

1 **Recapitulation of the evolution of biosynthetic gene clusters reveals hidden chemical**
2 **diversity on bacterial genomes**

3
4 Pablo Cruz-Morales^{1,*}, Christian E. Martínez-Guerrero¹, Marco A. Morales-Escalante¹, Luis
5 Yáñez-Guerra¹, Johannes Florian Kopp², Jörg Feldmann², Hilda E. Ramos-Aboites¹ & Francisco
6 Barona-Gómez^{1,*}

7
8
9 ¹ Evolution of Metabolic Diversity Laboratory, Langebio, Cinvestav-IPN. Irapuato, Guanajuato,
10 México.

11
12 ² Trace Element Speciation Laboratory (TESLA), College of Physical Sciences. Aberdeen,
13 Scotland, UK.

14
15 Authors for correspondence:

16 Francisco Barona-Gómez (fbarona@langebio.cinvestav.mx)

17 Pablo Cruz-Morales (pcruz@langebio.cinvestav.mx)

22 **Abstract**

23 Natural products have provided humans with antibiotics for millennia. However, a decline in the
24 pace of chemical discovery exerts pressure on human health as antibiotic resistance spreads. The
25 empirical nature of current genome mining approaches used for natural products research limits
26 the chemical space that is explored. By integration of evolutionary concepts related to emergence
27 of metabolism, we have gained fundamental insights that are translated into an alternative genome
28 mining approach, termed EvoMining. As the founding assumption of EvoMining is the evolution
29 of enzymes, we solved two milestone problems revealing unprecedented conversions. First, we
30 report the biosynthetic gene cluster of the ‘orphan’ metabolite leupeptin in *Streptomyces roseus*.
31 Second, we discover an enzyme involved in formation of an arsenic-carbon bond in *Streptomyces*
32 *coelicolor* and *Streptomyces lividans*. This work provides evidence that bacterial chemical
33 repertoire is underexploited, as well as an approach to accelerate the discovery of novel antibiotics
34 from bacterial genomes.

35

36

37 **Introduction**

38 The concept of genome mining can be defined as both an extension of the central dogma of
39 molecular biology, situating metabolites at the downstream end, and a novel approach that
40 promises to turn the discovery of natural products (NPs) drugs into a chance-free endeavor (1–3).
41 Within the context of increased antibiotic resistance and emergence of modern diseases, and given
42 the strong track of NPs (which are the result of specialized metabolism) in providing useful
43 molecules to human welfare (3, 4), genome mining has revitalized the investigation into NP
44 biosynthesis and their mechanisms of action (3, 5). Evidence of this resurgence has steadily
45 increased since the first NP that was discovered using genome mining approaches, i.e the
46 farnesylated benzodiazepinone ECO4601, entered into human clinical trials more than a decade
47 ago (6–8).

48 In contrast to experimental hurdles, which have been acknowledged elsewhere (9), *in silico*
49 genome mining has enjoyed a relatively higher success. Early genome mining approaches built up
50 from the merger between a wealth of genome sequences and an accumulated biosynthetic
51 empirical knowledge, mainly surrounding Polyketide Synthases (PKS) and Non-Ribosomal
52 Peptide Synthetases (NRPSs) (10, 11). These approaches, which rely on high quality genome
53 sequences due to the modularity and repetitive nature of PKSs and NRPSs, can be classified as: (i)
54 chemically-driven, where the structure of a metabolite is linked to potential enzymes, such that the
55 biosynthetic genes of an ‘orphan’ metabolite that has been isolated and structurally characterized,
56 are identified (12); or (ii) genetically-driven, where known sequences of protein domains (13) or
57 active-site motifs (14) help to identify putative biosynthetic gene clusters (BGCs) and their
58 products. The latter relates to the term ‘cryptic’ BGC, defined as a genetic locus that has been
59 predicted to direct the synthesis of a NP, but which remains to be experimentally confirmed (15).

60 Genome mining of NPs has also helped to prioritize strains and metabolites on which to
61 focus for further investigation. During this process, based on *a priori* biosynthetic insights,
62 educated guesses surrounding PKS and NRPS can be put forward, increasing the likelihood of
63 discovering interesting chemical and mechanistic variations. Moreover, biosynthetic logics for a
64 growing number of NP classes, such as phosphonates (16–18) and ribosomally synthesized post-
65 translationally modified peptides (RiPPs) (19) are complementing early NRPS/PKS-centric
66 approaches. Nevertheless, the current pace of deciphering novel biosynthetic logics, which can
67 only be achieved after long periods of research, hampers our ability to cope with the rate of
68 appearance of antibiotic resistance. Focusing on known chemical scaffolds with the concomitant
69 high rate of rediscovering the same classes of NPs, although useful under certain circumstances
70 and specific cases (20), seems also sub-optimal for the discovery of much-needed novel drugs.
71 From these observations it becomes apparent that more efficient approaches that will lead to the
72 discovery of novel chemistry are needed.

73 In this work we have developed an alternative method for current NP genome mining,
74 which is guided by evolutionary theory. By means of integrating three evolutionary concepts
75 related to emergence of specialized metabolism, we have gained fundamental insights that are
76 translated into the discovery of novel NPs. First, we embraced the concept that new enzymatic
77 functions evolve by retaining their reaction mechanisms, while expanding their substrate
78 specificities (21). In consequence, this process expands enzyme families. Second, evolution of
79 contemporary metabolic pathways frequently occurs through recruitment of existing enzyme
80 families to perform new metabolic functions (22). Indeed, in the context of NP biosynthesis, cases
81 of functional overlap driven by promiscuous enzymes that have been expanded and recruited have
82 been reported (23, 24). The correspondence of enzymes to either central or specialized metabolism

83 is typically solved through detailed experimental analyses, but we argue here that these could also
84 be achieved through phylogenomics. Third, BGCs are rapidly evolving metabolic systems,
85 consisting of smaller biochemical sub-systems or ‘sub-clusters’, which may have their origin in
86 central metabolism (25–27).

87 Integration of these three evolutionary principles was formalized as a bioinformatics
88 pipeline, termed EvoMining, which can be defined as a functional phylogenomics tool for
89 identification of expanded, repurposed enzyme families, with the potential to catalyze new
90 conversions in specialized metabolism (Figure 1A). As this process does not rely on sequence
91 similarity searches of previously identified NP biosynthetic enzymes, but rather on recapitulation
92 of an evolutionary process, the predictive power of evolutionary theory is fully embraced.
93 Moreover, given that predictions are done at the single-gene level, rather than looking at large
94 PKS, NRPS or BGC sequence assemblies, low quality draft genome and metagenome sequences
95 are compatible with this approach. Indeed, we demonstrate that EvoMining can predict
96 biosynthetic genes for orphan molecules, as well as new NP biosynthetic pathways in model
97 strains, both involving novel enzymatic conversions.

98 Experimentally, we focused in the phylum *Actinobacteria*, which includes renowned NP-
99 producing genera, such as *Streptomyces*, that have provided a plethora of useful NPs (3).
100 Experimental evidence for two key cases using *Streptomyces* species is reported to advance the
101 field of genome mining of NPs. First, the BGC for the biosynthesis of the orphan small peptide
102 aldehyde (SPA) leupeptin, of high economic importance, is identified in *Streptomyces roseus*
103 ATCC 31245. Second, a novel NP synthesized by the model organisms *Streptomyces coelicolor*
104 A3(2) and *Streptomyces lividans* 66 is experimentally characterized. Both of these case studies
105 include novel chemical conversions catalyzed by enzymes that were blindly targeted by

106 EvoMining. Therefore, these results validate EvoMining as an alternative and complementary
107 method for the discovery of potential drug leads. Moreover, the insights gained during this
108 integrative approach suggest that bacterial genomes encode a larger chemical diversity yet-to-be
109 discovered that can be untapped by means of using evolutionary theory.

110 **Results & Discussion**

111 The EvoMining bioinformatics pipeline, in its current version 1.0, uses three input
112 databases (green cylinders, **Figure 1A**). First, a genome database that contains the annotated
113 genomes of 230 members of the phylum *Actinobacteria*, as retrieved from the GenBank database
114 (**Table S1**). Second, a database containing the amino acid sequences of enzymes belonging to nine
115 ‘precursor supply central metabolic pathways’ (PSCP), defined as previously (28). This dataset
116 provides a universe of 103 enzyme families, to be used as query sequences (**Table S2**), which
117 were extracted from genome-scale metabolic network reconstructions (GSMR) of model
118 *Actinobacteria*. Third, a NP seed database consisting of 226 experimentally characterized BGCs
119 (mainly from *Actinobacteria*), including: (i) NRPSs and PKSs biosynthetic systems extracted from
120 specialized databases (10, 11); and (ii) other classes of well-described NPs biosynthetic pathways,
121 e.g. terpenes, phosphonates and RiPPs, extracted from the literature (**Table S3**).

122 The sequences in the PSCP database were used as queries to retrieve homologous
123 sequences contained in the genome database. The threshold used for defining homology was non-
124 conservative, such that expansion events resulting from both gene duplication and horizontal gene
125 transfer could be retrieved. When propagated through the genomes database, after homology
126 searches, these query sequences gave rise to an enzyme family internal database. After a heuristic
127 approach, an organism’s enzyme expansion was called by statistical measure when the number of
128 homologs on its genome was larger than the average of each enzyme family plus its standard

129 deviation. The enzymes that complied with this criterion were stored on the enzyme expansion
130 internal database (yellow cylinders, **Figure 1A**; **Table S4**). The expansion of each enzyme family
131 was sorted throughout a phylogenetic species tree (**Tree S1**), allowing taxonomic resolution of
132 expansions, as previously observed (28). With this approach we found that 98 enzyme families,
133 out of 103 enzymes from the PSCP database, had expansion events.

134 A critical function for the EvoMining approach is identification of enzyme families
135 expanded in concert with NP biosynthesis clusters. To accomplish this, the expanded enzyme
136 families were then mined for recruitment of their members within the context of NP biosynthesis.
137 In the cases where an expanded enzyme family could be connected via sequence homology to one
138 or more proteins within the NP seed database, their sequences were stored on the enzyme
139 recruitment internal database (yellow cylinders, **Figure 1A**, **Table S5**). It should be noted that
140 only a fraction of all sequences in the NP seed database have been characterized in the context of
141 NP biosynthesis. Often, not even functional annotation is provided. Therefore, the functional
142 association between recruited enzymes and this relatively large sequence space is supported by the
143 occurrence of the expanded enzymes within BGCs that have been linked to a known metabolite.
144 The enzyme recruitment internal database consisted of 23 enzyme families, including both known
145 recruitments, e.g. aconitase in phosphinothricin biosynthesis (16), and all related sequences
146 codified by the analyzed genomes. Thus, the functional potential of the NP seed and the genome
147 databases, together, is fully exploited.

148 The sequences of the recruited enzymes, together with those from the expanded enzyme
149 families, were used to make multiple sequence alignments (MSA) and Bayesian phylogenetic
150 reconstructions (**Figure 1A**). Moreover, in order to provide useful functional annotation for
151 interpretation of EvoMining phylogenetic trees (**Figure 1B**), a bidirectional best-hit analysis,

152 between the enzyme recruitment and the PSCP databases, was used for directing the labeling of
153 central metabolic orthologs (red branches). We adopted this simple strategy, as it is safe to assume
154 that NP biosynthetic enzymes will diverge significantly from central metabolic homologs.
155 Homologs related to the very few known recruitments (blue branches) were therefore considered
156 to be NP biosynthetic homologs identified by EvoMining. This proof-of-concept analysis provided
157 515 recruitment events, which are called EvoMining hits, and their gene identifiers (GIs) were
158 used as queries to retrieve contigs (12-109 Kbp, 71.3 Kbp in average, 19.9 Kbp standard
159 deviation). The retrieved contigs were then analyzed for putative NP BGCs using antiSMASH and
160 ClusterFinder (29, 30). When an NP positive hit was obtained after this process, this was also
161 noted in the phylogenetic tree (cyan branches).

