

1 iPRESTO: automated discovery of
2 biosynthetic sub-clusters linked to specific
3 natural product substructures
4

5 Joris J.R. Louwen¹, Satria A. Kautsar¹, Sven van der Burg², Marnix H. Medema^{1*}, Justin J.J. van der
6 Hooft^{1,3*}

7 1. Bioinformatics Group, Wageningen University, Wageningen, the Netherlands

8 2. Netherlands eScience Center, Amsterdam, the Netherlands

9 3. Department of Biochemistry, University of Johannesburg, Johannesburg, South Africa

10 * Corresponding authors

11 E-mail: marnix.medema@wur.nl, justin.vanderhooft@wur.nl

12 Abstract

13 Microbial specialised metabolism is full of valuable natural products that are applied clinically,
14 agriculturally, and industrially. The genes that encode their biosynthesis are often physically clustered on
15 the genome in biosynthetic gene clusters (BGCs). Many BGCs consist of multiple groups of co-evolving
16 genes called sub-clusters that are responsible for the biosynthesis of a specific chemical moiety in a
17 natural product. Sub-clusters therefore provide an important link between the structures of a natural
18 product and its BGC, which can be leveraged for predicting natural product structures from sequence, as
19 well as for linking chemical structures and metabolomics-derived mass features to BGCs.

20 While some initial computational methodologies have been devised for sub-cluster detection, current
21 approaches are not scalable, have only been run on small and outdated datasets, or produce an
22 impractically large number of possible sub-clusters to mine through.

23 Here, we constructed a scalable method for unsupervised sub-cluster detection, called iPRESTO, based
24 on topic modelling and statistical analysis of co-occurrence patterns of enzyme-coding protein families.
25 iPRESTO was used to mine sub-clusters across 150,000 prokaryotic BGCs from antiSMASH-DB. After
26 annotating a fraction of the resulting sub-cluster families, we could predict a substructure for 16% of the
27 antiSMASH-DB BGCs. Additionally, our method was able to confirm 83% of the experimentally
28 characterised sub-clusters in MIBiG reference BGCs. Based on iPRESTO-detected sub-clusters, we could
29 correctly identify the BGCs for xenorhabdin and salbostatin biosynthesis (which had not yet been
30 annotated in BGC databases), as well as propose a candidate BGC for akashin biosynthesis. Additionally,
31 we show for a collection of 145 actinobacteria how substructures can aid in linking BGCs to molecules by
32 correlating iPRESTO-detected sub-clusters to MS/MS-derived Mass2Motifs substructure patterns.

33 This work paves the way for deeper functional and structural annotation of microbial BGCs by improved
34 linking of orphan molecules to their cognate gene clusters, thus facilitating accelerated natural product
35 discovery.

36 Author summary

37 In this work, we introduce iPRESTO, a tool for scalable unsupervised sub-cluster detection in biosynthetic
38 gene clusters. This detection is important because these biosynthetic hotspots encode many products
39 useful for humanity, such as antibiotics, antitumor agents, or herbicides. Recent technological
40 developments have made identification of biosynthetic loci in genomes straightforward. Yet, methods to
41 connect these inferred biosynthetic genes to the final chemical structures of their cognate metabolites
42 are largely lacking. Being able to reliably predict parts of the final product would constitute a real step
43 forward in natural product genome mining. Therefore, we focussed on constructing a tool to
44 systematically detect and annotate small regions called sub-clusters, which code for the biosynthesis of
45 substructures in the final product, across all genomically inferred biosynthetic diversity. iPRESTO makes
46 it possible to query unknown biosynthetic regions and infer which substructures are present in their
47 metabolic products. This will facilitate more effective prioritization of chemical novelty, as well as linking
48 activities from bioassays and microbiome-associated phenotypes to the metabolites responsible for them.

49 Introduction

50 A considerable part of bacterial metabolism is dedicated to the biosynthesis of specialised metabolites.
51 These natural products (NPs) have many uses as pharmaceuticals, crop protection agents, and
52 ingredients for foods and cosmetics [1, 2]. NPs consist of a spectrum of different chemical classes, which
53 are often highly complex in structure [3]. Intriguingly, the genes necessary for the biosynthesis of NPs
54 cluster together physically in biosynthetic gene clusters (BGCs) [4]. The search and discovery of new
55 BGCs accelerates identification of new NPs, which is especially important in the field of antibiotics, as
56 antibiotic-resistant bacteria are becoming increasingly prevalent [5].

57 Due to the growing availability of genomic data, genome mining approaches have become more and
58 more useful for NP discovery. Currently, multiple algorithms exist that mine bacterial genomes for
59 putative BGCs, such as antiSMASH, ClusterFinder and PRISM [6-8]. These methods have provided a
60 better understanding of BGC diversity and the evolutionary mechanisms that govern BGC diversity.

61 Many classes of BGCs display a modular architecture [4]. As such, a BGC can be divided into multiple
62 modules or sub-clusters, where each sub-cluster is a group of co-evolving genes responsible for the
63 biosynthesis of a specific chemical moiety in the NP [4, 9, 10]. Sub-clusters therefore provide a direct
64 link between the substructures of an NP and its BGC. This makes information about sub-clusters and the
65 substructures they synthesise highly valuable for genome-based structure prediction, which would be a
66 great asset for tools like antiSMASH. Apart from enhancing structural predictions for existing BGC
67 classes, sub-cluster knowledge would facilitate predicting novel (partial) structures of currently
68 unclassified BGCs, such as the thousands of unclassified BGCs with yet unknown products in the
69 antiSMASH-DB [11].

70 Additionally, BGC modularity poses a great opportunity to connect metabolomics experiments to sub-
71 cluster data. Chemical moieties identified from fragments in mass spectrometry (MS) data could be
72 linked to sub-clusters responsible for their synthesis, as part of MS-guided genome mining strategies
73 [10, 12, 13]. Recent advances in substructure modelling [14] may aid such co-occurrence-based
74 metabologenomic approaches [15] by automating the identification of substructures from MS/MS data.

75 Recently, Del Carratore et al. [10] introduced an initial method for the detection of sub-clusters in BGCs.
76 By constructing Clusters of Orthologous Groups (COGs) and by using a statistical approach to group co-
77 occurring COGs in sub-clusters, they were able to detect several experimentally characterised sub-
78 clusters, as well as to discover novel ones. However, COG construction is not very scalable due to the all-
79 vs-all BLAST calculation required. As a result, their analysis was performed on a relatively small dataset
80 that is by now almost a decade old, and the chosen approach is hard to scale up to the massive amounts
81 of genomic data that have become available in recent years. Additionally, the proposed statistical
82 approach greatly overestimates the numbers of sub-clusters. This is due to the presence of redundant
83 BGCs, which leads to artificial sub-clusters spanning entire BGCs, and caused by the inherently nested
84 structure of the sub-clusters, where smaller, less specific sub-clusters are contained in larger, more
85 specific sub-clusters.

86 Here, we propose an improved scalable method for unsupervised sub-cluster detection which we called
87 the integrated Prediction and Rigorous Exploration of biosynthetic Sub-clusters Tool (iPRESTO). iPRESTO
88 is scalable to large datasets and takes phylogenetic bias into account by filtering the input in a more
89 advanced way. To detect sub-clusters, iPRESTO uses a statistical approach (PRESTO-STAT) as well as a
90 topic modelling algorithm (PRESTO-TOP). As a data source, we used the antiSMASH-DB, which is one of
91 the largest collections of BGCs that currently exists, and which has been scrutinized for underlying
92 genome assembly quality [11]; it contains over 150,000 BGCs from almost 25,000 bacterial species
93 selected to reduce taxonomical bias. These numbers represent a considerable improvement in
94 comparison with the previous method as it contains over ten times as many BGCs, while being less
95 redundant. After applying iPRESTO on this large collection of BGCs, we were able to annotate 45 sub-
96 cluster motifs based on occurrences in known BGCs from the MIBiG reference BGC database [16]. Using
97 these annotated sub-cluster motifs, we zoomed in on relevant sub-clusters, and showed direct usefulness
98 of our method by correctly predicting the BGCs for xenorhabdin and salbostatin biosynthesis (which have
99 been published but were missing from BGC databases) and identifying a candidate BGC for akashin
100 biosynthesis. Finally, as a starting point for the automated connection of BGCs to their NPs, we were able
101 to systematically link sub-clusters to substructures by using a metabologenomic correlation method in a
102 paired-genome-metabolome dataset of 145 actinobacteria.

103 Results & Discussion

104 Overview of iPRESTO

105 iPRESTO prepares each BGC for sub-cluster detection by tokenising each gene in a BGC as a combination
106 of Pfam domains (**Fig 1** and S1 Fig). If a pair of proteins share the same Pfam domains, this provides an
107 effective indication of (at least distant) sequence similarity, while Pfam detection is highly scalable. As
108 Pfams are quite broad sequence models (which would be a major disadvantage compared to using
109 COGs), we increased the resolution by splitting the 112 most abundant biosynthetic Pfams into a number
110 of subPfams, akin to the implementation in BiG-SLICE [17]. Each subPfams constitutes a narrower
111 domain model that covers a subset of a Pfam's sequence space. We only considered biosynthetic
112 domains (see Methods) to limit the search space and focus solely on finding biosynthetic sub-clusters.
113 With a graph-based filtering step, redundant BGCs are removed, after which iPRESTO detects sub-
114 clusters using PRESTO-STAT and PRESTO-TOP. PRESTO-STAT is based on the previously published
115 statistical method, which we expanded by partly removing nested sub-clusters, collapsing similar sub-
116 clusters into families, and joining similar families into clans.

117 **Fig 1. Outline of the iPRESTO workflow for the detection of sub-clusters.** All genes in BGCs are
118 converted into strings of Pfam domains, after which redundant BGCs are filtered out based on an Adjacency
119 Index of domains. Sub-clusters are detected using two methods: PRESTO-TOP (TOP) and PRESTO-STAT
120 (STAT). BGCs from the MIBiG database are used to annotate putative sub-clusters with sub-structures. These
121 annotations are used to predict sub-structures in unknown BGCs.

