

1 **Delving into the *Bacillus cereus* group biosynthetic gene clusters cosmos: a**
2 **comparative-genomics-based classification framework**

3 Hadj Ahmed Belaoui¹, Amine Yekkour¹, Abdelghani Zitouni¹, Atika Meklat¹

4 ¹ Laboratoire de Biologie des Systèmes Microbiens (LBSM), Ecole Normale Supérieure de
5 Kouba, Algiers, Algeria

6 Correspondence:

7 Hadj Ahmed Belaoui

8 E-mail: shooper5@yahoo.fr

9 **Abstract**

10 **Background:** In this study, the *Bacillus* sp. strain BH32 (a plant-beneficial bacterial
11 endophyte) and its closest non-type *Bacillus cereus* group strains were used to study the
12 organization, conservation, and diversity of biosynthetic gene clusters (BGCs) among this
13 group to propose a classification framework of gene cluster families (GCFs) among this
14 intricate group. A dataset consisting of 17 genomes was used in this study. Genomes were
15 annotated using PROKKA ver.1.14.5. The web tool antiSMASH ver. 5.1.2 was used to
16 predict the BGCs profiles of each strain, with a total number of 198 BGCs. The comparison
17 was made quantitatively based on a BGCs counts matrix comprising all the compared
18 genomes and visualized using the Morpheus tool. The constitution, distribution, and
19 evolutionary relationships of the detected BGCs were further analyzed using a manual
20 approach based on a BLASTp analysis (using BRIG ver. 0.95); a phylogenetic analysis of the
21 concatenated BGCs sequences to highlight the evolutionary relationships; and the
22 conservation, distribution and the genomic co-linearity of the studied BGCs using Mauve
23 aligner ver. 2.4.0. Finally, the BIG-SCAPE/CORASON automated pipeline was used as a
24 complementary strategy to investigate the gene cluster families (GCFs) among the *B. cereus*
25 group.

26 **Results:** Based on the manual approach, we identified BGCs conserved across the studied
27 strains with very low variation and interesting singletons BGCs. Moreover, we highlighted the
28 presence of two major BGCs synteny blocks (named “*synteny block A*” and “*synteny block*
29 *B*”), each composed of conserved homologous BGCs among the *B. cereus* group. For the
30 automatic approach, we identified 23 families among the different BGCs classes of the *B.*
31 *cereus* group, named using a rational basis. The proposed manual and automatic
32 approaches proved to be in harmony and complete each other, for the study of BGCs among
33 the selected genomes.

34 **Conclusion:** Ultimately, we propose a framework for an expanding classification of the *B.*
35 *cereus* group BGCs, based on a set of reference BGCs reported in this work.

36 **Keywords:** Comparative genomics, *Bacillus cereus* group, endophyte, Biosynthetic gene
37 clusters (BGCs), Synteny.

38

39 **Background**

40 The *Bacillus cereus* group of bacteria represents a homogeneous subdivision of the genus
41 *Bacillus* with closely related phylogeny within the Firmicutes phylum [1, 2]. The bacteria
42 constituting this group are Gram-positive, spore-forming, aerobic/facultative anaerobic, and
43 rod-shaped with low-GC content [2]. Numerous bacteria related to *B. cereus* group were
44 shown to produce several interesting compounds and enzymes, metabolize different kinds of
45 pollutants, and promote the growth of both plants and animals when used as biostimulators
46 [3–5]. The first described and most documented members of the group are *B. cereus*, *B.*
47 *thuringiensis*, and *B. anthracis* [2]. Many bacteria classified as *B. cereus* are ubiquitous in the
48 environment, with apparent soil origin, and present as commensal to intestines of insects or
49 foodborne opportunistic pathogens often related to human poisoning [6, 7]. Other members
50 of the group related to *B. thuringiensis* are insect pathogens widely used in agriculture for the
51 biocontrol of insect pests [7, 8]; while *B. anthracis* is the causative agent of anthrax [2].

52 Traditionally, pathogenic potential and virulence characteristics permit organism
53 differentiation within the group [1, 2, 7]: *B. cereus* by carrying biosynthetic gene cluster (BGC)
54 for cereulide cytotoxin, plasmids carrying the crystal insecticidal genes for *B. thuringiensis*
55 and presence of anthrax toxin and capsule genes for *B. anthracis* [2]. However, genetic
56 experiments have exhibited a high level of synteny and protein similarity, with limited
57 differences in gene content [1]. Conventional taxonomic markers, such as 16S and 23S
58 rRNA genes, as well as whole-genome DNA hybridization seemed essentially identical,
59 questioning the speciation of the group members [7]. Thus, extensive genomic similarities
60 have contributed to the suggestion that *B. anthracis*, *B. cereus*, and *B. thuringiensis* are
61 members of a single species, *B. cereus sensu lato* [1]. Moreover, the horizontal genetic
62 transfer of plasmid and chromosome DNA among the strains of the *B. cereus* group has
63 likely related to the diversity of this bacterial group, thus complicating the speciation [2].

64 In the present context of increased environmental screening with the generalization of whole
65 genome sequencing that revealed the presence of newly identified recombinant forms [2], it
66 is important to explore different approaches for understanding the evolution of the *B. cereus*
67 group that may contribute to a more accurate characterization of these organisms. With the
68 availability of a critical mass of compiled genomic data and the development of advanced
69 computational tools for genome mining, one such approach consists of the bioinformatics
70 exploration of an extensive range of potential biosynthetic gene clusters (BGCs) for the
71 presence among the *B. cereus* group.

72 BGCs are a physical grouping of all the genes required to produce secondary metabolites,
73 including pathway-specific regulatory genes [9]. BGCs mining is the process of identifying
74 and characterizing these clusters to understand the biosynthesis of natural products and to
75 discover new natural products and biosynthetic pathways [10]. BGCs are the source of many
76 important natural products, such as antibiotics, anti-cancer compounds, and enzymes [11].
77 BGCs mining can lead to the discovery of new natural products with unique properties, and
78 the development of new methods for natural product production [12]. BGCs can be used as

79 a source of new enzymes and biosynthetic pathways that can be exploited for biotechnology
80 applications such as the production of biofuels, fine chemicals, and enzymes for
81 bioremediation [13]. Many natural products produced by BGCs have medicinal properties
82 and can be used as leads for drug development. BGCs mining can lead to the discovery of
83 new natural products with potential therapeutic applications [14]. BGCs are often horizontally
84 transferred between bacteria [15]. To gain a genetic advantage, it is postulated that elements
85 in BGCs are horizontally acquired across species for quick adaptation to a new environment
86 [16]. Studying their evolution and distribution can provide insights into the evolution and
87 ecology of bacteria [17].

88 With recent developments in next-generation sequencing and advancements in genome
89 mining tools, it became possible to computationally identify thousands of BGCs and draw a
90 global map of BGCs within a group of bacteria that allow us to systematically explore those
91 of interest [9]. To overcome this challenge, researchers are increasingly using bioinformatics
92 tools such as ClusterFinder [18], antiSMASH [19], and Big-scape [20], which can help
93 automate the process of BGCs identification and classification. Additionally, efforts are being
94 made to establish a standardized nomenclature for BGCs and to create a comprehensive
95 database of BGCs, which lead to the establishment of the MIBiG database [21, 22] which
96 can help facilitate data sharing and comparison among researchers.

97 Big-scape is a bioinformatics tool that can be used to study BGCs in bacteria using an
98 automated streamlined pipeline [20]. It enables researchers to identify, classify, and compare
99 BGCs across different bacterial strains and can be used to infer the biosynthetic pathways
100 and natural products associated with each BGC, providing valuable insights into the
101 evolution and adaptation of bacteria to different environments, as well as the discovery of
102 new natural products and biosynthetic pathways. Automated bioinformatics pipelines and
103 manual bioinformatics are both useful methods for analyzing biological data, but they have
104 different advantages and disadvantages. Automated pipelines are more efficient and can be
105 run by researchers with minimal bioinformatics experience [23], while manual bioinformatics

106 is more flexible and customizable, but requires more expertise [24]. The choice of which
107 method to use will depend on the specific research question, the amount and type of data,
108 and the available resources.

109 Based on the known biosynthesis pathways potentially involved in the production of
110 specialized metabolites by *Bacillus* and closely related species, BGC predictions rely on the
111 development of bioinformatics tools and algorithms design to search for conserved motifs of
112 specific pathways; including peptide synthetases (NRPSs), polyketide synthases (PKSs),
113 and ribosomally synthesized and post-translationally modified peptides (RiPPs) pathways [25]

114 Despite BGCs being important among the *Bacillus cereus* group, there is currently limited
115 data on their conservation across the different strains of this group. Further research is
116 required to better understand the diversity and distribution of BGCs among this group.
117 Nevertheless, highlighting the synteny of BGCs in bacteria is one challenging yet beneficial
118 task.

119 Synteny refers to the preservation of gene order and chromosomal location among different
120 organisms [26]. One benefit of studying this phenomenon is to guide the discovery and
121 characterization of new natural products produced by these bacteria. By identifying the
122 synteny of BGCs among different strains, researchers can infer the presence of similar BGCs
123 in other strains, and can then use this information to guide their search for new natural
124 products. Another benefit is that it can provide insight into the evolutionary relationships
125 among different taxa. The conservation of BGCs in terms of chromosomal location and gene
126 organization among different strains can be used to infer evolutionary relationships among
127 these strains, and can also help to understand how these bacteria have adapted to different
128 environments over time.

129 Advances in sequencing technologies, high-throughput screening techniques, and improved
130 computational methods have led to a rapidly increasing number of BGCs being identified.
131 This has created a growing need for effective and standardized methods for BGC
132 classification [12]. There have been several efforts to establish a unified system for BGC

133 classification, such as the antiSMASH platform, but the field is still in its early stages and
134 much work remains to be done to develop a comprehensive and widely-accepted
135 classification system [19].

136 Moreover, during a screening for potentially beneficial endophytic bacteria, the strain *Bacillus*
137 sp. BH32, which belongs to the *B. cereus* group, was isolated from *Atriplex halimus* L., a
138 halophyte sampled from a continental hypersaline region (Sebkha) in Djelfa province, Algeria.
139 The strain, which was proven to help tomato and wheat seedlings tolerate salt stress at
140 various levels, was consecutively genome analyzed to determine putative mechanisms
141 involved in salt tolerance and plant promotion [27], and thus, is part of the genome dataset of
142 this study, along with its closely-related strains.

143 In the present study, we investigated the BGCs of the *Bacillus* sp. strain BH32 at a genomic
144 level along with its closest non-type strains to explore the conservation and putative evolution
145 patterns of BGCs among the *Bacillus cereus* group, and to highlight singletons. Based on a
146 combined strategy (manual and automatic), we aimed to establish the basis of a rational
147 classification of BGCs among the *B. cereus* group.

148 **Methods**

149 **1. Presenting the dataset**

150 The genomic dataset is composed of *Bacillus* sp. BH32 genomes and its closest non-type
151 strains genomes; all part of the *Bacillus cereus* group. *Bacillus* sp. BH32 is a beneficial
152 endophyte, isolated during a previous study from *Atriplex halimus* L., a halophyte sampled
153 from an Algerian continental Sebkha from the province of Djelfa. This strain was proven to
154 help tomato and wheat seedlings tolerate salt stress at various levels [27]. The selection of
155 the closest non-type strains of *Bacillus* sp. BH32 was performed with BLASTn 2.10.0+ [28],
156 using the whole genome of *Bacillus* sp. BH32 as a query against the NCBI's "Complete
157 prokaryotic genomes" database (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>, accessed on 02-04-
158 2020), targeting the first most significant 16 hits (high-quality complete genomes only) (**table**
159 **1**).

160 **2. BGCs counts and distribution**

161 BGCs from the dataset were predicted using antiSMASH ver. 5.1.2 [19]. BGCs counts were
162 compared quantitatively, after preparing a BGCs-types/counts matrix comprising all the
163 compared genomes, and visualized using Morpheus
164 (<https://software.broadinstitute.org/morpheus/>). The generated heatmap showed the counts
165 of BGCs (cyan-white-fuchsia, from 0 to 5) by type, and their distribution among the strains
166 (dendrograms were generated by hierarchical clustering using one minus Spearman rank
167 correlation) (**Fig. 1**).

168 **3. BGCs BLASTp comparisons**

169 A BLASTp (ver. 2.9.0+) comparison of *Bacillus* sp. BH32 BGCs against BGCs of the closest
170 non-type strains was done, using the BGCs amino-acid sequences (for each genome, a
171 multifasta file was used, containing the amino acid sequences of all the genes constituting
172 the whole predicted BGCs, prepared from the GenBank records generated by antiSMASH),
173 using BRIG ver. 0.95 [29]. The reference ring constitutes the detected BGCs in *Bacillus* sp.
174 BH32, while each ring represents BGCs of each genome of the closest non-type strains (the
175 color shades represent sequence identity, the greyer it gets; the lower the percentage
176 identity). The 3 outermost rings (from inside to outside) represent: genes composing BGCs in
177 *Bacillus* sp. BH32 according to their respective location in the genome; unique genes
178 (sequence identity < 30% with 100% of the strains) and rare genes (sequence identity < 30%
179 with 80% the strains) with their locus tag IDs; and the corresponding regions/BGCs types
180 (last ring) as detected by antiSMASH (**Fig. 2**). The annotation from antiSMASH and
181 PROKKA of each gene was retrieved using their corresponding locus tag IDs.

182 **4. BGCs gains, losses, and rearrangements**

183 **Manual approach**

184 First, fasta files were prepared; each containing concatenated amino acid sequences of all
185 BGCs for each strain (with conserved order of occurrence in the respective genome).

186 Sequences were aligned using MAFFT ver. 7.221.3 [30]. An ML tree was inferred from the
187 aligned sequences by using the Maximum Likelihood method based on the JTT matrix-based
188 model [31]. The tree harboring the highest log likelihood (-248721.7967) was selected. The
189 percentage of trees in which the associated sequences clustered together is displayed next
190 to the branches (branch-support values). Initial tree(s) for the heuristic search were
191 automatically obtained based on Neighbor-Join and BioNJ algorithms applied to a matrix of
192 pairwise distances (estimated by a JTT model), and then selecting the topology with the
193 highest log likelihood value. The tree was drawn to scale, with branch lengths reflecting the
194 number of substitutions per site. The analysis involved 17 amino acid sequences (full record
195 of BGCs from each strain). All positions containing gaps and missing data were eliminated.
196 There were a total of 28898 positions in the final dataset. Evolutionary analyses were
197 conducted in MEGA7 [32]. The final tree was drawn using Adobe Illustrator CC 18.1.1 (**Fig. 3**
198 **a**).

199 BGCs synteny among the strains was investigated with MAUVE ver. 2.4.0 [33–37], using the
200 alignment of concatenated Genbank sequences of the BGCs of each strain, according to
201 their order of appearance in the respective genome, where each row shows the conservation
202 and orientation of BGCs in the corresponding strain in the ML tree after alignment. *Bacillus*
203 *thuringiensis* c25 was set as a reference, according to the phylogenetic analysis (most
204 distant). Homologous segments indicating orthologous clusters (from MAUVE alignment
205 diagram) with a locally collinear block (LCB) weight ≥ 773 were confirmed by
206 annotation. BGCs relative orientation and order conservation were highlighted (**Fig. 3 b**).