162 The abovementioned functional annotation provides information that validates NP-related
163 phylogenetic clades that consist of EvoMining hits (**Figure 2B**). Subtraction of the known (blue
164 branches) and antiSMASH / ClusterFinder predicted (cyan branches) NP lineages, within the NP-
165 related clades, reveals putative BGCs coding for repurposed enzymes only accessible by
166 EvoMining (green branches). Henceforth, we refer to these homologs as EvoMining predictions,
167 which we define as unknown NP biosynthetic enzymes supported by phylogenetic evidence
168 encoded within previously undetected BGCs. Chemically, an implication of an EvoMining hit is
169 that it uncovers enzymatic conversions, mainly involving diverging substrate specificities (but
170 potentially also mechanistic variations), which in turn can lead to alternative biosynthetic logics
171 and therefore chemical scaffolds. Thus, an EvoMining prediction, as a concept, is here applied to
172 entire BGCs rather than to specific enzymes. In addition to **Figure 1**, EvoMining version 1.0 can
173 be explored as online supplementary material at
174 http://148.247.230.39/newevomining/new/evomining_web/index.html

175 **Evolutionary insights and performance of EvoMining**

176 Of the 515 EvoMining hits we successfully retrieved contigs containing their cognate
177 enzyme coding genes for 448 of them (71.3 Kbp on average; 87 %) (**Table S5**). Among these,
178 many EvoMining hits (20 %) were included in contigs with internal gaps, which hampers
179 sequence annotation. This subset, together with the remaining 13 % of the total hits whose contigs
180 could not be retrieved, account for one third of contigs that come from highly fragmented
181 genomes. So, the advantage of EvoMining in this respect is that predictions can be made, early on
182 during analysis, in draft genomes (and metagenomes) that can be further improved. This provides
183 an opportunity to prioritize in cost-effective manner large strain collections during genome-driven
184 drug discovery efforts. Focusing on an EvoMining hit related to the enolase enzyme family (**Tree**
185 **S2**), which was found in the genome of *Streptomyces sviveus* (1 scaffold of 9 Mbp with 552 gaps
186 and 8X coverage, GI: 297196766), illustrates the benefit of EvoMining in this respect. The contig
187 containing this recruited enolase (GI: 297146550) had 6 gaps including missing sequence at its 5'
188 end. After closing of these gaps by sequencing PCR products, the complete sequences for several
189 phosphonate-related enzymes, namely, alcohol dehydrogenase (*phpC*), phosphonopyruvate
190 decarboxylase (*ppd*), nicotinamide mononucleotide adenylyl transferase (*phpF*), carboxy-
191 phosphoenolpyruvate synthase (*phpH*, EvoMining hit) and aldehyde dehydrogenase (*phpJ*),
192 could be annotated. Further sequence analysis suggested that indeed this locus encodes for a
193 putative phosphonate BGC related to phosphinothricin (31–33) (**Text S1, Figures S1 and S2**).

194 We further asked the question of whether EvoMining enzymes are indeed encoded within
195 BGCs potentially directing the synthesis of NPs. The retrieved contigs were mined for BGCs of
196 known classes of NPs using antiSMASH (29) as well as for putative new BGCs using
197 ClusterFinder (30). From this analysis we found that these tools could predict BGCs for 62.5 %

198 and 10.5 %, respectively, of the contigs harboring EvoMining hits (73 % together). The remaining
199 27 % of contigs are unique EvoMining hits, and therefore potentially EvoMining predictions
200 related to emerging BGCs and chemical scaffolds (**Figure 2A**). The enzymes included in this 27
201 % represent the core of EvoMining, and due to the lack of functional information subsequent
202 analysis is experimentally and bioinformatically challenging. However, manual analysis of the
203 green branches within the NP-related clade provided as an example in **Figure 1B**, allowed us to
204 identify a highly conserved BGC-like locus present in the genus *Streptomyces* (see below
205 discussion related to **Figure S7, Table S7 and S8**).

206 As a prerequisite for embarking on characterization of completely unprecedented BGCs we
207 first determined whether EvoMining hits are associated to particular BGC classes. For this
208 purpose, we used antiSMASH to classify and count for the number of BGCs contained in each
209 contig. Globally, BGCs for 22 out of the 24 categories used by antiSMASH (34) could be
210 detected. Only aminoglycosides and indoles could not be detected, which may be due to the
211 limited enzymatic repertoire explored by our seed NP and PSCP input databases. Among the 22
212 antiSMASH categories, type I PKSs and NRPSs represent the most abundant classes (**Figure S3**),
213 confirming that EvoMining can identify well-known NP biosynthetic systems following a unique,
214 non-biased strategy. Despite this convergence, at least with this limited analysis, it was also found
215 that the number of EvoMining hits increased in parallel to the number of BGC classes (**Figure**
216 **2B**). The implication of this observation is that novel classes may be discovered for any given
217 enzyme recruitment, as long as enough sequence space is explored.

218 For instance, among the 23 recruited enzyme families, only that of indole-3-
219 glycerolphosphate synthase was linked to a single class of BGCs. It should be noted, however, that
220 this family has the smallest number of EvoMining hits. At the other end, the enzyme families with

221 the larger number of EvoMining hits showed the highest BGC diversity, namely, 17 BGC classes
222 for the asparagine synthetase enzyme family, followed by 10 BGC classes for both 3-
223 phosphoshikimate-1-carboxyvinyl-transferase and 2-dehydro-3-deoxyphosphoheptonate aldolase
224 families (**Figure 2B**). It seems, therefore, that some of the recruitments have high predictive
225 potential when used as beacons for detection of BGCs. The latter may be due to the evolvability of
226 enzymes towards specificities for common precursors of NPs, driven by enzyme promiscuity (23,
227 24), or because these enzymes catalyse recurrent reactions in NP biosynthesis with less
228 mechanistic restrictions (27).

229 After our *in silico* analysis, we aimed to demonstrate that EvoMining could indeed be used
230 to predict enzymes of unknown function involved in the synthesis of NPs. More specifically, we
231 focused in providing experimental evidence to support two critical cases, which could not be
232 solved with current methods: (i) the discovery of a BGC of an orphan metabolite that has been
233 extensively investigated; and (ii) the discovery of a BGC driving the synthesis of an NP produced
234 by well-studied model strains. Remarkably, and in agreement with the definition of an EvoMining
235 hit, the enzymes identified by this approach are proposed to catalyse unprecedented chemical
236 conversions.

237

238 **Discovery of leupeptin BGC in *Streptomyces roseus*: the orphan metabolite problem**

239 Leupeptin is the first member of a large family of NPs generically known as small
240 peptide aldehydes (SPA), and it is widely used in industry and bioresearch due to its potent anti-
241 proteolytic activity (35). Despite the fact that leupeptin was discovered during the golden age of
242 antibiotic research (35) its BGC remained elusive until now (**Figure 3**). The structure of leupeptin

243 includes a C-terminal aldehyde group in a peptide chain that includes an acyl group at its N-
244 terminal end. Early biochemical studies of leupeptin revealed that: (i) this peptide is produced
245 from acetate, leucine and arginine; (ii) an ATP-dependent synthetase is responsible for the
246 condensation of an acetyl-leucine-leucine intermediary; (iii) this intermediary is released from an
247 enzymatic complex before its condensation with an arginine residue; and (iv) the enzymatic
248 complex responsible for the synthesis of acetyl-leucine-leucine-arginine, called leupeptidic acid,
249 has a molecular mass of approximately 320 kDa (36–39).

250 A more recent study on the biosynthesis of the SPA flavopeptin produced by
251 *Streptomyces* sp. NRRL-F6652 (40), suggested that synthesis of leupeptin by producing strains,
252 such as *S. roseus* ATCC 31245 (35), may occur via a distinctive NRPS complex. Key features of
253 this putative NRPS were proposed to be: (i) an acyl transfer domain typically found as an N-
254 terminal starter C-domain, responsible for initial acylation of non-ribosomal peptides (41); (ii)
255 three complete modules for peptide synthesis, two for leucine and one for arginine residues; and
256 (iii) a reductase domain at its C-terminal end, responsible for reductive peptide release, leading to
257 an aldehyde. Based in this modern proposal, after sequencing the genome of *S. roseus* ATCC
258 31245 (BioProject: PRJNA287805, Accession: LFML00000000), we searched for flavopeptin-like
259 NRPSs. However, this approach proved to be unsuccessful.

260 Mining of the *S. roseus* genome sequence with EvoMining, by contrast, produced a hit to
261 an enzyme annotated as argininosuccinate lyase (ASL), typically involved in arginine
262 biosynthesis. ASLs condense succinate and arginine via an amidic bond at the guanidine group of
263 the arginine in a reversible reaction (42). Indeed, two ASL homologs sharing 25 % amino acid
264 sequence identity, which could be phylogenetically resolved, were found in the genome of *S.*
265 *roseus* (**Figure 3A** and **Tree S3**). The first homolog, as expected, is located within a clade that
266 includes central metabolic homologs (ASL1 or *argH* gene); whereas the second homolog is

267 located in a clade together with enzymes previously related to the biosynthesis of uridyl-peptides,
268 namely, napsamycin and pacidamycin (43, 44). Thus, the members of the latter clade, termed here
269 ASL2, are predicted to include recruited enzymes involved in NP biosynthesis, and that these
270 enzymes are performing a related but different chemistry than that catalyzed by ASL1.

271 After detailed annotation of the region surrounding the recruited ASL2 gene, a NRPS
272 gene was found (**Figure 3A**). The product of this gene was predicted to have an N-terminal
273 condensation domain (C₁), followed by an adenylation domain predicted to bind threonine (A₁); a
274 peptidyl carrier protein domain (PCP₁); a second condensation domain (C₂); an adenylation
275 domain predicted to bind serine (A₂); a second peptidyl carrier protein (PCP₂); and a thioesterase
276 domain (TE) (**Figure 3B**). On the basis of this annotation, which drastically differs from the
277 prediction based in the flavopeptin system, we predicted that this NRPS would produce an
278 acylated dipeptide, which is released upon the action of the thioesterase domain. This biosynthetic
279 logic is indeed consistent with the early biochemical data available for leupeptin (36–39).

280 We therefore renamed the EvoMining hit as *leupB*, and the NRPS as *leupA*, and a
281 functional association between LeupA and LeupB was assumed. Specifically, we speculated that
282 LeupB is capable of condensing the dipeptide produced by LeupA, possibly acyl-Leu-Leu, with an
283 arginine residue, leading to leupeptidic acid or acyl-Leu-Leu-Arg peptide (**Figure 3C**). In
284 accordance with this hypothesis the total predicted mass of LeupA and LeupB is 318 kDa, which
285 is strikingly close to that mass of 320 kDa early on estimated by Umezawa and co-workers for the
286 complex directing leupeptidic acid formation (39). Additional genes downstream of *leupB*,
287 transcribed in the same direction and therefore possibly involved in leupeptin biosynthesis include
288 *leupC* and *leupD*, are annotated as a threonine kinase and cysteine synthase, respectively (**Figure**
289 **3B**). The relevance of this functional annotation remains to be further investigated, but may have
290 to do with reduction of the leupeptidic acid.

291 To demonstrate that the predicted locus is involved in leupeptin biosynthesis, we used
292 insertional mutagenesis to disrupt the *leupA* gene. We chose this rather simple approach due to the
293 lack of genetic tools for manipulation of *S. roseus*, and because this gene could be considered to
294 be essential for synthesis of leupeptin. We did not target *leupB*, as the central ASL homolog may
295 complement its function via enzyme promiscuity, masking the expected phenotype (unless a
296 double mutant is obtained), as previously reported in other BGCs (23). The *S. roseus leupA*
297 mutant, termed PCMSr1, was obtained, and after comparative LC-MS analysis of this mutant with
298 the parental wild-type strain, it was found that mutation of *leupA* renders PCMSr1 unable to
299 produce leupeptin. The peak absent from the LC chromatograms obtained from PCMSr1, present
300 in the wild type strain, corresponds with leupeptin authentic standard, as confirmed after
301 electrospray ionization (ESI) mass spectrometry (m/z 427) and fragmentation pattern (ms^2 of 367
302 and 409) (**Figures S4A** and **S5**).

303 To further establish a link between leupeptin biosynthesis and the postulated locus,
304 currently involving *leupA-D* (**Figure 3**), we constructed a genomic library from which two clones,
305 containing at least these four genes, were isolated. Both constructs were introduced into *E. coli*
306 DH10B, and the resulting transformants were used for comparative fermentation experiments. LC-
307 MS analysis of these cultures revealed that extracts of supernatants from both strains presented
308 fractions with peaks with same retention times and expected mass as authentic leupeptin standard
309 (**Figures S4B** and **S5**). Therefore, we concluded that this locus, as predicted using EvoMining, is
310 indeed directing the synthesis of leupeptin in *S. roseus* ATCC31245. However, given that the
311 constructs bearing the *leupABCD* genes are 65 to 70 Kbp in size (clone 8_10B and 9_18N,
312 respectively), the involvement of other genes in leupeptin biosynthesis, which are included in
313 these constructs, cannot be ruled out at current time.

314

315 **Discovery of an arseno-organic metabolite in *Streptomyces lividans* and *Streptomyces***
316 ***coelicolor*: novel chemistry in model organisms.**

317 For the second case, we aimed to identify a novel NP in the model organisms *S. coelicolor*
318 A3(2) and *S. lividans* 66, which have been mined thoroughly and presumably most of their NP
319 repertoire has been elucidated (45). Furthermore, several methods for genetic manipulation for
320 these two strains are available (46), making these organisms ideal for these proof-of-concept
321 experiments. EvoMining hits for these two strains include recruitments belonging to the 3-
322 phosphoshikimate-1-carboxyvinyltransferase enzyme family, which was used as an example in
323 Figure 1B. This enzyme, or AroA, catalyzes the transfer of a vinyl group from
324 phosphoenolpyruvate (PEP) to 3-phosphoshikimate forming 5-enolpyruvylshikimate-3-phosphate
325 and releasing phosphate. The reaction is part of the shikimate pathway, a common pathway for the
326 biosynthesis of aromatic amino acids and other metabolites (47).