122 To enrich the discovery of sub-clusters with a method that does not produce nested sub-clusters, we
123 introduce PRESTO-TOP as a novel approach for sub-cluster detection. PRESTO-TOP is built on Latent
124 Dirichlet Allocation (LDA), which is used to model topics in text documents. LDA has already been used
125 successfully in genome and metabolome data analysis before [14, 18]. In the case of PRESTO-TOP, a
126 text document is a BGC, a word is a gene represented as a domain combination, and a topic can be
127 thought of as a sub-cluster motif. This highlights the use of PRESTO-TOP for sub-cluster detection, as we
128 assume that a BGC is a combination of multiple different sub-clusters, which consist of co-evolving genes
129 that co-occur in multiple BGCs. Another benefit of PRESTO-TOP is that a topic or sub-cluster motif will
130 usually consist of a set of core genes that encode the enzymes to synthesise the base of a substructure,
131 while various combinations of additional modifying genes can be found in PRESTO-STAT-detected
132 (nested) sub-clusters. In this way, the two iPRESTO methods can jointly capture substructure diversity,
133 by identifying the sub-cluster cores as well as their variants.

134 The resulting sub-clusters of both methods can be annotated with substructures and subsequently be
135 used to predict sub-structures in BGCs. iPRESTO is readily usable for anyone who wants to detect sub-
136 clusters in their own datasets, both by creating new sub-cluster models and by querying BGCs to the
137 collection of sub-clusters we detected in this study. iPRESTO can handle large amounts of BGCs:
138 tokenising and reducing redundancy in the 150,000 BGCs in the antiSMASH-DB dataset took around 48
139 hours each using 32 CPU cores on an Intel Xeon CPU E5-2670 v3. Detecting sub-clusters with PRESTO-
140 STAT and PRESTO-TOP completed in 24 and 8 hours, respectively. iPRESTO can query around 20 BGCs
141 per minute to the sub-clusters detected in this study including the tokenisation steps. iPRESTO also
142 contains a visualisation module to visualise the results of querying a BGC to PRESTO-STAT or PRESTO-
143 TOP output (see S2 Fig for an example of querying the rifamycin BGC).

144 PRESTO-STAT improves comprehensibility of existing statistical method

145 We applied iPRESTO to the antiSMASH-DB v2 dataset, which contained, after pre-processing, 60,028
146 BGCs with 10,539 domain combinations (Table A in S1 Text). Using the PRESTO-STAT method, we found
147 108,085 sub-clusters in the dataset. Over 80% of the statistical sub-clusters contain fewer than ten
148 genes, and 17% of the sub-clusters occur in more than 10 BGCs (S3 Fig). When comparing PRESTO-
149 STAT with the previous version of the method by Del Carratore et al. [10], we observed that PRESTO-
150 STAT produces on average roughly two sub-clusters per BGC, while the previous method resulted in
151 roughly fourteen sub-clusters per BGC. This indicates that we end up with fewer nested sub-cluster
152 structures, which is most likely due to our extended redundancy filtering that removed almost half of the

153 dataset (Table A in S1 Text). Even so, nested structures are still very apparent in our results (S2 Fig).
154 For example, thousands of BGCs have more than 30 sub-clusters, many of which overlap with one
155 another (S4A Fig). Not only do the nested structures inflate the results, but they also have the additional
156 disadvantage that their presence makes it harder to connect BGCs with similar yet distinct sub-clusters.

157 To facilitate the sub-cluster analysis, we connected related sub-clusters by clustering the statistical sub-
158 clusters into 10,000 sub-cluster families (SCFs) and the SCFs into 2,000 sub-cluster clans (SCCs). We
159 used K-means clustering and represented the statistical sub-clusters as a presence/absence matrix of the
160 tokenised genes. Although some SCCs grouped seemingly unrelated sub-clusters together that share
161 only one gene (based on having the same Pfam domain content), most SCCs (81%) provided groups of
162 related sub-clusters, sharing at least three genes.

163 Apart from the nested structures, the statistical method produces many sub-clusters of which only a
164 fraction probably provides meaningful information. This is illustrated by the fact that the PRESTO-STAT
165 results can be very noisy: in a group of BGCs sharing multiple sub-clusters, all combinations of these
166 shared sub-clusters could form new sub-clusters, which happens frequently (S2 Fig). Additionally, it is
167 rather difficult to query a BGC using the statistical sub-clusters while allowing inexact matching, as this
168 would quickly become very time consuming.

169 PRESTO-TOP identifies characterised and novel sub-clusters

170 The drawbacks of PRESTO-STAT present a clear reason as to why we chose to also develop PRESTO-TOP,
171 which can find multiple sub-clusters in a BGC and is able to capture sub-cluster diversity within sub-
172 cluster motifs. Furthermore, LDA, upon which PRESTO-TOP is built, allows for a scalable way to build and
173 query sub-cluster motifs.

174 We used PRESTO-TOP to train and query a model on the antiSMASH-DB dataset with 1,000 sub-cluster
175 motifs. In the Methods section, we provide information on (hyper)parameters used and reasoning for the
176 chosen settings. Over 80% of the BGCs in the dataset contained at least one sub-cluster motif (S4B Fig).
177 To assess the quality of the sub-cluster motifs, we visualised all sub-clusters individually, where each
178 sub-cluster is a group of genes matching against a sub-cluster motif (Fig 2A). For a sub-cluster to be
179 interesting, we would expect its size to be between 2-12 genes, as experimentally characterised sub-
180 clusters fall in this range [19]. Upon checking our results, most sub-clusters that were present across a
181 considerable number of BGCs were within this expected size range (Fig 2A), while some sub-clusters
182 were uninformative as they encompass (nearly) entire BGCs (Fig 2B). To validate the sub-cluster motifs,
183 we assessed whether we could detect a set of 109 experimentally verified sub-clusters, which are stored
184 in the SubClusterBlast module within the antiSMASH framework. The sub-cluster motifs from PRESTO-
185 TOP matched to 91 (83%) validated sub-clusters, where the methoxymalonate and AHBA sub-clusters of
186 macbecin are shown as examples (Fig C). Additionally, PRESTO-STAT was able to detect 78 of the
187 validated sub-clusters, of which 75 overlap with the sub-cluster motifs (S5 Fig). In general, we see that
188 PRESTO-TOP generates a more restricted amount of sub-cluster data, which might contain less
189 meaningful sub-clusters compared to PRESTO-STAT in absolute numbers but has a considerably higher
190 ratio of valid sub-cluster information.

191 **Fig 2. BGC length versus sub-cluster length.** (a) Scatterplot of the length of each BGC (number of non-
192 empty genes) from the antiSMASH-DB dataset versus the length of a match to a topic or sub-cluster motif,
193 representing a sub-cluster. The colour of each dot indicates how many times a BGC with a certain length
194 contains a sub-cluster with a certain length. (b) BGC for sipanmycin where the identified sub-cluster
195 encompasses the entire BGC, demonstrating an uninformative result. (c) BGC for macbecin where the two
196 characterised sub-clusters for AHBA (red) and methoxymalonyl (blue) are highlighted in the structure of
197 macbecin [20]. Sub-clusters from (b) and (c) are linked to their corresponding location in (a).

198 Our results provide clear examples of sub-cluster motifs that capture sub-cluster variety, by containing a
199 set of core genes responsible for synthesising the base of a substructure, and a set of modifying genes
200 that may not be present in all sub-clusters. For example, a motif like the sugar-related sub-cluster motif
201 680 is present in 134 MIBiG BGCs that represent different biosynthetic classes, such as different types of
202 polyketide synthases and nonribosomal peptide synthetases. This motif codes for the biosynthesis of
203 different (di)deoxy-sugars that are sometimes modified with amino or methyl-amino groups. However,

204 for some sub-cluster motifs, the biosynthetic context had an impact on shaping the motif. The sugar-
205 related sub-cluster motif 207, for example, contains several indolocarbazole biosynthesis genes as some
206 MIBiG BGCs matching to this motif encode the production of indolocarbazoles, and some of the
207 indolocarbazole-related genes ended up in this motif as weak features.

208 Exploring the sub-cluster motifs

209 Among the 90 identified characterised sub-clusters from the antiSMASH SubClusterBlast module, we
210 could readily annotate 23 sub-cluster motifs covering around 4,000 of the PRESTO-TOP-detected sub-
211 clusters. To extend on the sub-cluster knowledge stored in the SubClusterBlast module, we annotated
212 another 22 PRESTO-TOP-detected sub-cluster motifs for which sub-cluster instances were found inside
213 MIBiG BGCs. Together, these 45 annotations constitute 24 different types of substructures at different
214 levels of detail and allow us to explore the discovered sub-clusters more deeply (Fig 3 and S1 File). In
215 the non-redundant antiSMASH-DB dataset, around 9,500 (16%) putative BGCs contain at least one of
216 these annotated sub-cluster motifs. Through iPRESTO, we now gained relevant knowledge about these
217 putative BGCs that we can use to predict part of the structures of the products they encode.

218 **Fig 3. Sub-cluster motif annotations.** The pie chart visualises the annotations for the 45 sub-cluster motifs
219 divided into general substructure groups, where an example substructure is shown for several groups.
220 Additionally, examples of eight of the substructures are shown in the structures of apoptolidin, platencin,
221 fluvirucin b2 and pyralomicin 1a, where the colours of the substructures correspond to the sub-cluster motif
222 annotations in the pie chart. For these four metabolites, their respective BGCs are shown where the sub-cluster
223 motifs are highlighted in the same colour as the substructures they encode.