207 **BIGSCAPE automatic approach and GCFs nomenclature proposal**

208 The selected antiSMASH profiles of the *B. cereus* genomes were used to identify the gene
209 cluster families GCFs using BiG-SCAPE v.1.1.0. with default parameters [20]. The output of
210 the BiG-SCAPE was exploited to propose a classification of GCFs among the *B. cereus*
211 group, based on the respective BGCs classes/types, as organized by the unsupervised
212 Machine Learning clustering approach employed by BiG-SCAPE. Since the BiG-SCAPE

213 output consists of BGCs grouped into GCFs with arbitrary codes, we propose a systematic
214 and grounded nomenclature of the highlighted families. In our proposal, families names are
215 composed of: a “BC” prefix that stands for *Bacillus cereus* (group) (“BCS” in case of
216 singletons), followed by the type (e.g. bacteriocin), then an increasing number that should
217 follow the order of detection/characterization during this study (or future studies). The
218 product’s name should appear instead of merely the type, in case the cluster product is
219 known (from the MIBIG database). The proposed reference BGCs are based on ‘*exemplars*’,
220 according to the affinity propagation clustering approach used by BIG-SCAPE. Harmony
221 between the manual and automatic approaches was assessed manually based on the
222 correspondence between the BGCs region codes as given by antiSMASH and those
223 displayed among the BIG-SCAPE output. The similarity networks were arranged to fit their
224 corresponding figures, while the original BIG-SCAPE output, the used dataset (antiSMASH
225 profiles), and the proposed reference BGCs files (included in the antiSMASH profiles) are
226 available as a supplementary material. Final figures of the GCFs classes distribution (Fig.
227 5-8) were composed under Adobe Illustrator CC 18.1.1.

228

229 **Results**

230 **1. BGCs distribution and counts in the *B. cereus* group**

231 The BGCs were first compared in terms of distribution and counts among the compared
232 genomes (**Fig. 1**). All the strains have similar counts of NARPS-like, LAP-bacteriocin,
233 siderophore, and betalacton BGCs (one for each type). *B. cereus* JHU has the highest
234 number of bacteriocin BGCs (4), followed by *B. thuringiensis* L-76-01 and *B. thuringiensis*
235 c25 (3 each), while the other strains have 2 (except *Bacillus* sp. BH32, only one bacteriocin
236 BGC was detected). *B. cereus* ZB201708 is the only strain where antiSMASH failed to detect
237 a terpene BGC (all the other strains have one). *B. cereus* ZB201708, *B. thuringiensis* YGd
238 22-03, and *B. thuringiensis* c25 are the only strains to have a lanthipeptid BGC (one for each).
239 NRPS BGC-type has the highest counts among the selected strains, ranging from 3 to 5.

240 *Bacillus cereus* A1, *Bacillus bombysepticus* Wang, *Bacillus thuringiensis* c25 and *Bacillus*
241 *cereus* FORC087 have 3 NRPS BGCs ; *Bacillus* sp. BH32, *Bacillus wiedmannii* PL1, *Bacillus*
242 *thuringiensis* YGd22-03 and *Bacillus thuringiensis* serovar galleriae 4G5 have 4 ; and the
243 remaining strains have 5.

244 The strains *Bacillus thuringiensis* serovar morrisoni BGS, *Bacillus cereus* G9842, *Bacillus*
245 *thuringiensis* BT59, *Bacillus thuringiensis* HD1002, *Bacillus thuringiensis* JW-1, and *Bacillus*
246 *thuringiensis* HD 789 have the same BGC counts (NRPS = 5, NRPS-like = 1,
247 LAP/bacteriocin = 1, bacteriocin = 2, terpene = 1, betalactone = 1, siderophore = 1,
248 lanthipeptide = 0). Same thing for *Bacillus wiedmannii* PL1 and *Bacillus thuringiensis* serovar
249 galleriae 4G5 (NRPS = 4, NRPS-like = 1, LAP/bacteriocin = 1, bacteriocin = 2, terpene = 1,
250 betalactone = 1, siderophore = 1, lanthipeptide = 0). *Bacillus cereus* A1, *Bacillus*
251 *bombysepticus* Wang and *Bacillus cereus* FORC087 have the same profile as well (NRPS =
252 3, NRPS-like = 1, LAP/bacteriocin = 1, bacteriocin = 2, terpene = 1, betalactone = 1,
253 siderophore = 1, lanthipeptide = 0). The remaining strains have different BGC count profiles
254 **(Fig. 1)**.

255 [Insert Figure 1 here]

256 **2. Amino acid sequence conservation of the *B. cereus* group BGCs**

257 The BLASTp comparison of BGCs from *Bacillus* sp. BH32 vs. its closest non-type strains
258 **(Fig. 2)** shows heterogeneous content in terms of predicted protein sequences from detected
259 BGCs. The most conserved BGC seems to be the terpene BGC (region 2.1), while the less
260 conserved one seems to be the NRPS BGC (region 2.2).

261 Although the antiSMASH failed to detect the terpene biosynthetic genes, they seem to be
262 present in the BLASTp output of *Bacillus cereus* ZB201708 (14th ring, from the inner ring to
263 the outer ring), which could be due to assembly concerns, or else, the constituting genes of
264 this terpene cluster might be distributed in other BGCs from *Bacillus cereus* ZB201708.

265 NRPS BGC (region 2.2) from *Bacillus* sp. BH32 has very low sequence identity when
266 compared to its counterparts in *B. bombysepticus* Wang, *B. cereus* A1 and *B. thuringiensis*
267 c25 (<30%). Considering that the other NRPS BGCs from *Bacillus* sp. BH32 (regions 6.1, 8.1,
268 12.1, and 15.1) have higher sequence identity % with the 3 corresponding BGCs of the
269 aforementioned strains, which confirms the absence of a counterpart of the NRPS BGC from
270 region 2.2 in these strains as suggested by the counts (5 NRPS BGCs in *Bacillus* sp. BH32
271 against only 4 in *B. bombysepticus* Wang, *B. cereus* A1, and *B. thuringiensis* c25). This
272 cluster has two rare genes (identity below 30% for 80% of the strains): ctg2_300 =
273 unknown/hypothetical protein; ctg2_321 = unknown / IS200/IS605 family transposase ISAsp8.
274 The siderophore BGC (region 3.1) encoding petrobactin (100%) and the LAP/bacteriocin
275 BGC (region 5.1, unknown encoded product, closest BGCs with overall gene sequence
276 identity >70% from strains *Bacillus cereus* ZB201708, *Bacillus cereus* JHU and *Bacillus*
277 *thuringiensis* L-7601) seem to be well conserved, with a slight variation on the composing
278 genes. The NRPS BGC (region 6.1) is another example of a well-conserved cluster across
279 the compared strains (closest known cluster = bacillibactin, 46%).
280 The NRPS-like BGC (region 8.1) is present among all the strains. Nevertheless, *Bacillus* sp.
281 BH32 has one unique gene (sequence identity below 30% for 100% of the strains) in this
282 region (ctg8_81 = unknown/hypothetical protein), and a rare one (ctg8_90 =
283 unknown/hypothetical protein).
284 The NRPS BGC (region 12.1) shows variation in its first 19 genes (~ 1/3 of the total BGC),
285 including one unique gene (ctg12_31 = unknown / Spore germination protein B1), and one
286 rare gene (ctg12_35 = biosynthetic-additional; (smcogs) SMCOG1012:4'-
287 phosphopantetheinyl transferase / 4'-phosphopantetheinyl transferase Sfp).
288 The NRPS BGC (region 15.1) from *Bacillus* sp. BH32 has one unique gene (ctg15_85 =
289 biosynthetic-additional; (smcogs) SMCOG1028: crotonyl-CoA reductase - alcohol
290 dehydrogenase / Zinc-type alcohol dehydrogenase-like protein). This NRPS BGC doesn't

291 have a counterpart in *Bacillus cereus* FORC087, confirming the BGCs count (5 NRPS BGCs
292 in *Bacillus* sp. BH32 for only 4 in *Bacillus cereus* FORC087).

293 The bacteriocin BGC (region 24.1) from *Bacillus* sp. BH32 is another example of a well-
294 conserved cluster across its closely related non-type strains.

295 Finally, the betalacton BGC (region 27.1) from *Bacillus* sp. BH32 (fengycin, 40%) arbores
296 two unique genes (ctg27_13 = unknown/hypothetical protein; ctg27_20 =
297 unknown/hypothetical protein). Globally, all BGCs from *Bacillus* sp. BH32 are present in at
298 least 13 of its closest strains (>81%), either very well-conserved (regions) or less conserved,
299 with some rare/unique genes (**Fig. 2**).

300 *[Insert Figure 2 here]*

301 **3. Manual approach for *B. cereus* group BGCs synteny**

302 The ML tree based on BGCs amino acids sequences (**Fig. 3 a**) shows 4 clades:

303 Clade I: *B. thuringiensis* c25, *B. cereus* JHU, *B. thuringiensis* HD1002 and *B. thuringiensis*
304 BT-59; Clade II: *Bacillus* sp. BH32; Clade III: *B. wiedmannii* PL1, *B. cereus* FORC087, *B.*
305 *thuringiensis* serovar galleriae 4G5, *B. bombysepticus* Wang, *B. cereus* A1 and *B.*
306 *thuringiensis* YGd22-03; Clade IV: *B. thuringiensis* HD 789, *B. thuringiensis* JW1, *B.*
307 *thuringiensis* L-7601, *B. cereus* G9842, *B. thuringiensis* serovar morrisoni BGS and *B.*
308 *cereus* ZB201708.

309 From the ML tree (**Fig. 3 a**), *B. thuringiensis* c25 has been considered as a reference for
310 BGCs conservation (the most distant genome among the compared strains). *Bacillus* sp.
311 BH32 which usually clusters with *B. cereus* ZB201708 and *B. thuringiensis* serovar morrisoni
312 BGSC 4AA1 (GBDP and ANI analysis), doesn't seem to be close to these strains when it
313 comes to its BGCs, showing thus more complexity among this intricate group, in terms of
314 BGCs.

315 The analysis of the BGCs conservation is shown in (**Fig. 3 b**). Here, BGCs are mentioned
316 according to their attributed BGC tag number (**Fig. 3 b legends**). Considering *B. thuringiensis*

317 c25 as a reference, we observe that the orthologous BGCs group that seems to be best
318 conserved in order and appearance is the one composed of BGC3 (bacteriocin: unknown),
319 BGC4 (bacteriocin: unknown), BGC5 (betalacton, fengycin), BGC6 (NRPS: gramicidin in *B.*
320 *cereus* JHU; nostopeptolide A2 in *B. thuringiensis* serovar galleriae 4G5; unknown in the
321 remaining strains), BGC7 (NRPS: bacillibactin), BGC8 (siderophore: petrobactin) and BGC9
322 (linear azol(in)e-containing peptides “LAP” bacteriocin). This group was named “*synteny*
323 *block A*”, which appears in this order in all strains except *B. wiedmannii*, where BGC3
324 appears to be in the beginning, unusually separated from the other BGCs.

325 The second conserved group is composed of BGC1 (terpene: molybdenum co-factor) and
326 BGC2 (NRPS: polyoxypeptin, except for *B. thuringiensis* c25, unknown), named “*synteny*
327 *block B*”. BGC10 (bacteriocin: unknown) is only present in *B. thuringiensis* C25 (thus, will not
328 be mentioned here again).

329 In *B. cereus* JHU, BGC1 and BGC2 (*synteny block B*) are still in the same order, as well as
330 for the BGC3, BGC4, BGC5, BGC6, BGC7, BGC8 and BGC9 (*synteny block A*). Meanwhile,
331 4 new BGCs appeared: BGC13 (bacteriocin: unknown) and BGC14 (bacteriocin: unknown)
332 both in the beginning; BGC15 (NRPS-like 2: unknown) between the two synteny blocks A
333 and B; and BGC16 (NRPS-Polyketide: chejuenolide A / chejuenolide B) at the end. The
334 BGC11 (NRPS-like 1: unknown) moved to the beginning, right before the *synteny block B*.
335 BGC12 (lanthipeptid: cerecidin / cerecidin A1 / cerecidin A2 / cerecidin A3 / cerecidin A4 /
336 cerecidin A5 / cerecidin A6 / cerecidin A7) is lost.

337 BGC13 and BGC14 (both unknown bacteriocins) are only found in *B. cereus* JHU.

338 As for *B. cereus* JHU, *B. thuringiensis* HD-1002 and *B. thuringiensis* BT-59 have the two
339 synteny blocks with the BGC15 in between. BGC16 is at the right end of the *synteny block B*.
340 BGC12, BGC13, and BGC14 are lost and the BGC11 jumped to the end.

341 Since *Bacillus* sp. BH32 genome doesn't constitute a complete one, the observations for the
342 BGCs conservation for this strain concern only their presence, not their order (the use of a
343 reference genome for the reordering of the contigs doesn't take into count the possible

344 genomic rearrangements, and thus biases the synteny analysis, so we kept the contigs as
345 they resulted from the assembly, from the largest to the smallest one).

346 In *Bacillus* sp. BH32, the BGC4 (bacteriocin with no closest known cluster) seems to be
347 missing, while conserved across all the analyzed genomes of its closely related non-type
348 strains, suggesting this absence due to the incompleteness of the genome, and should be
349 found once the genomic gaps filled. The remaining BGCs from *synteny blocks A* and *B* are
350 all present in *Bacillus* sp. BH32.

351 In the BGCs clades III and IV, all BGCs are inversely oriented in comparison with the strains
352 of the first clade (including *B. thuringiensis* C25). In *B. wiedmannii*, BGC11 jumped between
353 *synteny blocks A* and *B*, while BGC3 and BGC15 are right before the *synteny block B*. This
354 strain doesn't have the BGC12.

355 *B. cereus* FORC087 has the same homologous BGCs from 1 to 9 (inversely oriented), with
356 BGC11 in the beginning. BGC15 and BGC16 are absent. *B. thuringiensis* serovar *galleriae*
357 4G5 has *synteny blocks A* and *B*, and the BGC15 reappears next to BGC2 from *synteny*
358 *block B*. As for *B. cereus* FORC087, BGC 11 is still in the beginning.

359 For *B. bombysepticus* Wang, BGC2 is missing. Despite the BGC15 (NRPS-like: unknown)
360 bears homologous segments linking it with the BGC2 (NRPS: polyoxypeptin), after manually
361 checking the annotation, it confirmed the assignment to its closest orthologous BGC, the
362 BGC15. The BGC1 is there, at the end. BGC11 is still in the beginning, before the *synteny*
363 *block A*.

364 *B. cereus* A1 has an arrangement similar to *B. bombysepticus* Wang, except for BGC11
365 jumping to the end. For *B. thuringiensis* YGd22-03, BGC2 reappears, in the *synteny block B*,
366 while the BGC12 lost since *B. thuringiensis* c25 reappears here.