327 The phylogenetic reconstruction of the actinobacterial AroA enzyme family (**Fig 1B** and
328 **Tree S4**), as expected, shows a major clade associated with central metabolism; this clade
329 includes SLI_5501 from *S. lividans* and SCO5212 from *S. coelicolor*. The phylogeny also has a
330 divergent clade that includes two family members linked to the BGCs of the polyketide
331 asukamycin (48) and phenazines (49), as well as AroA homologs from 26 % of the genomes in the
332 database. In *S. coelicolor* and *S. lividans* these recruited homologs are encoded by SLI_1096 and
333 SCO6819, respectively. Moreover, in *S. lividans*, these orthologous genes are located within a
334 large genomic island, SliGI-1, which has been functionally linked to metal homeostasis (50), and
335 they are situated only six genes upstream of a two-gene PKS system spanning SCO6826-7 and
336 SLI_1088-9, respectively. This PKS was identified in *S. coelicolor* since the early days of the
337 genome mining of this organism, but often it was referred to as a cryptic BGC (51, 52). Indeed,

338 the divergent *AroA* homologs have not been associated with it. Furthermore, these homologs are
339 classified as “other genes” when both genomes are mined using antiSMASH.

340 The gene neighborhood of these *aroA* genes is highly conserved between the genomes of
341 *S. lividans* and *S. coelicolor*. Thus, from this point onwards we will refer to the *S. lividans* genes
342 only. The syntenic region spans from SLI_1077 to SLI_1103, including several biosynthetic
343 enzymes, regulators, transporters and the PKS, suggesting that these genes, together with other
344 biosynthetic genes in this locus, are functionally linked and form a single BGC (**Figure 4A**).
345 Detailed annotation of this BGC (**Table S5**) revealed the presence of a 2,3-bisphosphoglycerate-
346 independent phosphoenolpyruvate mutase enzyme (SLI_1097; PPM), downstream and possibly
347 transcriptionally coupled to the *aroA* homolog. Thus, a functional link between these genes, as
348 well as with phosphonopyruvate decarboxylase gene (PPD; SLI_1091) encoded in this BGC, was
349 proposed. The combination of mutase-decarboxylase enzymes is a conserved biosynthetic feature
350 of NPs containing Carbon-Phosphate bonds (16).

351 Other non-enzymatic functions were found encoded within this BGC, including a set of
352 ABC transporters, originally annotated as phosphonate transporters (SLI_1100 and SLI_1101).
353 Four arsenic tolerance-related genes (SLI_1077-1080) located upstream of the PKS could also be
354 annotated. These genes are paralogous to the main arsenic tolerance system encoded by the *ars*
355 operon (53), which is located at the core of the *S. lividans* chromosome (SLI_3946-50). This BGC
356 also codes for regulatory proteins, mainly arsenic responsive repressor proteins (SLI_1078,
357 SLI_1092, SLI_1102 and SLI_1103). Thus, overall, our detailed annotation suggests a link
358 between arsenic and phosphonate biosynthetic chemistry. Accordingly, in order to reconcile the
359 presence of phosphonate-like biosynthetic, transporter and arsenic resistance genes, within a BGC,

360 we postulated a biosynthetic logic analogous to that of phosphonate biosynthesis, but involving
361 arsenate as the driving chemical moiety.

362 The abovementioned hypothesis was further supported by the three following observations.
363 First, arsenate and phosphate are similar in their chemical and thermodynamic properties, causing
364 phosphate and arsenate utilizing enzymes to have overlapping affinities and kinetic parameters
365 (54, 55). Second, previous studies have demonstrated that AroA is able to inefficiently catalyze a
366 reaction in the opposite direction to the biosynthesis of aromatic amino acids, namely, the
367 formation of PEP and 3-phosphoshikimate from enolpyruvyl shikimate 3-phosphate and
368 phosphate (47). Indeed, arsenate and enolpyruvyl shikimate 3-phosphate can react to produce
369 arsenoenolpyruvate (AEP), a labile analog of PEP, which is spontaneously broken down into
370 pyruvate and arsenate (47). Third, it has been demonstrated that the phosphoenolpyruvate mutase,
371 PPM, an enzyme responsible for the isomerization of PEP to produce phosphonopyruvate, is
372 capable of recognizing AEP as a substrate. Although at low catalytic efficiency, the formation of
373 3-arsenopyruvate by this enzyme, a product analog of the phosphonopyruvate intermediate in
374 phosphonate NPs biosynthesis (16), has been reported (56).

375 The previous evidence was used to postulate a novel biosynthetic pathway encoded by
376 SLI_1077-SLI_1103. A putative arseno-organic product synthesized by this pathway may
377 resemble the structural characteristics and properties of a phospholipid. A detailed functional
378 annotation, and biosynthetic proposal, is provided as supplementary information (**Figure S6** and
379 **Table S6**). To determine the product of the predicted BGC we used expression analysis, as well as
380 comparative metabolic profiling of wild type and mutant strains, in both *S. lividans* and *S.*
381 *coelicolor*. Using RT-PCR analysis, we first determined the transcriptional expression profiles of
382 the PKS (SLI_1088), *aroA* (SLI_1096), one of the *arsR*-like regulator (SLI_1103), and the

383 periplasmic-binding protein of the ABC-type transporter (SLI_1099). As expected for a cryptic
384 BGC, the results of these experiments demonstrate that the proposed pathway is repressed under
385 standard laboratory conditions. We then analyzed the potential role of arsenate as an inducer of the
386 expression of this BGC, either alone or in combination with phosphate deprivation. Indeed, we
387 found that the analyzed genes were induced when *S. lividans* 66 was grown in the presence of 500
388 μM of arsenic and 3 μM of phosphate (**Figure 4B**).

389 In parallel, we used PCR-targeted gene replacement to produce mutants of the SLI_1096
390 and SCO6819 genes, and analyzed the phenotypes of the mutant and wild type strains on a
391 combined arsenate/phosphate gradient, i.e. low phosphate and high arsenate, and *vice versa*. After
392 this, we cultivated the wild type and mutant strains in liquid cultures, with and without arsenic,
393 during 14 days to obtain enough biomass for chemical analysis. Organic extracts from the pellets
394 of these cultures were analyzed using HPLC coupled with an ICP-MS calibrated to detect arsenic-
395 containing molecular species. Simultaneously, a high-resolution mass spectrometer determined the
396 mass over charge of the ions detected by the ICP. This set up allows for high-resolution detection
397 of arseno-organic metabolites (57). Using this approach, we detected the presence of arseno-
398 organic metabolites in the organic extracts from the pellets of both wild-type *S. coelicolor* and *S.*
399 *lividans*, with m/z of 331.1248, 333.1041, and 351.1147 (**Figure 4C** and **Figure S6**). These
400 metabolites could not be detected in either identical extracts from wild type strains grown in the
401 absence of arsenate or in the mutant strains deficient for the SLI_1096/SCO6819 genes. Thus, the
402 product of this pathway may be a relatively polar arsenolipid (**Figure S6**). The actual structures of
403 these products are still subject to further investigation and will be discussed in detail in a future
404 publication.

405

406

407 **EvoMining and its future impact into NP genome mining**

408 On one hand, confirmation of a link between SLI_1096/SCO6819 and the synthesis of an
409 arseno-organic metabolite provides an example on how genome-mining efforts, based in novel
410 enzyme sequences, can be advanced. For instance, co-occurrence of divergent SLI_1096
411 orthologs, now called arsenoenolpyruvate synthases (AEPS); arsenopyruvate mutase (APM) and
412 arsonopyruvate decarboxylase (APD), can be used as beacons to mine publically available
413 bacterial genomes. Indeed, thirteen BGCs with the potential to synthesize arseno-organic
414 metabolites, all of them encoded in genomes of myceliated *Actinobacteria*, were identified after
415 sequence similarity searches using the non-redundant GenBank database. The divergence and
416 potential chemical diversity within these arseno-related BGCs was characterized after a
417 phylogenetic analysis, using as matrix all conserved genes of these BGCs (**Figure 5**). This
418 analysis suggests three possible sub-classes with distinctive features, PKS-related, PKS-
419 independent and PKS/NRPS hybrid, which warrant further investigation.

420 On the other hand, once an EvoMining hit is validated as an enzyme with an
421 unprecedented function in the context of a hidden BGC, i.e. an EvoMining prediction, this could
422 lead to identification of novel classes of conserved BGCs. To illustrate this, we focused in the
423 AroA EvoMining tree (Figure 1B). The green branches within the NP-related clade of this tree
424 were manually curated in search for a conserved BGC (**Figure 1B**). One particular case was found
425 to appear frequently, and thus its genes were annotated in detail. This BGC was found to be
426 conserved in at least sixty-three *Streptmyces* genomes, and in the genome of *Microtetraspora*
427 *glauca* NRRL B-3735, included in the non-redundant GenBank database (**Table S8**).

428 For these analyzes, fifteen genes upstream and downstream the EvoMining AroA hit were
429 extracted as before, and annotated on the basis of the locus from *S. griseolus* NRRL B-2925, as
430 this organism provides a condensed version of this locus (**Table S7**). Indeed, this locus has some
431 of the expected features for an NP BGC, including: (i) gene organization suggesting an assembly
432 of operons, most of them transcribed in the same direction; (ii) genes encoding for enzymes,
433 regulators and potential resistance mechanisms; and (iii) enzymes that have been found in other
434 known NP BGCs. The latter observation, which actually includes seven genes out of thirty-one,
435 present in an equal number of NP BGCs, actually confirms the NP nature of this locus. The reason
436 why EvoMining did not lead to these homologs as expanded and recruited enzymes has to do with
437 the fact that none of these enzyme families are included within the limited space explored by our
438 PSCP database.

439 Thus, we conclude that EvoMining has great potential to ease natural product and drug
440 discovery by means of accelerating the conceptual genome-mining loop that goes from novel
441 enzymatic conversions, their sequences, and propagation after homology searches.

442 **Methods**

443 **Bioinformatics**

444 *Seed NP database*: NRPS and PKS BGCs were obtained from the DoBISCUIT and ClusterMine
445 360 databases (10, 11). BGCs for other NP classes were collected from available literature (**Table**
446 **S1**). The database included amino acid fasta sequences from GenBank, DoBISCUIT and
447 ClusterMine360. Annotated GenBank formatted files were downloaded from the GenBank
448 database to assemble a database that included 226 BGCs. *Genome database*: Complete and draft
449 genomes of 230 members of the *Actinobacteria* family (**Table S1**) were retrieved from the

450 GenBank either as single contigs or as groups of contigs in GenBank format, amino acid and DNA
451 sequences were extracted from these files using in-house made scripts. *Precursor supply central*
452 *pathway (PSCP) database*: the amino acid sequences from the proteins involved in central
453 metabolism were obtained from a database that we assembled for a previous enzyme expansion
454 assessment, published elsewhere (28). The final database for this work included a total of 339
455 queries for nine pathways, including amino acid biosynthesis, glycolysis, pentose phosphate
456 pathway and tricarboxylic acids cycle (**Table S2**).

457 *EvoMining pipeline*: The PSCPs database was used as query to retrieve PSCP enzyme families
458 from the genome database using BlastP (65), with an e-value cutoff of 0.0001 and a score cutoff
459 of 100. At least three query sequences, representing each of the GSMRs used as sources of these
460 sequences (28), were used for Blast searches. The average number of homologs of each enzyme
461 family per genome and the standard deviation were calculated to establish a cutoff to identify and
462 highlight significant expansion events (**Table S4**). An enzyme family expansion was scored if the
463 number of homologs in at least one genome was higher than the average number of homologs plus
464 a standard deviation unit. Enzyme families with expansions were used for the next step of the
465 analysis. To identify enzyme recruitments, the amino acid sequences of the seed NP database were
466 used as queries for BlastP searches against expanded enzyme families identified in the previous
467 step using an e-value cutoff of 0.0001 and a bit score cutoff of 100. These parameters were
468 consistent with the approach used by EvoMining, as confirmed heuristically.

469 The homologs found in known BGCs were added as seeds to the sequences from the
470 expanded enzyme families with recruitments for future clade identification and labeling. These
471 sets of sequences were aligned using Muscle version 3.8.31 (58). The alignments were inspected
472 and curated manually using JalView (59). The curated alignments were used for phylogenetic
473 reconstructions, which were estimated using MrBayes (60) with the following parameters:

474 aamodelpr=mixed, samplefreq=100, burnfrac=0.25 in four chains and for 1000000 generations.
475 The bidirectional best hits with the sequences in the PSCP database were identified and tagged
476 using in-house scripts to distinguish PSCP orthologs from other homologs that result from
477 expansion events. The gene identifiers (GIs) of the EvoMining hits were used as queries to
478 retrieve their genome context as DNA regions of approximately 80 Kbs including the EvoMining
479 hit coding genes. These contigs were retrieved from the genomes in GenBank format using an in-
480 house script. These contigs were annotated using antiSMASH (29) through its web interface. The
481 whole process was executed semi-automatically using in-house scripts written in Perl.