224 On average, an annotated sub-cluster motif occurs in 239 non-redundant BGCs, ranging from 19 BGCs
225 for sub-cluster motif 190, to 873 BGCs for sub-cluster motif 220, which encode the biosynthesis of
226 caprazol and dihydroxybenzoic acid moieties, respectively (S6 Fig). Some of the annotated sub-cluster
227 motifs are mainly present in one BGC class, while others occur in diverse BGC classes (S6 and S7 Figs).
228 An example of the latter is sub-cluster motif 773, which occurs in 153 BGCs mostly encoding
229 nonribosomal peptide synthetases and type I polyketide synthases. This sub-cluster motif encodes the
230 production of a 3-amino-2-methylpropionyl starter unit that appears in the known gene cluster
231 BGC0001597 (fluvirucin b2) (Fig 3). Interestingly, the motif also occurs in some BGCs of the class
232 "Other", meaning they cannot be classified by antiSMASH, like two BGCs from *Amycolatopsis alba* DSM
233 44262 (NZ_KB913032.1.cluster021; AMYAL_RS0129245 - AMYAL_RS0129610) and *Bradyrhizobium sp.*
234 *Ec3.3* (NZ_AXAS01000001.cluster006; YUU_RS0100020 - YUU_RS49645). This does not only provide
235 interesting leads for these BGCs with previously unknown structural predictions, but it also adds to their
236 validity. In total, 6.5% of the 10,000 "Other" class BGCs in the antiSMASH-DB contain one of the
237 annotated sub-cluster motifs.

238 iPRESTO can identify BGCs of orphan metabolites through sub-cluster 239 presence

240 Information about the sub-clusters present in a BGC is not only useful to predict the product of a BGC,
241 but it could also be used as a tool to identify BGCs for 'orphan' known metabolites. To demonstrate this,
242 we searched NPAtlas [21] with substructures that are encoded by our annotated sub-cluster motifs and
243 looked for metabolites without a MIBiG BGC that are found in one of the strains in the antiSMASH-DB
244 dataset. We first searched for metabolites that contain the dithiopyrrolone substructure for which the
245 biosynthesis is encoded by sub-cluster motif 517, annotated as such based on the MIBiG BGCs encoding
246 thiomarinol, holomycin and thiolutin [22-24]. In doing so, we found xenorhabdins 1-6, produced by
247 many *Xenorhabdus* strains that are also present in the antiSMASH-DB [25]. By searching for BGCs in
248 those strains that contain a match to the dithiopyrrolone sub-cluster motif, we found 12 *Xenorhabdus*
249 strains that contain such a BGC (Fig 4). In one of those strains, *X. doucetiae*, the BGC for xenorhabdin
250 biosynthesis has recently been described, corroborating that we accurately identified BGCs for
251 xenorhabdin biosynthesis based on iPRESTO-detected sub-clusters [26]. Next, we searched NPAtlas for
252 metabolites with the valienol moiety present in validamycin and pyralomicins, which is encoded by sub-
253 cluster motif 940 [27, 28]. As a result, we found salbostatin, which is produced by *Streptomyces*
254 *albus* ATCC 21838 in our dataset [29]. By investigating BGCs in that strain, we identified a BGC that
255 contains sub-cluster motif 940 and should therefore be responsible for salbostatin biosynthesis (Fig 4).
256 Indeed, it turned out that this BGC has already been described in 2008 to encode the production of

257 salbostatin [30], but it has been lacking from the MIBiG database [16]. This valienol sub-cluster motif
258 encoding C7-cyclitol-like substructures is an interesting example of a sub-cluster motif that can be found
259 in different biosynthetic contexts, *i.e.*, PKS-NRPS-like pyralomicins and different kinds of saccharides like
260 validomycin and salbostatin. This analysis highlights that iPRESTO allows identifying correct links
261 between BGCs and molecules that are published but were yet missing in public BGC databases (and
262 which can thus be added to these resources).

263 **Fig 4. Connecting non-MIBiG BGCs to their metabolic products through iPRESTO-detected sub-**
264 **clusters.** (a) Phylogenetic tree made with CORASON of 12 *Xenorhabdus* BGCs and 3 MIBiG BGCs, that contain
265 an iPRESTO-predicted sub-cluster for dithiopyrrolone biosynthesis [31]. The A-domain containing gene of
266 NZ_F0704550.1.cluster001 was used as query for CORASON. Structures of thiomarinol (1), thiolutin (2) and
267 holomycin (3) are linked to their MIBiG BGCs. *Xenorhabdus* (4-9) are encoded by *X. doucetiae* str. FRM16 as
268 indicated by the asterisk, while we infer based on sub-cluster presence that the other *Xenorhabdus* BGCs are
269 also responsible for xenorhabdin biosynthesis. (b) Phylogenetic tree made with CORASON
270 NZ_CP010519.1.cluster004 from *S. albus* ATCC 21838 and 4 MIBiG BGCs, that contain an iPRESTO-predicted
271 sub-cluster for C7 cyclitol biosynthesis. The predicted 2-epi-5-epi-valiolone synthase from
272 NZ_CP010519.1.cluster004 was used as query for CORASON. Structures of validomycin A (10) and pyralomycin
273 1A (11) are linked to their MIBiG BGCs. Salbostatin (12) is encoded by *S. albus* ATCC 21838 as indicated by the
274 hash symbol.

275 By searching in NPAtlas for chlorinated indoles, we found the orphan metabolites akashin A-C produced
276 by the diazaquinomycins producer *Streptomyces* sp. F001 [32]. The BGC of akashins has not been
277 described before in literature. As this strain was not present in the antiSMASH-DB, we ran antiSMASH 6
278 on the genome of this strain and used iPRESTO to infer sub-clusters in the predicted BGCs. As akashins
279 have chlorinated-indole moieties and are glycosylated, we sought for such sub-cluster motifs in the BGCs
280 of *S. sp. F001*. Interestingly, we identified the genomic region in QZWF01000007.1.region003
281 (StrepF001_25985 - StrepF001_26130) directly upstream of the diazaquinomycin BGC, based on the
282 presence of sub-cluster motifs 194, 607 and 680 that were annotated as methylaminosugar, halogenated
283 aromatic ring, and (amino)deoxysugar, respectively (Fig 5). The formation of the indigo-derived
284 backbone of akashins could potentially be formed by the two p450 enzymes, akin to CYP102G4, a
285 recently described p450 enzyme from *S. cattleya* [33]. This p450 enzyme can catalyse the reaction from
286 indole to 3-hydroxyindole after which spontaneous oxidation forms indigo. CYP102G4 was even shown to
287 accept chloro-indole as substrate, in the case that chlorination occurs before indole formation in akashin
288 biosynthesis. This shows that iPRESTO can aid in generating meaningful hypotheses about the
289 biosynthesis of orphan metabolites.

290 **Fig 5. Putative BGC for akashin A biosynthesis.** The antiSMASH-predicted BGC
291 QZWF01000007.1.region003 is shown (StrepF001_26130-StrepF001_26145), which is hypothetically
292 responsible for akashin A biosynthesis in *S. sp. F001*. Genes are coloured by their iPRESTO-predicted sub-
293 clusters or predicted function based on Pfam domains.

294 Correlation analysis in substructure-based integrative omics mining

295 To automatically link unknown molecules to BGCs at a larger scale, correlating substructures predicted
296 from metabolomics data to sub-clusters from genome data would potentially be of great added value
297 [12, 13]. To test such an approach, we used a previously defined correlation score which assumes that a
298 BGC is needed to synthesise a product, but that a BGC may be cryptic and not synthesise anything [15].
299 Ernst et al. [34] used the MS2LDA tool to discover substructure mass patterns, called Mass2Motifs, from
300 metabolomics data of 145 *Salinispora* and *Streptomyces* species for all of which (except one) genomic
301 data and BGC predictions are also available (the 'Streptomyces/Salinispora dataset') [14]. To identify
302 sub-clusters in these, we used iPRESTO to query all Streptomyces/Salinispora BGCs on the sub-cluster
303 motifs and sub-cluster clans (SCCs) of the antiSMASH-DB dataset. For each of the 107,590 pairs of
304 Mass2Motifs and sub-cluster motifs, we used the correlation score from Doroghazi et al. [15] to calculate
305 how frequently they co-occur across the Streptomyces/Salinispora strains, while we did the same for the
306 122,404 pairs of Mass2Motifs and SCCs (S8 Fig). To prioritise interesting substructure-sub-cluster pairs,
307 we performed permutation tests for all pairs to assess the likelihood of a high scoring pair arising by
308 chance. This was especially needed as the Streptomyces/Salinispora dataset includes highly related
309 strains, in which many BGCs and compounds are shared. Abundant sub-clusters and substructures

310 therefore get high correlation scores by default. Permutation testing resulted in 3,230 and 1,939
311 'significant' pairs of Mass2Motifs and sub-cluster motifs or SCCs, respectively (S8 Fig). As an example of
312 how such an approach connects substructure information inferred from genome and metabolome mining,
313 we identified 5 high correlation scores with low p-values between two staurosporine-related mass2motifs
314 and both sub-cluster motifs and SCCs constituting the amino-sugar moiety of staurosporine (Fig 6).
315 Since currently only a fraction of the Mass2Motifs, sub-cluster motifs and SCCs are annotated, our
316 analysis serves as an illustration of how such an approach could help to link metabolome and genome
317 data in the future.

318 This correlation method generally results in a lot of noise, as sub-clusters and substructures that occur in
319 a shared subset of strains will all correlate to each other. Such co-correlating structures make the
320 identification of the actual correlating pair therefore difficult, especially with limited annotations.
321 Identifying clusters of co-correlating pairs could therefore provide a way to make the interpretation of
322 this analysis easier. Additionally, the correlation analysis is not perfect in our case, as multiple different
323 sub-clusters are often responsible for synthesising the same kind of substructure. For example, we
324 identified multiple sub-cluster motifs that can encode methylated aminosugars, while only one
325 mass2motif is annotated as a methylated aminosugar. In future approaches, such mismatches between
326 genome and metabolome could be overcome by finding ways to group sub-cluster motifs together that
327 encode similar structures before running such metabologenomic correlation analyses. Combining such
328 solutions with the integration of more diverse species, new annotations, and improved correlation scoring
329 methods like the one developed in Hjörleifsson Eldjárn et al. [35] would improve such analyses
330 drastically. Furthermore, we expect that combining co-occurrence based scores (such as standardized
331 Metcalf) with feature-based scores, such as NPClassScore [36], and the here developed iPRESTO, will
332 further help to prioritize plausible BGC-MS/MS spectral links [12, 13]. Indeed, we expect that tools like
333 iPRESTO could in the future be built into frameworks like NPLinker [35]. As our current contribution
334 represents a first step in linking substructure-and sub-cluster models with rather limited (annotated)
335 information, we expect that analyses like these will have great impact in the future to facilitate
336 metabologenomics experiments that use integrative omics mining.

