

1 Automating Assessment of the Undiscovered 2 Biosynthetic Potential of Actinobacteria

3 Bogdan Tokovenko^{1*}, Yuriy Rebets¹, Andriy Luzhetskyy^{1,2*}

4 ¹ Actinomycetes Metabolic Engineering Group, Helmholtz Institute for Pharmaceutical Research
5 Saarland, Saarbrücken, Germany

6 ² Department of Pharmaceutical Biotechnology, Faculty of Natural Sciences and Technology, University of
7 Saarland, Saarbrücken, Germany

8 * Corresponding author

9 E-mail: Andriy.Luzhetskyy@helmholtz-hzi.de (AL), to.bogdan@gmail.com (BT)

1 Abstract

2 **Background.** Biosynthetic potential of Actinobacteria has long been the subject of theoretical estimates.
3 Such an estimate is indeed important as a test of further exploitability of a taxon or group of taxa for new
4 therapeutics. As neither a set of available genomes nor a set of bacterial cultivation methods are static, it
5 makes sense to simplify as much as possible and to improve reproducibility of biosynthetic gene clusters
6 similarity, diversity, and abundance estimations.

7 **Results.** We have developed a command-line computational pipeline (available at
8 <https://bitbucket.org/qmentis/clusterscluster/>) that assists in performing empirical (genome-based)
9 assessment of microbial secondary metabolite gene clusters similarity and abundance, and applied it to a
10 set of 208 complete and de-duplicated Actinobacteria genomes. After a brief overview of Actinobacteria
11 biosynthetic potential as compared to other bacterial taxa, we use similarity thresholds derived from 4
12 pairs of known similar gene clusters to identify up to 40-48% of 3247 gene clusters in our set of genomes
13 as unique. There is no saturation of the cumulative unique gene clusters curve within the examined
14 dataset, and Heap's alpha is 0.129, suggesting an open pan-clustome. We identify and highlight pitfalls
15 and possible improvements of genome-based gene cluster similarity measurements.

16 Introduction

17 Theoretical estimates of Actinobacteria biosynthetic potential often point at a huge possible number of
18 unique antibiotics which are still to be discovered [1–3]. Such an estimate is indeed important as a test of
19 further exploitability of a taxon for new therapeutics. Focus on the *Streptomyces* genus is quite
20 understandable: of the circa 12000 secondary metabolites with antibiotic activity known in 1995, 55%
21 were produced by *Streptomyces* and additional 11% by other Actinobacteria [4].

22 While Actinobacteria were *the* sources of new antibiotics in the 1960-70s, progressively fewer and fewer
23 new antibiotics were found since then, feeding a growing disappointment and disbelief in the ability to find
24 more novel compounds. However, it is now recognized [5,6] that most of the BGCs (biosynthetic gene
25 clusters) are silent and require activation to produce compounds. Thus, genome-based estimate, which
26 also includes silent BGCs, is important for the revival of interest and belief in the potential of
27 Actinobacteria.

1 Currently dominating antibiotics discovery strategy is to search wide – that is, characterize and sequence
2 more strains (for example, marine Actinobacteria), looking for promising compounds. Quite apparently
3 such a strategy is not cheap, and normally it does not overcome the issue of silent BGCs, as these
4 remain silent unless specific steps are taken to activate them. BGC activation is a laborious process.
5 Therefore, an important step preceding BGC activation experiments is prioritisation of the most promising
6 genomes and clusters.

7 Moreover, based on our own experience of understanding and working with BGCs, we felt that there is a
8 need to quantify a persistent feeling of finding highly similar BGCs in different genomes. Thus, we
9 decided to perform medium-scale genomes-based estimation of the number of unique BGCs discovered
10 with each new genome, as well as develop and make public the tools for comparing BGCs en masse.

11 With this study we prepare an analysis-driven strategy for estimating BGCs diversity/similarity in a sample
12 of multiple genomes, and also (to a lesser degree) a strategy for BGCs (and genomes) prioritisation
13 before their activation. Ideally, this strategy would start with an automated exploration of a large corpus of
14 BGCs (looking for similarities and/or specific genes), then proceed through manual examination (and,
15 possibly, adjusting annotations) of the candidate BGC(s), and culminate in the experimental activation of
16 the candidate BGC(s). Based on our genome set, we answer a number of questions about the secondary
17 metabolites genomics of Actinobacteria. Do highly similar BGCs occur in different genomes? How often?
18 How many previously unseen, novel BGCs one may expect when sequencing new strain's genome?
19 Which BGC types are the least and the most unique? What should we expect in terms of novel BGCs if
20 we extrapolate from the current set of genomes?

21 There are several significant benefits of automated estimation compared to infrequent individual studies.
22 First of all, estimation results become more easily reproducible. Secondly, comparison results become
23 readily available as a simple table, which only requires re-calculations when new genomes/clusters are
24 added. Finally, it is much easier (and faster) to update estimates and comparisons as new genomes
25 become available, including discoveries of new species, or development of new bacteria cultivation
26 methods, significantly extending the range of microorganisms which can be grown and studied in the lab.

27 To the best of our knowledge, this is the first attempt at automating empirical (genome-based)
28 quantification of Actinobacteria potential. Presented approach and software can likely be applied to other
29 natural sources of BGCs.

1 Results and discussion

2 Overview of potential

3 Why do we focus on Actinobacteria? Empirically, Actinobacteria were the largest natural product source
4 since 1960s. But does this taxon also stand out at the genomic level, taking into account silent clusters?
5 To answer this question, we performed an overview study of 2775 bacterial genomes. We had also
6 performed the same overview separately for 285 genomes of Actinobacteria.

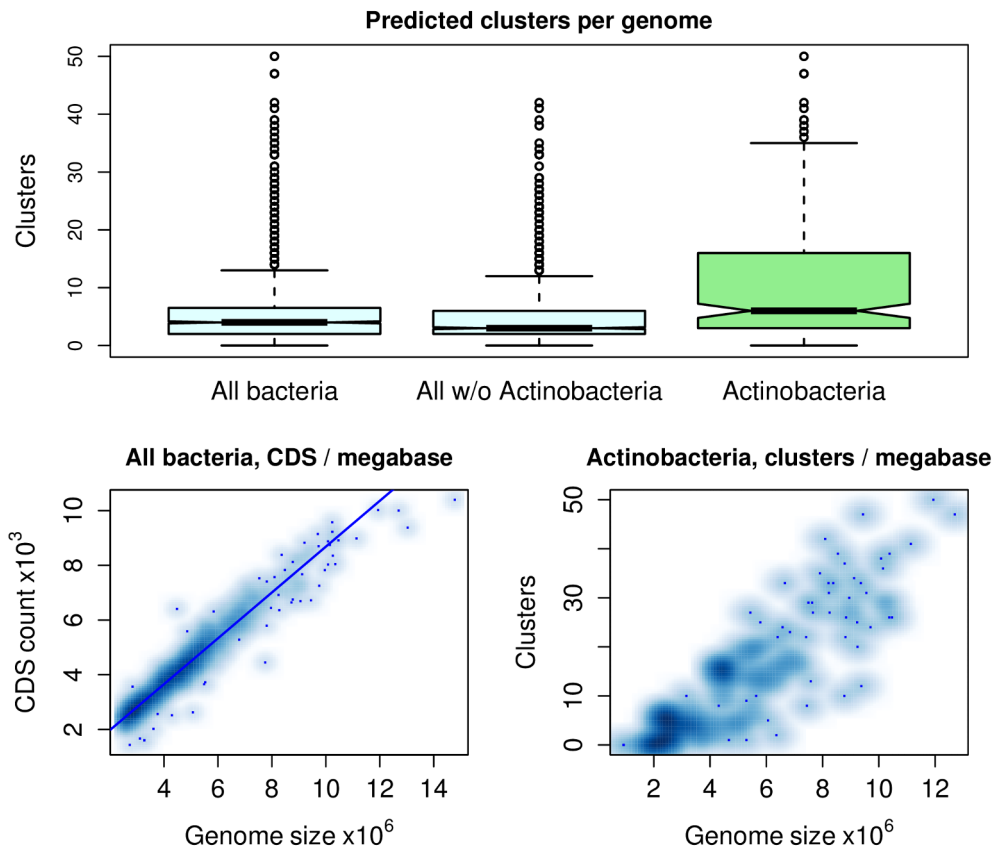
7 We wanted to know how many putative BGCs (biosynthetic gene clusters) can be identified using
8 antismash2 in the genomes, and then group this data by genus/family and other taxonomic levels to
9 reveal highest-cluster-count groups.