482

483 **Mutagenesis analysis**

484 *S. coelicolor* (SCO6819) and *S. lividans* 66 (SLI_1096) knock-out mutants were constructed using
485 in-frame PCR-targeted gene replacement of their coding sequences with an apramycin resistance
486 cassette (*acc(3)IV*) (61). The plasmid pIJ773 was used as template to obtain a mutagenic cassette
487 containing the apramycin resistance marker by PCR amplification with the primers reported in
488 **Table S9**. The mutagenic cassettes were used to disrupt the coding sequences of the genes of
489 interest from the cosmid clone 1A2 that spans from SCO6971 to SCO6824 (62). Given the high
490 sequence identity between the regions covered by cosmid 1A2 with the orthologous region in *S.*
491 *lividans*, this cosmid clone was also used for disruption of SLI_1096. The gene disruptions were
492 performed using the Redirect system reported elsewhere (61). Double cross-over ex-conjugants
493 were selected using apramycin resistance and kanamycin sensitivity as phenotypic markers. The
494 genotype of the clones was confirmed by PCR. The strains and plasmids of the Redirect system
495 were obtained from the John Innes Centre (Norwich, UK).

496 The *S. roseus leupA* mutant was constructed following an insertional mutagenesis strategy.
497 For this purpose, a fragment of 640 bp was amplified by PCR and cloned in the vector pCR2.1-
498 TOPO (ampicillin/ kanamycin resistance) using a TA cloning kit from Invitrogen (Carlsbad,
499 USA), to produce the suicide plasmid pLEUPA that cannot be replicated in *S. roseus*. This
500 plasmid was introduced into *S. roseus* via protoplasts, generated following standard protocols. The
501 transformants were selected using kanamycin (50 μ /mL) and the genotype of the insertional
502 mutants was confirmed by PCR.

503

504 **Transcriptional analysis**

505 The *S. lividans* 66 wild type strain was grown on 0 and 3; 0 and 300; 500 and 3; 500 and 300 μ M
506 of Arsenate and KH_2PO_4 respectively in solid modified R5 media for eight days. The complete
507 culture conditions used in this work are further detailed in a following section. Mycelium
508 collected from plates was used for RNA extraction with a NucleoSpin RNA II kit (Macherey-
509 Nagel). The RNA samples were used as template for RT-PCR using the one step RT-PCR kit
510 (Qiagen) (2ng RNA template for each 40 μ l reaction). The housekeeping sigma factor *hrdB*
511 (SLI_6088) was used as a control.

512

513 ***S. roseus* genome sequencing and library construction**

514 *S. roseus* ATCC31245 was obtained from the ATCC collection, and its genomic DNA was
515 extracted using common protocols (46) and sequenced at the genomic sequencing facilities of
516 Langebio, Cinvestav-IPN (Irapuato, Mexico), using an Illumina MiSeq platform in paired-end
517 format with read lengths of 250 bases and insert length of 800 bases. In total, 721 Mbp of

518 sequence was obtained. The raw reads were filtered using Trimmomatic (63) and assembled with
519 velvet (64), obtaining a 7.8 Mb assembly in 165 contigs with a coverage of 95 X and a GC content
520 of 72 %. This assembly was annotated using RAST (65) antiSMASH (29) and EvoMining. A
521 genomic library of *S. roseus* ATCC31245 was further obtained for cloning into the pESAC13A
522 vector with an average insert length of 70 Kbps (Bio S&T, Montreal, Canada). pESCA13A is a
523 derivative from pPAC-S1 (66), which has an apramycin resistance as selection marker. This
524 library was screened for the *leup* locus by PCR, leading two clones named 9_18N and 8_10B,
525 containing the desired region. As described further, these constructs were used for heterologous
526 expression in *E. coli*.

527

528 **LC-MS metabolite profile analysis**

529 The SLI_1096 and SCO6818 minus mutants were grown on modified R5 medium (K_2SO_4
530 0.25 gr; $MgCl_2 \cdot 6H_2O$ 10.12 gr; glucose 10 gr; casamino acids 0.1 gr; TES buffer 5.73 gr; trace
531 element solution (46) 2 ml; agar 20gr) supplemented with a gradient of KH_2PO_4 and Na_3AsO_4
532 ranging from 3 to 300 μM and 0 to 500 μM , respectively. Induction of the arseno-organic BGC in
533 both strains was detected in the condition where phosphate is limited and arsenic is available.
534 Therefore, modified R5 liquid media supplemented with 3 μM KH_2PO_4 and 500 μM Na_3AsO_4 was
535 used for production of arseno-organic metabolites, and the cultures were incubated for 14 days in
536 shaken flasks with metal springs for mycelium dispersion at 30 C. The mycelium was obtained by
537 filtration, and the filtered mycelium was washed thoroughly with deionized water and freeze-
538 dried. The samples were extracted overnight twice with MeOH/DCM (1:2). The extracts were
539 combined and evaporated to dryness, and the dry residues were re-dissolved in 1 mL of MeOH
540 (HPLC-Grade) and injected to the HPLC. The detection of organic arsenic species was achieved

541 by online-splitting of the HPLC-eluent with 75% going to ESI-Orbitrap MS (Thermo Orbitrap
542 Discovery) for accurate mass analysis and 25% to ICP-QQQ-MS (Agilent 8800) for the detection
543 of arsenic. For HPLC, an Agilent Eclipse XDB-C18 reversed phase column was used with a H₂O
544 /MeOH gradient (0-20 min: 0-100% MeOH; 20-45 min: 100% MeOH; 45-50 min: 100% H₂O).
545 The ICP was set to oxygen mode and the transition $^{75}\text{As}^+ \rightarrow (^{75}\text{As}^{16}\text{O})^+$ (Q1: m/z = 75, Q2: m/z =
546 91) was observed. The correction for carbon enhancement from the gradient was achieved using a
547 mathematical approach as described previously (57). The ESI-Orbitrap-MS was set to positive ion
548 mode in a scan range from 250-1100 amu. Also, MS²-spectra for the major occurring ions were
549 generated.

550 For native leupeptin production wild type *S. roseus* and *LeupA*, mutants were grown in
551 leupeptin production media (67) containing: Glucose 3gr; NH₄NO₃ 0.5gr; MgSO₄ (7H₂O) 0.5gr;
552 KCl 0.05 gr; L-leucine 0.75 gr; L-arginine 0.75 gr; glycine 0.75 gr; casaminoacids 0.1 gr; yeast
553 extract 0.4 gr per liter. These cultures were set up in shaken flasks with metal springs for 48 hours
554 at 30 C. For heterologous production of leupeptins, *E. coli* DH10B transformants carrying the
555 9_18N and 8_10B PAC clones were grown in Luria-Bertani media (LB) in shaken flasks with 50
556 µg per mL of apramycinm at 37 C for 48 hours. The cultures were centrifuged and the
557 supernatants freeze-dried to obtain 10X concentrates. The crude extracts were analyzed using a
558 C18-218TP vydac column (Grace Healthcare; Columbia, USA) or Restek C18 column (Restek
559 Chromatography; Bellefonte, US), with a 0-100% gradient of [trifluoroacetic acid 0.01% in
560 water]-acetonitrile, and detected by diode array (DAD) at λ=210 nm. Leupeptin authentic standard
561 (L2884, Sigma-Aldrich, St Louis, USA) was used as reference and the peaks with equivalent
562 retention times from the extracts were collected for MS analysis performed on an ion trap LTQ-

563 VELOS equipment in positive mode (Thermo scientific, Waltham, USA) at the MS Unit of
564 Unidad Irapuato Cinvestav-IPN (Irapuato, Mexico).

565 **Acknowledgements**

566 We are indebted with Angélica Cibrián-Jaramillo, Marnix Medema, Paul Straight and Sean
567 Rovito, for useful discussions and critical reading of the manuscript, as well as with Alicia
568 Chagolla and Yolanda Rodriguez of the MS Unit of unidad Irapuato, Cinvestav, for analytical
569 services. This work was funded by Conacyt Mexico (grants No. 179290 and 177568) and
570 FINNOVA Mexico (grant No. 214716) to FBG. PCM was funded by Conacyt scholarship (No.
571 28830) and a Cinvestav postdoctoral fellowship. JF and JFK acknowledge funding from the
572 College of Physical Sciences, University of Aberdeen, UK.

573

574 **Competing interests**

575 PCM, CEMG, MAME, HERA and FBG have filed patent applications related to this work. The
576 other authors declare that no competing interests exist.

577 **References**

- 578 1. Schreiber S (2005) Small molecules: the missing link in the central dogma. *Nature Chemical*
579 *Biology* 1:64–66. DOI: 10.1038/nchembio0705-64
- 580 2. Bachmann BO, Van Lanen SG, Baltz RH (2014) Microbial genome mining for accelerated
581 natural products discovery: is a renaissance in the making? *J Ind Microbiol Biotechnol*
582 41:175–84. DOI: 10.1007/s10295-013-1389-9
- 583 3. Demain AL (2014) Importance of microbial natural products and the need to revitalize their
584 discovery. *J Ind Microbiol Biotechnol* 41:185–201. DOI: 10.1007/s10295-013-1325-z
- 585 4. *Antimicrobial resistance: global report on surveillance 2014* (World Health Organization-
586 UN). ISBN: 9789241564748.
- 587 5. Harvey A, Edrada-Ebel R, Quinn R (2015) The re-emergence of natural products for drug
588 discovery in the genomics era. *Nature Reviews Drug Discovery* 14:111–129.
589 DOI:10.1038/nrd4510
- 590 6. Bachmann BO, McAlpine JB, Zazopoulos E (2006) Farnesyl dibenzodiazepinone, and processes for
591 its production. *US Patent CA2466340 A*.
- 592 7. McAlpine J et al. (2008) Biosynthesis of diazepinomicin/ECO-4601, a *Micromonospora*
593 secondary metabolite with a novel ring system. *Journal of natural products* 71:1585–90. DOI:
594 10.1021/np800376n
- 595 8. Gourdeau H et al. (2007) Identification, characterization and potent antitumor activity of
596 ECO-4601, a novel peripheral benzodiazepine receptor ligand. *Cancer Chemotherapy and*
597 *Pharmacology* 61:911921. DOI: 10.1007/s00280-007-0544-2
- 598 9. Jensen PR, Chavarria KL, Fenical W, Moore BS, Ziemert N (2014) Challenges and triumphs
599 to genomics-based natural product discovery. *J Ind Microbiol Biotechnol* 41:203–9. DOI:
600 10.1007/s10295-013-1353-8
- 601 10. Conway K, Boddy C (2013) ClusterMine360: a database of microbial PKS/NRPS
602 biosynthesis. *Nucleic Acids Research* 41:D402–D407. DOI: 10.1093/nar/gks993
- 603 11. Ichikawa N et al. (2013) DoBISCUIT: a database of secondary metabolite biosynthetic gene
604 clusters. *Nucleic acids research* 41:D408–14. DOI: 10.1093/nar/gks1177
- 605 12. Barona-Gómez F, Wong U, Giannakopoulos A, Derrick P, Challis G (2004) Identification of a
606 cluster of genes that directs desferrioxamine biosynthesis in *Streptomyces coelicolor* M145.
607 *Journal of the American Chemical Society* 126:1628216283. DOI: 10.1021/ja045774k

- 608 13. Lautru S, Deeth R, Bailey L, Challis G (2005) Discovery of a new peptide natural product by
609 Streptomyces coelicolor genome mining. *Nature Chemical Biology* 1:265–269. DOI:
610 10.1038/nchembio731
- 611 14. Udvary D et al. (2007) Genome sequencing reveals complex secondary metabolome in the
612 marine actinomycete *Salinispora tropica*. *Proceedings of the National Academy of Sciences*
613 104:10376–10381. DOI: 10.1073/pnas.0700962104
- 614 15. Challis G (2008) Mining microbial genomes for new natural products and biosynthetic
615 pathways. *Microbiology (Reading, England)* 154:1555–69. DOI:
616 10.1099/mic.0.2008/018523-0
- 617 16. Metcalf W, Donk W (2009) Biosynthesis of phosphonic and phosphinic acid natural products.
618 *Biochemistry* 78:65–94. DOI: 10.1146/annurev.biochem.78.091707.100215
- 619 17. Yu X et al. (2013) Diversity and abundance of phosphonate biosynthetic genes in nature.
620 *Proceedings of the National Academy of Sciences* 110:20759–20764. DOI:
621 10.1073/pnas.1315107110
- 622 18. Ju K-S, Doroghazi J, Metcalf W (2013) Genomics-enabled discovery of phosphonate natural
623 products and their biosynthetic pathways. *Journal of Industrial Microbiology &*
624 *Biotechnology*. 41:345-56 doi: 10.1007/s10295-013-1375-2
- 625 19. Arnison P et al. (2012) Ribosomally synthesized and post-translationally modified peptide
626 natural products: overview and recommendations for a universal nomenclature. *Natural*
627 *Product Reports* 30:108–160. DOI: 10.1039/c2np20085f
- 628 20. Thaker M et al. (2013) Identifying producers of antibacterial compounds by screening for
629 antibiotic resistance. *Nature Biotechnology* 31:922–927. DOI: 10.1038/nbt.2685
- 630 21. Gerlt JA, Babbitt PC (2001) Divergent evolution of enzymatic function: mechanistically
631 diverse superfamilies and functionally distinct suprafamilies. *Annual review of biochemistry*
632 70:209–46. DOI: 10.1146/annurev.biochem.70.1.209
- 633 22. Caetano-Anollés G et al. (2009) The origin and evolution of modern metabolism. *Int J*
634 *Biochem Cell Biol* 41:285–97. DOI: 10.1016/j.biocel.2008.08.022
- 635 23. Tahlan K, Park HU, Wong A, Beatty PH, Jensen SE (2004) Two sets of paralogous genes
636 encode the enzymes involved in the early stages of clavulanic acid and clavam metabolite
637 biosynthesis in *Streptomyces clavuligerus*. *Antimicrob Agents Chemother* 48:930–9. DOI:
638 10.1128/AAC.48.3.930-939.2004
- 639 24. Verdel-Aranda K, López-Cortina ST, Hodgson DA, Barona-Gómez F (2015) Molecular
640 annotation of ketol-acid reductoisomerases from *Streptomyces* reveals a novel amino acid
641 biosynthesis interlock mediated by enzyme promiscuity. *Microb Biotechnol* 8:239–52. DOI:
642 10.1111/1751-7915.12175