337 **Fig 6. Metabologenomic correlation scores between sub-clusters and mass2motifs.** Stacked histogram
338 of the correlation scores across the *Streptomyces*/*Salinispora* strains between the mass2motifs paired with
339 either the SCCs or sub-cluster motifs with a p-value below 0.1. Highlighted with their scores are the pairs
340 mass2motif_108 with SSC_452, SSC_1010, sub-cluster_motif_207 and sub-cluster_motif_680, and the pair
341 mass2motif_8 with SSC_452. The aforementioned sub-cluster motifs (blue) and SCCs (brown) are responsible
342 for sugar synthesis in staurosporine, while both mass2motifs (red) are staurosporine related.

343 Conclusion and future perspectives

344 This study introduces the iPRESTO concept and makes it available as a command line tool. We plan to
345 include iPRESTO in one of the future releases of antiSMASH, so the collection of sub-clusters we
346 generated in this study can be used to detect and visualize them in antiSMASH-predicted BGCs. We
347 anticipate that this will enhance the current scope of sub-cluster detection, as antiSMASH's current sub-
348 cluster predictor SubClusterBlast offers a limited amount of sub-cluster data, whereas our sub-cluster set
349 will allow making more connections between predicted BGCs and MIBiG reference BGCs. This will
350 accelerate NP discovery by linking structural information from genome and metabolome data.

351 Due to the above discussed limitations of PRESTO-STAT, we plan to use PRESTO-TOP as the main
352 method for sub-cluster detection in the antiSMASH implementation, as it also captures sub-cluster
353 variety in the sub-cluster motifs and yet can be used easily to query BGCs for sub-cluster motifs.
354 PRESTO-STAT could still be used to identify the sub-cluster boundaries better, by for example linking
355 groups of related PRESTO-STAT sub-clusters to 'parent' PRESTO-TOP sub-cluster motifs, and by using
356 the PRESTO-STAT modules to more specifically identify the sub-cluster variant found in a given BGC. The
357 drawback of the statistical method that it produces highly nested and variable sub-clusters could as such
358 be used as a strength. A way to further improve PRESTO-TOP would be to apply PRESTO-TOP in a semi-
359 supervised manner, which constitutes a major potential benefit of this approach. Before training an LDA
360 model, certain motifs could be seeded beforehand, which allows accurate sub-cluster motifs to be reused
361 in new analyses, analogous to the metabolomics database MotifDB, in which annotated Mass2Motifs are

362 stored in MotifSets [37]. Such semi-supervised approaches would allow for noise to be eliminated from
363 sub-cluster motifs and sub-cluster motifs to be finetuned. Another way to reduce noise and to identify
364 the more robust sub-cluster motifs would be to train multiple PRESTO-TOP models on the same dataset. Sub-cluster motifs
365 arise through chance would be filtered out, as they would only occur in one or a few of the many LDA
366 models. Noisy genes in accurate sub-cluster motifs could be filtered out by taking intersects of multiple
367 similar sub-cluster motifs. As another option, each BGC could be represented multiple times in training to
368 increase the observations of less frequently occurring sub-clusters. This could lead to better estimation of
369 the sub-cluster motif distributions over the data and cause less erroneous mixed sub-cluster motifs. We
370 have attempted this for a small subset and noticed that the overlap with SubClusterBlast increased
371 slightly, making this an interesting avenue to continue PRESTO-TOP sub-cluster algorithmic
372 developments.

373 Using iPRESTO, in our current study we were able to characterise 45 different sub-cluster motifs present
374 in diverse BGC classes. The remaining 955 sub-cluster motifs remain largely unexplored, of which many
375 are likely to encode useful substructures. We expect that, in the future, more annotations will increase
376 the value of our results even more, which will be aided by the inclusion of updated (expanded) versions
377 of the MIBiG database. Using one of the characterised sub-cluster motifs, we showed a direct practical
378 application of our method by hypothesising a putative BGC for akashin A production. Additionally, we
379 provided the initial step for linking sub-clusters to substructures in a systematic way, which in the future
380 could facilitate the automated connection of BGCs to their NPs.

381 Methods

382 Data and Code availability

383 iPRESTO is available as a command-line tool at <https://git.wageningenur.nl/bioinformatics/iPRESTO/>. The
384 annotated sub-cluster motifs and other relevant data can be found at
385 <https://doi.org/10.5281/zenodo.6953657>. The following sections describe the most important steps of
386 this project, while the Supplementary methods in S1 Text provide more detailed explanations.

387 Data selection

388 The antiSMASH-DB dataset consisted of three data sources: the MIBiG database, the
389 Streptomyces/Salinispora dataset and the antiSMASH-DB. Version 1.4 of the MIBiG database was used
390 which contains 1,819 BGCs (https://dl.secondarymetabolites.org/mibig/mibig_gbk_1.4.tar.gz). The
391 Streptomyces/Salinispora dataset consists of 5,927 BGCs that originate from the 146 *Streptomyces* and
392 *Salinispora* strains investigated by Crüsemann et al. [38]. antiSMASH 3.0 was used for the detection of
393 BGCs in the Streptomyces/Salinispora dataset. The antiSMASH-DB version 2 is comprised of 152,122
394 BGCs detected with antiSMASH 4.0, where we included BGCs from draft genomes (Table A in S1 Text;
395 https://dl.secondarymetabolites.org/database/2.0/asdb_20180828_all_results.tar.xz). BGCs were
396 discarded if they were flagged by antiSMASH as lying on a contig-edge, as these BGCs are probably
397 incomplete (fragmented) and less accurate. Additionally, BGC class information was included in the
398 analysis, by using the assigned antiSMASH biosynthetic classes.

399 Data pre-processing

400 BGCs were tokenised by converting each gene into a string of (sub)Pfam domains. To detect (sub)Pfams,
401 the HMMER3 tool hmmscan was used with a custom profile hidden Markov model (pHMM) database
402 consisting of Pfam database version 32.0, where 112 Pfams were replaced by corresponding subPfams
403 [39, 40]. These 112 Pfams were selected as they are the most abundant biosynthetic Pfams in the
404 antiSMASH-DB (S2 File). To create subPfams, the multiple sequence alignment of a Pfam is split into
405 clades, after which a new pHMM is built for each clade, each of which constitutes a subPfam (S1A Fig and
406 https://github.com/satriaphd/build_subpfam).

407 Redundant BGCs were removed from the analysis using a similarity network of BGCs, where BGCs were
408 connected based on an Adjacency Index of domains higher than 0.95 or if BGCs were fully contained
409 within one another. From each maximal clique in the network, only the BGC with the most domains was

410 chosen to remain in the analysis (Table A in S1 Text and S9 Fig) [41]. After redundancy filtering, all non-
411 biosynthetic domains were removed from all BGCs. To select biosynthetic domains, EC-associated Pfams
412 were collected with ECDomainMiner, from which Pfams were selected if they occurred in pre-calculated
413 BGCs [42]. After manual curation, this resulted in a list of 1,839 biosynthetic Pfams (S3 File).
414 Additionally, Pfams that occurred less than three times in the dataset were removed as well as BGCs that
415 contained less than two non-empty genes (S4 File).

416 PRESTO-STAT

417 The statistical method for sub-cluster detection was re-implemented in Python based on Del Carratore et
418 al. [10] with some alterations, resulting in PRESTO-STAT. Instead of representing genes as COGs as in
419 the previous method, we represent each gene as a combination of its domains. First, all possible
420 adjacency and co-localisation interactions between each pair of genes are counted. To assess whether an
421 observed interaction between two genes occurs more than by random chance, one needs to distribute
422 such a pair of genes randomly through the dataset and calculate the probability of the observed
423 interaction. To reduce the computational burden of a permutation-based approach, for each pair of genes
424 one gene is kept fixed while the other is being randomly distributed throughout the data. For an
425 adjacency interaction this gives a hypergeometric equation describing all available positions of one gene
426 while the other is fixed (Table B1 in S1 Text). This follows from the fact that there are three options for
427 the position of gene B while keeping gene A fixed: not adjacent to gene A (B_1), adjacent to gene A (B_2),
428 or adjacent to gene A on both sides (B_3). N_1 , N_2 and N_3 represent all available positions in these three
429 categories, while N_{tot} represents all positions and B_{tot} all occurrences of gene B. For a co-localisation
430 interaction the same applies, except for the fact that gene B can be co-localised with n_{max} genes A, where
431 n_{max} is the number of genes A co-localised with gene B (Table B2 in S1 Text). When n_{max} is large this
432 becomes computationally hard, which is why we replaced duplicate genes with an empty gene (a dash)
433 and placed one copy of the duplicate gene at the end of the cluster separated by an empty gene. This
434 simplifies the equation as only two types of co-localisations need to be counted: co-localisation and no
435 co-localisation (Table B3 in S1 Text). A p-value can be calculated by summing all probabilities in the
436 hypergeometric distribution that correspond to several interactions higher or equal to the observed
437 number of interactions. Or, to make it easier, by subtracting the sum of all possible interactions smaller
438 than the observed interaction from one (Table B4 in S1 Text).

439 Calculating an interaction between each pair of genes results in two p-values, one coming from gene A
440 and one coming from gene B. Only the largest p-value for both the co-localisation, and the adjacency
441 interactions is considered, to be conservative. To control false discovery rate under dependency we used
442 the Benjamini–Yekutieli method on both the co-localisation and adjacency p-values [43].

443 To group interacting pairs of genes into sub-clusters, undirected graphs are constructed, where each
444 gene is a node. An edge is made between two genes if they have an adjacency or co-localisation p-value
445 below a threshold of 0.1. All maximal cliques are selected as sub-clusters, while changing the threshold
446 iteratively to all the p-values in the dataset smaller than the original threshold of 0.1. To reduce false
447 positives, we removed putative sub-clusters if they contained fewer than three genes and if they only
448 occurred in one BGC. Next, we grouped similar sub-clusters together using K-means clustering into sub-
449 cluster families and sub-cluster clans and removed redundant sub-clusters (Supplementary methods in
450 S1 Text) [44, 45].