10 Prior to the overview analysis highly similar genomes and genomes shorter than 2.5 megabases were
11 removed as described in the *Methods*. Among the smallest genomes, the following phyla were highly
12 represented (roughly in the order of increasing genome size): Proteobacteria, Tenericutes, Bacteroidetes,
13 Spirochaetes, Chlamydiae, Chloroflexi, Firmicutes, Euryarchaeota, Crenarchaeota, Korarchaeota,
14 Thaumarchaeota, Aquificae, Cyanobacteria, Thermotogae. There were also a few Actinobacteria (e.g.
15 *Tropheryma*, *Cryptobacterium*, *Bifidobacteria*, etc). The highest number of BGCs among genomes smaller
16 than 2 megabases was observed in Cyanobacteria – up to 13 gene clusters in *Prochlorococcus marinus*
17 *str.* MIT 9313. Table S1 contains the spreadsheet summary of small genome characteristics.

18 The two largest genomes were *Sorangium cellulosum* So0157-2 (NC_021658) and *Sorangium*
19 *cellulosum* So ce56 (NC_010162) at 14.8 and 13 megabases long, respectively. The next-largest genome
20 was *Streptomyces rapamycinicus* NRRL 5491 (NC_022785) at 12.7 megabases.

21 Among Actinobacteria, the Actinomycetales order (Actinobacteriadae subclass) is the most well-
22 represented (232 genomes), followed by Bifidobacteriales (33 genomes, the entire Bifidobacteriaceae
23 suborder of this dataset). The most well-represented suborder is Corynebacterineae (127 genomes),
24 followed by Bifidobacteriaceae (33 genomes), Micrococccineae (30 genomes), Streptomycineae (19
25 genomes) and others. *Mycobacterium* (63 genomes), *Corynebacterium* (50 genomes), and
26 *Bifidobacterium* (30 genomes) dominate at the genus taxonomical level.

1 Actinobacteria had an average of 10.7 BGCs per genome (median 6), while the complete set of genomes
2 had mean 5.6 (median 4). Maximal number of BGCs identified was 50 in *Streptomyces bingchenggensis*
3 BCW-1 (NC_016582). Fig. 1 (top) shows boxplots of BGC counts per genome for the full genome set and
4 Actinobacteria genomes.



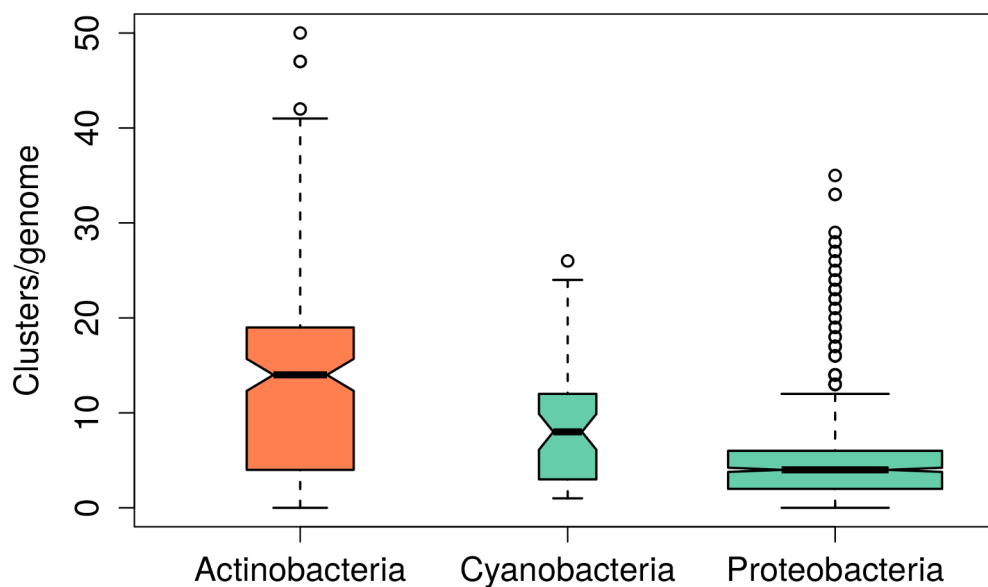
5 **Fig. 1. Predicted gene clusters per genome. Top:** boxplots of predicted gene cluster counts per
6 genome for 1651 bacterial genomes larger than 2.5 megabases, 1508 non-Actinobacteria large genomes,
7 and for 278 Actinobacteria genomes of all sizes. **Bottom left:** smoothed scatterplot of CDS counts and
8 genome sizes for 1651 bacterial genomes. **Bottom right:** smoothed scatterplot of putative gene cluster
9 counts and genome sizes for 278 Actinobacteria genomes.

10 There is a clear linear relationship between genome size and the number of annotated CDS (Fig. 1,
11 bottom left; Pearson correlation 0.97, p-value < 10⁻⁵). Expressed as a linear model, $CDS_count =$
12 $304.3 + 837.8 \times genome_size_in_megabases$. Linear relationship holds for all studied subsets of
13 bacterial genomes.

14 In Actinobacteria the count of putative BGCs appears to correlate better with genome size (Fig. 1, bottom
15 right; Pearson correlation coefficient 0.88, p-value < 10⁻⁵) than in all bacteria (not shown). Expressed as a

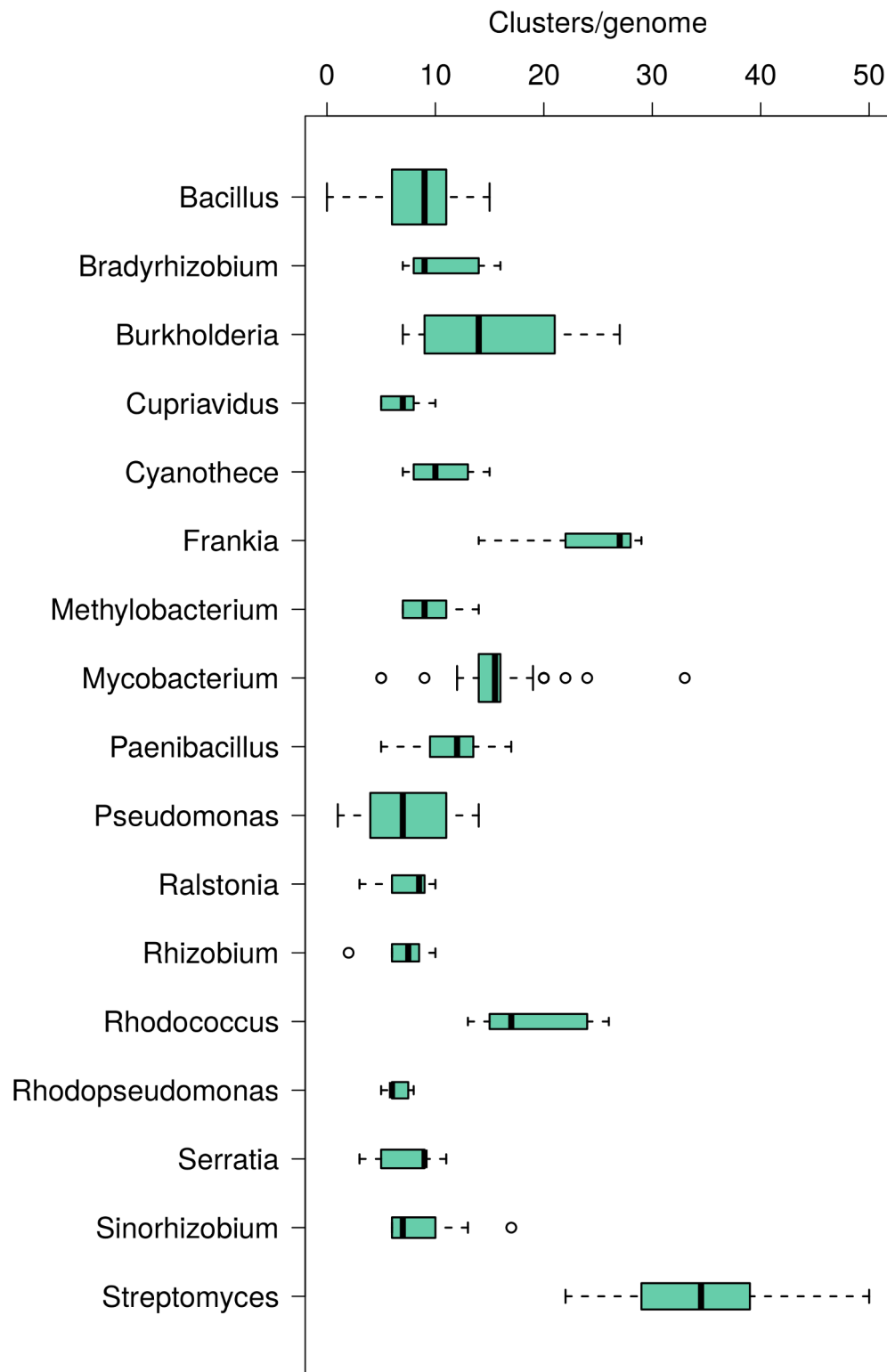
1 linear model, $\text{clusters_count} = -6.8 + 3.9 \times \text{actino_genome_size_in_megabases}$ (almost
2 4 clusters per megabase, except for the first 1.74 megabases).