367 For clade IV, the *synteny blocks A*, *B*, and BGC15 hold the same position in all its strains. In
368 *B. thuringiensis* HD-789, BGC11 is in the end, next to the *synteny block B*. BGC16 reappears,
369 to stand in the beginning.

370 *B. thuringiensis* HD-789, *B. thuringiensis* JW-1, *B. thuringiensis* serovar morrisoni BGSC
371 4AA1 *B. cereus* G9842, *B. thuringiensis* L76-01 and *B. cereus* ZB201708 have the BGC11
372 jumping again to the front position, followed by the BGC16, except *B. cereus* ZB201708
373 where the appearing BGC18 (lanthipeptid: surfactin) appeared between BGC11 and BGC16;
374 and the *synteny blocks A* and *B* with BGC15 in between. *B. thuringiensis* L76-01 and *B.*
375 *cereus* ZB201708 are characterized with a new BGC for each one: BGC17 (bacteriocin:
376 unknown) and BGC19 (bacteriocin: unknown) respectively, at the end.

377 BGCs of the *synteny blocks A* and *B* being consistently present in all the strains, will naturally
378 tend to occur in chromosomic locations among the analyzed genomes. BGC11 (NRPS-like 1:
379 unknown) is present across all strains, at various positions, suggesting its possible jumping
380 between chromosomic and plasmidic locations.

381 BGC15 (NRPS-like 2: unknown) is largely present (except in *B. thuringiensis* c25 and *B.*
382 *cereus* FORC087), and is located most of the time between synteny blocks A and B,
383 suggesting a more likely chromosomic location as well. BGC12 (lanthipeptid: cerecidin /
384 cerecidin A1 / cerecidin A2 / cerecidin A3 / cerecidin A4 / cerecidin A5 / cerecidin A6 /
385 cerecidin A7), when present, is next to BGC11.

386 The similarity plots in the homologous segments are confirming the BLASTp results from the
387 BRIG ring generator (**Fig. 2**). The present synteny analysis can give direction to the
388 prediction of the presence of certain BGCs, as we showed that they are usually conserved in
389 presence and order. It can also be used to detect the evolution and distribution of
390 homologous BGCs across the complex *B. cereus* group. The observed BGCs positions
391 should be further investigated, in their genomic context, by studying synteny not only for
392 BGCs, but extending the analysis to the whole genomic content.

393 *[Insert Figure 3 here]*

394 **4. BIG-SCAPE results**

395 The BIG-SCAPE pipeline was used as well, in line with the manual approach, as a
396 complementary strategy to investigate the gene cluster families (GCFs) among the *B. cereus*
397 group.

398 There was a total of 198 BGCs, organized into 4 main classes (Fig. 4, table 2):

- 399 ▪ NRPS: consisting of a total of 89 BGCs, it is the biggest class representing 44.94% of
400 the total BGCs;
- 401 ▪ RiPPs: 58 BGCs fall into this class, representing 29.29% of the sum of the analyzed
402 BGCs;
- 403 ▪ Terpenes: represent 5.58% of the total BGCs (with 17 BGCs identified as coding for a
404 terpene);
- 405 ▪ Other (remaining classes): composed of 34 BGCs, representing 17.17% of the
406 studied BGCs.

407 We suggest a framework for an expanding classification of the *B. cereus* group BGCs, based
408 on a set of reference BGCs described in this work (tables 3-4; supplementary figures S1, S2,
409 S3, and S4), as anchoring points to affiliate unknown query BGCs to the proposed
410 families/clans accordingly. Such reference BGCs should be included in future attempts to
411 assign a set of BGCs (from genomes belonging to the *B. cereus* group) using the BIG-
412 SCAPE pipeline into one of the proposed families in the current proposal. Moreover, the
413 same strategy can be reproduced for other bacterial groups.

414 *[Insert Figure 4 here]*

415 Hence, we propose the following families:

416 **The RiPPs gene cluster families of the *Bacillus cereus* group**

417 The BIG-SCAPE output yielded 13 RiPPs GCFs (Fig. 5), with an average of 4 BGCs by
418 family. It is the 2nd largest group of BGCs with 58 clusters (including 7 singletons, table 4),
419 without any known reference BGC in the MIBIG database (Fig. 5a). Following the proposed

420 nomenclature, the identified families were named as follows (GCF name/ reference BGC
421 strain and genomic region) (Fig. 5b; table 3; supp. Fig. S1):

- 422 • BC-Bacteriocin 1: *B. cereus* ZB201708, NZ_CP030982.1.region009
- 423 • BC-Bacteriocin 2: *B. thuringiensis* HD1002, NZ_CP009351.1.region005
- 424 • BC-Bacteriocin 3: *B. thuringiensis* serovar *galleriae* 4G5, NZ_CP010089.1.region007
- 425 • BC-Lasso peptide 1: *B. cereus* A1, NZ_CP015727.1.region001
- 426 • BC-Bacteriocin 4: *B. thuringiensis* serovar *morrisoni* BGSC 4AA1,
427 NZ_CP010577.1.region009
- 428 • BC-Bacteriocin 5: *B. cereus* ZB201708, NZ_CP030982.1.region010
- 429 • BC-Cerecidin: *B. thuringiensis* YGd22-03, NZ_CP019230.1.region011

430 Families BC-Bacteriocin 1, BC-Bacteriocin 2 and the singleton BC-Bacteriocin 3 are part of
431 the RiPPs Clan I, while the families BC-Bacteriocin 4 and BC-Bacteriocin 5 are part of the
432 RiPPs Clan II (Fig. 5c).

433 The BC-Cerecidin family was named based on the cerecidin compounds (lanthipeptides),
434 with the closest known BGC (MIBIG database) being that of the cerecidin / cerecidin A1 /
435 cerecidin A2 / cerecidin A3 / cerecidin A4 / cerecidin A5 / cerecidin A6 / cerecidin A7
436 (lanthipeptides putative class II), with a similarity of 94%.

437 The BC-Lasso peptide 1 was named based on its similarity with a LAP–bacteriocin BGC
438 (sactipeptide/lassopeptide).

439 Meanwhile the singletons are (Fig. 5c, table 4; supp. Fig. S1):

- 440 • BCS-Bacteriocin 1: *B. cereus* JHU, CP046511.1.region001
- 441 • BCS-Surfactin-like 1: *B. cereus* ZB201708, NZ_CP030982.1.region002
- 442 • BCS-Bacteriocin 2: *B. cereus* JHU, CP046511.1.region002
- 443 • BCS-Bacteriocin 3: *B. thuringiensis* c25, NZ_CP022345.1.region010
- 444 • BCS-Bacteriocin 4: *B. cereus* ZB201708, NZ_CP030982.1.region014
- 445 • BCS-Bacteriocin 5: *B. thuringiensis* L-7601, NZ_CP020002.1.region013

446 The largest family was named BC-Lasso peptide 1 (with 17 analogous BGCs present in all of
447 the studied genomes), followed by BC-Bacteriocin 4 (with 14 analogous BGCs). The
448 remaining families have somewhere between 7 and 1 BGC; while the smallest family is
449 composed of only 2 analogous BGCs (the BC-Cerecidin family) (Fig. 5b,c).

450 Furthermore, the RiPPs class is characterized by 7 singleton BGCs (Fig. 5c), displaying the
451 highest number of singletons among the analyzed BGCs classes, which makes the RiPPs
452 class the best niche for peculiar and likely to be unique secondary metabolites in the *B.*
453 *cereus* group, in addition to pinpointing to possible horizontal gene transfer (HGT) at a BGC-
454 level.

455 The singleton BCS-Surfactin-like 1 was named this way, due to its similarity (8%) with a
456 known lanthipeptide BGC coding for the surfactin lipopeptide.

457 *[Insert Figure 5 here]*

458 **The NRPS gene cluster families of the *Bacillus cereus* group**

459 For the NRPS class (Fig. 6), the BIG-SCAPE approach highlighted 6 families, without any
460 reference BGC being known, which pinpoints further the future possible discovery of
461 numerous novel metabolites in this intricate group. The NRPS class has the most
462 represented families among the studied genomes, composed of 89 BGCs with an average
463 number of 15 BGCs by family, and a staggering number of links (634 links) (Fig. 6a). No
464 clans were proposed for this class.

465 We proposed the following names for the described NRPS families (Fig. 6b, table 3; supp.
466 Fig. S2):

- 467 • BC-NRPS 1: *Bacillus*. sp. BH32, SJAS00000000.2.2.region002
- 468 • BC-NRPS 2: *Bacillus* sp. BH32, SJAS00000000.2.15.region001
- 469 • BC-Bacillibactin-like: *B. cereus* A1, NZ_CP015727.1.region003
- 470 • BC-NRPS-like 1: *B. cereus* G9842, NC_011772.1.region001

- 471 • BC-NRPS 3: *B. cereus* G9842, NC_011772.1.region006
- 472 • BC-NRPS 4: *B. cereus* ZB201708, NZ_CP030982.1.region003

473 The BC-NRPS 3 (showing low similarity with known BGCs coding for nostocyclopeptide A2:
474 28%, gramicidine: 16% from the MIBIG database; and the kurstakin C12 with a score of
475 0.639 according to the norine database) and the BC-Bacillibactin-like (showing a genes
476 similarity of 46% with the known Bacillibactin BGC from the MIBIG database) families are
477 represented in all of the studied genomes (17 analogous BGCs for both of them) (Fig. 6b,c).

478 It is noteworthy to mention that the BC-NRPS 4 family has a reference BGC showing little
479 similarity with the known BGCs coding for polyketides chejuenolide A / chejuenolide B (7%)
480 from the MIBIG database (table 3; supp. Fig. S2).

481 *[Insert Figure 6 here]*

482 **The terpene gene cluster family of the *Bacillus cereus* group**

483 For the terpene BGCs (Fig. 7), they all belong to the same family, named BC-Terpene 1,
484 which bears similarity with the molybdenum cofactor (reference BGC:
485 NZ_CP020002.1.region012 from *B. thuringiensis* L-7601, with a similarity of 17% with the
486 molybdenum cofactor from the MIBIG database) (table 3; supp. Fig. S3).

487 It is the most conserved GCF (137 link in the corresponding similarity network) (Fig. 7 b)
488 among the studied *B. cereus* genomes. The phylogenetic analysis (Fig. 7 c) confirms this
489 conservation, showing two almost identical clades, with the exception of two different genes
490 bearing 2 different pfam domains: the PF08445 FR47-like protein domain (Fig. 7 c, clade I)
491 and the PF00583 Acetyl transferase GNAT family domain (Fig. 7 c, clade II).

492 *[Insert Figure 7 here]*

493 **The remaining siderophores and betalactone gene cluster families of the *Bacillus*** 494 ***cereus* group**

495 Regarding the remaining BGC classes (the 3 GCFs gathered under the section “others” in
496 the BIG-SCAPE output)(Fig. 8), namely the BC-Petrobactin 1, BC-Petrobactin 2 families

497 (both part of the BC-Petrobactin clan, coding for petrobactin siderophores, with a genes
498 similarity of 100% with the known petrobactin BGC) and the BC-Fengycin-like 1 family
499 (coding for betalactone products, with a genes similarity of 40% with known fengycin BGC
500 from the MIBIG database), are another example of well-conserved families across the *B.*
501 *cereus* group genomes, with an average number of BGCs by family of 11 (Fig. 8a).

502 The proposed families and their respective reference BGCs are (Fig. 8, table 3; supp. Fig.
503 S4):

- 504 • BC-Petrobactin 2: *B. cereus* A1, NZ_CP015727.1.region002
- 505 • BC-Petrobactin 1: *Bacillus* sp. BH32, SJAS000000000.2.3.region001
- 506 • BC-Fengycin-like 1: *B. cereus* G9842, NC_011772.1.region007

507 Indeed, for the BC-Petrobactin clan (Fig. 8c), it has analogous BGCs across all of the studied
508 genomes, being part of either the BC-Petrobactin 1 or BC-Petrobactin 2 families; and the BC-
509 Fengycin-like 1 is represented in all of the 17 genomes (Fig. 8b), which is in line with the
510 broad prevalence of Fengycin-like compounds reported in numerous strains belonging to the
511 *B. cereus* group and their closely related taxa [38].

512 *[Insert Figure 8 here]*

513 **5. Harmony between the manual and the automatic approach**

514 Based on the manual approach (**Fig. 3**), we highlighted 2 BGCs synteny blocks “*synteny*
515 *block A*” and “*synteny block B*”.

516 After performing a BIG-SCAPE analysis of the antiSMASH profiles of the studied genomes
517 (Fig. 5, 6, 7, and 8), we could link the BGCs of the highlighted synteny blocks with their
518 corresponding proposed families as follow:

519 For the **synteny block A**:

- 520 • BGC3 (bacteriocin: unknown) belongs to the **BC-Bacteriocin 4** family;

- 521 • BGC4 (bacteriocin: unknown) belongs to the **BC-Bacteriocin 1-2-3** families (BC-
522 Bacteriocin clan I);
- 523 • BGC5 (betalacton, fengycin) belongs to the **BC-Fengycin-like 1** family;
- 524 • BGC6 (NRPS: gramicidin in *B. cereus* JHU; nostopeptolide A2 in *B. thuringiensis*
525 serovar galleriae 4G5 ; unknown in the remaining strains) belongs to the **BC-NRPS 3**
526 family;
- 527 • BGC7 (NRPS: bacillibactin) belongs to the **BC-Bacillibactin-like** family;
- 528 • BGC8 (siderophore: petrobactin) belongs to the **BC-Petrobactin** clan; and
- 529 • BGC9 (linear azol(in)e-containing peptides “LAP” bacteriocin) belongs to the **BC-**
530 **Lasso peptide 1** family.

531 For the second conserved group named “**synteny block B**”:

- 532 • BGC1 (terpene: molybdenum co-factor) belongs to the **BC-Terpene 1** family; and
- 533 • BGC2 (NRPS: polyoxypeptin, except for *B. thuringiensis* c25, unknown) belongs to the
534 **BC-NRPS 1** family.

535 That is to say, the synteny block A is a series of BGCs belonging (in order) to the families
536 /clans (with the corresponding number of analogous BGCs found in each family):

537 BC-Bacteriocin 4 (14 BGCs), BC-Bacteriocin clan I (7+8+1 BGCs of the BC-Bacteriocin 1-2-3
538 families, respectively), BC-Fengycin-like 1 (17 BGCs), BC-NRPS 3 (17 BGCs), BC-
539 Bacillibactin-like (17 BGCs), BC-Petrobactin clan (families 1 and 2, 17 BGCs), BC-Lasso
540 peptide 1 (17 BGCs);

541 And the synteny block B is composed of BGCs belonging (in order) to the families:

542 BC-Terpene 1 (17 BGCs), and BC-NRPS 1 (16 BGCs).

543 The remaining families are either under-represented across the *B. cereus* group or
544 considered as singletons, which confirms the harmony and complementarities between both
545 manual and automated approaches with BIG-SCAPE.