- 643 25. Vining LC (1992) Secondary metabolism, inventive evolution and biochemical diversity--a
644 review. *Gene* 115:135–40. DOI: 10.1016/0378-1119(92)90551-Y
- 645 26. Firn R, Jones C (2009) A Darwinian view of metabolism: molecular properties determine
646 fitness. *Journal of Experimental Botany* 60:719–726. DOI: 10.1093/jxb/erp002
- 647 27. Medema MH, Cimermancic P, Sali A, Takano E, Fischbach MA (2014) A systematic
648 computational analysis of biosynthetic gene cluster evolution: lessons for engineering
649 biosynthesis. *PLoS Comput Biol* 10:e1004016.
- 650 28. Barona-Gómez F, Cruz-Morales P, Noda-García L (2012) What can genome-scale metabolic
651 network reconstructions do for prokaryotic systematics? *Antonie Van Leeuwenhoek* 101:35–
652 43. DOI: 10.1007/s10482-011-9655-1
- 653 29. Medema M et al. (2011) antiSMASH: rapid identification, annotation and analysis of
654 secondary metabolite biosynthesis gene clusters in bacterial and fungal genome sequences.
655 *Nucleic Acids Research* 39:W339–W346. DOI: 10.1093/nar/gkr466
- 656 30. Cimermancic P et al. (2014) Insights into secondary metabolism from a global analysis of
657 prokaryotic biosynthetic gene clusters. *Cell* 158:412–21. DOI: 10.1016/j.cell.2014.06.034
- 658 31. Blodgett JA et al. (2007) Unusual transformations in the biosynthesis of the antibiotic
659 phosphinothricin tripeptide. *Nat Chem Biol* 3:480–5. DOI:10.1038/nchembio.2007.9
- 660 32. Blodgett J, Zhang J, Metcalf W (2005) Molecular Cloning, Sequence Analysis, and
661 Heterologous Expression of the Phosphinothricin Tripeptide Biosynthetic Gene Cluster from
662 *Streptomyces viridochromogenes* DSM 40736. *Antimicrobial Agents and Chemotherapy*
663 49:230–240. DOI: 10.1128/AAC.49.1.230-240.2005
- 664 33. Schwartz D et al. (2004) Biosynthetic Gene Cluster of the Herbicide Phosphinothricin
665 Tripeptide from *Streptomyces viridochromogenes* Tü494. *Applied and Environmental*
666 *Microbiology* 70:7093–7102. DOI: 10.1128/AEM.70.12.7093-7102.2004
- 667 34. Blin K et al. (2013) antiSMASH 2.0-a versatile platform for genome mining of secondary
668 metabolite producers. *Nucleic acids research* 41:W204–12. DOI: 10.1093/nar/gkt449
- 669 35. Aoyagi T, Takeuchi T, Matsuzaki A, Kawamura K, Kondo S (1969) Leupeptins, new protease
670 inhibitors from Actinomycetes. *The Journal of antibiotics* 22:283–6. DOI:
671 <http://doi.org/10.7164/antibiotics.22.283>
- 672 36. Suzukake K, Fujiyama T, Hayashi H, Hori M (1979) Biosynthesis of leupeptin. II. Purification and
673 properties of leupeptin acid synthetase. *The Journal of antibiotics* 32:523-30. DOI:
674 <http://doi.org/10.7164/antibiotics.22.283>
- 675 37. Hori M, Hemmi H, Suzukake K, Hayashi H (1978) Biosynthesis of leupeptin. *The Journal of*
676 *antibiotics* 31:95-8. DOI: <http://doi.org/10.7164/antibiotics.31.95>

- 677 38. Suzukake K, Hori M, Tamemasa O, Umezawa H (1981) Purification and properties of an enzyme
678 reducing leupeptin acid to leupeptin. *Biochim Biophys Acta*. 661:175-81. DOI: 10.1016/0005-
679 2744(81)90001-2
- 680 39. Suzukake K, Hayashi H, Hori M (1980) Biosynthesis of leupeptin. III. Isolation and properties of an
681 enzyme synthesizing acetyl-L-leucine. *The Journal of antibiotics* 33:857-62. DOI:
682 <http://doi.org/10.7164/antibiotics.33.857>
- 683 40. Chen Y, McClure R, Zheng Y, Thomson R, Kelleher N (2013) Proteomics guided discovery
684 of flavopeptins: anti-proliferative aldehydes synthesized by a reductase domain-containing
685 non-ribosomal peptide synthetase. *Journal of the American Chemical Society* 135:10449–56.
686 DOI: 10.1021/ja4031193
- 687 41. Rausch C, Hoof I, Weber T, Wohlleben W, Huson DH (2007) Phylogenetic analysis of
688 condensation domains in NRPS sheds light on their functional evolution. *BMC evolutionary*
689 *biology* 7:78. DOI: 10.1186/1471-2148-7-78
- 690 42. Sampaleanu LM, Vallée F, Thompson GD, Howell PL (2002) Three-dimensional structure of
691 the argininosuccinate lyase frequently complementing allele Q286R. *Biochemistry* 40:15570–
692 80. DOI: 10.1021/bi011525m
- 693 43. Kaysser L et al. (2011) Identification of a napsamycin biosynthesis gene cluster by genome
694 mining. *Chembiochem* 12:477–87. DOI: 10.1002/cbic.201000460
- 695 44. Zhang W, Ostash B, Walsh CT (2010) Identification of the biosynthetic gene cluster for the
696 pacidamycin group of peptidyl nucleoside antibiotics. *Proceedings of the National Academy*
697 *of Sciences of the United States of America* 107:16828–33. DOI: 10.1073/pnas.1011557107
- 698 45. Challis GL (2014) Exploitation of the *Streptomyces coelicolor* A3(2) genome sequence for
699 discovery of new natural products and biosynthetic pathways. *J Ind Microbiol Biotechnol*
700 41:219–32. DOI: 10.1007/s10295-013-1383-2
- 701 46. Kieser T, Bibb MJ, Buttner MJ, Chater KF, Hopwood DA (2000) *Practical Streptomyces*
702 *genetics* (The John Innes Foundation, Norwich UK). ISBN 0-7084-0623-8
- 703 47. Zhang F, Berti PJ (2006) Phosphate analogues as probes of the catalytic mechanisms of MurA
704 and AroA, two carboxyvinyl transferases. *Biochemistry* 45:6027–37. DOI:
705 10.1021/bi0601914
- 706 48. Rui Z et al. (2010) Biochemical and genetic insights into asukamycin biosynthesis. *The*
707 *Journal of biological chemistry* 285:24915–24. DOI: 10.1074/jbc.M110.128850
- 708 49. Seeger K et al. (2011) The biosynthetic genes for prenylated phenazines are located at two
709 different chromosomal loci of *Streptomyces cinnamomensis* DSM 1042. *Microbial*
710 *biotechnology* 4:252–62. DOI: 10.1111/j.1751-7915.2010.00234.x

- 711 50. Cruz-Morales P et al. (2013) The genome sequence of *Streptomyces lividans* 66 reveals a
712 novel tRNA-dependent peptide biosynthetic system within a metal-related genomic island.
713 *Genome Biology and Evolution* 5:1165–1175. DOI: 10.1093/gbe/evt082
- 714 51. Nett M, Ikeda H, Moore BS (2009) Genomic basis for natural product biosynthetic diversity
715 in the actinomycetes. *Nat Prod Rep* 26:1362–84. DOI: 10.1039/B817069J
- 716 52. Bentley SD et al. (2002) Complete genome sequence of the model actinomycete
717 *Streptomyces coelicolor* A3(2). *Nature* 417:141–7. DOI: 10.1038/417141a
- 718 53. Wang L et al. (2006) arsRBOCT arsenic resistance system encoded by linear plasmid pHZ227
719 in *Streptomyces* sp. strain FR-008. *Applied and environmental microbiology* 72:3738–42.
720 DOI: 10.1128/AEM.72.5.3738-3742.2006
- 721 54. Elias M et al. (2012) The molecular basis of phosphate discrimination in arsenate-rich
722 environments. *Nature* 491:134–137. DOI: 10.1038/nature11517
- 723 55. Tawfik D, Viola R (2011) Arsenate replacing phosphate: alternative life chemistries and ion
724 promiscuity. *Biochemistry* 50:1128–34. DOI: 10.1021/bi200002a
- 725 56. Chawla S, Mutenda EK, Dixon HB, Freeman S, Smith AW (1995) Synthesis of 3-
726 arsonopyruvate and its interaction with phosphoenolpyruvate mutase. *The Biochemical*
727 *journal* 308 (Pt 3):931–5.
- 728 57. Amayo KO, Raab A, Krupp EM, Gunnlaugsdottir H, Feldmann J (2013) Novel identification
729 of arsenolipids using chemical derivatizations in conjunction with RP-HPLC-ICPMS/ESMS.
730 *Analytical chemistry* 85:9321–7. DOI: 10.1021/ac4020935
- 731 58. Edgar RC (2004) MUSCLE: multiple sequence alignment with high accuracy and high
732 throughput. *Nucleic acids research* 32:1792–7. DOI: 10.1093/nar/gkh340
- 733 59. Waterhouse AM, Procter JB, Martin DM, Clamp M, Barton GJ (2009) Jalview Version 2--a
734 multiple sequence alignment editor and analysis workbench. *Bioinformatics (Oxford,*
735 *England)* 25:1189–91. DOI: 10.1093/bioinformatics/btp033
- 736 60. Ronquist F et al. (2012) MrBayes 3.2: efficient Bayesian phylogenetic inference and model
737 choice across a large model space. *Systematic biology* 61:539–42. DOI:
738 10.1093/sysbio/sys029
- 739 61. Gust B, Challis G, Fowler K, Kieser T, Chater K (2003) PCR-targeted *Streptomyces* gene
740 replacement identifies a protein domain needed for biosynthesis of the sesquiterpene soil odor
741 geosmin. *Proceedings of the National Academy of Sciences* 100:1541–1546. DOI:
742 10.1073/pnas.0337542100

- 743 62. Redenbach M et al. (1996) A set of ordered cosmids and a detailed genetic and physical map
744 for the 8 Mb *Streptomyces coelicolor* A3(2) chromosome. *Molecular microbiology* 21:77–96.
745 DOI: 10.1046/j.1365-2958.1996.6191336.x
- 746 63. Bolger AM, Lohse M, Usadel B (2014) Trimmomatic: a flexible trimmer for Illumina
747 sequence data. *Bioinformatics (Oxford, England)* 30:2114–20. DOI:
748 10.1093/bioinformatics/btu170
- 749 64. Zerbino DR, Birney E (2008) Velvet: algorithms for de novo short read assembly using de
750 Bruijn graphs. *Genome research* 18:821–9. DOI: 10.1101/gr.074492.107
- 751 65. Aziz RK et al. (2008) The RAST Server: rapid annotations using subsystems technology.
752 *BMC genomics* 9:75. DOI: 10.1186/1471-2164-9-75
- 753 66. Sosio M et al. (2000) Artificial chromosomes for antibiotic-producing actinomycetes. *Nat*
754 *Biotechnol* 18:343–5. DOI:10.1038/73810
- 755 67. Takagi K, Yamamoto Y, Yamazaki T, Yamaguchi H, Umezawa H (1978) Process for
756 producing l-leupeptins. *US patent 4066507 A*

757 **Figures**

758 **Figure 1. EvoMining pipeline for the recapitulation of the evolution of NP biosynthesis. A.**

759 Bioinformatic workflow: the three input databases, as discussed in the text, are shown in green.

760 Internal databases are shown in yellow, whereas grey boxes depict processes. **B.** An example of a

761 typical EvoMining phylogenetic tree (**Tree S4**) using the case of 3-carboxyvinyl-

762 phosphoshikimate synthase family. Red branches include homologs related to central metabolism

763 and their topology resembles that of a species guide tree (**Tree S1**), while blue branches have been

764 recruited into known BGCs. Cyan braches are EvoMining hits found within regions recognised as

765 NP-related also by antiSMASH or ClusterFinder. Green branches are not classifiable by other

766 methods and thus represent EvoMining predictions that may form part of BGCs for novel classes

767 of NPs (see **Figure 2A** for further details).