451 PRESTO-TOP

452 PRESTO-TOP uses Latent Dirichlet Allocation (LDA) latent sub-cluster composition in BGCs [46]. LDA
453 assumes a bag-of-words representation, where each BGC is depicted as a frequency vector of its domain
454 combinations, not taking gene order into account. We used the multicore LDA implementation from
455 Gensim, that makes use of online variational Bayes [47, 48]. In this implementation, an LDA model is
456 trained by updating it with mini-batches from the data, which has low time and memory complexity. We
457 chose the chunk size of each mini-batch to be 5% of the data with a minimum chunk size of 2,000,
458 which is loosely based on testing different chunk sizes by Hoffman et al. [48]. We considered that using
459 500 iterations to train a model was enough after assessing that the log-likelihood converged sufficiently

460 (S10 Fig). For the sake of computational resources, we did limited hyperparameter optimisation for the
461 number of sub-cluster motifs (topics) N , α , and β . To test the performance of the different models, we
462 considered the coherence score as measured with the `u_mass` method [49] and the overlap with
463 validated sub-clusters from SubClusterBlast (Supplementary methods in S1 Text). Based on the
464 coherence score of the different models, choosing 250 sub-cluster motifs seemed optimal (S11A Fig).
465 However, upon manual inspection of some of the motifs, it turned out that many motifs are hard to
466 annotate with a single substructure due to the presence of many noisy features. This is corroborated by
467 the fact that choosing 250 sub-cluster motifs does not produce the highest overlap with SubClusterBlast
468 (S11B Fig). Instead, the model with 1000 sub-cluster motifs produced the highest overlap with
469 SubClusterBlast while having a similar coherence score to the model with 250 motifs, which is why we
470 chose 1000 sub-cluster motifs. We chose the default setting of a symmetric $1/N$ for hyperparameters α
471 and β , as we could not find better SubClusterBlast overlap when setting α and β to symmetric,
472 asymmetric, auto, or 1.

473 Each sub-cluster motif in an LDA model consists of a probability vector of domain combinations,
474 representing the contribution of each domain combination to a sub-cluster motif. To filter out noise, we
475 sorted this vector from high to low probability, summed the probabilities and included all domain
476 combinations until 0.95 was reached. When a group of genes from a BGC match to a sub-cluster motif,
477 each gene is assigned a feature probability describing how well it fits in the sub-cluster motif, for which
478 we set a cut-off of 0.3. For a sub-cluster to be considered it needs to consist of more than one gene, for
479 which we set a cut-off of 1.1 on the summed feature probabilities. Additionally, we calculated an overlap
480 score for each match, which we computed by summing the domain combination probabilities from the
481 sub-cluster motif present in the match [50]. We set a threshold of 0.15 on the overlap score, as this was
482 the highest threshold that did not remove manually validated SubClusterBlast sub-clusters from the
483 analysis.

484 Acknowledgements

485 We thank Dr Dick de Ridder and Dr Simon Rogers for useful comments and discussions.

486 Supporting information

487 **S1 Text. Supplementary information for *iPRESTO: automated discovery of biosynthetic sub-clusters***
488 ***linked to specific natural product substructures.***

489 **S1 Fig. Schematic depiction of BGC tokenisation.** (A) subPfam are constructed for the 112 most
490 frequent Pfam domains in the antiSMASH-DB by dividing the multiple sequence alignment of a Pfam into clades
491 and converting each clade into a new pHMM. (B) The BGCs predicted by antiSMASH are tokenised by detecting
492 (sub)Pfam in each gene, where non-biosynthetic Pfams are removed. After tokenising the BGCs, sub-cluster
493 can be detected with the statistical method (Stat), where the tokenised genes are represented in their original
494 order, or by LDA, which assumes a bag of words model where original gene order is not considered.

495 **S2 Fig. Result of querying rifamycin (BGC000373) to the PRESTO-TOP and PRESTO-STAT sub-**
496 **clusters generated in this project.** Only around 25% of the PRESTO-STAT sub-clusters are shown. Each
497 gene is depicted as a token, where all (sub)Pfam domains are coloured. The visualisation of the BGC, the
498 PRESTO-TOP and PRESTO-STAT output are separated by a dashed line, respectively. All PRESTO-STAT sub-
499 clusters clearly exhibit a nested structure, where all combinations of genes in an actual sub-cluster are detected
500 as individual sub-clusters. The PRESTO-STAT sub-clusters shown here are also examples of noisy sub-clusters
501 comprised of combinations of genes from different actual sub-clusters, like detected PRESTO-STAT sub-clusters
502 that are combinations of genes responsible for the biosynthesis of AHBA (green), sugars (blue) and the
503 polyketide scaffold (purple).

504 **S3 Fig. Information about the PRESTO-STAT sub-clusters.** (A) The distribution of the number of genes
505 per PRESTO-STAT sub-cluster in the antiSMASH-DB dataset. (B) The distribution of the \log_{10} transformed
506 PRESTO-STAT sub-cluster occurrences in the antiSMASH-DB dataset.

507 **S4 Fig. Number of PRESTO-STAT and PRESTO-TOP sub-clusters per BGC.** (A) Distribution of the log₁₀
508 transformed number of PRESTO-STAT sub-clusters per BGC in the non-redundant antiSMASH-DB dataset,
509 where the bin with the seemingly negative value represents BGCs without any PRESTO-STAT sub-cluster. (B)
510 The number of topics or sub-cluster motifs per BGC in the non-redundant antiSMASH-DB dataset, not counting
511 sub-clusters of length one as these are almost definitely noise (see Methods). (C) All BGCs with at least one
512 annotated sub-cluster motif grouped by how many annotated sub-cluster motifs they have. In total there are
513 9,425 putative BGCs with at least one annotated sub-cluster motif, and 350 MIBiG BGCs.

514 **S5 Fig. PRESTO-STAT and PRESTO-TOP overlap with validated sub-clusters from SubClusterBlast.**
515 Overlap between detected SubClusterBlast sub-clusters and output of both sub-cluster detection methods
516 applied on the antiSMASH-DB dataset according to different overlap cut-offs. The overlap expresses the fraction
517 of genes from the original SubClusterBlast sub-cluster that is found in the iPRESTO-detected sub-cluster. We
518 considered an overlap of 0.6 sufficient for having detected a sub-cluster (see Supplementary methods in S1
519 Text).

520 **S6 Fig. Degrees (occurrences) of the annotated sub-cluster motifs within the antiSMASH-DB dataset**
521 **(non-redundant).**

522 **S7 Fig. BGC class distribution across sub-cluster motifs.** Relative abundance of antiSMASH classes when
523 querying the non-redundant antiSMASH-DB dataset on the 45 annotated sub-cluster motifs. Matches of length
524 1 are ignored and hybrid class BGCs are counted for all classes they contain. RIPPs classes are grouped
525 together.

526 **S8 Fig. Correlation scores between Mass2Motifs and sub-clusters.** (A) Correlation scores between
527 Mass2Motifs and SCCs. (B) Correlation scores between Mass2Motifs and sub-cluster motifs. In both panels the
528 significant pairs are highlighted.

529 **S9 Fig. Graphical representation of graph-based filtering for the small dataset: MIBiG-and**
530 **Streptomyces/Salinispora BGCs.** Each node represents a BGC and an edge represents an adjacency index
531 (AI) of 0.95 or higher. In blue are the BGCs chosen as representatives, while BGCs that are filtered out are
532 shown in black. We show the small dataset here as it was difficult to visualize this process for the antiSMASH-
533 DB dataset.

534 **S10 Fig. LDA model convergence.** Convergence of the log-likelihood of an LDA model with 1,000 topics/sub-
535 cluster motifs trained on the non-redundant 60,028 BGCs from the antiSMASH-DB dataset, which also contains
536 the Streptomyces/Salinispora dataset and the MIBiG database, using 2,000 iterations of chunk size 3,000. Log-
537 likelihood based on 28 held out BGCs.

538 **S11 Fig. Coherence scores and overlap with SubClusterBlast sub-clusters for different LDA models.**
539 (A) Coherence scores of different LDA models trained using PRESTO-TOP on the non-redundant antiSMASH-DB
540 dataset with different number of topics. (B) Number of validated SubClusterBlast sub-clusters found with
541 different LDA models trained using PRESTO-TOP on the non-redundant antiSMASH-DB dataset with different
542 number of topics.

543 **S1 File. Excel sheet containing the current information about the 45 annotated sub-cluster motifs.**

544 **S2 File. The 112 domains for which we created subPfam.**

545 **S3 File. The biosynthetic domains we considered in this study.**

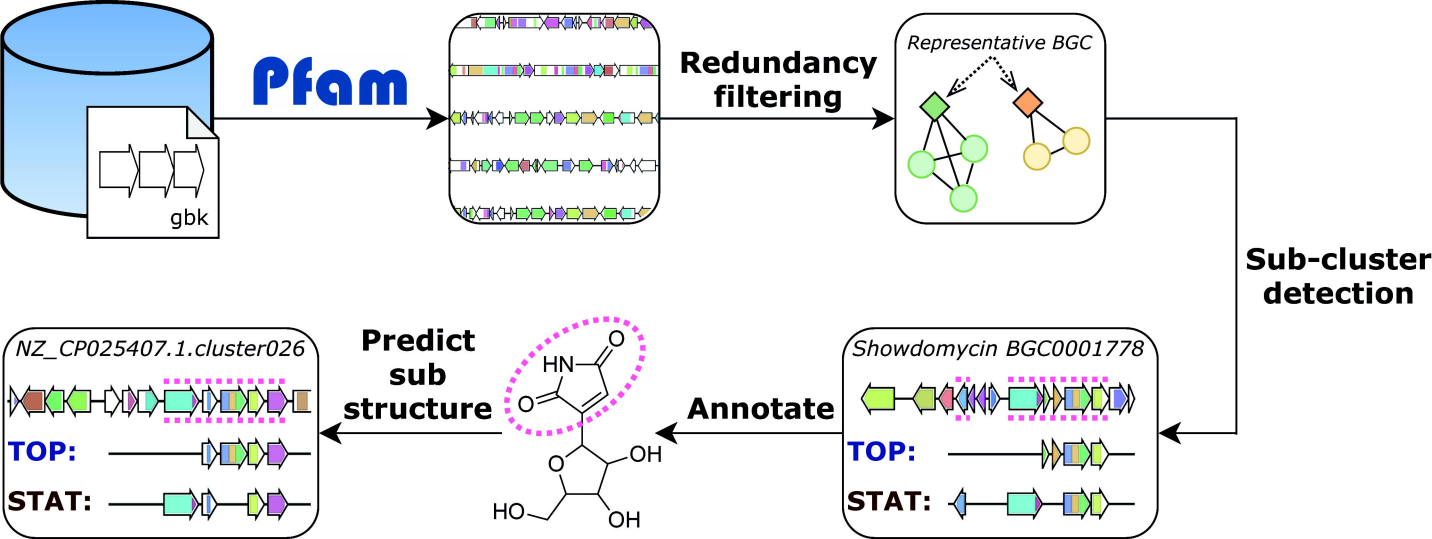
546 **S4 File. All used domain-combinations present in the antiSMASH-DB dataset after filtering.**

