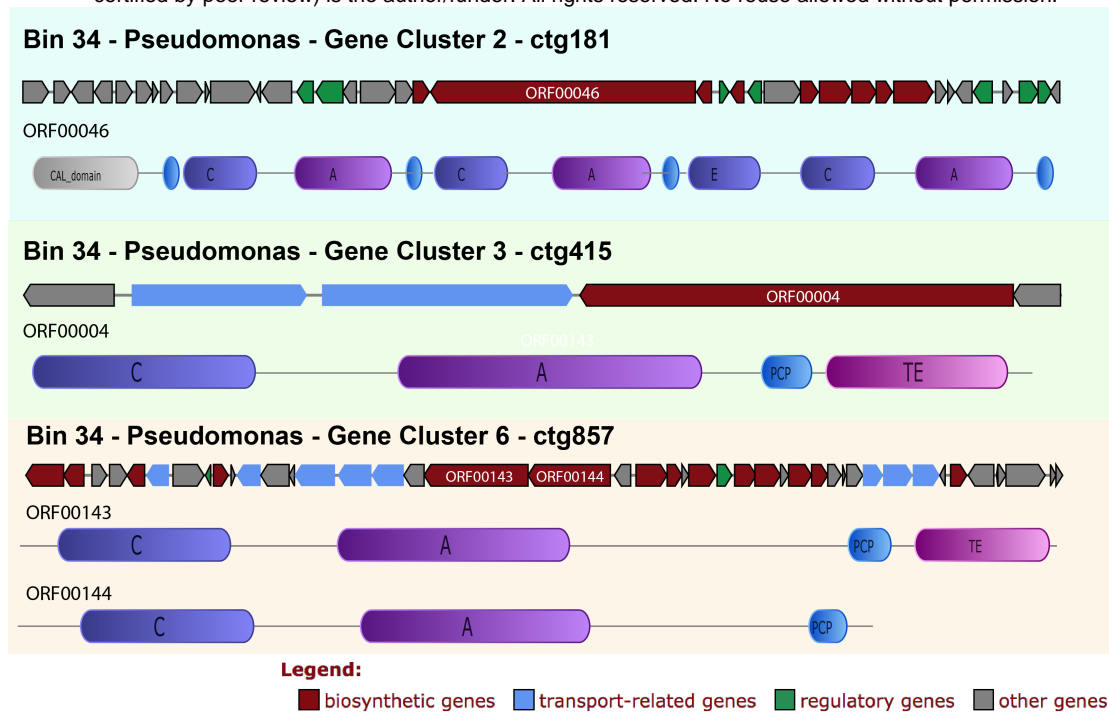


**phylogenetic tree of C domains (environmental domains, the top 3 blast results on RefSeq and the NAPDOS reference sequences).** The shadow colours represent the domain classifications (LCL, CYC, Start domains, EPIM, ModAA, Dual, and DCL). The sidebars represent phyla (Proteobacteria, Cyanobacteria, Firmicutes, Actinobacteria, Verrucomicrobia). All sequences from environmental bins are in red.

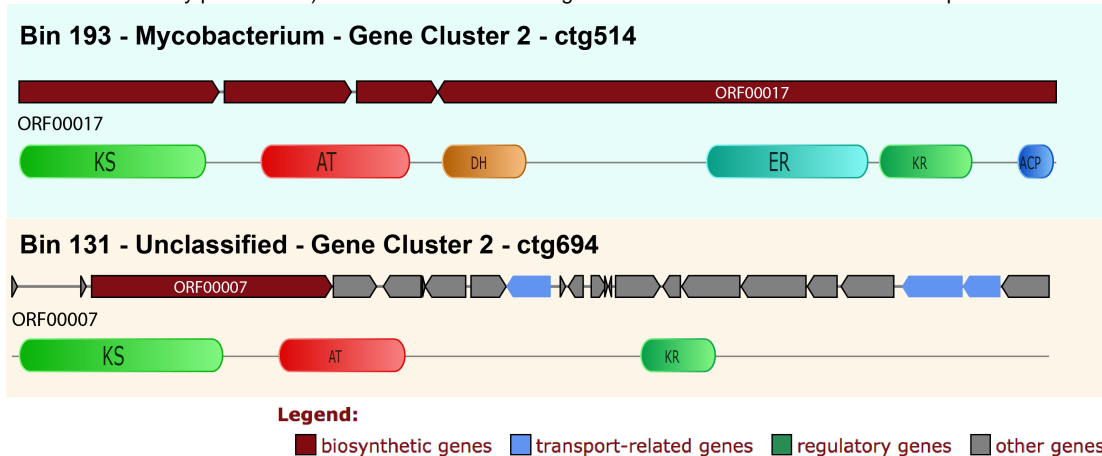




**Figure 4: NAPDOS tree of KS domains (environmental domains, the top 3 blast results on RefSeq and the NAPDOS reference sequences).** The shadow colours represent the domain classifications (Modular, KS1, Iterative, Trans-AT, Hybrid, PUFA, Enediyenes, Type II and Fabs – Fatty acid synthase). The sidebars represent phyla (Cyanobacteria and Actinobacteria). All the sequences from environmental bins are in red.



**Figure 5: Bin 34, NRPS clusters detailed annotation and synteny.** The synteny of the clusters with a functional classification for each ORF is given. In addition, for the NPRS biosynthetic ORFS the domain annotations are given. CAL: Co-enzyme A ligase domain, C: condensation, A: adenylation, E: epimerization, TE: Termination, KR: Ketoreductase domain, and ECH: Enoyl-CoA hydratase



**Figure 6: Bin 193 and 131 type I PKS clusters detailed annotation and synteny.** It is possible to see the synteny of the cluster with the functional classification for each ORF. In addition, for the PKS biosynthetic ORF the domain-specific annotations can be seen. KS: keto-synthase, AT: acyltransferase, KR: ketoreductase, E: epimerization, DH: dehydratase, and ER: enoylreductase.