768 **Figure 2. Analysis of EvoMining Hits. A.** Pie chart of the whole set of EvoMining hits as

769 annotated using antiSMASH and ClusterFinder. **B.** Diversity of BGCs per recruited enzyme

770 family. Top panel, the number of hits and BGC classes per family are compared. As the number of

771 hits increases, more BGC classes are found. Bottom panel, a diversity plot for each enzyme

772 family, showing the proportion and number of BGC classes as defined by AntiSMaSH. The label

773 “Detected by ClusterFinder” means EvoMining hits that are also found by this algorithm, whereas

774 the label “EvoMining predictions” includes all hits that could not be detected by antiSMASH or

775 Cluster Finder.

776

777

778

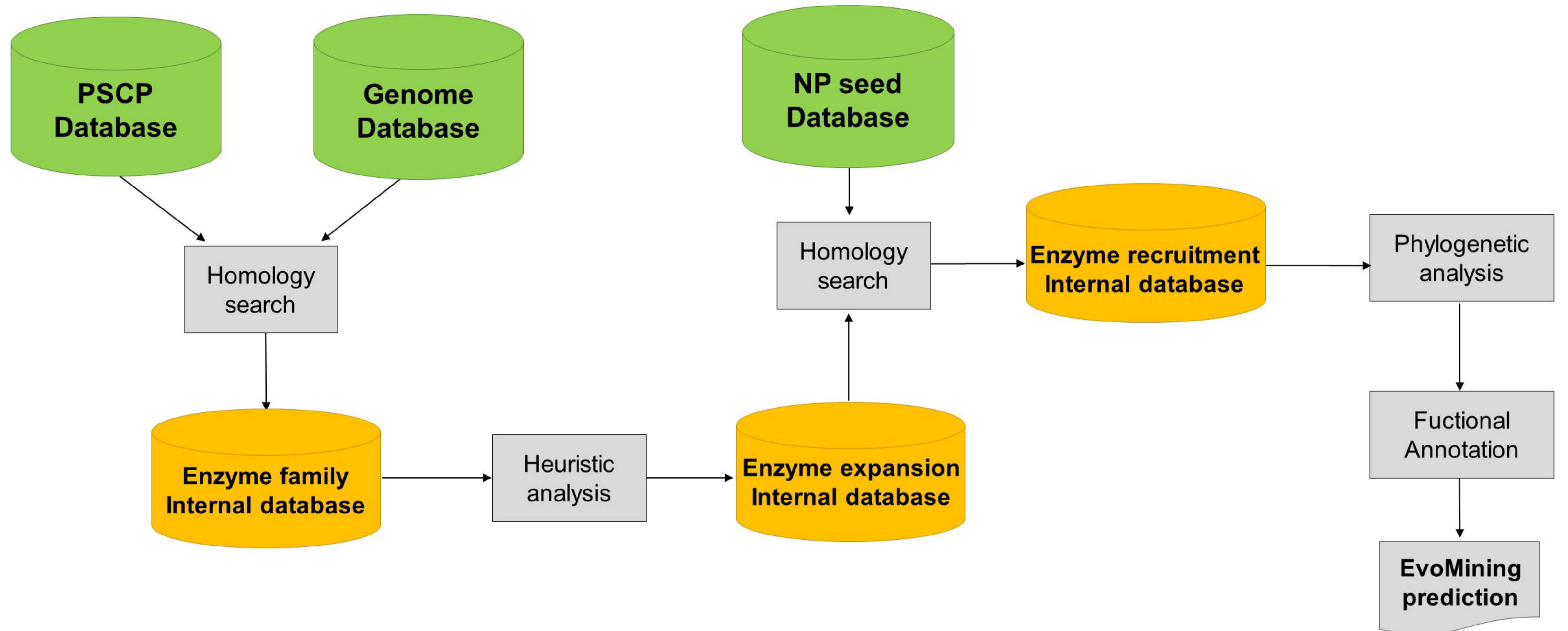
779 **Figure 3. Discovery of the BGC for leupeptin.** **A.** Left panel, phylogenetic reconstruction of the
780 actinobacterial argininosuccinate lyase enzyme family (**Tree S3**). Homologs related to central
781 metabolism are shown in red branches. A clade including recruited homologs is shown in cyan.
782 The LeupB homolog from *S. roseus*, a leupeptin producer, together with the known recruitments
783 for Pacidamycin and Napsamycin (blue branches), are indicated. **B.** Genome context of *leupB*,
784 including *leupA* (novel NRPS), *leupC* (annotated as threonine kinase) and *leupD* (annotated as
785 cysteine synthase). **C.** Biosynthetic proposal for leupeptin, based in the EvoMining prediction and
786 earlier biochemical data, as discussed within the text. LeupA is proposed to produce Acyl-Leu-
787 Leu, which is used by LeupB, together with arginine, for the synthesis of leupeptidic acid. A
788 reductase activity, that remains to be identified, is required for formation of the characteristic
789 aldehyde group.

790

791 **Figure 4. Discovery of a BGC for arseno-organic NPs in *S. coelicolor* and *S. lividans*.** **A.** The
792 conserved arseno-organic BGC in *S. lividans* 66 and *S. coelicolor* is shown. The proposed
793 biosynthetic logic for early intermediates in the biosynthesis of arseno-organic metabolites is
794 shown. **B.** Transcriptional analysis of selected genes within the arseno-organic BGC, showing that
795 gene expression is repressed under standard conditions, but induced upon the presence of arsenate.
796 **C.** HPLC-Orbitrap/QQQ-MS trace of organic extracts from mycelium of wild type and the
797 SLI_1096 mutant showing the detection of arsenic-containing species. Three m/z signals were
798 detected within the two peaks found in the trace from the wild type strain grown on the presence
799 of arsenate. These m/z signals are absent from the wild type strain grown without arsenate, and
800 from the SLI_1096 mutant strain grown on phosphate limitation and the presence of arsenate.
801 Identical results were obtained for *S. coelicolor* and the SCO6819 mutant

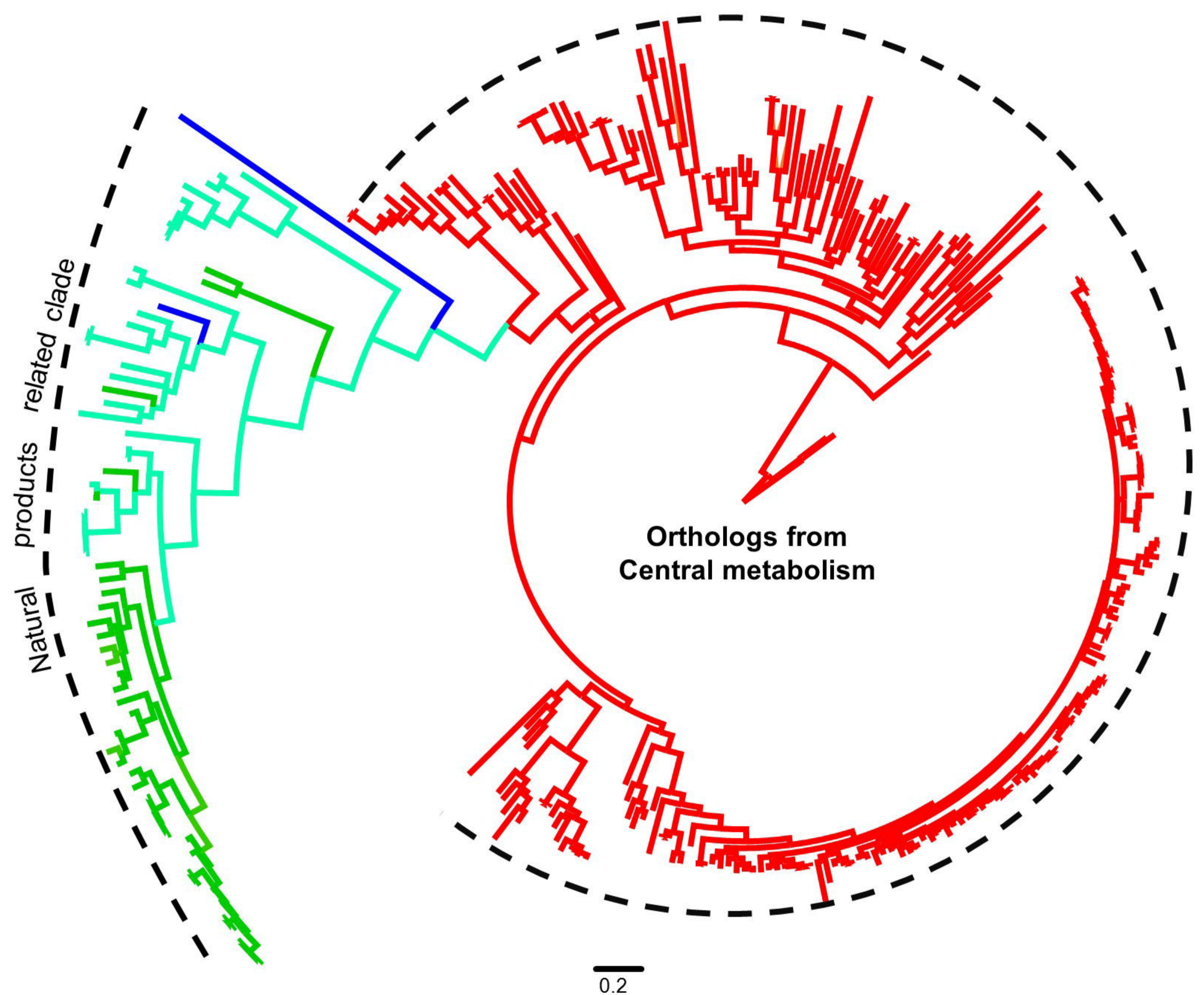
802

803 **Figure 5. Novel BGCs for arseno-organic metabolites found in *Actinobacteria*.** The BGCs
804 were found by mining for the co-occurrence of arsenoenol pyruvate synthase (AEPS),
805 arsenopyruvate mutase (APM), arsenoenolpyruvate decarboxylase (APD), in available bacterial
806 genomes from the GenBank, as of November 2014. Related BGCs were only found in
807 actinomycetes. The phylogeny was constructed with a concatenated matrix of conserved enzymes
808 among the BGCs that included AEPS, APM, APD (purple arrows), plus CTP synthase and
809 anaerobic dehydrogenase (shown as blue arrows together with other enzymes). Variations in the
810 functional content of the BGCs are accounted by PKSs (red arrows), hybrid PKS-NRPS (yellow
811 genes), arsenic regulation and metabolism proteins (brown arrows), and other regulators and
812 transporters (green arrows). Three main classes of arseno-organic BGCs could be expected from
813 this analysis: PKS-independent, PKS-NRPS-dependent and PKS-dependent biosynthetic systems.
814 Dotted lines indicate sequence gaps, and an asterisk marks the sequence from *Nocardiopsis*
815 *lucentensis*, which is assumed that have a missing PKS gene in one of the sequence gaps.

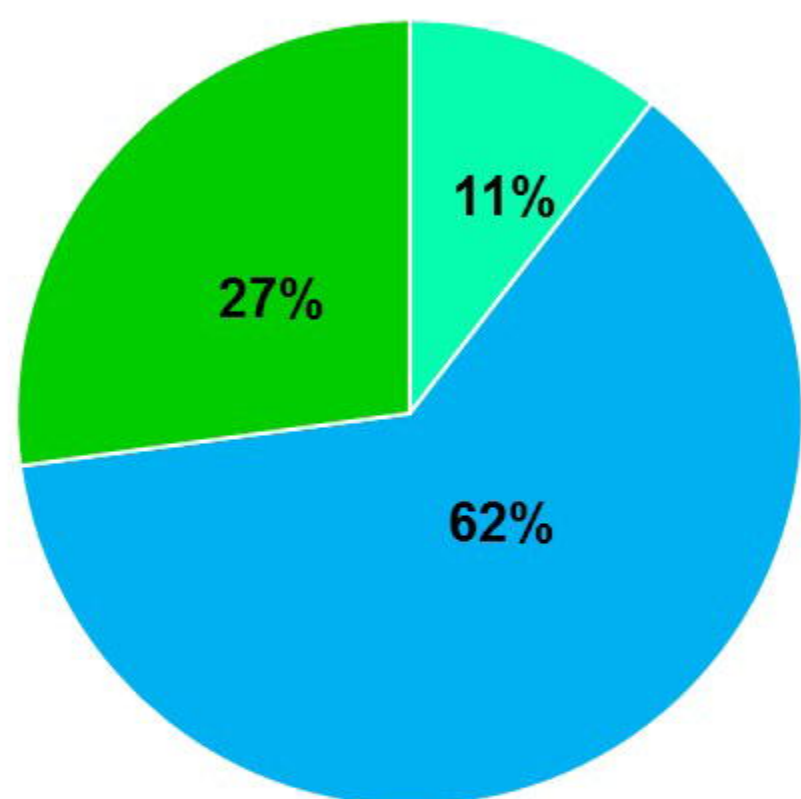
A**B**

bioRxiv preprint doi: <https://doi.org/10.1101/020503>; this version posted July 6, 2015. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

- **Not involved in NP biosynthesis**
- **Known recruitments**
- **EvoMining hits
(Detected by AntiSMASH/ClusterFinder)**
- **EvoMining predictions**