547 References

- 548 1. Dayan FE, Cantrell CL, Duke SO. Natural products in crop protection. *Bioorganic & medicinal chemistry*.
549 2009;17(12):4022-34.
- 550 2. Li JWH, Vederas JC. Drug Discovery and Natural Products: End of an Era or an Endless Frontier?
551 *Science*. 2009;325(5937):161. doi: 10.1126/science.1168243.
- 552 3. Pye CR, Bertin MJ, Lokey RS, Gerwick WH, Linington RG. Retrospective analysis of natural products
553 provides insights for future discovery trends. *Proc Natl Acad Sci U S A*. 2017;114(22):5601-6. Epub

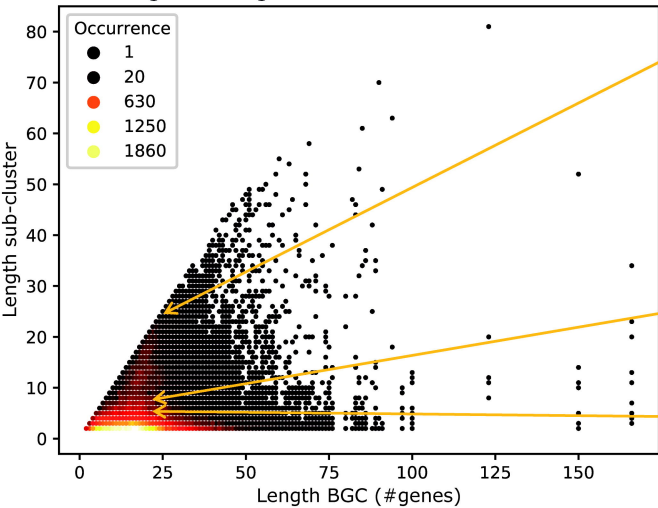
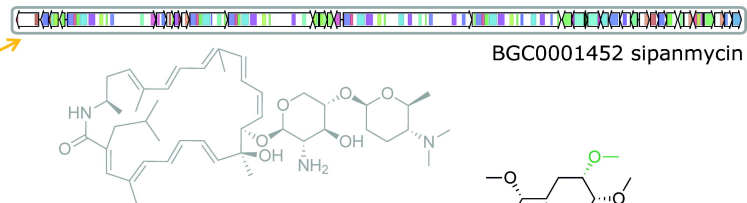
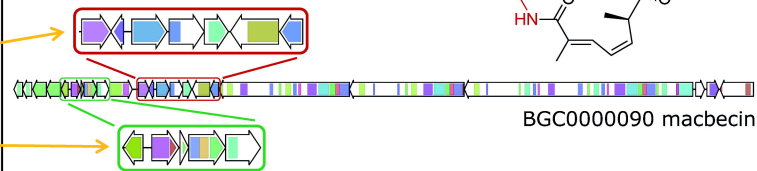
- 554 2017/05/04. doi: 10.1073/pnas.1614680114. PubMed PMID: 28461474; PubMed Central PMCID:
555 PMCPCMC5465889.
- 556 4. Medema MH, Cimermancic P, Sali A, Takano E, Fischbach MA. A systematic computational analysis of
557 biosynthetic gene cluster evolution: lessons for engineering biosynthesis. *PLoS Comput Biol.*
558 2014;10(12):e1004016. Epub 2014/12/05. doi: 10.1371/journal.pcbi.1004016. PubMed PMID: 25474254;
559 PubMed Central PMCID: PMCPCMC4256081.
- 560 5. Chevrette MG, Currie CR. Emerging evolutionary paradigms in antibiotic discovery. *J Ind Microbiol*
561 *Biotechnol.* 2018. Epub 2018/10/01. doi: 10.1007/s10295-018-2085-6. PubMed PMID: 30269177.
- 562 6. Cimermancic P, Medema MH, Claesen J, Kurita K, Wieland Brown LC, Mavrommatis K, et al. Insights
563 into secondary metabolism from a global analysis of prokaryotic biosynthetic gene clusters. *Cell.*
564 2014;158(2):412-21. Epub 2014/07/19. doi: 10.1016/j.cell.2014.06.034. PubMed PMID: 25036635; PubMed
565 Central PMCID: PMCPCMC4123684.
- 566 7. Blin K, Shaw S, Kloosterman AM, Charlop-Powers Z, van Wezel GP, Medema Marnix H, et al.
567 antiSMASH 6.0: improving cluster detection and comparison capabilities. *Nucleic Acids Research.*
568 2021;49(W1):W29-W35. doi: 10.1093/nar/gkab335.
- 569 8. Skinnider MA, Johnston CW, Gunabalasingam M, Merwin NJ, Kieliszek AM, MacLellan RJ, et al.
570 Comprehensive prediction of secondary metabolite structure and biological activity from microbial genome
571 sequences. *Nature Communications.* 2020;11(1):6058. doi: 10.1038/s41467-020-19986-1.
- 572 9. Fischbach MA, Walsh CT, Clardy J. The evolution of gene collectives: How natural selection drives
573 chemical innovation. *Proceedings of the National Academy of Sciences.* 2008;105(12):4601. doi:
574 10.1073/pnas.0709132105.
- 575 10. Del Carratore F, Zych K, Cummings M, Takano E, Medema MH, Breitling R. Computational identification
576 of co-evolving multi-gene modules in microbial biosynthetic gene clusters. *Communications Biology.* 2019;2(1).
577 doi: 10.1038/s42003-019-0333-6.
- 578 11. Blin K, Shaw S, Kautsar SA, Medema MH, Weber T. The antiSMASH database version 3: increased
579 taxonomic coverage and new query features for modular enzymes. *Nucleic Acids Research.* 2020;49(D1):D639-
580 D43. doi: 10.1093/nar/gkaa978.
- 581 12. Louwen JJ, Van Der Hooft JJ. Comprehensive large-scale integrative analysis of omics data to
582 accelerate specialized metabolite discovery. *Msystems.* 2021;6(4):e00726-21.
- 583 13. van der Hooft JJJ, Mohimani H, Bauermeister A, Dorrestein PC, Duncan KR, Medema MH. Linking
584 genomics and metabolomics to chart specialized metabolic diversity. *Chemical Society Reviews.*
585 2020;49(11):3297-314. doi: 10.1039/D0CS00162G.
- 586 14. van der Hooft JJ, Wandy J, Barrett MP, Burgess KE, Rogers S. Topic modeling for untargeted
587 substructure exploration in metabolomics. *Proc Natl Acad Sci U S A.* 2016;113(48):13738-43. Epub
588 2016/11/20. doi: 10.1073/pnas.1608041113. PubMed PMID: 27856765; PubMed Central PMCID:
589 PMCPCMC5137707.
- 590 15. Doroghazi JR, Albright JC, Goering AW, Ju KS, Haines RR, Tchalukov KA, et al. A roadmap for natural
591 product discovery based on large-scale genomics and metabolomics. *Nat Chem Biol.* 2014;10(11):963-8. doi:
592 10.1038/nchembio.1659. PubMed PMID: 25262415; PubMed Central PMCID: PMCPCMC4201863.
- 593 16. Kautsar SA, Blin K, Shaw S, Navarro-Muñoz JC, Terlouw BR, van der Hooft JJJ, et al. MIBiG 2.0: a
594 repository for biosynthetic gene clusters of known function. *Nucleic Acids Research.* 2019;48(D1):D454-D8.
595 doi: 10.1093/nar/gkz882.
- 596 17. Kautsar SA, van der Hooft JJJ, de Ridder D, Medema MH. BiG-SLiCE: A highly scalable tool maps the
597 diversity of 1.2 million biosynthetic gene clusters. *GigaScience.* 2021;10(1). doi: 10.1093/gigascience/giaa154.
- 598 18. Chen X, Hu X, Shen X, Rosen G, editors. Probabilistic topic modeling for genomic data interpretation.
599 2010 IEEE International Conference on Bioinformatics and Biomedicine (BIBM); 2010: IEEE.
- 600 19. Blin K, Medema MH, Kazempour D, Fischbach MA, Breitling R, Takano E, et al. antiSMASH 2.0--a
601 versatile platform for genome mining of secondary metabolite producers. *Nucleic Acids Res.* 2013;41(Web
602 Server issue):W204-12. Epub 2013/06/06. doi: 10.1093/nar/gkt449. PubMed PMID: 23737449; PubMed
603 Central PMCID: PMCPCMC3692088.
- 604 20. Zhang M-Q, Gaisser S, Nur-E-Alam M, Sheehan LS, Vousden WA, Gaitatzis N, et al. Optimizing Natural
605 Products by Biosynthetic Engineering: Discovery of Nonquinone Hsp90 Inhibitors. *Journal of Medicinal*
606 *Chemistry.* 2008;51(18):5494-7. doi: 10.1021/jm8006068.
- 607 21. van Santen JA, Jacob G, Singh AL, Aniebok V, Balunas MJ, Bunsko D, et al. The Natural Products Atlas:
608 An Open Access Knowledge Base for Microbial Natural Products Discovery. *ACS Central Science.*
609 2019;5(11):1824-33. doi: 10.1021/acscentsci.9b00806.
- 610 22. Li B, Walsh CT. Identification of the gene cluster for the dithiopyrrolone antibiotic holomycin in
611 *Streptomyces clavuligerus*. *Proceedings of the National Academy of Sciences.* 2010;107(46):19731-5. doi:
612 doi:10.1073/pnas.1014140107.
- 613 23. Fukuda D, Haines AS, Song Z, Murphy AC, Hothersall J, Stephens ER, et al. A Natural Plasmid Uniquely
614 Encodes Two Biosynthetic Pathways Creating a Potent Anti-MRSA Antibiotic. *PLOS ONE.* 2011;6(3):e18031. doi:
615 10.1371/journal.pone.0018031.
- 616 24. Huang S, Him Tong M, Qin Z, Deng Z, Deng H, Yu Y. Identification and characterization of the
617 biosynthetic gene cluster of thiolutin, a tumor angiogenesis inhibitor, in *Saccharothrix algeriensis* NRRL B-
618 24137. *Anti-Cancer Agents in Medicinal Chemistry (Formerly Current Medicinal Chemistry-Anti-Cancer Agents).*
619 2015;15(3):277-84.
- 620 25. McInerney BV, Gregson RP, Lacey MJ, Akhurst RJ, Lyons GR, Rhodes SH, et al. Biologically Active
621 Metabolites from *Xenorhabdus* Spp., Part 1. Dithiopyrrolone Derivatives with Antibiotic Activity. *Journal of*
622 *Natural Products.* 1991;54(3):774-84. doi: 10.1021/np50075a005.
- 623 26. Bode E, Brachmann AO, Kegler C, Simsek R, Dauth C, Zhou Q, et al. Simple "On-Demand" Production
624 of Bioactive Natural Products. *ChemBioChem.* 2015;16(7):1115-9. doi:
625 <https://doi.org/10.1002/cbic.201500094>.