3 At the phylum taxonomy level, the two leaders are Actinobacteria and Proteobacteria. Cyanobacteria are
4 also promising: this phylum has at least 3 genomes with over 20 putative BGCs each (Fig. 2). Fig. 2
5 shows that the majority of Proteobacteria have fewer than 10 BGCs (mean is 4.9 clusters/genome), but
6 there is also a group of high BGC count outliers which make this phylum stand out. Cyanobacteria mean
7 is higher than that of Proteobacteria (8.5 clusters/genome). Going down to the class level, we see that
8 specifically Beta- and to a lesser extent Delta-Proteobacteria have the most promising BGCs per genome
9 count (not shown).



10 **Fig. 2. Boxplot of putative gene clusters count per genome by phyla.** Only phyla with at least 9
11 genomes, and with mean gene clusters per genome larger than or equal to 4 were included. Width of
12 each boxplot is proportional to the square root of the number of genomes in the phylum.

13 At the genus level (Fig. 3), *Streptomyces* clearly dominate. Runners-up are *Frankia* and *Rhodococcus*
14 (both with low genome counts), also *Burkholderia* and *Mycobacteria*. 60 *Mycobacteria* genomes appear
15 to be fairly heterogeneous with respect to their BGCs count per genome, as can be seen by the highest
16 number of outliers this genus has in Fig. 3.



- 1 **Fig. 3. Boxplot of putative gene clusters count per genome by genus.** Only genii with at least 5
- 2 genomes, and with mean gene clusters per genome larger than or equal to 6 were included. Width of
- 3 each boxplot is proportional to the square root of the number of genomes in the genus.
- 4 For a more in-depth analysis of BGCs by their types and proportions within Actinobacteria please see [7].

1 Automating biosynthetic potential assessment

2 After confirming that Actinobacteria indeed have the highest genomic potential to produce bioactive
3 compounds, the next step is to assess BGCs similarity. We developed our method based on how a
4 human expert would approach the task of measuring BGCs similarity. A pair of (potentially similar) BGCs
5 is the basic analysis unit of our method. A normalized BGC similarity score enables multiple types of
6 downstream analysis. An ability to calculate this similarity score from different metrics ensures some
7 flexibility of the developed pipeline to handle multiple and changing scenarios.

8 The software pipeline comes in two parts: the primary genomes processing Python script (with multiple
9 external dependencies), which produces a CSV file for the second part – R scripts for CSV data analysis
10 and summarization (Fig. S10 and dataset S11). The software has command-line interface and was tested
11 on Linux only. However, program options and the source code are well-documented. Source code is
12 available at <https://bitbucket.org/qmentis/clusterscluster>.

13 Python program accepts a list of (possibly multi-locus) GenBank (.gb or .gbk) files, 1 per genome. As long
14 as these GenBank files can be read with BioPython [8], they will be compatible with the program.

15 GenBank was chosen as a sufficiently common sequence and annotation format, though annotations are
16 currently not required for processing. In fact, only genomic sequence(s) and sequence meta-data (ID,
17 species/strain information) are used. Exact and ID-only duplicates are detected and reported by the
18 program.

19 The BGC prediction software that we used, AntiSMASH 2 [9], can either use existing gene annotations or
20 quickly predict CDS using Glimmer. As gene annotations can be quite different, affecting BGC prediction,
21 we decided to strip any existing gene annotations and allow AntiSMASH 2 to use Glimmer for all
22 genomes. This likely had a negative effect on the quality of BGC prediction, but also removed any
23 annotation quality biases.

24 AntiSMASH 2 works by identifying core secondary metabolite genes, and then adds variable-size
25 extensions (5-20kbp, depending on BGC type) up- and down-stream from the core genes. These
26 extensions contain both cluster genes and some non-cluster genes, possibly hampering cluster similarity
27 analysis. We provide an optional patch which allows using AntiSMASH 2 without these extensions.
28 However, as tailoring enzymes are important, by default BGC extensions are included.

1 All the CDS sequences are then translated to a multi-fasta protein file, which is then used for all protein-
2 level comparisons.

3 Our software can identify cluster genes from different genomes which belong to the same group of
4 orthologs across all the genomes in the analysis. We use InParanoid [10] and QuickParanoid
5 (<http://pl.postech.ac.kr/QuickParanoid/>) for orthologous groups construction across all the analyzed
6 genomes. InParanoid was slightly modified to allow incremental addition of genomes without recalculating
7 existing genome pair similarities. This processing step generates a single table of all orthologous groups
8 for the analyzed genomes. Genes from different clusters belonging to a common orthologous group form
9 a single “link” between the BGCs; the number of such orthologous links can then be used to quantify
10 cluster similarity.

11 In our earlier pilot study of BGCs similarity (not published) we found that gene orthology is not powerful
12 enough for BGC comparisons. For example, many NRPS proteins despite performing different synthesis
13 steps have high (>70%) identities, which assigns them all to a handful of orthologous groups. Orthology
14 analysis, which requires 4 full-genome blasts per one pair of genomes, is also quite computationally
15 intensive. To partially remedy this, we provide a separate script dependent only on a local InParanoid
16 directory to allow splitting computations to multiple computers/nodes/GRID workers. Although we keep
17 the orthologous groups functionality, our focus is now on full-length protein-protein alignments, which
18 provide better resolution (lower granularity) of gene similarity measurements.

19 Full length protein alignments are performed in parallel with usearch [11], an extremely fast heuristic
20 protein alignment/search tool. It is possible to avoid heuristics by enabling the full dynamic programming
21 mode of usearch, which is significantly slower. In this initial version, all possible pairs of proteins from
22 BGCs are aligned. Computing time for this step can be significantly improved in future versions by only
23 aligning reasonable pairs (e.g. those with similar lengths). All protein alignment results are saved to a file-
24 based cache, enabling process restarts and incremental addition of genomes to the analysis.

25 For all genes similar between any two clusters, their average protein identity is calculated, and gene order
26 in both clusters is compared. Finally, all collected and computed results are written to a single CSV file
27 (one similar BGC pair per row, dataset S11) for further analysis. See Fig. S10 for pseudo-code of BGCs
28 comparison, as a part of genomes processing flowchart.

1 Further analysis and figure plotting are performed with several R scripts, which were used for the
2 preparation of this manuscript.

3 **Gene cluster similarity definition**

4 Ideally, two BGCs should be referred to as similar if the compounds they produce are similar. If one had
5 compound structures for both BGCs, one could then use Tanimoto compound similarity index with a
6 certain threshold for this purpose. However, it is not immediately obvious what should the value of such a
7 threshold be: indeed, what is the minimal difference between similar and dissimilar molecules? For
8 example, would methylation make a new molecule dissimilar? What about one more methylation?
9 Moreover, predicting compound structure purely from gene composition of a cluster is neither precise nor
10 reliable. This further complicates the definition of what “similar gene clusters” really are (except for
11 obvious cases, like our control BGCs). Although we thus realize that both the metric and the threshold we
12 are using are not perfect, we believe they provide useful proxies for identifying similar and unique BGCs,
13 responsible for biosynthesis of small molecules with new chemotypes. Making our toolbox available
14 allows interested parties to modify and create new metrics and thresholds, adapting them to specific
15 tasks.

16 **Gene cluster similarity data table**

17 After developing the software and being motivated to test it, we sought applying it and getting BGCs
18 similarity estimate for a subset of Actinobacteria genomes. In this and the following sections we describe
19 some of the analyses of the data generated with our software.