A

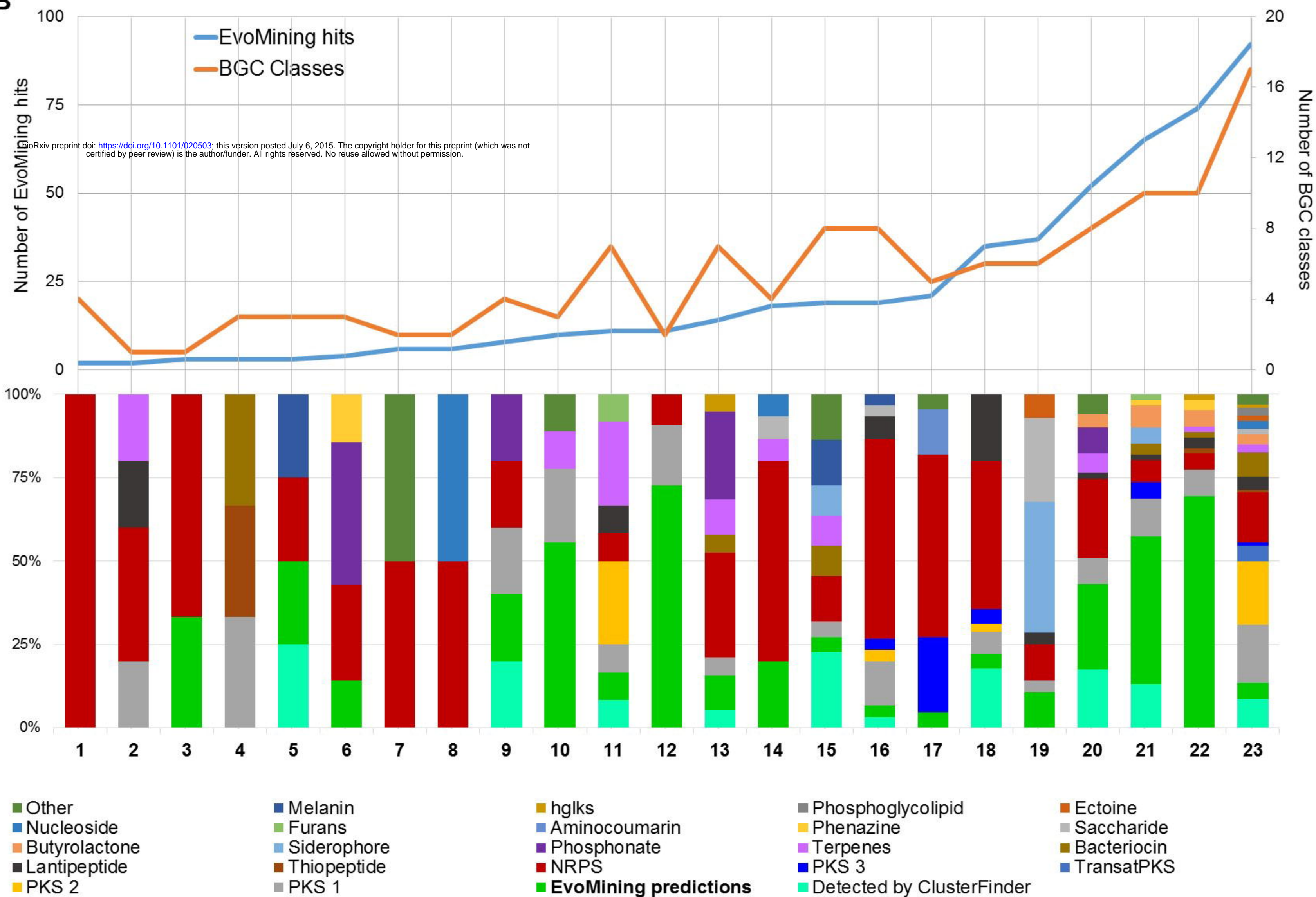


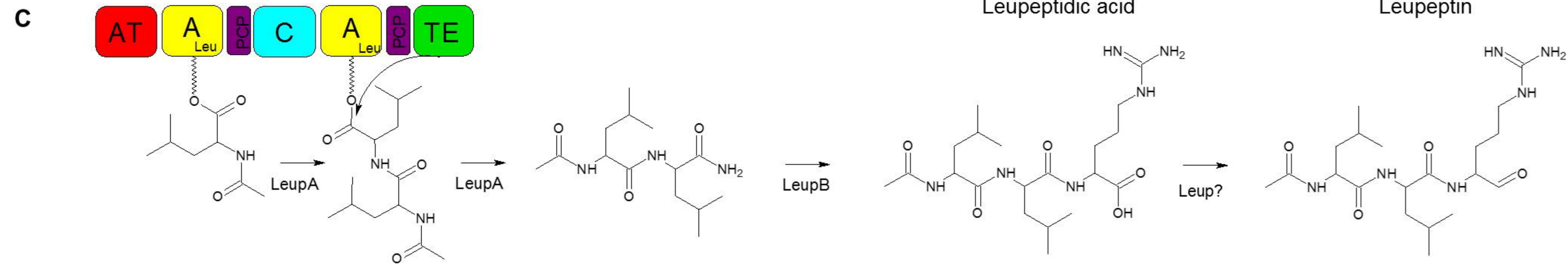
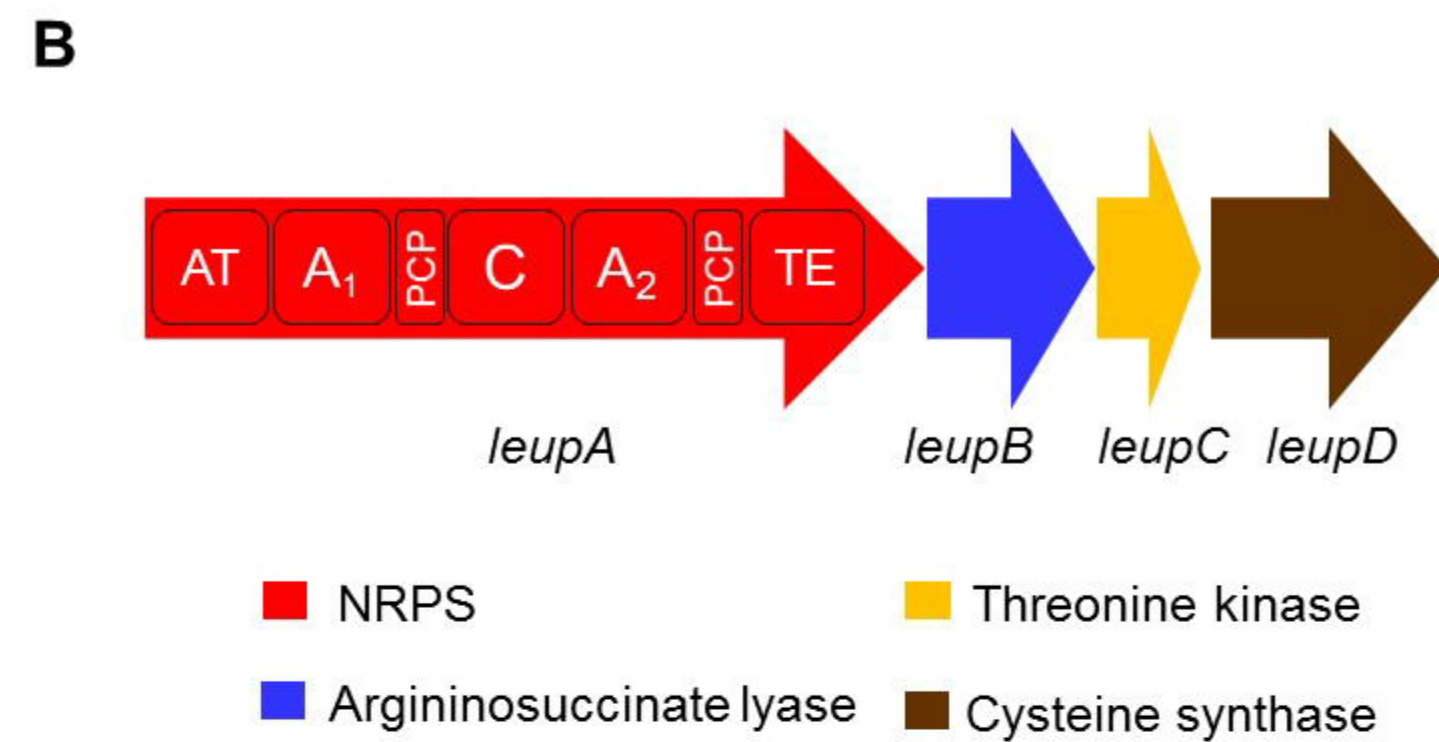
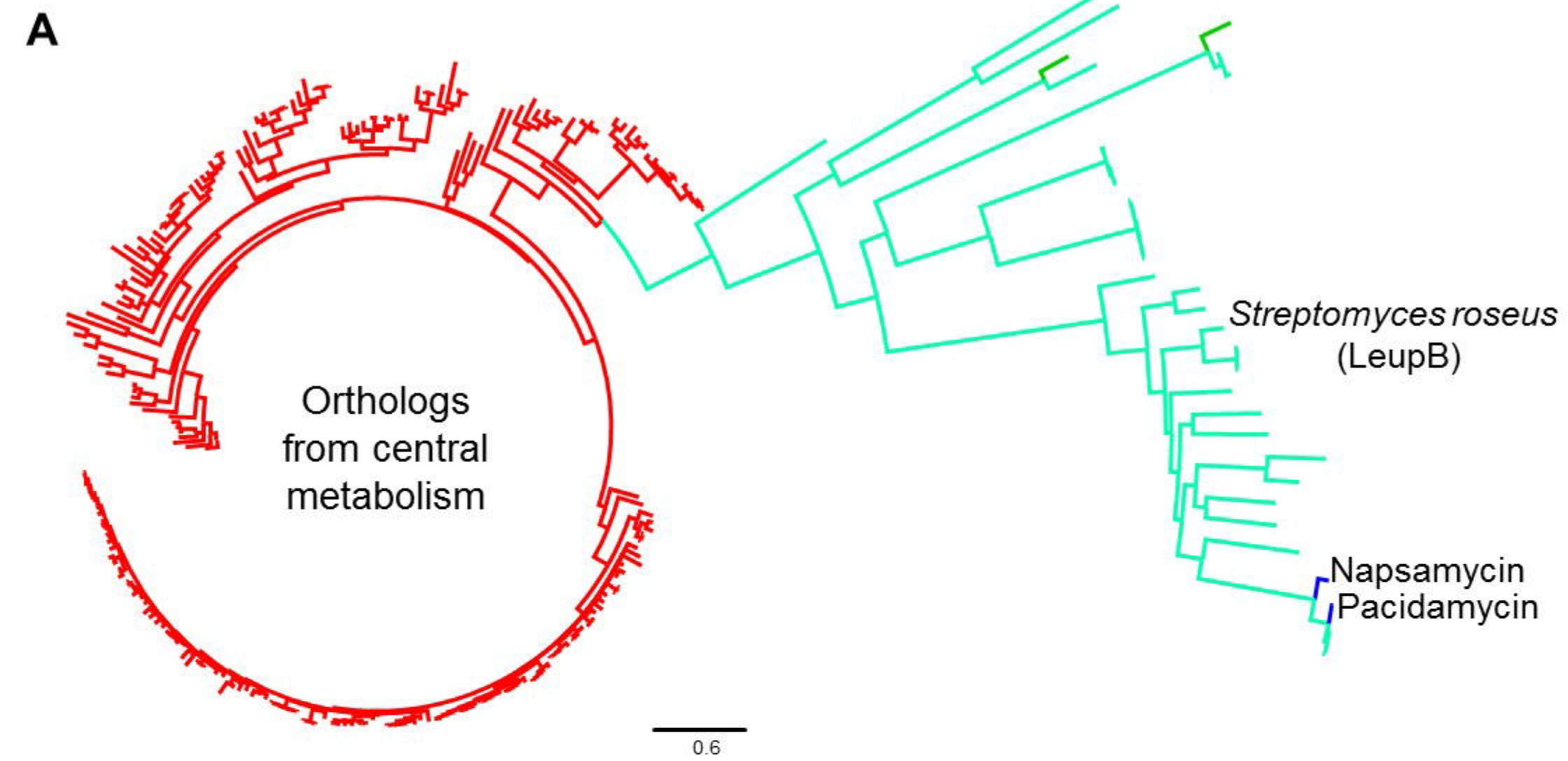
1. Indole-3-Glycerolphosphate synthase
2. *N*-acetyl-gamma-glutamyl-phosphate reductase
3. Aspartate transaminase
4. Histidinolphosphate aminotransferase
5. Homoserine-*O*-succinyl transferase
6. Enolase
7. Anthranilate phosphoribosyltransferase
8. Histidinol phosphatase
9. Citrate synthase
10. Acetolactate synthase
11. Glyceraldehyde-3-phosphate dehydrogenase
12. Phosphoglycerate dehydrogenase
13. Aconitate hydratase
14. Acetyl glutamate kinase
15. Aspartate kinase
16. Cysteine synthase
17. Prephenate dehydrogenase
18. Argininosuccinate lyase
19. Acetyl ornithine aminotransferase
20. Isopropylmalate dehydrogenase
21. 3-Phosphoshikimate-1-carboxyvinyl-transferase
22. 2-dehydro-3-deoxyphosphoheptonate aldolase
23. Asparagine synthase