- 626 27. Bai L, Li L, Xu H, Minagawa K, Yu Y, Zhang Y, et al. Functional analysis of the validamycin biosynthetic
627 gene cluster and engineered production of validoxylamine A. *Chemistry & biology*. 2006;13(4):387-97.
628 28. Flatt PM, Wu X, Perry S, Mahmud T. Genetic Insights into Pyralomicin Biosynthesis in *Nonomuraea*
629 *spiralis* IMC A-0156. *Journal of Natural Products*. 2013;76(5):939-46. doi: 10.1021/np400159a.
630 29. Vértesy L, Fehlihaber H-W, Schulz A. The Trehalase Inhibitor Salbostatin, a Novel Metabolite from
631 *Streptomyces albus*, ATCC21838. *Angewandte Chemie International Edition in English*. 1994;33(18):1844-6.
632 doi: <https://doi.org/10.1002/anie.199418441>.
633 30. Choi WS, Wu X, Choeng Y-H, Mahmud T, Jeong BC, Lee SH, et al. Genetic organization of the putative
634 salbostatin biosynthetic gene cluster including the 2-epi-5-epi-valiolone synthase gene in *Streptomyces albus*
635 ATCC 21838. *Applied Microbiology and Biotechnology*. 2008;80(4):637-45. doi: 10.1007/s00253-008-1591-2.
636 31. Navarro-Muñoz JC, Selem-Mojica N, Mullowney MW, Kautsar SA, Tryon JH, Parkinson EI, et al. A
637 computational framework to explore large-scale biosynthetic diversity. *Nature Chemical Biology*.
638 2020;16(1):60-8. doi: 10.1038/s41589-019-0400-9.
639 32. Braesel J, Clark CM, Kunstman KJ, Green SJ, Maienschein-Cline M, Murphy BT, et al. Genome
640 Sequence of Marine-Derived *Streptomyces* sp. Strain F001, a Producer of Akashin A and Diazaquinomycins.
641 *Microbiology Resource Announcements*. 2019;8(19):e00165-19. doi: doi:10.1128/MRA.00165-19.
642 33. Kim J, Lee P-g, Jung E-o, Kim B-G. In vitro characterization of CYP102G4 from *Streptomyces cattleya*:
643 A self-sufficient P450 naturally producing indigo. *Biochimica et Biophysica Acta (BBA) - Proteins and*
644 *Proteomics*. 2018;1866(1):60-7. doi: <https://doi.org/10.1016/j.bbapap.2017.08.002>.
645 34. Ernst M, Kang KB, Caraballo-Rodríguez AM, Nothias L-F, Wandy J, Chen C, et al. MolNetEnhancer:
646 Enhanced Molecular Networks by Integrating Metabolome Mining and Annotation Tools. *Metabolites*.
647 2019;9(7):144. PubMed PMID: doi:10.3390/metabo9070144.
648 35. Hjörleifsson Eldjárn G, Ramsay A, van der Hooft JJJ, Duncan KR, Soldatou S, Rousu J, et al. Ranking
649 microbial metabolomic and genomic links in the NPLinker framework using complementary scoring functions.
650 *PLOS Computational Biology*. 2021;17(5):e1008920. doi: 10.1371/journal.pcbi.1008920.
651 36. Louwen JJ, Medema MH, van der Hooft JJ. Enhanced correlation-based linking of biosynthetic gene
652 clusters to their metabolic products through chemical class matching. 2022. doi:
653 <https://doi.org/10.21203/rs.3.rs-1391827/v2>.
654 37. Rogers S, Ong CW, Wandy J, Ernst M, Ridder L, van der Hooft JJJ. Deciphering complex metabolite
655 mixtures by unsupervised and supervised substructure discovery and semi-automated annotation from MS/MS
656 spectra. *Faraday Discussions*. 2019;218(0):284-302. doi: 10.1039/C8FD00235E.
657 38. Crüsemann M, O'Neill EC, Larson CB, Melnik AV, Floros DJ, da Silva RR, et al. Prioritizing Natural
658 Product Diversity in a Collection of 146 Bacterial Strains Based on Growth and Extraction Protocols. *J Nat Prod*.
659 2017;80(3):588-97. doi: 10.1021/acs.jnatprod.6b00722. PubMed PMID: 28335604; PubMed Central PMCID:
660 PMC5367486.
661 39. Bateman A, Smart A, Luciani A, Salazar GA, Mistry J, Richardson LJ, et al. The Pfam protein families
662 database in 2019. *Nucleic Acids Research*. 2018;47(D1):D427-D32. doi: 10.1093/nar/gky995.
663 40. Mistry J, Finn RD, Eddy SR, Bateman A, Punta M. Challenges in homology search: HMMER3 and
664 convergent evolution of coiled-coil regions. *Nucleic acids research*. 2013;41(12):e121-e.
665 41. Bron C, Kerbosch J. Algorithm 457: finding all cliques of an undirected graph. *Commun ACM*.
666 1973;16(9):575-7. doi: 10.1145/362342.362367.
667 42. Alborzi SZ, Devignes M-D, Ritchie DW. ECDomainMiner: discovering hidden associations between
668 enzyme commission numbers and Pfam domains. *BMC Bioinformatics*. 2017;18(1):107. doi: 10.1186/s12859-
669 017-1519-x.
670 43. Benjamini Y, Yekutieli D. The control of the false discovery rate in multiple testing under dependency.
671 *The annals of statistics*. 2001;29(4):1165-88.
672 44. Arthur D, Vassilvitskii S, editors. k-means++: The advantages of careful seeding. *Proceedings of the*
673 *eighteenth annual ACM-SIAM symposium on Discrete algorithms*; 2007: Society for Industrial and Applied
674 *Mathematics*.
675 45. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine
676 learning in Python. *the Journal of machine Learning research*. 2011;12:2825-30.
677 46. Blei DM, Ng AY, Jordan MI. Latent dirichlet allocation. *Journal of machine Learning research*.
678 2003;3(Jan):993-1022.
679 47. Rehurek R, Sojka P, editors. *Software framework for topic modelling with large corpora*. In *Proceedings*
680 *of the LREC 2010 Workshop on New Challenges for NLP Frameworks*; 2010: Citeseer.
681 48. Hoffman M, Bach FR, Blei DM, editors. *Online learning for latent dirichlet allocation*. *advances in neural*
682 *information processing systems*; 2010.
683 49. Röder M, Both A, Hinneburg A, editors. *Exploring the space of topic coherence measures*. *Proceedings*
684 *of the eighth ACM international conference on Web search and data mining*; 2015.
685 50. van der Hooft JJJ, Wandy J, Young F, Padmanabhan S, Gerasimidis K, Burgess KEV, et al.
686 Unsupervised Discovery and Comparison of Structural Families Across Multiple Samples in Untargeted
687 Metabolomics. *Anal Chem*. 2017;89(14):7569-77. Epub 2017/06/18. doi: 10.1021/acs.analchem.7b01391.
688 PubMed PMID: 28621528; PubMed Central PMCID: PMC5524435.



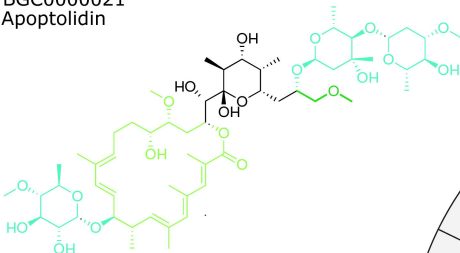
a

BGC length vs length of a match to a sub-cluster motif

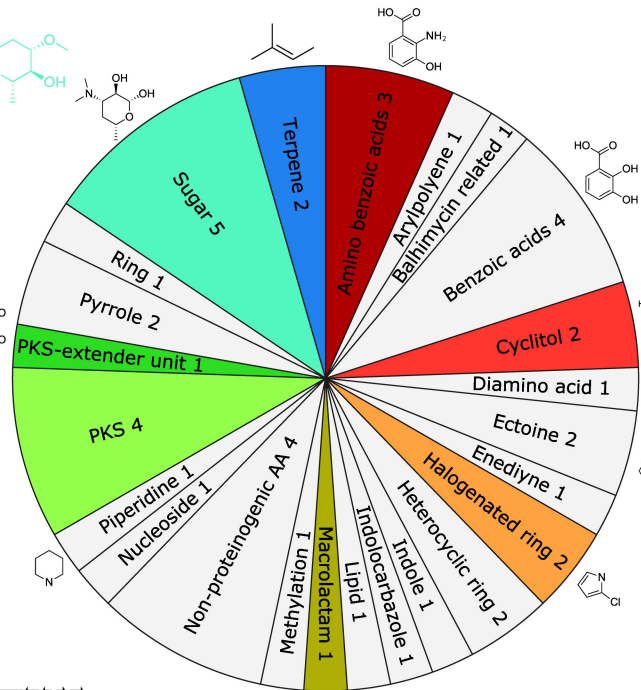
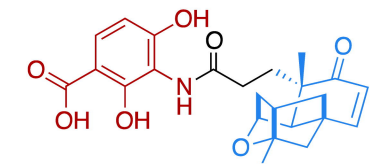
**b****c**