546

547 **Discussion**

548 **1. The importance of BGCs investigations**

549 Deciphering beneficial features in plant growth-promoting bacteria requires research into the
550 encoded parvome (the secondary metabolome inferred from the genome) [39]. Genes
551 responsible for the production of secondary metabolites (SMs) are typically grouped into
552 often quite large and complex BGCs [40]. BGCs are self-contained sets of co-located genes
553 that accomplish the coordinated and regulated biosynthesis of a single set of SM congeners,
554 with some exceptions, including BGCs that lack genes for required modifying enzymes that
555 are located in different parts of the genome; distributed BGCs where two or more sub-
556 clusters located in different parts of the genome collaborate during convergent biosynthesis
557 of a single set of SM congeners; and superclusters that contain intertwined genes for the
558 biosynthesis of more than one set of SM scaffolds [39, 41, 42].

559 BGCs are organized around genes that encode biosynthetic enzymes that yield the SM
560 carbon skeleton (“backbone” enzymes), such as NRPSs, PKSs, PKS–NRPS hybrids, etc.
561 The BGCs also feature genes encoding various enzymes that further modify the SM carbon
562 skeleton (such as cytochrome P450 monooxygenases, various other oxidoreductases, etc)
563 together with genes for transporters, regulators, and self-resistance determinants. In addition,
564 many BGCs also harbor genes for enzymes that synthesize specialized monomers for the
565 corresponding pathway [43].

566 **2. BGCs screening and mining approaches**

567 A biological approach in BGCs screening is considered the best prospect for resolving the
568 potential of microbial parvomes. However, **Baltz** [44] estimated that 10^7 strains would need
569 to be examined to discover a novel class of antibiotics [44]. Thus platforms combining
570 different approaches are the appropriate strategy [45]. Moreover, the apparent failure to
571 uncover the full potential of natural product-producing microorganisms is likely due to the

572 lack of understanding that is required to activate the expression of their BGCs in the
573 laboratory [46]. Hence, new strategies for microbial SMs discovery are being employed,
574 comprising genomics, metabolomics, and analytical tools that permit the study of complex
575 systems [47].

576 Advances in genomics have unveiled a vast reservoir of BGCs in microbial genomes [46].
577 There are two main approaches to predicting BGCs. The first is based on pre-computed
578 pHMMs derived from a set of genes known to participate in SM metabolism to identify
579 sequences of interest [48, 49]. The second uses some function-agnostic criteria, such as
580 synteny conservation or shared evolutionary history, to implicate genes as part of a gene
581 cluster [50] which highlights the importance of synteny analysis among BGCs. Moreover, due
582 to common metabolic functions across distant taxa, approaches employed by SMURF and
583 antiSMASH are enormously successful [51]. Given the size of genomic data, studying BGCs
584 on a case-by-case basis is no longer interesting. Hence, sequence similarity networking
585 approaches can automatically relate predicted BGCs to gene clusters of known function and
586 group them into gene cluster families (GCFs) [12, 46, 52, 53]. “Old school” methods gave
587 way to new workflows that entail: (1) sequencing of the whole genome of strains that produce
588 interesting SMs; (2) bioinformatic prediction of all BGCs; (3) comparison of the predicted
589 BGCs with the retro-biosynthetic assessment of the structure of the target SM; (4)
590 comparative genome analysis with organisms that produce similar SMs and with taxa that
591 are phylogenetically close but not known to yield the target SM; and (5) comparative analysis
592 of transcriptomes under SM producing vs. nonproducing conditions [39]. This led to the
593 differentiation between BGCs that are likely to produce known compounds and BGCs that
594 may encode novel chemistry. However, the number of GCFs to which no known functions
595 can be linked is so great that it is difficult to know which of the BGCs encode the most
596 interesting molecules [46].

597 **3. Original BGC exploration methods**

598 Out-of-the-box approaches are encouraged in targeting underexplored environmental niches
599 and bacterial phyla [54]. Rare environments have rendered an interesting source of SMs.
600 The chemical diversity of natural products correlates with the diversity of source
601 microorganisms. This is probably due to the evolution of organism-specific biosynthetic
602 machinery selected based on the adaptation of the microbe to the habitat, where beneficial
603 SMs play an important part [47]. Applying an isolate-based genome-mining approach on
604 bacteria obtained from unusual environments -such as deserts and arid regions like the
605 sampling sites of the present study- can be used for screening unique BGCs that may govern
606 the biosynthesis of novel natural products [54].

607 **4. Automated bioinformatics pipelines vs. manual bioinformatics**

608 Automated bioinformatics pipelines and manual bioinformatics are both methods used to
609 analyze biological data, but they have some key differences. Automated bioinformatics
610 pipelines are computer programs that are designed to perform specific bioinformatics tasks,
611 such as sequence alignment, gene annotation, and variant calling, in a pre-defined, step-by-
612 step manner. They are typically used to analyze large amounts of data and can be run on
613 high-performance computing clusters, allowing for efficient and rapid data processing [55].

614 Automated bioinformatics pipelines are also designed to be easy to use and can be run by
615 researchers with minimal bioinformatics experience [56]. Manual bioinformatics, on the other
616 hand, is the process of analyzing biological data using manual methods, such as web-based
617 tools, command-line programs, or custom scripts [57]. This method is typically used for
618 smaller data sets, or when an automated pipeline is not available or not suitable for the task
619 at hand. Manual bioinformatics requires a higher level of bioinformatics expertise and can be
620 more time-consuming, but it can also be more flexible and customizable [58]. Both
621 automated bioinformatics pipelines and manual bioinformatics have their advantages and
622 disadvantages. Automated pipelines can be faster and more efficient at processing large
623 amounts of data, but they may not be as flexible or customizable as manual methods.

624 Manual bioinformatics requires more expertise and can be more time-consuming, but it can
625 also be more tailored to specific research questions.

626 During our study, Big-scape has proven to be useful in the exploration of BGCs in the *B.*
627 *cereus* group thanks to its ability to identify and analyze large numbers of BGCs
628 comprehensively and efficiently.

629 As Big-scape allowed the comparison of BGCs across different strains of *Bacillus cereus*, it
630 helped us identify conserved and divergent BGCs among the *B. cereus* group, which will
631 ultimately provide insights into the evolution and adaptation of *Bacillus cereus* to different
632 environments.

633 **5. The lack of BGCs classification**

634 Despite the importance of BGCs, they are often difficult to classify and identify due to their
635 complex genetic organization and a large number of different types of BGCs that can be
636 found in different bacterial species [9]. One of the main challenges in BGCs classification is
637 the lack of consensus on the criteria and methods used to classify BGCs [59]. Different
638 researchers may use different criteria, such as gene content, gene organization, or
639 evolutionary relationships, to classify BGCs, which can lead to inconsistent and conflicting
640 results. Another challenge is that many BGCs are not well-annotated, making it difficult to
641 identify the genes and functions associated with each BGC [60]. This can be particularly
642 difficult in the case of novel BGCs, which may not have been previously described in the
643 literature. Additionally, the lack of a standardized nomenclature for BGCs can further
644 complicate the classification process [52].

645 **6. Highlighting the hidden BGC potential and synteny in the *B. cereus* group**

646 Many BGCs are not accounted for in the corresponding parvomes, referred to as “orphan”
647 BGCs. These are either ‘cryptic’ (cannot yet be linked to a product, activity, or phenotype) or
648 ‘silent’ (the compound is known in another organism, but experimental validation is required)
649 [46]. The development of facile strategies to induce the expression of these silent BGCs, and

650 to assign SM structures to orphan clusters will allow us to elucidate the microbial parvomes
651 [39], including those of the *Bacillus cereus* group.

652 **7. Lack of data about BGCs conservation in *Bacillus cereus* sensu lato**

653 There is currently limited data available about the conservation of biosynthetic gene clusters
654 (BGCs) in the *Bacillus cereus* group [61]. One reason for this is that the *Bacillus cereus*
655 group is a large and diverse group of bacteria, comprising several closely related species
656 and subspecies. This diversity can make it challenging to identify and study BGCs across the
657 different strains of *Bacillus cereus*. Additionally, the identification and characterization of
658 BGCs often require sophisticated techniques such as genome sequencing, transcriptomics,
659 and metabolomics, which can be time-consuming and expensive [62].

660 Therefore, not all strains of *Bacillus cereus* have been fully characterized in terms of their
661 BGCs. Furthermore, the majority of studies focus on certain strains of *Bacillus cereus*, such
662 as *Bacillus anthracis*, *Bacillus thuringiensis*, and *Bacillus cereus* sensu stricto, that are
663 known to produce important natural products [63–69], leaving out other strains.

664 The availability of genomes allowed the functional validation of identified BGCs based on the
665 structures of known SMs. However, comparative genomics can also be used to derive
666 hypotheses for the structures of the products of orphan BGCs. Thus, *de novo* sequenced
667 BGCs may be assigned to known SM structural families if the core genes and the constituent
668 tailoring genes are orthologous (or even syntenic) to functionally characterized BGCs, as
669 those revealed in the present work. The synteny and the tight phylogenetic distances
670 observed in our study support the conclusion that the BGCs in the *B. cereus* group arose
671 dependently with the acquisition of conserved core component genes.

672 In the case of the *B. cereus* group, we observed that many of the biosynthetic gene clusters
673 responsible for the production of natural products are syntenic, meaning that they are located
674 in the same chromosomal position and have a similar gene organization among different
675 strains of *B. cereus*. This suggests that these biosynthetic gene clusters have been inherited
676 through common ancestry, and have been conserved over time through selective pressures.

677 The synteny of these gene clusters can be used to help understand the evolutionary
678 relationships among different strains of *B. cereus*, and can also aid in the identification of
679 new natural products. The growing number of genomes in databases revealed lineage-
680 specific conservation of certain orthologous BGCs [39], as for the synteny block A and B of
681 the *B. cereus* group highlighted by the present study, or by pointing out the absence of
682 certain widely present BGCs in some species, thereby allowing the generation of
683 evolutionary hypotheses correlating the production of a given SM with the lifestyle and the
684 evolutionary history of the producer [39].

685 Additionally, the synteny of BGCs can also aid in the development of new antibiotics and
686 other bioactive compounds [10]. By identifying conserved regions of BGCs among different
687 strains of *B. cereus*, researchers can infer the presence of similar biosynthetic pathways and
688 enzymes among these strains, and can then use this information to guide the development of
689 new antibiotics and other bioactive compounds. Overall, highlighting the synteny of BGCs
690 can provide a valuable tool for understanding the biology, evolution, and biotechnology of the
691 *B. cereus* group [70].

692 For instance, a syntenic BGC with orthologous genes to those of the *B. bassiana* oosporein
693 BGC has recently been identified in the genome of *Cordyceps cycadae*, and the production
694 of oosporein was confirmed by HPLC [71], and comparative genomics of *C. militaris* and *A.*
695 *nidulans* revealed a syntenic BGC with four orthologous genes each in these fungi [72].

696 **8. The RiPPs in the *B. cereus* group**

697 RiPPs are an important class of natural products produced by the *Bacillus cereus* group. We
698 highlighted a consequent diversity among this BGC class, as it appeared to be the most
699 diversified compared to the remaining classes, with 13 RiPPs GCFs, and 7 singletons. RiPPs
700 have a wide range of biological activities, and they have potential applications in medicine,
701 agriculture, and industry [73]. RiPPs have a wide range of structural diversity and chemical
702 complexity, making them an attractive target for natural product discovery. Additionally, the
703 lack of resistance to RiPPs amongst pathogenic bacteria makes them an attractive target for

704 the development of new antibiotics. The study of RiPPs biosynthesis and the enzymes
705 responsible for the ribosomal synthesis and post-translational modification of these peptides
706 is an active area of research, which has the potential to lead to the discovery of new RiPPs
707 with unique properties and the development of new methods for RiPPs production [74].

708 We suggested that some BGCs could be part of mobile elements (BGC11, from the manual
709 approach). One common method of BGCs acquisition is through integrative and conjugative
710 elements (ICEs). The prevalence of these ICEs seems to be partially dictated by their
711 ecological background: bacteria originating from soil, plants, or aquatic environments contain
712 a greater number of ICEs than species from other environments [75].

713 **9. BGCs in endophytes**

714 Endophytes developed a variety of ways to successfully colonize plants; subdue their
715 immunity and their physiology as a nutrient source, and defend the plant host from
716 pathogens and opportunists. BGCs were also reported to mediate crucial functions in plant
717 colonization by beneficial endophytes. These functions are often facilitated by the vast array
718 of SMs produced by root-associated bacteria, which play a key role in inter- and intra-species
719 interactions [76, 77]. To date, a handful of studies have explored the diversity and
720 composition of bacterial SM-encoding BGC in soil [78].

721 Bioactive metabolites mediate important ecological functions, which are as diverse as their
722 chemical structures. Siderophores enhance iron uptake in environments where the
723 bioavailability of iron is limited [79], pigments protect against ultraviolet radiation and have
724 antioxidant activity [80] and compatible solutes protect against osmotic stress [81]. Besides,
725 it is possible to find in nature several examples of mutualistic relationships that have
726 coevolved whereby the microorganisms are actively cultured in exchange for producing
727 bioactive small molecules [47].

728 For instance, NRPS and PKS BGCs are responsible for the synthesis of a wide array of
729 siderophores, toxins, pigments, and antimicrobial compounds [82] that are believed to play a
730 pivotal role in bacterial adaptation to soil and rhizosphere ecosystems, and in plant health

731 and development [83]. However, little is known regarding the distribution of these gene
732 families in the root microbiome, and their functional role in the complex community
733 interactions in root ecosystems remains an enigma [84].

734 Studying the encoded secondary metabolome of endophytes will amplify our understanding
735 of the multiple roles that SMs play in the biotic and abiotic interactions in plants, leading to
736 the unveiling of natural products that can be used for various applications.

737 **Conclusions**

738 The classification of Biosynthetic Gene Clusters (BGCs) is an evolving field that has gained
739 significant attention in recent years. While some initial efforts have been made to classify
740 BGCs, the field is still relatively new and the classification methodologies are constantly
741 being refined and improved [22, 85]. Our work is an objective proposal for a consistent and
742 standardized approach to BGCs classification among the *Bacillus cereus* group, based on a
743 reproducible strategy that can be extended to other taxa, allowing comparison and
744 integration of data from different studies to expand the initial classification scheme that we
745 proposed. The current investigation is a substantive contribution to the discovery and
746 characterization of new natural products and biosynthetic pathways, based on BGCs analogy
747 and synteny.

748 **Supplementary Information**

749 The online version contains available supplementary material.

750 **Additional file 1.** Supplementary figures.

751 **Additional file 2.** antiSMASH profiles.

752 **Additional file 3.** BIG-SCAPE output files.

753 **Funding**

754 This work didn't receive any funding.

755 **Authors' contributions**

756 HAB conceived, designed the study; and analyzed the data. HAB and AY drafted the
757 manuscript; HAB, AY, AZ, and AM reviewed the manuscript. All authors read and approved
758 the final manuscript.