20 Using 208 thoroughly de-duplicated complete Actinobacteria genomes as inputs, we obtained a data table
21 of BGC pairs and their similarity. This table is the starting point for all the following analyses. The resulting
22 spreadsheet is available as dataset S11. Columns have the following meanings:

- 23 • **is_intra**: “1” means that both BGCs of the pair belong to the same genome (intra-genome
24 similarity); “0” means inter-genome similarity;
- 25 • **genome1_ID, genome2_ID**: GenBank accession number for the longest nucleotide sequence of
26 the genome (i.e. chromosome); for a subset of non-GenBank genomes – an internally assigned
27 genome ID;

- 1 • **species1, species2**: species information extracted from the original GenBank files; this is
- 2 missing for some of the non-GenBank genomes;
- 3 • **cluster1, cluster2**: numeric identifiers of BGCs within the genome, as assigned by the cluster
- 4 prediction method;
- 5 • **type1, type2**: generic BGC type, as assigned by the cluster prediction method;
- 6 • **genes1_count, genes2_count**: number of genes in the 1st and 2nd clusters of the pair;
- 7 • **size1_kb, size2_kb**: BGC sizes in kilobases;
- 8 • **ortholinks_count1, ortholinks_count2**: number of orthology links between the BGCs (see
- 9 *Methods* for details); we have not calculated these for our test set;
- 10 • **similar_genes_count**: a number of unique gene pairs between the clusters with protein identity
- 11 above 60%;
- 12 • **avg_protein_identity**: actual average protein identity of all similar genes;
- 13 • **P, K, S**: Pearson, Kendall, and Spearman correlations of two integer vectors of similar gene
- 14 positions in both clusters;
- 15 • **ratio1, ratio2**: ratio of similar_genes_count to genes1_count and genes2_count, respectively;
- 16 • **ratioa**: an average of ratio1 and ratio2;
- 17 • **clusterid1, clusterid2**: string concatenation of genome1_ID with cluster1 and genome2_ID with
- 18 cluster2, respectively; this serves as a unique BGC identifier.

19 Of these columns, the last five (ratio1, ratio2, ratioa, clusterid1, clusterid2) were added by the additional
20 preprocess.R script.

21 **Gene-level synteny**

22 To quantify local gene-level synteny for BGC pairs, we tested as measures Spearman and Kendall rank-
23 correlation coefficients, as well as Pearson's correlation coefficient and Kendall tau distance for two
24 vectors of similar gene positions within their respective gene clusters. None of these methods was perfect
25 for gene synteny. Rank-correlation measures ignore gene deletions/insertions, while Pearson's correlation
26 is not sensitive to gene rearrangements in longer gene clusters. We are working on a more sensitive
27 gene-level synteny metric, which will be added to the software repository after testing.

1 **Thresholds and unique gene clusters**

2 As already mentioned in *Gene cluster similarity definition*, defining a single perfect threshold might not be
3 possible, at least for genome-based gene cluster similarity comparison.

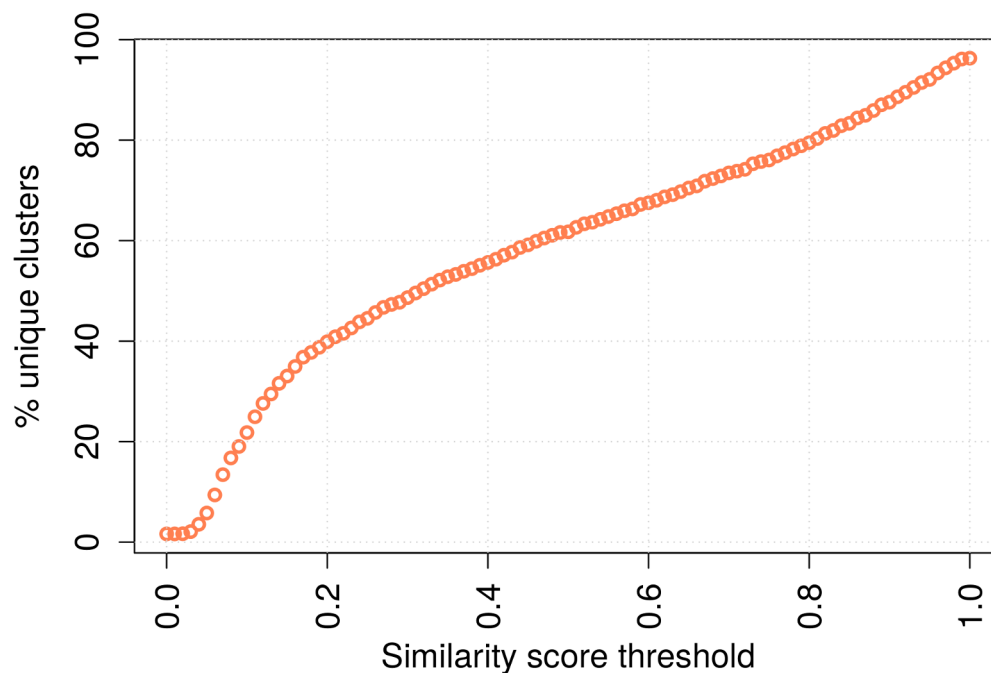
4 Similarity thresholds that we used are primarily based on values obtained for 4 pairs of known BGCs from
5 different genomes responsible for the production of the same or highly similar compound: landomycin,
6 microsclerodermin, moenomycin, nonactin (Table 1). First, only BGC sequences were compared, out of
7 their genomic contexts – that is, with well-defined boundaries (manually annotated, no extensions with
8 unrelated genes; top 4 rows of Table 1). For moenomycin, complete genome sequences were also
9 available, so the manually annotated BGCs were first compared to their respective automatically-
10 annotated counterparts in the genomic context (yielding lower similarity scores), and then both *S.*
11 *Ghanaensis* and *S. Clavuligerus* complete genomes were analysed, yielding the lowest control similarity
12 score of 0.29 for automatically annotated moenomycin BGCs in their genomic context (last 3 rows of
13 Table 1). Based on these controls, the lowest acceptable BGC similarity score was set to 0.5 – sufficient
14 for identifying all but one of control cases as similar.

15 **Table 1. Pairs of previously characterized similar gene clusters.**

Gene cluster	Species	Gene pairs	Similarity score
Landomycin	<i>S. cyanogenus</i> / <i>S. Globisporus</i>	26	0.79
Microsclerodermin	2 <i>Streptosporangium</i> species	9	0.53
Moenomycin	<i>S. ghanaensis</i> / <i>S. Clavuligerus</i>	12	0.73
Nonactin	<i>S. sp. CcalMP-8W</i> / <i>S. Fulvissimus</i>	24	0.84
Moenomycin	<i>S. clavuligerus</i> cluster/genome	14	0.55
Moenomycin	<i>S. ghanaensis</i> cluster/genome	15	0.67
Moenomycin	<i>S. ghanaensis</i> / <i>S. Clavuligerus</i> whole genomes	13	0.29

16 Looking into the 0.29 similarity score for moenomycin in genomic context, we found that a nearby NRPS-
17 T1PKS BGC and a single T1PKS gene in *S. Clavuligerus* were appended to the predicted moenomycin
18 BGC, resulting in a single huge composite cluster (97 kbp and 85 genes versus 35 kbp and 30 genes for
19 predicted moenomycin in *S. Ghanaensis*). Thus, even despite 13 similar genes with average protein
20 identity 68.3%, these two gene clusters received a low similarity score. In our method, the estimated

1 percentage of unique BGCs will be sensitive (proportional) to the frequency of predicting a single
2 combined/composite BGC instead of two or more distinct clusters.
3 Despite using control BGCs, the 0.5 similarity score would benefit from more supporting evidence. Thus,
4 we deemed it necessary to estimate the variability of unique BGCs percentage over the full range of
5 threshold values (Fig. 4). Between similarity score threshold values 0.2 and 1.0 the percentage of unique
6 BGCs grows slower than the threshold (slope < 1), indicating relative stability of the unique percentage
7 estimate, and suggesting that minimal percentage of unique BGCs in our dataset is approximately 40 (at
8 threshold value 0.2).



9 **Fig. 4. Effects of changing similarity score threshold on the percentage of unique gene clusters.**

10 To summarize, we provide the following support for the 0.5 score threshold used throughout this
11 manuscript:

- 12 • it is the lowest threshold allowing detection of 4 pairs of control BGCs (outside their genomic
13 context and after manually refining BGC boundaries);
- 14 • it is higher than the minimal threshold 0.2 suggested by diagnostic plot for thresholds (Fig. 4).