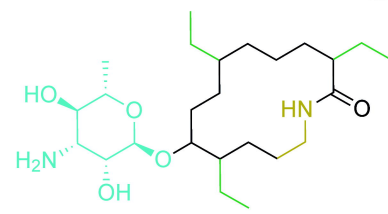
BGC0000021
Apoptolidin



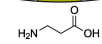
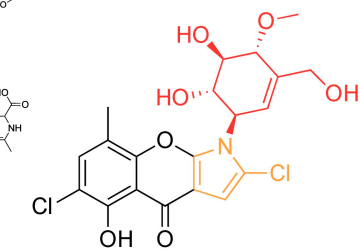
BGC0001140
Platencin

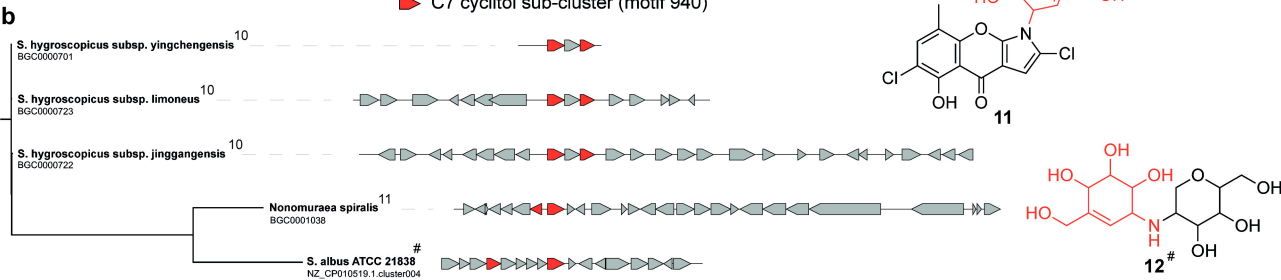
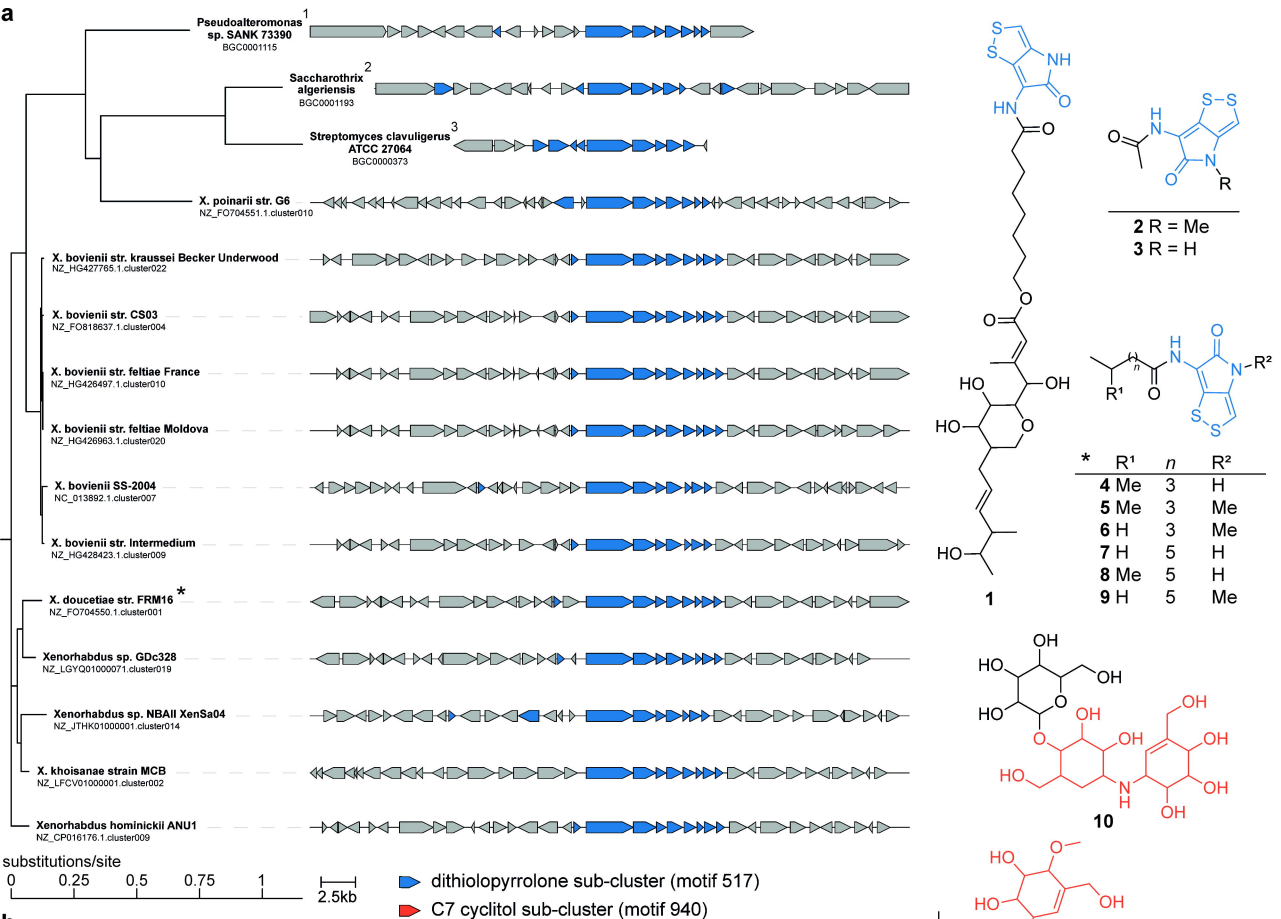


BGC0001597
Fluvirucin b2

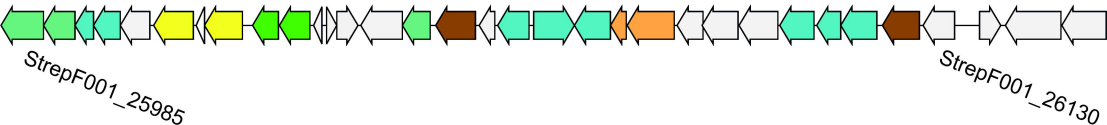


BGC0001038
Pyralomycin 1a

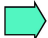





QZWF01000007.1.region003



iPRESTO predictions

 deoxy-aminosugar (motif 194 & 680)

 halogenated ring (motif 607)

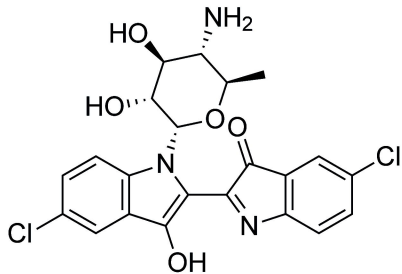
Other features

 sugar-related

 transport

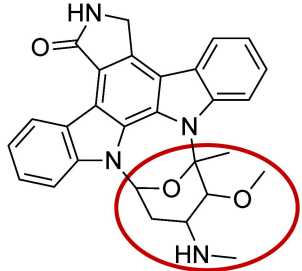
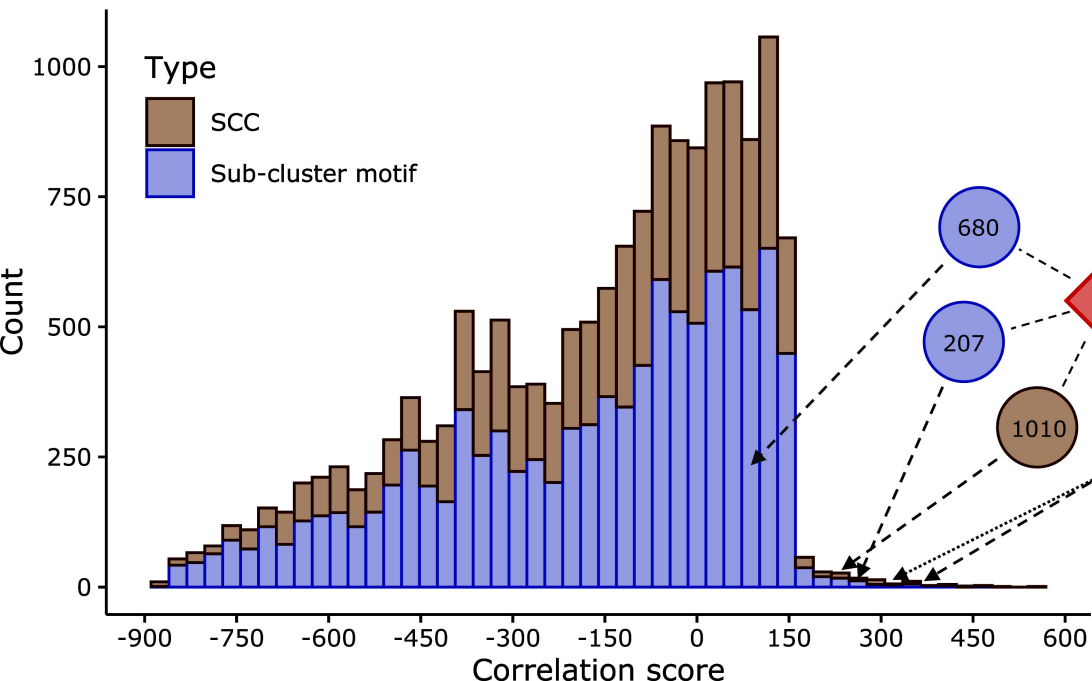
 p450

 oxidoreductase



akashin a

Correlation scores between mass2motifs and the two sub-cluster types



Staurosporine related mass2motifs

680

207

108

8

1010

452