759 **Acknowledgments**

760 We gratefully acknowledge the inspirational impact of Dr. Livio Antonielli. Furthermore, our
761 warmest appreciation for 2019's EMBO genomics and bioinformatics course organizers, with
762 special thanks to Dr. Fredj Tekaiia.

763 **Availability of data and materials**

764 The datasets supporting the conclusions of this article are included within the article (and its
765 additional files). The genomic dataset can be accessed through the corresponding accession
766 numbers. The antiSMASH profiles and the BIG-SCAPE output are available as
767 supplementary material, from which the reference BGCs can be fetched (see table 3).

768 **Declarations**

769 **Ethics approval and consent to participate**

770 Not applicable.

771 **Consent for publication**

772 Not applicable.

773 **Competing interests**

774 The authors declare that they have no competing interest.

775 **References**

- 776 1. Rasko DA, Altherr MR, Han CS, Ravel J. Genomics of the *Bacillus cereus* group of
777 organisms. *FEMS Microbiology Reviews*. 2005.
- 778 2. Ehling-Schulz M, Lereclus D, Koehler TM. The *Bacillus cereus* Group: *Bacillus* Species
779 with Pathogenic Potential . *Microbiol Spectr*. 2019.
- 780 3. Hong HA, Le HD, Cutting SM. The use of bacterial spore formers as probiotics. *FEMS*

- 781 Microbiology Reviews. 2005.
- 782 4. Guinebretière MH, Auger S, Galleron N, Contzen M, de Sarrau B, de Buyser ML, et al.
783 *Bacillus cytotoxicus* sp. nov. is a novel thermotolerant species of the *Bacillus cereus* group
784 occasionally associated with food poisoning. *Int J Syst Evol Microbiol*. 2013.
- 785 5. Liu Y, Lai Q, Göker M, Meier-Kolthoff JP, Wang M, Sun Y, et al. Genomic insights into the
786 taxonomic status of the *Bacillus cereus* group. *Sci Rep*. 2015.
- 787 6. Margulis L, Jorgensen JZ, Dolan S, Kolchinsky R, Rainey FA, Lo SC. The Arthromitus
788 stage of *Bacillus cereus*: Intestinal symbionts of animals. *Proc Natl Acad Sci U S A*. 1998.
- 789 7. Priest FG, Barker M, Baillie LWJ, Holmes EC, Maiden MCJ. Population structure and
790 evolution of the *Bacillus cereus* group. *J Bacteriol*. 2004.
- 791 8. Schnepf E, Crickmore N, Van Rie J, Lereclus D, Baum J, Feitelson J, et al. *Bacillus*
792 *thuringiensis* and Its Pesticidal Crystal Proteins . *Microbiol Mol Biol Rev*. 1998.
- 793 9. Chen R, Wong H, Burns B. New Approaches to Detect Biosynthetic Gene Clusters in the
794 Environment. *Medicines*. 2019.
- 795 10. Ziemert N, Alanjary M, Weber T. The evolution of genome mining in microbes-a review.
796 *Natural Product Reports*. 2016.
- 797 11. Albarano L, Esposito R, Ruocco N, Costantini M. Genome mining as new challenge in
798 natural products discovery. *Marine Drugs*. 2020.
- 799 12. Cimermancic P, Medema MH, Claesen J, Kurita K, Wieland Brown LC, Mavrommatis K,
800 et al. Insights into secondary metabolism from a global analysis of prokaryotic biosynthetic
801 gene clusters. *Cell*. 2014.
- 802 13. Ju KS, Gao J, Doroghazi JR, Wang KKA, Thibodeaux CJ, Li S, et al. Discovery of
803 phosphonic acid natural products by mining the genomes of 10,000 actinomycetes. *Proc Natl*
804 *Acad Sci U S A*. 2015.
- 805 14. Cane DE, Walsh CT. The parallel and convergent universes of polyketide synthases and
806 nonribosomal peptide synthetases. *Chem Biol*. 1999.

- 807 15. Romero D, Traxler MF, López D, Kolter R. Antibiotics as signal molecules. *Chemical*
808 *Reviews*. 2011.
- 809 16. Tran PN, Yen MR, Chiang CY, Lin HC, Chen PY. Detecting and prioritizing biosynthetic
810 gene clusters for bioactive compounds in bacteria and fungi. *Applied Microbiology and*
811 *Biotechnology*. 2019;103:3277–87.
- 812 17. Flórez L V., Biedermann PHW, Engl T, Kaltenpoth M. Defensive symbioses of animals
813 with prokaryotic and eukaryotic microorganisms. *Nat Prod Rep*. 2015.
- 814 18. Carey J, Nguyen T, Korchak J, Beecher C, De Jong F, Lane AL. An isotopic ratio outlier
815 analysis approach for global metabolomics of biosynthetically talented actinomycetes.
816 *Metabolites*. 2019.
- 817 19. Blin K, Shaw S, Steinke K, Villebro R, Ziemert N, Lee SY, et al. AntiSMASH 5.0: Updates
818 to the secondary metabolite genome mining pipeline. *Nucleic Acids Res*. 2019;47:W81–7.
- 819 20. Navarro-Muñoz JC, Selem-Mojica N, Mullowney MW, Kautsar SA, Tryon JH, Parkinson
820 EI, et al. A computational framework to explore large-scale biosynthetic diversity. *Nat Chem*
821 *Biol*. 2020.
- 822 21. Kautsar SA, Blin K, Shaw S, Navarro-Muñoz JC, Terlouw BR, Van Der Hooft JJJ, et al.
823 MIBiG 2.0: A repository for biosynthetic gene clusters of known function. *Nucleic Acids Res*.
824 2020.
- 825 22. Terlouw BR, Blin K, Navarro-Muñoz JC, Avalon NE, Chevrette MG, Egbert S, et al.
826 MIBiG 3.0: a community-driven effort to annotate experimentally validated biosynthetic gene
827 clusters. *Nucleic Acids Res*. 2023;51:D603–10.
- 828 23. Liu Z, Ma H, Goryanin I. A semi-automated genome annotation comparison and
829 integration scheme. *BMC Bioinformatics*. 2013.
- 830 24. Pfeiffer F, Oesterhelt D. A manual curation strategy to improve genome annotation:
831 Application to a set of haloarchael genomes. *Life*. 2015.
- 832 25. Sánchez-Hidalgo M, Martín J, Genilloud O. Identification and heterologous expression of

- 833 the biosynthetic gene cluster encoding the lasso peptide humidimycin, a caspofungin activity
834 potentiator. *Antibiotics*. 2020;9.
- 835 26. Makałowski W, Boguski MS. Evolutionary parameters of the transcribed mammalian
836 genome: An analysis of 2,820 orthologous rodent and human sequences. *Proc Natl Acad Sci*
837 U S A. 1998.
- 838 27. Belaouni HA, Compant S, Antonielli L, Nikolic B, Zitouni A, Sessitsch A. In-depth genome
839 analysis of *Bacillus* sp. BH32, a salt stress-tolerant endophyte obtained from a halophyte in a
840 semiarid region. *Appl Microbiol Biotechnol*. 2022;106:3113–37. doi:10.1007/s00253-022-
841 11907-0.
- 842 28. Zhang Z, Schwartz S, Wagner L, Miller W. A greedy algorithm for aligning DNA
843 sequences. *Journal of Computational Biology*. 2000;7:203–14.
- 844 29. Alikhan NF, Petty NK, Ben Zakour NL, Beatson SA. BLAST Ring Image Generator
845 (BRIG): Simple prokaryote genome comparisons. *BMC Genomics*. 2011;12.
- 846 30. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7:
847 Improvements in performance and usability. *Mol Biol Evol*. 2013.
- 848 31. Jones DT, Taylor WR, Thornton JM. The rapid generation of mutation data matrices from
849 protein sequences. *Bioinformatics*. 1992;8:275–82.
- 850 32. Kumar S, Stecher G, Tamura K. MEGA7: Molecular Evolutionary Genetics Analysis
851 Version 7.0 for Bigger Datasets. *Mol Biol Evol*. 2016;33:1870–4.
- 852 33. Darling AE, Mau B, Perna NT. Progressivemauve: Multiple genome alignment with gene
853 gain, loss and rearrangement. *PLoS One*. 2010.
- 854 34. Deshpande NP, Kaakoush NO, Wilkins MR, Mitchell HM. Comparative genomics of
855 *Campylobacter concisus* isolates reveals genetic diversity and provides insights into disease
856 association. *BMC Genomics*. 2013.
- 857 35. Johnson MC, Tatum KB, Lynn JS, Brewer TE, Lu S, Washburn BK, et al. *Sinorhizobium*
858 *meliloti* Phage Φ M9 Defines a New Group of T4 Superfamily Phages with Unusual Genomic

- 859 Features but a Common T=16 Capsid. *J Virol*. 2015.
- 860 36. Martínez-Carranza E, Ponce-Soto GY, Servín-González L, Alcaraz LD, Soberón-Chávez
861 G. Evolution of bacteria seen through their essential genes: The case of *Pseudomonas*
862 *aeruginosa* and *Azotobacter vinelandii*. *Microbiol (United Kingdom)*. 2019.
- 863 37. de la Haba RR, López-Hermoso C, Sánchez-Porro C, Konstantinidis KT, Ventosa A.
864 Comparative Genomics and Phylogenomic Analysis of the Genus *Salinivibrio*. *Front*
865 *Microbiol*. 2019.
- 866 38. Daas MS, Acedo JZ, Rosana ARR, Orata FD, Reiz B, Zheng J, et al. *Bacillus*
867 *amyloliquefaciens* ssp. *plantarum* F11 isolated from Algerian salty lake as a source of
868 biosurfactants and bioactive lipopeptides. *FEMS Microbiol Lett*. 2018.
- 869 39. Zhang L, Yue Q, Wang C, Xu Y, Molnár I. Secondary metabolites from hypocrealean
870 entomopathogenic fungi: Genomics as a tool to elucidate the encoded parvome. *Nat Prod*
871 *Rep*. 2020;37:1164–80.
- 872 40. Scherlach K, Hertweck C. Triggering cryptic natural product biosynthesis in
873 microorganisms. *Org Biomol Chem*. 2009.
- 874 41. Wiemann P, Guo CJ, Palmer JM, Sekonyela R, Wang CCC, Keller NP. Prototype of an
875 intertwined secondarymetabolite supercluster. *Proc Natl Acad Sci U S A*. 2013.
- 876 42. Rokas A, Wisecaver JH, Lind AL. The birth, evolution and death of metabolic gene
877 clusters in fungi. *Nature Reviews Microbiology*. 2018.
- 878 43. Wang B, Kang Q, Lu Y, Bai L, Wang C. Unveiling the biosynthetic puzzle of destruxins in
879 *Metarhizium* species. *Proc Natl Acad Sci U S A*. 2012.
- 880 44. Baltz RH. Function of MbtH homologs in nonribosomal peptide biosynthesis and
881 applications in secondary metabolite discovery. *Journal of Industrial Microbiology and*
882 *Biotechnology*. 2011.
- 883 45. Lewis K. Platforms for antibiotic discovery. *Nature Reviews Drug Discovery*.
884 2013;12:371–87.

- 885 46. van Bergeijk DA, Terlouw BR, Medema MH, van Wezel GP. Ecology and genomics of
886 Actinobacteria: new concepts for natural product discovery. *Nature Reviews Microbiology*.
887 2020.
- 888 47. Granada Garc'ia SD. Metabolitos secundarios microbianos como alternativa de control
889 frente a fitopatógenos del aguacate (*Persea americana* Mill.). *Esc Biociencias*. 2019.
- 890 48. Khaldi N, Seifuddin FT, Turner G, Haft D, Nierman WC, Wolfe KH, et al. SMURF:
891 Genomic mapping of fungal secondary metabolite clusters. *Fungal Genet Biol*. 2010;47:736–
892 41.
- 893 49. Blin K, Wolf T, Chevrette MG, Lu X, Schwalen CJ, Kautsar SA, et al. AntiSMASH 4.0 -
894 improvements in chemistry prediction and gene cluster boundary identification. *Nucleic Acids*
895 *Res*. 2017;45:W36–41.
- 896 50. Gluck-Thaler E, Slot JC. Specialized plant biochemistry drives gene clustering in fungi.
897 *ISME J*. 2018.
- 898 51. Gluck-Thaler E, Haridas S, Binder M, Grigoriev I V., Crous PW, Spatafora JW, et al. The
899 architecture of metabolism maximizes biosynthetic diversity in the largest class of Fungi. *Mol*
900 *Biol Evol*. 2020;37:2838–56.
- 901 52. Ziemert N, Lechner A, Wietz M, Millañ-Aguiñaga N, Chavarria KL, Jensen PR. Diversity
902 and evolution of secondary metabolism in the marine actinomycete genus *Salinispora*. *Proc*
903 *Natl Acad Sci U S A*. 2014.
- 904 53. Doroghazi JR, Albright JC, Goering AW, Ju K-S, Haines RR, Tchalukov KA, et al. A
905 Roadmap for Natural Product Discovery Based on Large-Scale Genomics and Metabolomics
906 HHS Public Access Author manuscript. *Nat Chem Biol*. 2014.
- 907 54. Sekurova ON, Schneider O, Zotchev SB. Novel bioactive natural products from bacteria
908 via bioprospecting, genome mining and metabolic engineering. *Microbial Biotechnology*.
909 2019;12:828–44.
- 910 55. Roy S, LaFramboise WA, Nikiforov YE, Nikiforova MN, Routbort MJ, Pfeifer J, et al. Next-

- 911 generation sequencing informatics: Challenges and strategies for implementation in a clinical
912 environment. *Archives of Pathology and Laboratory Medicine*. 2016.
- 913 56. Schmid K, Dohmen H, Ritschel N, Selignow C, Zohner J, Sehring J, et al. SangeR: the
914 high-throughput Sanger sequencing analysis pipeline. *Bioinforma Adv*. 2022;2.
- 915 57. Neves M, Ševa J. An extensive review of tools for manual annotation of documents.
916 *Briefings in Bioinformatics*. 2021;22:146–63.
- 917 58. Cherkasov A. *Bioinformatics: A practical guide to the analysis of genes and proteins*. Am
918 *J Hum Biol*. 2005.
- 919 59. Hannigan GD, Prihoda D, Palicka A, Soukup J, Klempir O, Rampula L, et al. A deep
920 learning genome-mining strategy for biosynthetic gene cluster prediction. *Nucleic Acids Res*.
921 2019.
- 922 60. Nguyen DD, Wu CH, Moree WJ, Lamsa A, Medema MH, Zhao X, et al. MS/MS
923 networking guided analysis of molecule and gene cluster families. *Proc Natl Acad Sci U S A*.
924 2013.
- 925 61. Xia L, Miao Y, Cao A, Liu Y, Liu Z, Sun X, et al. Biosynthetic gene cluster profiling
926 predicts the positive association between antagonism and phylogeny in *Bacillus*. *Nat*
927 *Commun*. 2022.
- 928 62. Medema MH, Trefzer A, Kovalchuk A, Van Den Berg M, Müller U, Heijne W, et al. The
929 sequence of a 1.8-Mb bacterial linear plasmid reveals a rich evolutionary reservoir of
930 secondary metabolic pathways. *Genome Biol Evol*. 2010.
- 931 63. Stohl EA, Milner JL, Handelsman J. Zwittermicin A biosynthetic cluster. *Gene*. 1999.
- 932 64. Hoton FM, Fornelos N, N'Guessan E, Hu X, Swiecicka I, Dierick K, et al. Family portrait
933 of *Bacillus cereus* and *Bacillus weihenstephanensis* cereulide-producing strains. *Environ*
934 *Microbiol Rep*. 2009.
- 935 65. Arias AA, Ongena M, Devreese B, Terrak M, Joris B, Fickers P. Characterization of
936 amylolysin, a novel lantibiotic from *Bacillus amyloliquefaciens* GA1. *PLoS One*. 2013.