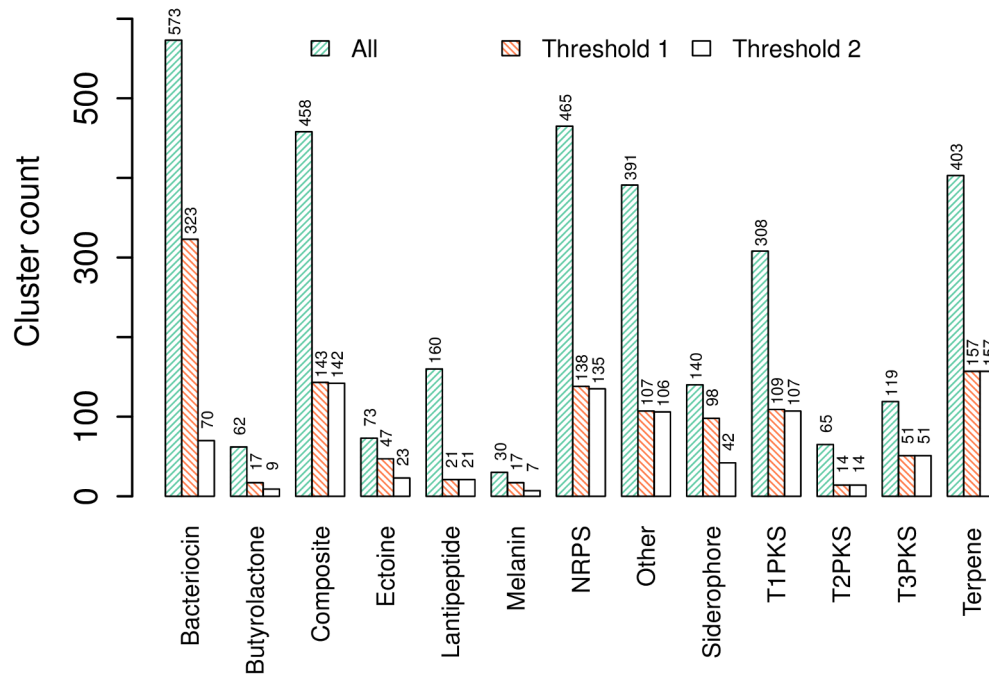
15 The 0.5 threshold is too high to detect the only control BGC tested in its genomic context (without manual
16 cluster boundaries adjustment). This suggests that our estimates of unique BGCs may have upward bias.

1 **Gene clusters diversity: percentage of unique clusters**

2 After building the BGC similarity table, we focused on finding the percentage of unique BGCs – that is,
3 clusters which do not have above the threshold similarities to other clusters. We decided to use the *ratioa*
4 column as the primary measure of BGC similarity (hereafter referred to as similarity score). It is of course
5 possible to use other dataset columns and their combinations as measures of similarity; however, existing
6 analysis R scripts will need to be updated for that.

7 Based on our control BGCs (see *Thresholds and unique gene clusters* above), we started with cluster
8 similarity score 0.5 or higher (threshold 1), identifying 61.7% BGCs as unique. This threshold seems low,
9 but 4 pairs of control similar BGCs would need such a low score to be identified as similar pairs. Low
10 score may identify some of the longer unique BGCs as “similar”, but given the presence of up- and down-
11 stream extensions around putative BGCs, it is more likely that even such a low threshold would fail to
12 detect truly similar BGCs, as exemplified by the low similarity score of 0.29 for the control moenomycin
13 BGC in 2 different genomic contexts (Table 1 in *Thresholds and unique gene clusters*).

14 If, in addition to threshold 1, we require at least 10 genes to be similar (protein identity > 60%) between
15 the BGCs (threshold 2), then the percentage of unique clusters increases to 73%. The number of genes
16 for threshold 2 was also based on our control clusters (Table 1) and diagnostic plots for thresholds
17 (Appendix S2). This additional threshold has two major effects: first, clusters shorter than 10 genes are
18 automatically considered “unique”; second, shorter spurious similar BGCs obtained at similarity score
19 threshold 0.5 are filtered out. Based on this threshold increase, we can conclude that clusters shorter
20 than 10 genes constitute less than 11% of all the BGCs, and that false-positive detections of shorter
21 similar BGCs with a score threshold 0.5 are smaller than 11%. Examining distribution of putative BGCs by
22 predicted type in 3 sets – all, unique, and similar – shows that additional 10-gene threshold affects the
23 most bacteriocin, butyrolactone, ectoine, melanin, and siderophore BGCs, which are indeed normally
24 shorter than 10 genes (Fig. 5). Taken together, decrease of the count of these BGC types explains almost
25 all of the estimate change. In other words, requiring at least 10 genes to be similar does not change the
26 unique BGCs estimate based on similarity score.



1 **Fig. 5. Distribution of secondary metabolite gene clusters by type in 3 datasets:** all 3247 gene
 2 clusters, non-unique at threshold 1 (similarity score above 0.5), non-unique at threshold 2 (similarity score
 3 above 0.5 and at least 10 similar genes).

4 Interestingly, if we only use the existence of 10 similar genes as a threshold (disregarding all the other
 5 metrics of BGC pair similarity), we get 63.35% of unique BGCs (see Appendix S2 for a diagnostic plot of
 6 similar genes threshold versus percent unique BGCs). Adjusting by 11% of the clusters known to be
 7 shorter than 10 genes, we arrive at the base estimate of 52% unique BGCs.

8 As we learn from diagnostic plots, around 4% of all BGCs with near-perfect similarities might come from
 9 the undetected similar genomes, decreasing our base estimate from 52% to 48%.

10 To summarize, we think that the 40-48% range of unique BGCs is a good estimate of uniqueness for the
 11 3247 BGCs examined in 208 de-duplicated genomes.

12 **Highest counts of unique gene clusters per genome**

13 The next question we wanted to answer is which genomes have the highest unique BGC counts (as
 14 defined in the previous section). After sorting all the genomes by the descending count of unique BGCs,
 15 we obtain a list of soil-dwelling Streptomycetacea (Table 2). The only surprising entry here is *Nocardia*
 16 *brasilensis*, which can be pathogenic. As we know, the larger the automatically predicted BGC is –

1 possibly including several unrelated smaller clusters – the lower is the chance that it will be properly
 2 recognized as similar to something else. Thus, this table of the most unique BGCs-rich genomes is likely
 3 to be biased towards bigger or closely located but unrelated BGCs.

4 **Table 2. Top 5 genomes by absolute counts of unique secondary metabolite gene clusters.**

Genome ID	Name	Family	Habitat	Unique gene clusters	Total gene clusters
CP007155.1	<i>Kutzneria albida</i> DSM 43870	Pseudonocardiaceae	Soil	53	54
NC_016582.1	<i>Streptomyces bingchenggensis</i> BCW 1	Streptomycetaceae	Soil	46	53
NC_016109.1	<i>Kitasatospora setae</i> KM-6054	Streptomycetaceae	Soil	43	49
NC_018681.1	<i>Nocardia brasiliensis</i> ATCC 700358	Nocardiacea	Soil	40	51
CM001015.1	<i>Streptomyces clavuligerus</i> ATCC 27064	Streptomycetaceae	Soil	40	52

5 We had previously reported that secondary metabolism genes occupy a significantly bigger proportion of
 6 *K. albida* genome compared to other actinomycetes, with 9.6% of all genes located inside manually-
 7 curated secondary metabolite BGCs [12]. It is thus not too surprising to find *K. albida* genome at the top
 8 of Table 2. But does it really have 53 unique BGCs? Expert examination of 46 BGCs reported previously
 9 [12] identified at least 9 BGCs which are not unique (a *whiE*-like type II PKS, several siderophores and
 10 terpenes, and also bacteriocin and ectoine BGCs). It is thus clear that 53 unique BGCs is an
 11 overestimate. Correspondingly, the above-derived 40-48% of unique BGCs estimate should be
 12 interpreted as an upper bound estimate.