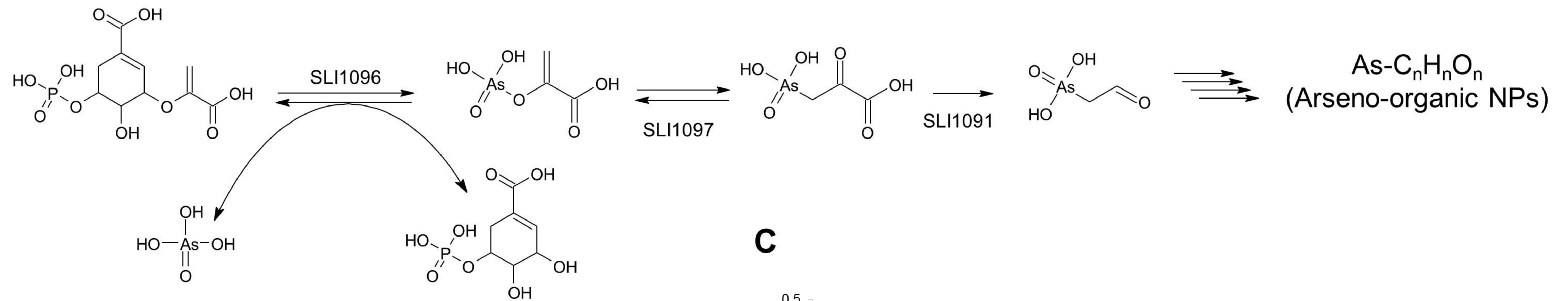
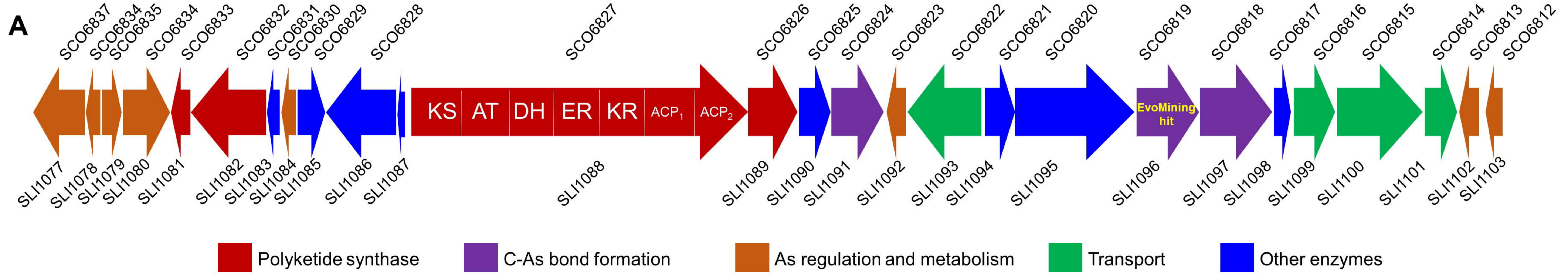
EvoMining hits

- Detected by ClusterFinder
- Detected by AntiSMASH
- EvoMining predictions

B

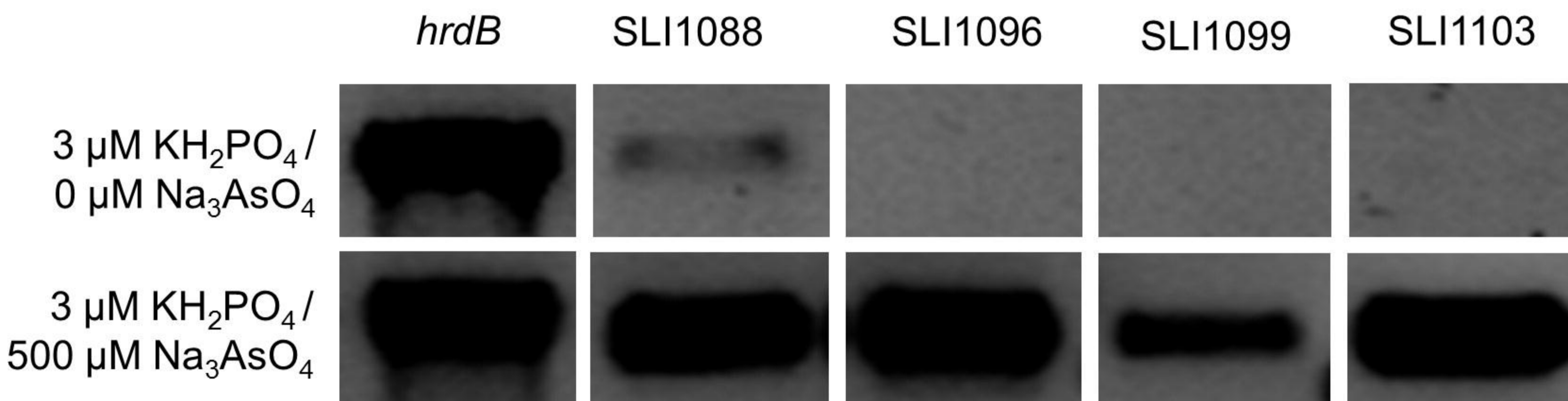




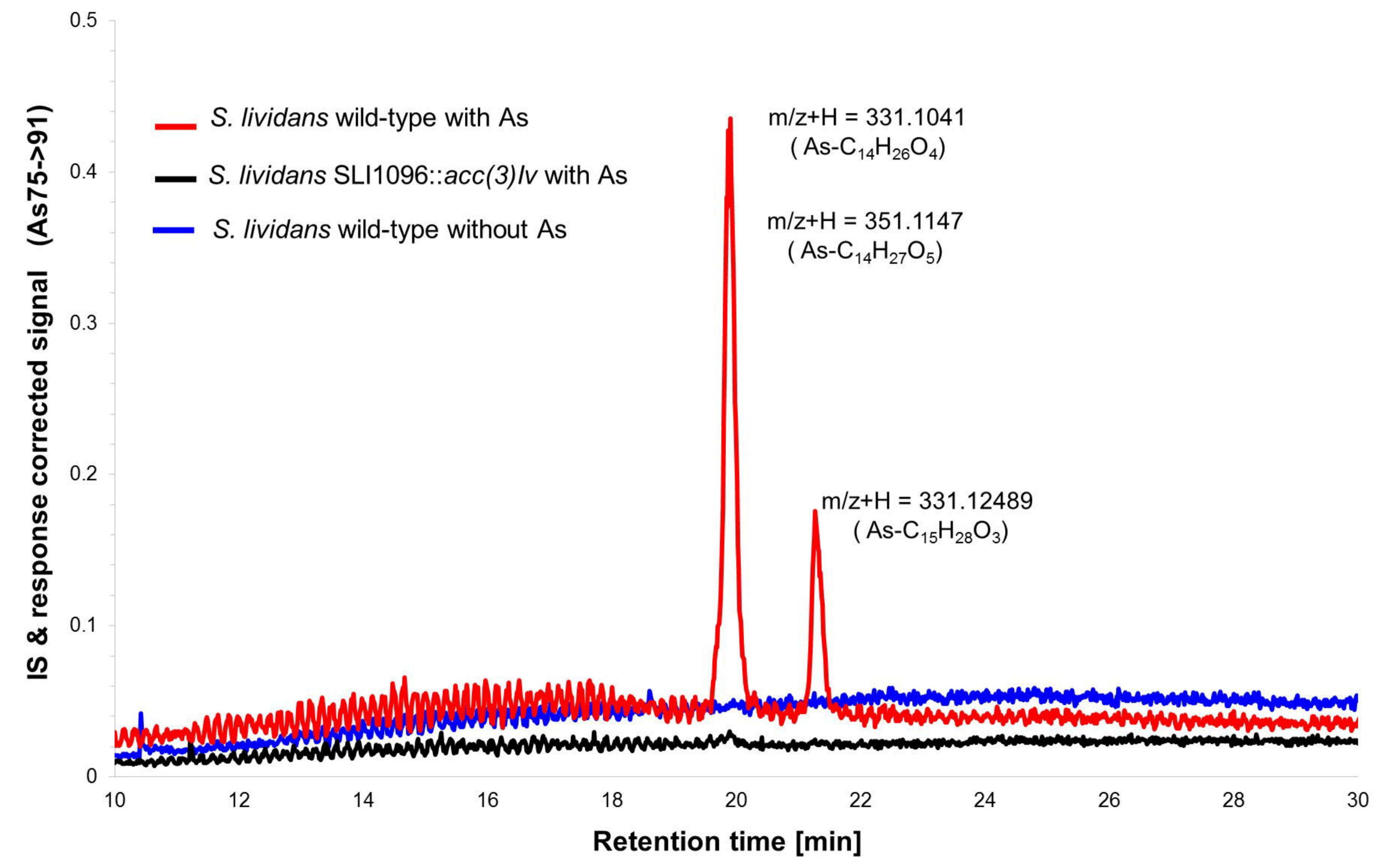


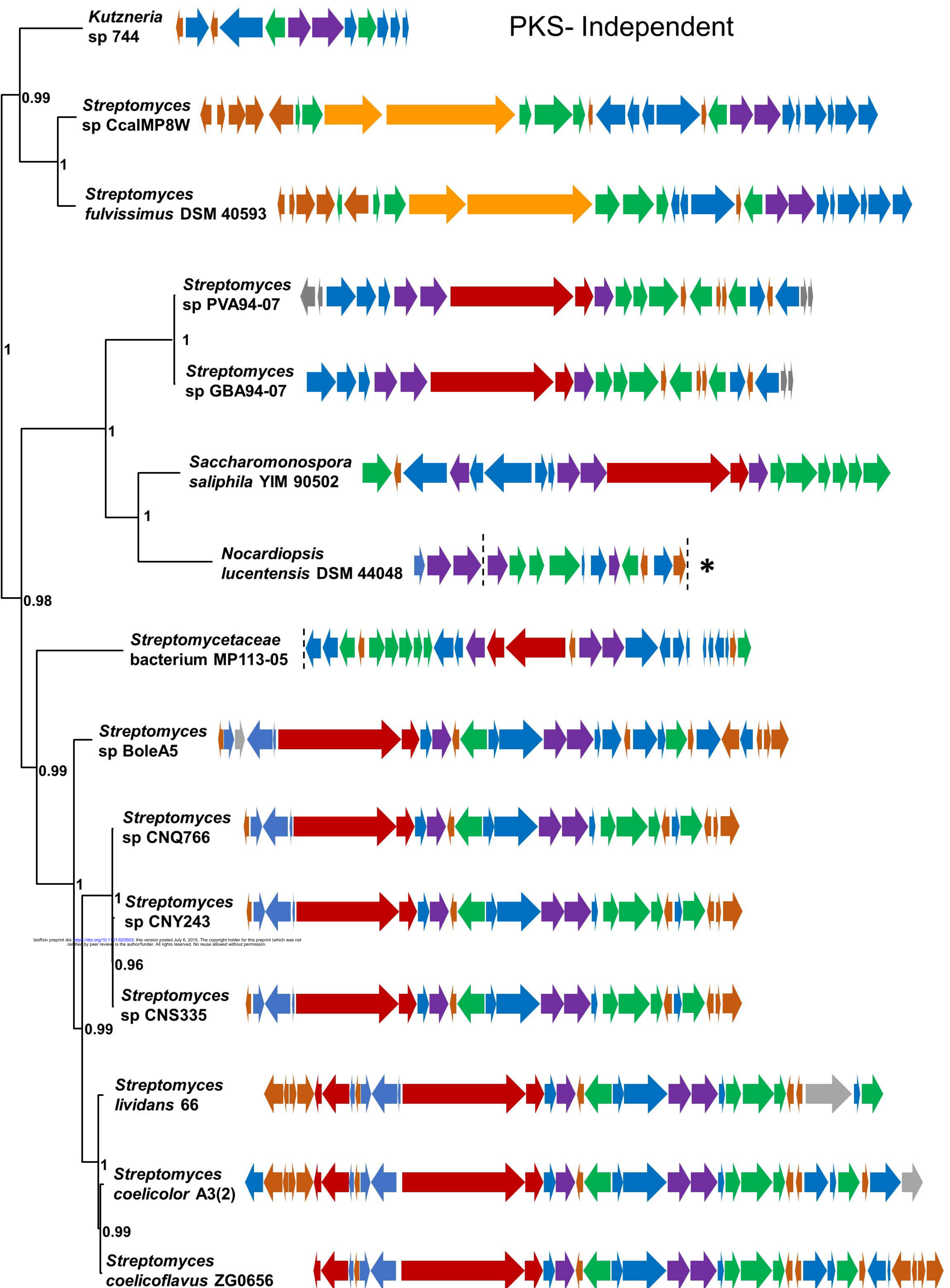
bioRxiv preprint doi: <https://doi.org/10.1101/020503>; this version posted July 6, 2015. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

B



C





Hybrid PKS-NRPS dependent

PKS- dependent