- 937 66. Mei X, Xu K, Yang L, Yuan Z, Mahillon J, Hu X. The genetic diversity of cereulide
938 biosynthesis gene cluster indicates a composite transposon Tnces in emetic *Bacillus*
939 *weihenstephanensis*. *BMC Microbiol.* 2014.
- 940 67. Zhang L, Teng K, Wang J, Zhang Z, Zhang J, Sun S, et al. CerR, a single-domain
941 regulatory protein of the LuxR family, promotes cerecidin production and immunity in *Bacillus*
942 *cereus*. *Appl Environ Microbiol.* 2018.
- 943 68. Farlow J, Bolkvadze D, Leshkasheli L, Kusradze I, Kotorashvili A, Kotaria N, et al.
944 Genomic characterization of three novel Basilisk-like phages infecting *Bacillus anthracis* 06
945 Biological Sciences 0604 Genetics. *BMC Genomics.* 2018.
- 946 69. Prasertanan T, Palmer DRJ. The kanosamine biosynthetic pathway in *Bacillus cereus*
947 UW85: Functional and kinetic characterization of KabA, KabB, and KabC. *Arch Biochem*
948 *Biophys.* 2019.
- 949 70. Steinke K, Mohite OS, Weber T, Kovács ÁT. Phylogenetic Distribution of Secondary
950 Metabolites in the *Bacillus subtilis* Species Complex. *mSystems.* 2021.
- 951 71. Lu Y, Luo F, Cen K, Xiao G, Yin Y, Li C, et al. Omics data reveal the unusual asexual-
952 fruiting nature and secondary metabolic potentials of the medicinal fungus *Cordyceps*
953 *cicadae*. *BMC Genomics.* 2017.
- 954 72. Asai T, Chung YM, Sakurai H, Ozeki T, Chang FR, Wu YC, et al. Highly oxidized
955 ergosterols and isariotin analogs from an entomopathogenic fungus, *Gibellula formosana*,
956 cultivated in the presence of epigenetic modifying agents. *Tetrahedron.* 2012.
- 957 73. Inoue M. Total synthesis and functional analysis of non-ribosomal peptides. *Chem Rec.*
958 2011.
- 959 74. Ortega MA, Van Der Donk WA. New Insights into the Biosynthetic Logic of Ribosomally
960 Synthesized and Post-translationally Modified Peptide Natural Products. *Cell Chemical*
961 *Biology.* 2016.
- 962 75. Ghinet MG, Bordeleau E, Beaudin J, Brzezinski R, Roy S, Burrus V. Uncovering the

- 963 prevalence and diversity of integrating conjugative elements in actinobacteria. PLoS One.
964 2011.
- 965 76. Tyc O, Song C, Dickschat JS, Vos M, Garbeva P. The Ecological Role of Volatile and
966 Soluble Secondary Metabolites Produced by Soil Bacteria. Trends in Microbiology. 2017.
- 967 77. Philippot L, Raaijmakers JM, Lemanceau P, Van Der Putten WH. Going back to the roots:
968 The microbial ecology of the rhizosphere. Nature Reviews Microbiology. 2013.
- 969 78. Dror B, Wang Z, Brady SF, Jurkevitch E, Cytryn E. Elucidating the Diversity and Potential
970 Function of Nonribosomal Peptide and Polyketide Biosynthetic Gene Clusters in the Root
971 Microbiome. mSystems. 2020.
- 972 79. Guerinot ML. Microbial iron transport. Annual Review of Microbiology. 1994.
- 973 80. Li C, Ji C, Tang B. Purification, characterisation and biological activity of melanin from
974 Streptomyces sp. FEMS Microbiol Lett. 2018.
- 975 81. Sadeghi A, Soltani BM, Nekouei MK, Jouzani GS, Mirzaei HH, Sadeghizadeh M.
976 Diversity of the ectoines biosynthesis genes in the salt tolerant Streptomyces and evidence
977 for inductive effect of ectoines on their accumulation. Microbiol Res. 2014.
- 978 82. Newman DJ, Cragg GM. Natural Products as Sources of New Drugs from 1981 to 2014.
979 Journal of Natural Products. 2016.
- 980 83. Chowdhury SP, Hartmann A, Gao XW, Borriss R. Biocontrol mechanism by root-
981 associated Bacillus amyloliquefaciens FZB42 - A review. Frontiers in Microbiology. 2015.
- 982 84. Charlop-Powers Z, Pregitzer CC, Lemetre C, Ternei MA, Maniko J, Hover BM, et al.
983 Urban park soil microbiomes are a rich reservoir of natural product biosynthetic diversity.
984 Proc Natl Acad Sci U S A. 2016.
- 985 85. Medema MH, Kottmann R, Yilmaz P, Cummings M, Biggins JB, Blin K, et al. Minimum
986 Information about a Biosynthetic Gene cluster. Nature Chemical Biology. 2015.
- 987

988 **Tables**

989 **Table 1.** Dataset for the comparative analysis

Preferred name	Deposit	Base pairs	% G+C	Proteins	Bioproject accession	Biosample accession	Assembly accession
<i>Bacillus wiedmannii</i>	PL1	5,309,441	35.26	5272	PRJDB9286	SAMD00204526	GCA_011405335.1
<i>Bacillus</i> sp.	BH32	5,661,597	34.86	5682	PRJNA523918	SAMN10992522	GCF_004367825.2
<i>Bacillus cereus</i>	JHU	5,323,903	35.26	5489	PRJNA591929	SAMN13391868	GCA_009738575.1
<i>Bacillus cereus</i>	G9842	5,387,334	35.26	5379	PRJNA224116	SAMN02604060	GCF_000021305.1
<i>Bacillus thuringiensis</i>	HD-789	5,495,278	35.26	5550	PRJNA171844	SAMN02603564	GCA_000292705.1
<i>Bacillus bombysepticus</i>	Wang	5,295,783	35.25	5265	PRJNA242213	SAMN02691825	GCF_000831065.1
<i>Bacillus thuringiensis</i>	HD1002	5,491,311	35.25	5548	PRJNA236049	SAMN03010437	GCF_000835025.1
<i>Bacillus thuringiensis</i> serovar <i>galleriae</i>	4G5	5,701,188	35.29	5770	PRJNA224116	SAMN03074947	GCF_000803665.1
<i>Bacillus thuringiensis</i> serovar <i>morrisoni</i>	BGSC 4AA1	5,652,292	35.33	5693	PRJNA224116	SAMN03274640	GCF_000940785.1
<i>Bacillus cereus</i>	A1	5,352,307	35.28	5273	PRJNA242371	SAMN02693464	GCA_000635895.2
<i>Bacillus thuringiensis</i>	YGd22-03	5,420,545	35.23	5455	PRJNA224116	SAMN06197294	GCF_002184245.1
<i>Bacillus thuringiensis</i>	L7601	5,790,408	35.18	5846	PRJNA224116	SAMN06473083	GCF_002025105.1
<i>Bacillus thuringiensis</i>	c25	5,334,660	35.32	5300	PRJNA38828 7	SAMN0731742 5	GCA_002222555 .1

<i>Bacillus cereus</i>	FORC087	5,271,204	35.27	5272	PRJNA47081	SAMN0911197	GCA_006384875
					8	7	.1
<i>Bacillus cereus</i>	ZB201708	5,466,652	35.22	5423	PRJNA48089	SAMN0965200	GCA_004006495
					9	6	.1
<i>Bacillus thuringiensis</i>	BT-59	5,500,615	35.26	5549	PRJNA53418	SAMN1147929	GCA_005155285
					9	9	.1
<i>Bacillus thuringiensis</i>	JW-1	5,500,376	35.26	5552	PRJNA57420	SAMN1286033	GCA_009025915
					1	8	.1

990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007

1008 **Table 2. BGC classes/families overview**

	RiPPS	NRPS	Terpene	others	total
# of families:	13	6	1	3	23
Average # of BGCs per family:	4	15	17	11	8
Max # of BGCs in a family:	17	17	17	17	-
Families with MIBiG Reference BGCs:	0	0	0	0	0

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037 **Table 3. Reference BGCs per class/family**

BGC families	strain	region	genomic location (nt)	most similar known cluster	similarity
RIPPs					
BC-Bacteriocin 1	<i>B. cereus</i> ZB201708	NZ_CP030982.1.region009	2,468,894–2,477,159	-	-
BC-Bacteriocin 2	<i>B. thuringiensis</i> HD1002	NZ_CP009351.1.region005	1,379,714–1,389,169	-	-
BC-Bacteriocin 3	<i>B. thuringiensis</i> 4G5	NZ_CP010089.1.region007	2,741,425–2,750,241	-	-
BC-Lasso peptide 1	<i>B. cereus</i> A1	NZ_CP015727.1.region001	568,125–591,631	-	-
BC-Bacteriocin 4	<i>B. thuringiensis</i> BGSC 4AA1	NZ_CP010577.1.region009	2,725,350–2,735,610	-	-
BC-Bacteriocin 5	<i>B. cereus</i> ZB201708	NZ_CP030982.1.region010	2,580,757–2,589,106	-	-
BC-Cerecidin	<i>B. thuringiensis</i> YGd22-03	NZ_CP019230.1.region011	4,472,628–4,495,786	cerecidin / cerecidin A1 / cerecidin A2 / cerecidin A3 / cerecidin A4 / cerecidin A5 / cerecidin A6 / cerecidin A7	94%
NRPS families					
BC- NRPS 1	<i>Bacillus</i> . sp. BH32	SJAS00000000.2.2.region002	286,950–346,871	polyoxypeptin	5%
BC- NRPS 2	<i>Bacillus</i> sp. BH32	SJAS00000000.2.15.region001	61,643–91,585	-	-
BC- Bacillibactin-like	<i>B. cereus</i> A1	NZ_CP015727.1.region003	1,501,835–1,550,463	bacillibactin	46%
BC- NRPS-like 1	<i>B. cereus</i> G9842	NC_011772.1.region001	386,831–430,076	-	-
BC-NRPS 3	<i>B. cereus</i> G9842	NC_011772.1.region006	2,287,034–2,352,013	-	-
BC-NRPS 4	<i>B. cereus</i> ZB201708	NZ_CP030982.1.region003	791,630–836,873	chejuenolide A / chejuenolide B	7%
Terpene family					
BC- terpene 1	<i>B. thuringiensis</i> L-7601	NZ_CP020002.1.region012	3,779,510–3,801,363	molybdenum cofactor	17%
Siderophore families					
BC- Petrobactin 1	<i>B. cereus</i> G9842	NC_011772.1.region004	1,852,404–1,866,111	petrobactin	100%
BC- Petrobactin 2	<i>B. cereus</i> A1	NZ_CP015727.1.region002	1,182,329–1,196,036	petrobactin	100%
Betalactone family					
BC- Fengycin-like 1	<i>B. cereus</i> G9842	NC_011772.1.region007	2,387,411–2,412,649	fengycin	40%

1038

1039

1040

1041

1042

1043

1044

1045

1046 **Table 4. Reference BGCs for RiPPs singletons**

RiPPs singletons	strain	region	genomic location (nt)	most similar known cluster	similarity
BCS-Bacteriocin 1	<i>B. cereus</i> JHU	CP046511.1.region001	24,944 – 34,346	-	-
BCS-Surfactin-like 1	<i>B. cereus</i> ZB201708	NZ_CP030982.1.region002	748,698 – 771,271	-	-
BCS-Bacteriocin 2	<i>B. cereus</i> JHU	CP046511.1.region002	128,885 – 139,342	-	-
BCS-Bacteriocin 3	<i>B. thuringiensis</i> c25	NZ_CP022345.1.region010	4,376,042 – 4,388,249	-	-
BCS-Bacteriocin 4	<i>B. cereus</i> ZB201708	NZ_CP030982.1.region014	5,325,429 – 5,338,733	-	-
BCS-Bacteriocin 5	<i>B. thuringiensis</i> L-7601	NZ_CP020002.1.region013	5,353,484 – 5,365,667	-	-

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075 **Figure captions**

1076 **Fig. 1** Distribution of BGCs classes and counts. The heatmap shows the number of BGCs by
1077 type (from 0 to 5), and their distribution among the strains (dendrograms were generated by
1078 hierarchical clustering using *one minus Spearman* rank correlation)

1079 **Fig. 2** BLASTp comparisons of *Bacillus* sp. BH32 BGCs against BGCs of the closest non-
1080 type strains. The reference ring constitutes the detected BGCs in *Bacillus* sp. BH32
1081 (translated to amino acid sequence), while individual rings represent BGCs of the closest
1082 non-type strains (the color represents sequence identity on a sliding scale, the greyer it gets;
1083 the lower the percentage identity). The 3 outermost rings depict (from inside to outside):
1084 genes composing BGCs in *Bacillus* sp. BH32 (in blue); *unique genes* (sequence identity
1085 below 30% for 100% of the strains, in green) and *rare genes* (sequence identity below 30%
1086 for 80% of the strains, in light velvet) with their locus tag IDs; and the *regions/BGCs types* (as
1087 detected by antiSMASH). Unique and rare genes have the following annotation (from
1088 antiSMASH / prokka, with the same location, or overlapping locations): ctg8_81 =
1089 unknown/hypothetical protein; ctg12_31 = unknown / Spore germination protein
1090 B1; ctg15_85 = biosynthetic-additional (smcogs) SMCOG1028:crotonyl-CoA reductase -
1091 alcohol dehydrogenase / Zinc-type alcohol dehydrogenase-like protein; ctg27_13 =
1092 unknown/hypothetical protein; ctg27_20 = unknown/hypothetical protein; ctg2_300 =
1093 unknown/hypothetical protein; ctg2_321 = unknown / IS200/IS605 family transposase
1094 ISAsp8 ; ctg8_90 = unknown / hypothetical protein; ctg12_35 = biosynthetic-additional
1095 (smcogs) SMCOG1012:4'-phosphopantetheinyl transferase / 4'-phosphopantetheinyl
1096 transferase Sfp

1097 **Fig. 3** BGCs phylogeny and conservation. **a)** Maximum Likelihood (ML) phylogenetic tree
1098 generated from concatenated BGCs amino acid sequences. Numbers on branches represent
1099 bootstrap values (average value: 92.21%). The tree is drawn to scale, with branch lengths
1100 measured in the number of substitutions per site. **b)** BGCs conservation among the strains
1101 (MAUVE alignment of concatenated genebank sequences of the BGCs of each strain). Each

1102 row represents the conservation and orientation of BGCs of the corresponding strain in the
1103 ML tree (left) after alignment, in comparison to *Bacillus thuringiensis* c25 (bottom row, set as
1104 reference). Red bars represent BGC sequence limits. Colored blocks refer to homologous
1105 segments with similarity plots (from MAUVE alignment diagram), with a *locally collinear block*
1106 (LCB) weight ≥ 773 . Each arrow symbolizes the BGC relative orientation with its tag number.
1107 BGCs with conserved order are framed in purple. BGC tag number, type, and most similar
1108 known clusters are shown on legends. The distance scale is shown under the alignment
1109 diagram (in amino acids)

1110 **Fig. 4** BGCs fractions (counts and percentages by class, according to BIG-SCAPE output).