1 **Gene clusters diversity: groups of gene clusters**

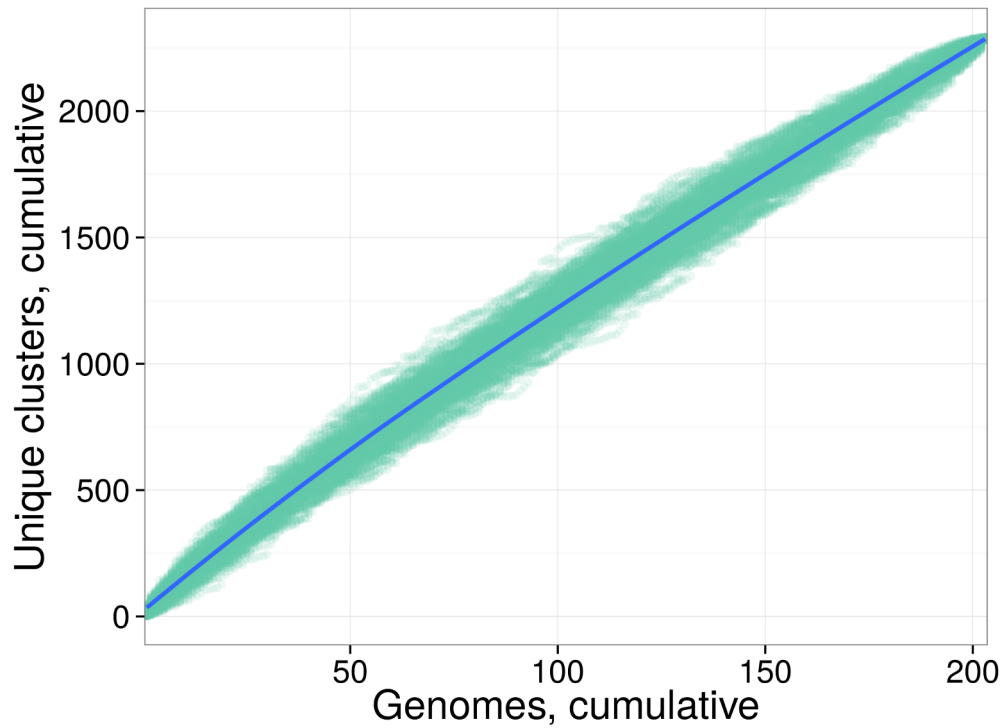
2 Another way to look at gene clusters diversity is to group/cluster them by similarity. We provide an R
3 script which does this based on the similarity score as defined above (see *Methods/Grouping gene*
4 *clusters by similarity*). Applying this approach at similarity score threshold 0.5, we identified 2393 groups
5 of 3247 BGCs.

6 Doroghazi *et al* [2] used a clustering-based approach to identify groups of similar BGCs. In our study a
7 single group of clusters contains on average 1.4 clusters, while groups identified by Doroghazi and
8 colleagues have 2.8 clusters per group (average calculation takes into account singleton groups).
9 Apparently, the similarity threshold that we use is more stringent than in the aforementioned study. In
10 addition, Doroghazi *et al* used their own method for finding clusters; despite having much higher fraction
11 of clusters-rich *Streptomyces* (approximately 40% of all genomes examined, compared to our 15%),
12 Doroghazi *et al* found on average 13.8 clusters per genome, compared to 15.9 found with AntiSMASH 2
13 in this study.

14 **Extrapolating the number of unique gene clusters**

15 Estimation of gene cluster similarity would be incomplete without extrapolating to more genomes. A
16 popular approach is to apply rarefaction analysis. However, rarefaction analysis was originally designed
17 for interpolating to a smaller sample size, not extrapolating to a larger sample size (at least without
18 assuming a specific underlying distribution), and it is not clear if all the assumptions are met (such as
19 sufficiently large sample size and random distribution of “individuals” in samples), in addition to the
20 underlying biological limitations (see *Rarefaction analysis* below).

21 Thus, we had applied a more appropriate resampling approach to look for saturation of the total unique
22 BGCs curve (Fig. 6). With a little over 200 genomes, the curve is slightly bent, but there is definitely no
23 saturation. We thus refrained from extrapolating this curve to higher numbers of genomes. See Appendix
24 S3 for 2 more scatterplots, at thresholds 0.2 and 0.8.



1 **Fig. 6. Total unique gene clusters smoothed scatterplot with genome order resampling.** Similarity
2 score threshold used: 0.5. Solid line is a loess fit to the scatterplot. The final point is always the same in
3 all resamplings (the total number of unique gene clusters), causing scatterplot to narrow down at the last
4 examined genome.

5 **Heap's alpha**

6 In studies of bacterial core and pan-genomes, the Heap's alpha parameter is often used to estimate how
7 many of the genes from a specific taxonomic group we haven't seen yet. For pan-genomes, Heap's alpha
8 smaller than 1 means an open pan-genome (i.e. there are many more genes we haven't seen yet), while
9 alpha larger than 1 means a closed pan-genome (there is only a small/limited number of genes we
10 haven't seen yet). Rephrasing for our case, each genome has one or more groups of gene clusters (as
11 described earlier), including duplicates. We applied the micropan R package [13] and obtained alpha =
12 0.129 for our set of genomes, indicating an open pan-clustome.

1 **Rarefaction analysis**

2 Application of rarefaction analysis to extrapolate expected groups of gene clusters to higher genome
3 counts has, in our opinion, argumentative value, because several assumptions of this analysis are likely
4 not satisfied.

5 The first assumption is that “individuals are randomly distributed”. For BGC diversity, this implies that
6 BGC types (and cluster similarity groups) are randomly distributed across genomes – which is not the
7 case, as cluster composition is strongly defined by each bacteria’s environment, distinguishing, for
8 example, clustomes of soil and water-dwelling bacteria. One has to control for the environment variable to
9 satisfy this requirement.

10 The second assumption is that sample size is sufficiently large, which may or may not be the case.

11 Finally, to extrapolate from rarefaction analysis, one has to assume a certain underlying parent
12 distribution.

13 On top of the technical limitations, rarefaction analysis does not account for several important biological
14 peculiarities of BGCs. First of all, BGCs consist of biosynthetically-active sub-units (genes and gene
15 modules), which can be re-arranged to produce different compounds, resulting in a combinatorial
16 expansion of the BGCs diversity. Secondly, genome-based extrapolations rely solely on the information
17 obtained from cultured strains. In other words, we might be missing a significant portion of BGCs diversity
18 until methods are developed to grow hard to culture strains. Finally, genome-based analysis relies on the
19 sequenced subset of bacterial genomes, inheriting biases which influenced strain selection for
20 sequencing. As the number of genomes grows, this final bias should decrease.

21 That said, rarefaction analysis is still useful, and can be found in Appendix S4. The same analysis was
22 performed on a small subset of 38 *Streptomyces*-only genomes; extrapolation indicated that
23 *Streptomyces* alone might be responsible for approximately half of all the groups of BGCs we have yet to
24 discover.

1 **Other types of analysis**

2 There are several more types of analysis which can be performed on the table of similar gene cluster
3 pairs generated by our software. Probably the first one to mention is looking for relationships between
4 16S rDNA phylogeny and BGC similarity, and/or percentage of unique clusters per genome. Succeeding
5 in establishing a link between the two will provide more proof for the horizontal transfer of entire
6 secondary metabolite BGCs. Another interesting (and related) question is whether fungal and bacterial
7 BGCs have similar elements, and if yes – to what extent, and which types of BGCs? Our software should
8 be applicable to fungal genomes as well, as long as BGC prediction will function for fungal genomes.

9 Another interesting type of analysis to perform is a graph of BGC similarity links (weighted edges)
10 between all the genomes (nodes). This graph can be examined for links density, thus grouping genomes
11 by BGC similarity. Another variation could add taxonomy information, e.g. by creating graph nodes from
12 genus-level aggregates of genomes, with edges representing a sum of similarities between genomes of
13 the connected genii. Described graph analysis may help identify biosynthetic hubs and evolutionary links
14 between BGCs.

15 **AntiSMASH 3 known gene clusters report**

16 AntiSMASH 3 (which has the functionality to report similarity to known gene clusters) was not yet
17 available at the time of performing this study. As soon as it was published [14], we briefly compared our
18 similarity approach to the known clusters reports of AntiSMASH 3.

19 The major difference is that AntiSMASH 3 is better suited for analysing individual genomes (or small
20 groups thereof): BGCs of the given genome(s) are compared against a set of curated clusters.
21 Conversely, our pipeline is better suited for the analysis of larger groups of genomes either with or without
22 any reference BGC datasets, to produce overview BGC similarity statistics. Incremental analysis
23 (comparing groups of new genomes to those analysed previously) is also possible.

24 There are also some other differences – for example, our protein similarity estimates are based on full-
25 length alignments instead of protein BLAST results, and we use more cluster similarity metrics – but the
26 most important difference is in the intended use.

1 **Limitations of the study**

2 As secondary metabolite gene clusters were predicted using AntiSMASH 2, this study is clearly limited to
3 the range of similarity-detectable compounds. However, assuming the same biosynthetic logic of possible
4 new secondary metabolite BGCs predicted using other/newer methods, our approach to estimate BGC
5 similarity and uniqueness should stay valid, and – thanks to the publicly available source code – can be
6 applied to new prediction results as well.