1111 **Fig. 5** RiPPs BGC families/clans. **a)** RiPPs BGCs families statistics, **b)** proposed RiPPs
1112 families distribution matrix (presence: 1, absence: 0) among the *B. cereus* group genomes, **c)**
1113 RiPPs BGCs families/clans cluster networks (generated by the highest cutoff selected) and
1114 singletons (separate dots).

1115 **Fig. 6** NRPS BGCs families. **a)** NRPS BGCs families statistics, **b)** proposed NRPS families
1116 distribution matrix (presence: 1, absence: 0) among the *B. cereus* group genomes, **c)** NRPS
1117 BGCs families cluster networks (generated by the highest cutoff selected).

1118 **Fig. 7** BC-Terpene 1 family features. **a)** BC-Terpene 1 BGCs family statistics, **b)** BC-Terpene
1119 1 BGCs cluster network generated by the highest cutoff selected, **c)** CORASON-like tree
1120 generated for the BC-Terpene 1 GCF. This tree was created using the sequences of the
1121 Core Domains in the BC-Terpene 1 gene cluster family. These are defined as the domain
1122 type(s) that (1) appeared with the highest frequency in the BC-Terpene 1 gene cluster family
1123 and (2) were detected in the exemplar cluster (defined by the affinity propagation cluster),
1124 which is, in this case, that of *B. thuringiensis* L-7601 (in red). All copies of the Core Domains
1125 in the exemplar were automatically concatenated, as well as those from the best-matching
1126 domains of the rest of the BGCs in the BC-Terpene 1 gene cluster family (aligned domain
1127 sequences were used). The tree was inferred using FastTree43 with default parameters.
1128 Visual alignment was attempted based on the position of the 'longest common information'

1129 from the distance calculation step (between the exemplar BGCs vs. each of the remaining
1130 clusters). The Pfam domains of the 2 main clades are reported at the bottom part of this
1131 figure. The Pfam domains' color significance can be retraced from a list, available here:
1132 https://git.wageningenur.nl/medema-group/BiG-SCAPE/blob/master/domains_color_file.tsv

1133 **Fig. 8** BC-Petrobactin 1, BC-Petrobactin 2, and BC-Fengycin BGCs families/clan. **a)** the
1134 remaining BGCs families statistics, **b)** proposed BC-Petrobactin 1, BC-Petrobactin 2, and
1135 BC-Fengycin families distribution matrix (presence: 1, absence: 0) among the *B. cereus*
1136 group genomes, **c)** BC-Petrobactin 1, BC-Petrobactin 2 and BC-Fengycin BGCs families/clan
1137 cluster networks (generated by the highest cutoff selected).

1138 **Fig. S1** RiPPs reference BGCs for each proposed family. For each reference BGC, from top
1139 to bottom: reference BGC info (from left to right are mentioned: the name of the proposed
1140 RiPPs family; the strain bearing the reference BGC; the genomic region; and the most similar
1141 known cluster with the similarity %); a depiction of the BGCs regions/distribution among the
1142 genome; and the BGC organization (from antiSMASH output).

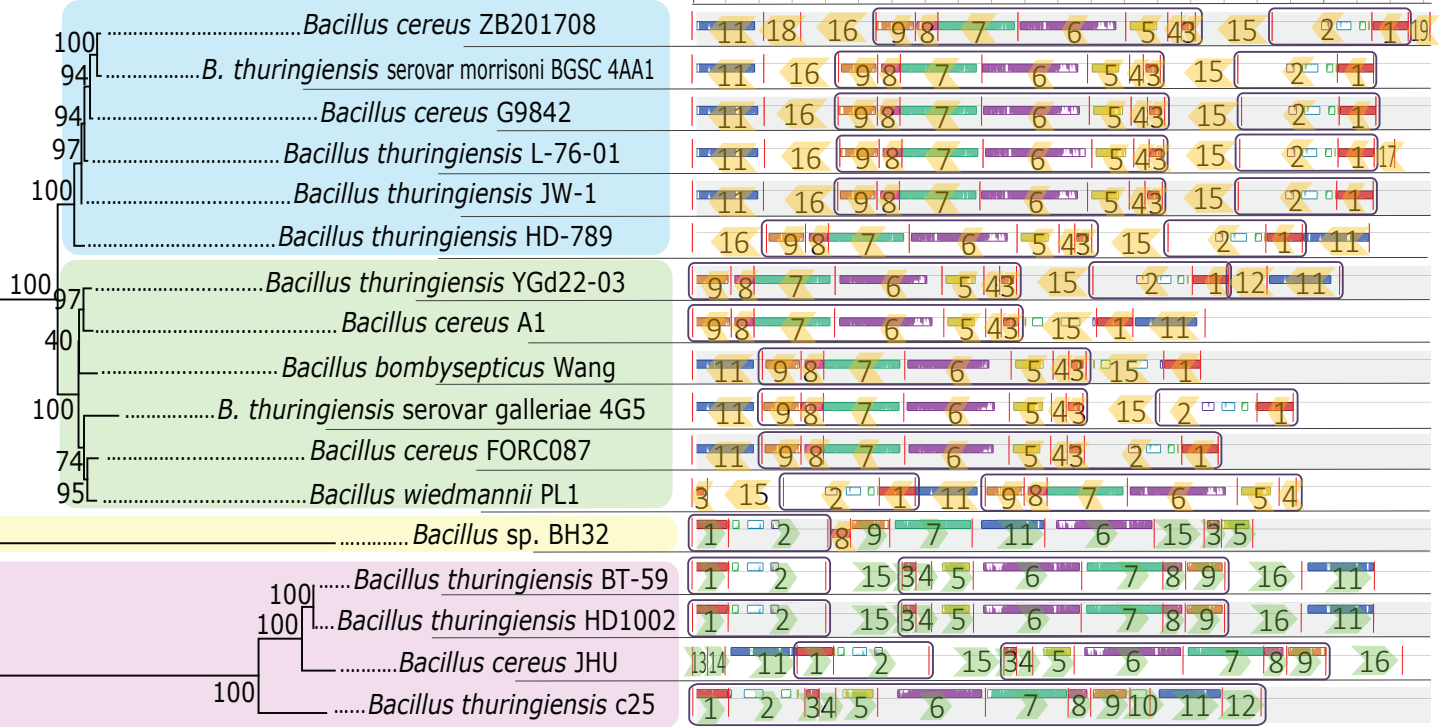
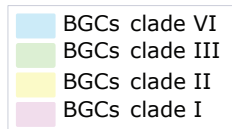
1143 **Fig. S2** NRPS reference BGCs for each proposed family. For each reference BGC, from top
1144 to bottom: reference BGC info (from left to right are mentioned: the name of the proposed
1145 NRPS family; the strain bearing the reference BGC; the genomic region; and the most similar
1146 known cluster with the similarity %); a depiction of the BGCs regions/distribution among the
1147 genome; and the BGC organization (from antiSMASH output).

1148 **Fig. S3** BC-Terpene 1 reference BGC. From top to bottom: reference BGC info (from left to
1149 right are mentioned: the name of the proposed family; the strain bearing the reference BGC;
1150 the genomic region; and the most similar known cluster with the similarity %); a depiction of
1151 the BGCs regions/distribution among the genome; and the BGC organization (from
1152 antiSMASH output).

1153 **Fig. S4** Siderophore/Betalactone reference BGCs for each proposed family. For each
1154 reference BGC, from top to bottom: reference BGC info (from left to right are mentioned: the
1155 name of the proposed family; the strain bearing the reference BGC; the genomic region; and

1156 the most similar known cluster with the similarity %); a depiction of the BGCs
1157 regions/distribution among the genome; and the BGC organization (from antiSMASH output).

a)



0.050

400000 AA

b)

0 20 40 60 80 100 120 140 160 180 200 220 240 260 280 300 320 340 360 380 400 420 440 k AA



BGC sequence limits



BGC orientation same as reference



BGC orientation adverse to reference



BGCs with conserved order



BGC homologous segment with similarity plots (segments with same colour are homologous, with Locally Collinear Blocks "LCB" weight threshold of 733)

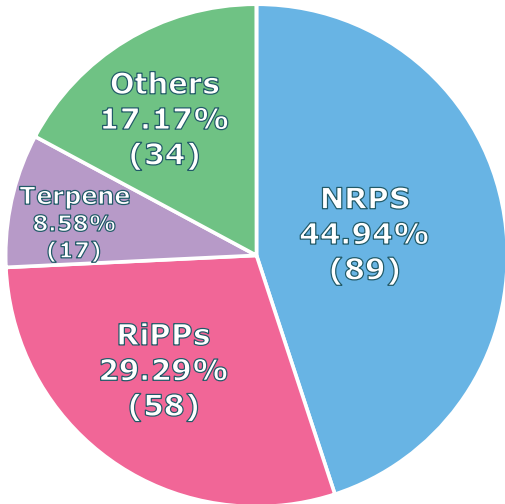
1, 2, 3, ...19

BGC tag number, according to appearance in genomes, from the reference strain *B. thuringiensis* c25 to *B. cereus* ZB201708 BGC with same tag numbers are considered as homologous.

BGC tag number, type and most similar known clusters:

- 1:terpen (molybdenum co-factor)
- 2: NRPS 1(polyoxypeptin*)
- 3:bacteriocin 1 (unknown)
- 4:bacteriocin 2 (unknown)
- 5: betalacton (fengycin)
- 6: NRPS 2 (unknown, gramicidin**, nostopeptolide A2***)
- 7: NRPS 3 (bacillibactin)
- 8:siderophore (petrobactin)
- 9: linear azol(in)e-containing peptides "LAP" bacteriocin
- 10:bacteriocin 3 (unknown)
- 11: NRPS-like 1 (unknown)
- 12:lanthipeptid 1(cerecidin/cerecidin A1/cerecidin A2/cerecidin A3 / cerecidin A4/cerecidin A5/cerecidin A6/cerecidin A7)
- 13:bacteriocin 4 (unknown)
- 14:bacteriocin 5 (unknown)
- 15: NRPS-like 2 (unknown)
- 16: NRPS Polyketide (chejuenolide A / chejuenolide B)
- 17:bacteriocin 6 (unknown)
- 18:lanthipeptid 2 (surfactin)
- 19:bacteriocin 7 (unknown)

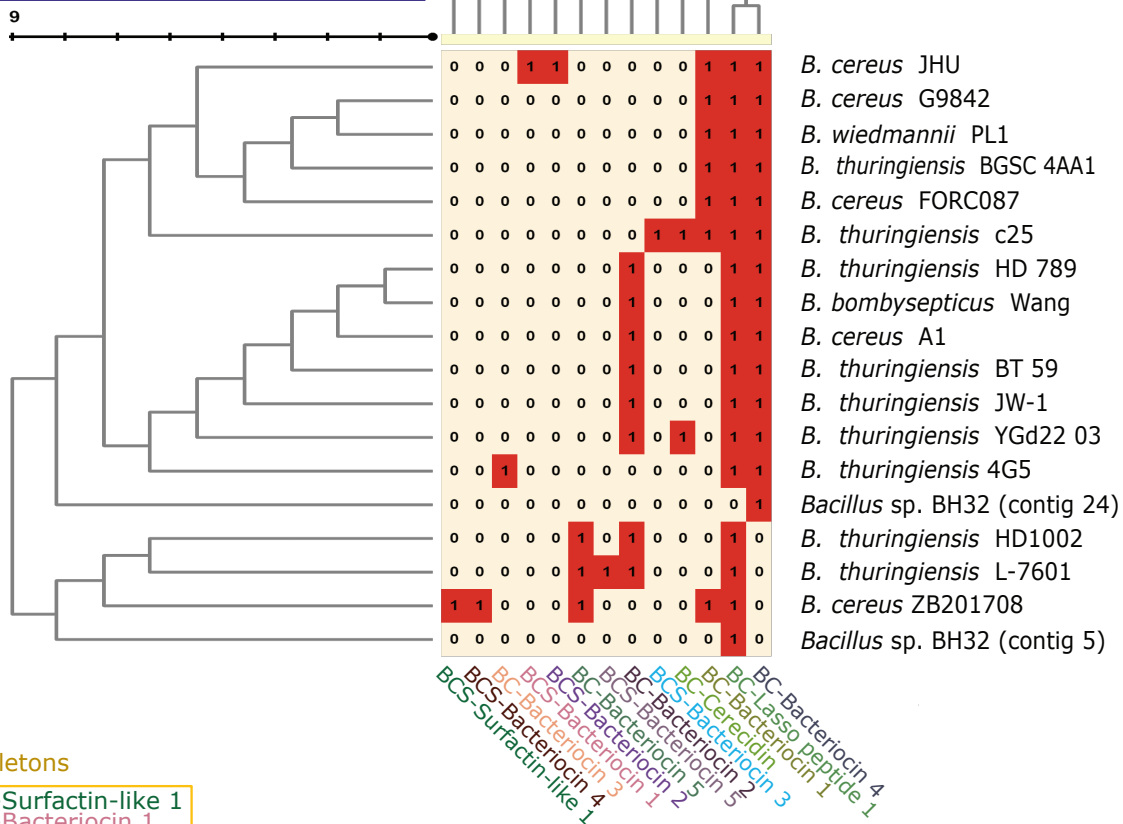
* Closest known cluster in all the strains, except *B. thuringiensis* c25** Closest known cluster in *B. cereus* JHU*** Closest known cluster in *B. thuringiensis* serovar galleriae4G5