7 Another potential pitfall was high internal similarity of the ~300 genomes set downloaded from GenBank.
8 However, we have invested significant effort into genome de-duplication (discarding ~100 highly-similar
9 genomes), and accounted for any remaining similar genome pairs during subsequent data analysis.

10 Finally, the number of currently available *draft* or otherwise *incomplete* Actinobacteria genomes is much
11 higher than 300, but we intentionally limited the used genomes quality to *complete*, as otherwise, due to
12 genome fragmentation, the number of predicted BGCs would be artificially inflated due to BGC fragment
13 predictions.

14 **Future work**

15 The approach (and software) that we present can be further improved with:

- 16 • better synteny measures, as discussed in *Gene-level synteny*;
- 17 • taking into account the order of biosynthetic domains (using a method similar to proper gene-level
18 synteny); this should significantly increase precision of similarity measures;
- 19 • better gene cluster boundaries by the gene cluster prediction software; as discussed above,
20 correct boundaries are crucial for proper similarity measurements;
- 21 • common genome annotation pipeline such as prokka [15] instead of stripping existing gene/CDS
22 annotations and then performing basic re-annotation with Glimmer;
- 23 • possibly, predicted structure similarities as an additional metric; this improvement is hard to
24 achieve, unless compound structure prediction from genes significantly improves;
- 25 • a graphical user interface.

26 We welcome interested groups to submit comments, suggestions, and patches.

1 **Applications**

2 In addition to answering questions about gene clusters diversity, our software can be used to catalogue
3 BGCs and create databases – either local, for one or several labs, or global, as a web-service.
4 Implementing such a database, especially with incremental genome additions (which is already possible
5 thanks to caching computationally-intensive steps) will allow further strengthening and improvement of
6 genome-based BGC selection and activation approaches. An obvious use would be to focus on (extract
7 and examine) only unique BGCs from a genome newly added to such a database. Another obvious use
8 would be to look for similar BGCs (and additional producers) of the BGC under investigation. Even low
9 similarity scores might be beneficial, as they help discovering BGCs with only some of the biosynthetic
10 modules being similar, indicating a potential modified compound.

11 We have created and keep extending an internal database of similar BGCs for the purposes listed above.
12 Further development should enable providing such a database service to the wider research community.

13 **Conclusions**

14 We found that Actinobacteria (on average) have more predicted secondary metabolite gene clusters per
15 genome than other bacteria. Cyanobacteria and Proteobacteria are runners-up in terms of the number of
16 predicted gene clusters per genome.

17 The number of predicted gene clusters per genome in Actinobacteria correlates with genome size,
18 increasing by almost 4 gene clusters per megabase beyond the first 1.7 megabases.

19 At the Genus level, Streptomyces, Frankia, Rhodococcus and Burkholderia have the highest average
20 gene clusters per genome counts.

21 In the examined set of 3247 predicted secondary metabolite gene clusters from 208 de-duplicated
22 Actinobacteria genomes we identified 40-48% as the upper boundary of the unique gene clusters
23 percentage. Defining gene cluster similarity threshold (as well as compound [dis-]similarity) is a
24 conceptually hard task – would a single gene make a difference, or two? What if this different gene results
25 in a major compound modification? To formulate the similarity threshold, we used 4 pairs of known similar
26 gene clusters and diagnostic plots across the full range of threshold values. We found that imprecise
27 identification of gene cluster boundaries is the major hurdle for precise gene cluster similarity evaluation.

1 On our dataset, Heap's alpha is 0.129, and there is no saturation of the cumulative unique gene clusters
2 curve for 208 genomes. This implies that Actinobacteria pan-clustome is open, and has many new gene
3 clusters to be discovered.

4 Command-line software pipeline that we developed automates generation of the dataset which served as
5 input for this article's analysis. We also include scripts used for the analysis in the open git repository. Our
6 software pipeline can be used to create and maintain databases of genomes and gene clusters, and to
7 compare new genomes/clusters to such databases. As the next step, we plan developing such a
8 database as an open service to the scientific community.

9 **Methods**

10 **Secondary metabolite gene clusters**

11 We use the term "secondary metabolite gene clusters" as it is used when publishing genome
12 announcements. However, some of the secondary metabolite gene cluster (GC) types (siderophores
13 being the most prominent example) are essential for bacteria survival. We decided to keep these BGC
14 types in the study, so as not to artificially inflate the percentage of unique BGCs based on the currently-
15 accepted use of "secondary".

16 **Overview of bacterial biosynthetic potential**

17 Genomes were downloaded on the 4th of February, 2014, from
18 <ftp://ftp.ncbi.nlm.nih.gov/genomes/Bacteria/all.gbк.tar.gz>. The archive contained 2773 sub-directories,
19 normally one directory per species. However, 9 directories contained multiple strain genomes of the same
20 species, or even genomes of other species (S5 has a full list of such directories). For each of the
21 genomes, multiple GenBank files (representing chromosomes and plasmids) were mechanically merged
22 into a single multi-locus GenBank file. Table S6 lists all 5165 sequence accession numbers together with
23 descriptions and sequence sizes. All 2775 genomes were then analysed with antiSMASH 2.0 [9]. For
24 Actinobacteria-only biosynthetic potential overview, all 285 Actinobacteria results were copied out into a
25 separate directory. For both the complete 2775-genome and 285-genome sets, results were collected into
26 a CSV file using a custom summarize_antismash_results.py script

1 (<https://bitbucket.org/qmentis/bioinformatics->
2 [scripts/src/b3faa96eeac2359a60016ded830ef3dab7c4f228/summarize_antismash_results.py](https://bitbucket.org/qmentis/bioinformatics-scripts/src/b3faa96eeac2359a60016ded830ef3dab7c4f228/summarize_antismash_results.py)), and then
3 loaded into R [16] for further analysis.

4 Both genome lists were searched for same species with different accessions. The highest number of such
5 cases was observed for *Staphylococcus aureus subsp. aureus* ST228, from the genome evolution study
6 [17]. After cleaning out highly similar genomes (see Appendix S7 for a table of accessions), 2724
7 genomes (278 for the Actinobacteria subset) were left for further analysis.

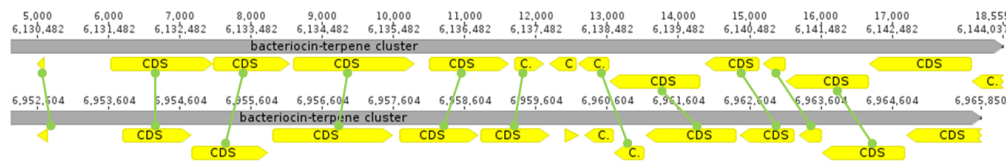
8 Genomes shorter than 2.5 megabases (Table S1) were very likely to have no or very few predicted BGCs;
9 short genomes were removed from the all-bacteria overview, leaving 1651 genomes (Table S8) for
10 analysis.

11 **Assessment of biosynthetic gene clusters diversity**

12 285 complete Actinobacteria genomes from GenBank (see above) were supplemented with 25 genomes
13 from other sources (all are Actinobacteria), then pruned from duplicates and highly-similar genomes
14 (manual pruning, relying on the initial BGC comparison runs and on whole-genome Mauve [18]
15 alignments), resulting in a list of 208 genomes. Secondary metabolite BGCs were predicted in these
16 genomes using AntiSMASH 2 [9]. 5 of 208 genomes had zero predicted BGCs and were excluded from
17 further analysis. Four pairs of known similar BGCs were then added to the analysis as controls and the
18 basis for estimating BGC similarity thresholds (see Table S9 for a complete list of genomes used). See
19 Fig. S10 for a flowchart of genomic data processing.

20 GC pair similarity is defined as a set of metrics, including count of similar gene pairs, average protein-
21 level identity, synteny of cluster genes, and possibly others (Fig. 7). For simplicity, in this study overall
22 BGC similarity score was defined as a ratio of genes similar between the clusters to the count of genes in
23 those clusters, $score = 0.5 * (similar_genes/cluster1_genes + similar_genes/cluster2_genes)$. Each BGC
24 is categorized as "unique" or "similar to other clusters" based on the thresholds defined for each of the
25 metrics. Minimal required protein-level identity was 60%; similar gene pairs were not reported below this
26 identity. Synteny coefficient was calculated as gene order correlation coefficient: similar genes were
27 sequentially numbered in both BGCs, followed by calculating Pearson correlation for both integer vectors.