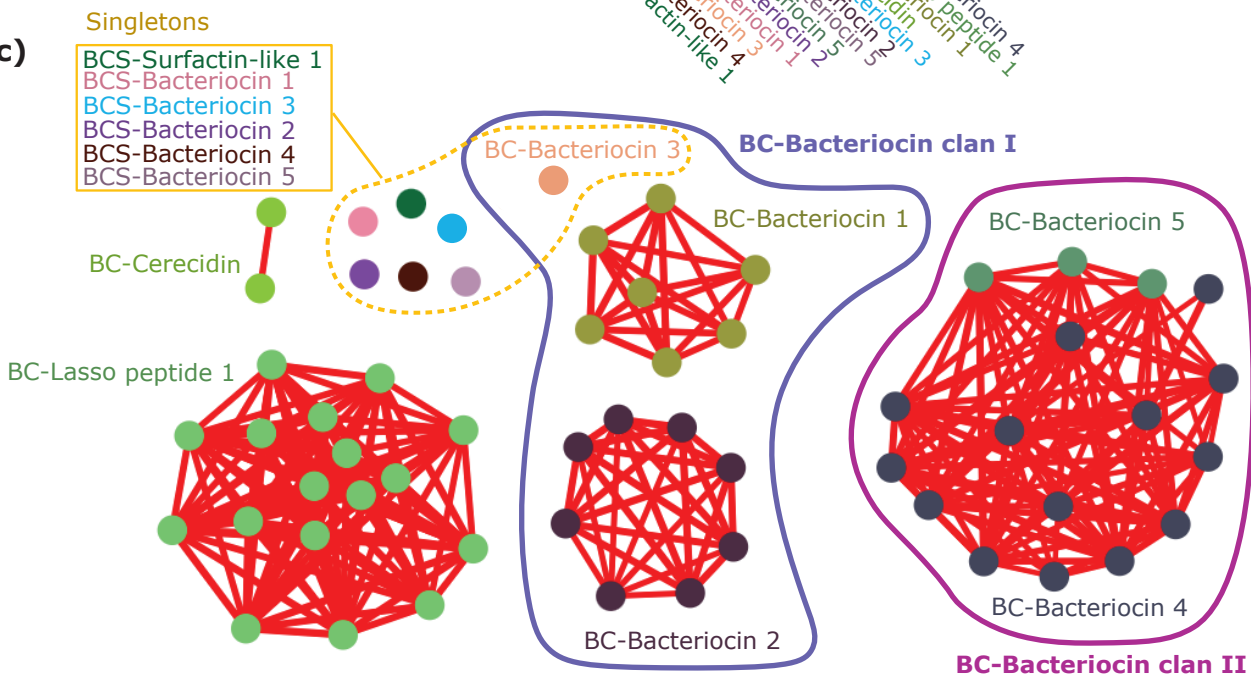
a)

# of families:	13
Average # of BGCs per family:	4
Max # of BGCs in a family:	17
Families with MIBiG Reference BGCs:	0
Total BGCs: 58 (7 singleton/s), links:	293

b)

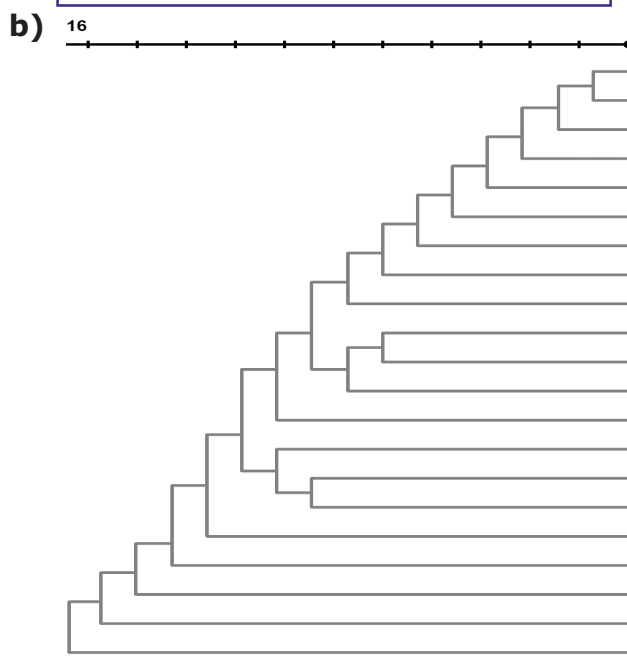
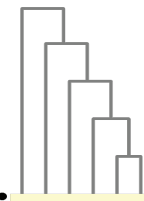


c)



a)

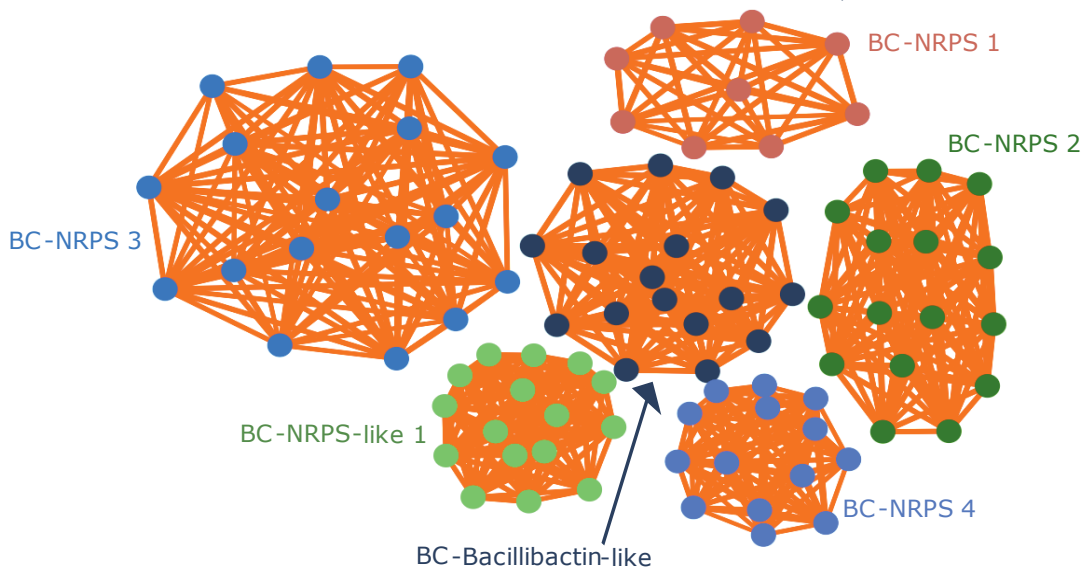
# of families:	6
Average # of BGCs per family:	15
Max # of BGCs in a family:	17
Families with MIBiG Reference BGCs:	0
Total BGCs: 89 (0 singleton/s), links:	634



1	1	1	1	1	1	<i>B. cereus</i> G9842
1	1	1	1	1	1	<i>B. cereus</i> JHU
1	1	1	1	1	1	<i>B. thuringiensis</i> HD 789
1	1	1	1	1	1	<i>B. thuringiensis</i> HD1002
1	1	1	1	1	1	<i>B. thuringiensis</i> BGSC 4AA1
1	1	1	1	1	1	<i>B. thuringiensis</i> L-7601
1	1	1	1	1	1	<i>B. cereus</i> ZB201708
1	1	1	1	1	1	<i>B. thuringiensis</i> BT 59
1	1	1	1	1	1	<i>B. thuringiensis</i> JW-1
0	1	1	1	1	1	<i>B. thuringiensis</i> 4G5
0	1	1	1	1	1	<i>B. wiedmannii</i> PL1
0	1	1	1	1	1	<i>B. thuringiensis</i> YGd22 03
0	1	1	0	1	1	<i>B. cereus</i> FORC087
0	0	0	1	1	1	<i>B. bombysepticus</i> Wang
0	0	1	1	1	1	<i>B. cereus</i> A1
0	0	1	1	1	1	<i>B. thuringiensis</i> c25
0	0	0	0	0	1	<i>Bacillus</i> sp. BH32 (contig 12)
0	0	0	1	0	0	<i>Bacillus</i> sp. BH32 (contig 15)
0	1	0	0	0	0	<i>Bacillus</i> sp. BH32 (contig 2)
0	0	0	0	1	0	<i>Bacillus</i> sp. BH32 (contig 6)
0	0	1	0	0	0	<i>Bacillus</i> sp. BH32 (contig 8)

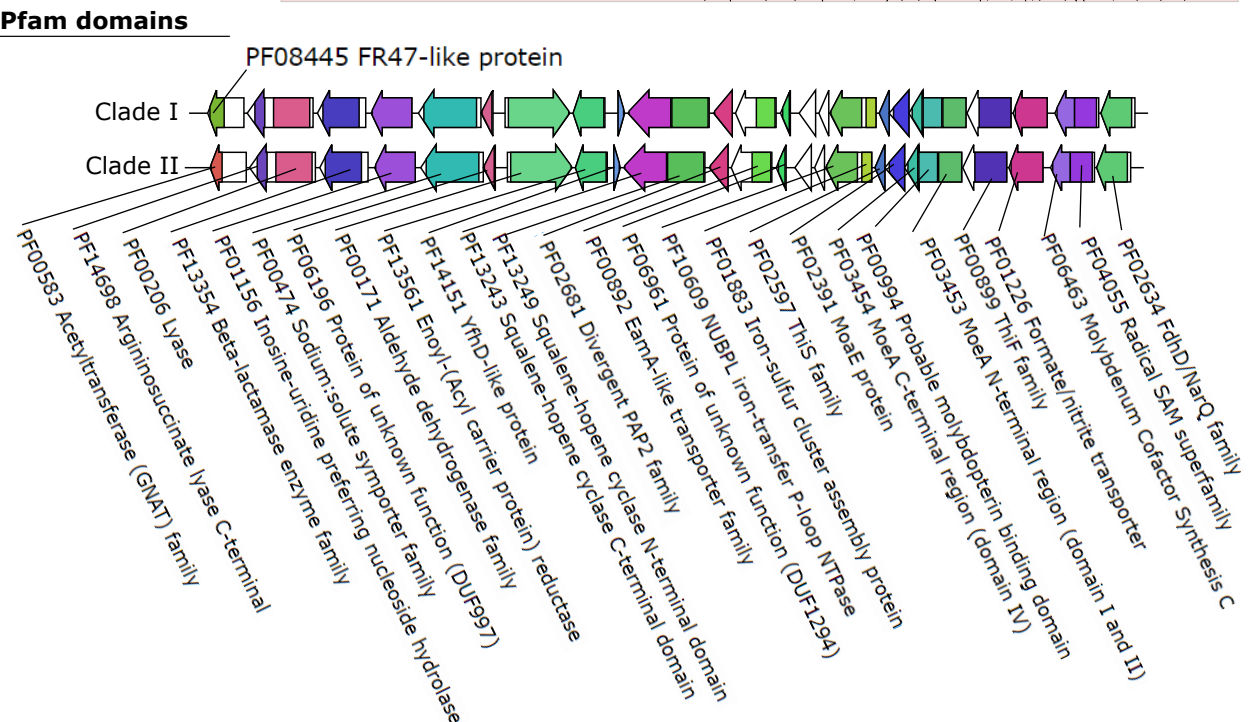
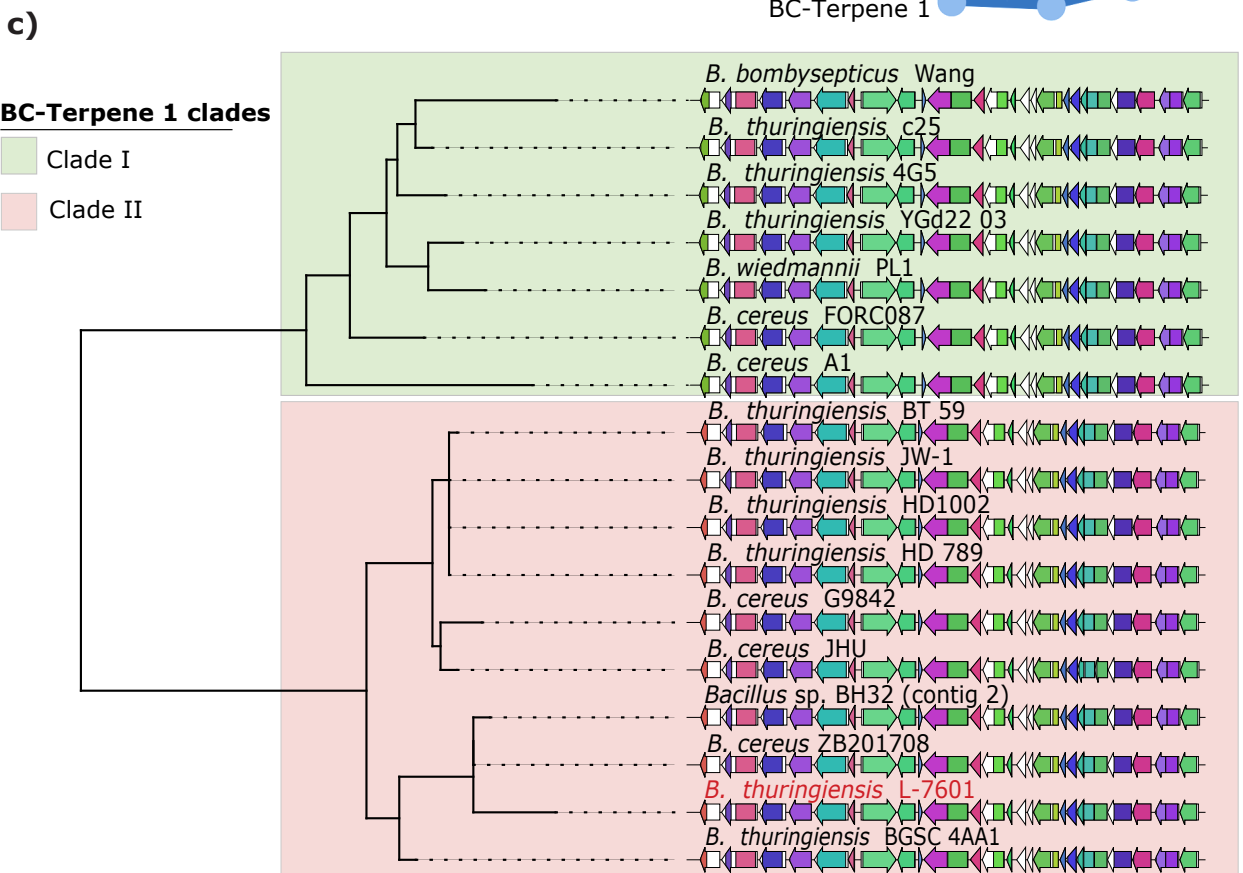
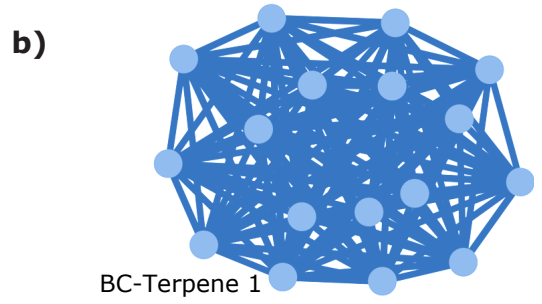
BC-NRPS 3
 BC-NRPS 1
 BC-NRPS 2
 BC-NRPS 4
 BC-NRPS-like 1
 BC-Bacillibactin-like

c)



a)

# of families:	1
Average # of BGCs per family:	17
Max # of BGCs in a family:	17
Families with MIBiG Reference BGCs:	0
Total BGCs: 17 (0 singleton/s), links:	136



a) # of families: 3
 Average # of BGCs per family: 11
 Max # of BGCs in a family: 17
 Families with MIBiG Reference BGCs: 0
 Total BGCs: 34 (0 singleton/s), links: 272