1 Downstream results analysis uncovered weak filtering capability of synteny coefficient, so despite being
2 calculated it was not used for discerning similar and unique BGCs.



3 **Fig. 7. Example of a comparison of bacteriocin-terpene gene clusters from *Actinosynnema mirum***
4 **DSM 43827 (top) and *Saccharothrix espanaensis* DSM 44229 (bottom).** Rounded-end handles
5 between the CDS of the two genomes mark 11 gene pairs, where protein alignment identity was higher
6 than 60%. For these 11 gene pairs, synteny coefficient (gene order correlation) is 1.0, and average actual
7 protein identity is 82%. Gene cluster similarity score is 0.76.

8 Four pairs of known BGCs from different genomes responsible for the production of the same compound
9 (landomycin, microsclerodermin, moenomycin, nonactin) were used as controls.

10 **Cumulative unique gene clusters curve with genome order** 11 **resampling**

12 To look for possible saturation of the total number of unique gene clusters from examined genomes, we
13 performed the following steps:

- 14 1. Randomize genome order.
- 15 2. Take the 1st genome. As it is the first seen genome, all of its BGCs are unique. Plot the number of
16 observed BGCs (X: genome number; Y: total unique BGCs found).
- 17 3. Take the 2nd genome, compare its BGCs to those seen before. Sum unseen (unique) BGCs from
18 the 2nd genome with the unique BGCs from the 1st, and plot the cumulative value.
- 19 4. Repeat until all the genomes are examined, every time comparing genome's clusters only to
20 those we had seen in previously examined genomes. At every genome, plot cumulative unique
21 BGCs.
- 22 5. Reset seen BGCs to an empty set. Repeat steps 1-4 200 times.

1 **Grouping gene clusters by similarity**

2 Sample gene clusters A, B, and C will be put into a single group if they have pairwise links with similarity
3 scores above the specified threshold. For example, high-scoring links A-B and B-C are sufficient to put all
4 three into a single group, even if A-C happens to be below the threshold. Alternatively, A-B and A-C are
5 sufficient to group all three, even if B-C is below the threshold. Unique BGCs become singleton groups.

6 **Rarefaction analysis**

7 Rarefaction analysis was performed in both incidence and abundance modes using the iNEXT R package
8 [19]. Abundance analysis was performed on counts of gene clusters (“individuals”) per group of gene
9 clusters (“species”, see *Grouping gene clusters by similarity* above), extrapolating to 45 and 60 thousand
10 examined gene clusters in total. Incidence analysis was performed on a vector of genome counts
11 (“samples”) per group of gene clusters (“species”), extrapolating to 15 thousand genomes.

12 **References**

- 13 1. Watve MG, Tickoo R, Jog MM, Bhole BD. How many antibiotics are produced by the genus
14 Streptomyces? Arch Microbiol. 2001;176: 386–390. doi:10.1007/s002030100345
- 15 2. Doroghazi JR, Albright JC, Goering AW, Ju K-S, Haines RR, Tchalukov KA, et al. A roadmap for
16 natural product discovery based on large-scale genomics and metabolomics. Nat Chem Biol.
17 2014;10: 963–968. doi:10.1038/nchembio.1659
- 18 3. Cimermancic P, Medema MH, Claesen J, Kurita K, Wieland Brown LC, Mavrommatis K, et al.
19 Insights into Secondary Metabolism from a Global Analysis of Prokaryotic Biosynthetic Gene
20 Clusters. Cell. 2014;158: 412–421. doi:10.1016/j.cell.2014.06.034
- 21 4. Demain AL. Prescription for an ailing pharmaceutical industry. Nat Biotechnol. 2002;20: 331–331.
22 doi:10.1038/nbt0402-331
- 23 5. Challis GL. Mining microbial genomes for new natural products and biosynthetic pathways.
24 Microbiology. 2008;154: 1555–1569. doi:10.1099/mic.0.2008/018523-0
- 25 6. Rebets Y, Brötz E, Tokovenko B, Luzhetskyy A. Actinomycetes biosynthetic potential: how to bridge
26 in silico and in vivo? J Ind Microbiol Biotechnol. 2013;41: 387–402. doi:10.1007/s10295-013-1352-9
- 27 7. Doroghazi JR, Metcalf WW. Comparative genomics of actinomycetes with a focus on natural
28 product biosynthetic genes. BMC Genomics. 2013;14: 611. doi:10.1186/1471-2164-14-611
- 29 8. Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, et al. Biopython: freely available
30 Python tools for computational molecular biology and. Bioinformatics. 2009;25: 1422–1423.
31 doi:10.1093/bioinformatics/btp163

- 1 9. Blin K, Medema MH, Kazempour D, Fischbach MA, Breitling R, Takano E, et al. antiSMASH 2.0—a
2 versatile platform for genome mining of secondary metabolite producers. *Nucleic Acids Res.*
3 2013;41: W204–W212. doi:10.1093/nar/gkt449
- 4 10. Ostlund G, Schmitt T, Forslund K, Kostler T, Messina DN, Roopra S, et al. InParanoid 7: new
5 algorithms and tools for eukaryotic orthology analysis. *Nucleic Acids Res.* 2010;38: D196–D203.
6 doi:10.1093/nar/gkp931
- 7 11. Edgar RC. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics.* 2010;26:
8 2460–2461. doi:10.1093/bioinformatics/btq461
- 9 12. Rebets Y, Tokovenko B, Lushchik I, Rückert C, Zaburannyi N, Bechthold A, et al. Complete genome
10 sequence of producer of the glycopeptide antibiotic Aculeximycin *Kutzneria albida* DSM 43870T, a
11 representative of minor genus of Pseudonocardiaceae. *BMC Genomics.* 2014;15: 885.
12 doi:10.1186/1471-2164-15-885
- 13 13. Liland LS and KH. micropan: Microbial Pan-genome Analysis [Internet]. 2014. Available:
14 <http://cran.r-project.org/web/packages/micropan/index.html>
- 15 14. Weber T, Blin K, Duddela S, Krug D, Kim HU, Bruccoleri R, et al. antiSMASH 3.0—a comprehensive
16 resource for the genome mining of biosynthetic gene clusters. *Nucleic Acids Res.* 2015;43: W237–
17 W243. doi:10.1093/nar/gkv437
- 18 15. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics.* 2014; btu153.
19 doi:10.1093/bioinformatics/btu153
- 20 16. R Core Team. R: A Language and Environment for Statistical Computing [Internet]. Vienna, Austria:
21 R Foundation for Statistical Computing; 2014. Available: <http://www.R-project.org>
- 22 17. Vogel V, Falquet L, Calderon-Copete SP, Basset P, Blanc DS. Short Term Evolution of a Highly
23 Transmissible Methicillin-Resistant *Staphylococcus aureus* Clone (ST228) in a Tertiary Care
24 Hospital. *PLoS ONE.* 2012;7: e38969. doi:10.1371/journal.pone.0038969
- 25 18. Darling ACE, Mau B, Blattner FR, Perna NT. Mauve: Multiple Alignment of Conserved Genomic
26 Sequence With Rearrangements. *Genome Res.* 2004;14: 1394–1403. doi:10.1101/gr.2289704
- 27 19. Chao A, Gotelli NJ, Hsieh TC, Sander EL, Ma KH, Colwell RK, et al. Rarefaction and extrapolation
28 with Hill numbers: a framework for sampling and estimation in species diversity studies. *Ecol*
29 *Monogr.* 2014;84: 45–67. doi:10.1890/13-0133.1

30

31 **Supporting information captions**

32 **S1 Table. Summary of 1074 deduplicated small genomes shorter than 2.5 megabases.**

33 **S2 Appendix. Diagnostic plots for ranges of threshold values.**

34 **S3 Appendix. Total unique gene clusters resampling scatterplots.**

35 **S4 Appendix. Rarefaction plots.**

36 **S5 Text. Multigenome directories downloaded from GenBank FTP.**

37 **S6 Table. Table of 2775 genomes used for overview analysis.**

38 **S7 Appendix. List of highly similar genomes removed prior to analysis.**

39 **S8 Table. 1651 deduplicated genomes larger than 2.5 megabases used for overview analysis.**

40 **S9 table. All 318 genomes and sequences used for similarity analysis.**

41 **S10 Figure. Similarity analysis genomic data processing flowchart.**

42 **S11 Dataset. Full table (tab-separated values) of similar gene cluster pairs.**