

# A global survey of specialized metabolic diversity encoded in bacterial genomes

Athina Gavriilidou<sup>\*1</sup>, Satria A. Kautsar<sup>\*2</sup>, Nestor Zaburanyi<sup>3,4</sup>, Daniel Krug<sup>3,4</sup>, Rolf Müller<sup>3,4</sup>, Marnix H. Medema<sup>\*\*\*2</sup>, Nadine Ziemert<sup>\*\*\*1,5,6</sup>

\*These authors contributed equally: Athina Gavriilidou, Satria A. Kautsar.

\*\*These authors jointly supervised this work: Marnix H. Medema, Nadine Ziemert.

<sup>†</sup>e-mail: [marnix.medema@wur.nl](mailto:marnix.medema@wur.nl); [nadine.ziemert@uni-tuebingen.de](mailto:nadine.ziemert@uni-tuebingen.de)

<sup>1</sup>Translational Genome Mining for Natural Products, Interfaculty Institute of Microbiology and Infection Medicine Tübingen (IMIT), Interfaculty Institute for Biomedical Informatics (IBMI), University of Tübingen, Auf der Morgenstelle 28, 72076 Tübingen, Germany.

<sup>2</sup>Bioinformatics Group, Wageningen University, the Netherlands

<sup>3</sup>Helmholtz Institute for Pharmaceutical Research Saarland (HIPS)—Helmholtz Centre for Infection Research (HZI), Campus E8 1, 66123 Saarbrücken, Germany.

<sup>4</sup>German Center for Infection Research (DZIF), Partner site Hannover-Braunschweig, Inhoffenstr. 7, 38124 Braunschweig, Germany.

<sup>5</sup>Cluster of Excellence ‘Controlling Microbes to Fight Infections’ (CMFI), University of Tübingen, Tübingen, Germany.

<sup>6</sup>German Centre for Infection Research (DZIF), Partnersite Tübingen, Tübingen, Germany.

## Abstract

Bacterial secondary metabolites have been studied for decades for their usefulness as drugs, such as antibiotics. However, the identification of new structures has been decelerating, in part due to rediscovery of known compounds. Meanwhile, multi-resistant pathogens continue to emerge, urging the need for new antibiotics. It is unclear how much chemical diversity exists in Nature and whether discovery efforts should be focused on established antibiotic producers or rather on understudied taxa. Here, we surveyed around 170,000 bacterial genomes as well as several thousands of Metagenome Assembled Genomes (MAGs) for their diversity in Biosynthetic Gene Clusters (BGCs) known to encode the biosynthetic machinery for producing secondary metabolites. We used two distinct algorithms to provide a global overview of the biosynthetic diversity present in the sequenced part of the bacterial kingdom. Our results indicate that only 3% of genomic potential for natural products has been experimentally discovered. We connect the emergence of most biosynthetic diversity in evolutionary history close to the taxonomic rank of genus. Despite enormous differences in potential among taxa, we identify *Streptomyces* as by far the most biosynthetically diverse based on currently available data. Simultaneously, our analysis highlights multiple promising high-producing taxa that have thus far escaped investigation.

## Background

Secondary metabolites (also called specialized metabolites) are biomolecules that are not essential for life but rather offer specific ecological or physiological advantages to their producers allowing them to thrive in particular niches. These Natural Products (NPs) are more chemically diverse than the molecules of primary metabolism, varying in both structure and mode of action among different organisms<sup>1</sup>. Historically, secondary metabolites have been the major source of human medicine and continue to contribute a substantial part of new chemical entities brought to the clinic<sup>2-5</sup>. Microbial NPs and their derivatives are especially dominant among anticancer compounds and antibiotics<sup>2,5-7</sup>. Regrettably, the emergence of antimicrobial-resistant pathogens has escalated<sup>4,5,8-10</sup> at a time when little attention is given to antibiotic development<sup>2-4</sup>. Combined with the fact that the rate of discovery of novel compounds has been slowing down, this is now leading to a global public health crisis<sup>4,5,8-10</sup>.

Nonetheless, genomics-based approaches to NP discovery<sup>11-13</sup> have revealed that an untapped source of biosynthetic potential lies hidden in the genomes, which displays much greater diversity than the compounds in use so far<sup>5,14,15</sup>. These findings were possible due to the discovery that bacterial genes encoding the biosynthesis of secondary metabolites are usually located in close proximity to each other, forming recognizable Biosynthetic Gene Clusters (BGCs). However, there are large differences in the numbers as well as the kinds of BGCs found in microbial genomes<sup>14,16</sup> and, while metabolomic data indicate that some biosynthetic pathways

are unique to specific taxa<sup>17</sup>, a systematic analysis of the taxonomic distribution of BGCs has not yet been performed. Similarly, while useful estimates of the chemical diversity of specific taxa have been provided<sup>16</sup>, systematic comparisons across taxa are lacking. Because of this, the scientific community appears divided on the best strategy for natural product discovery: should the established known NP producers be studied further or should the community be investigating underexplored taxa<sup>14,18</sup>?

Here, we harnessed recent advances in computational genomic analysis of BGCs to survey the enormous amount of genome data accumulated by the scientific community so far. Using a global approach based on more than 170,000 publicly available genomes, we created a comprehensive overview of the biosynthetic diversity found across the entire bacterial kingdom. We clustered 1,094,877 BGCs into 53,927 Gene Cluster Families (GCFs), and calibrated the granularity of the clustering to make it directly comparable to chemical classes as defined in NP Atlas<sup>19</sup>. This facilitated an analysis of the variance of diversity across major taxonomic ranks, which showed the genus rank to be the most appropriate to compare biosynthetic diversity across homogeneous groups. This finding allowed us to conduct comparisons within the bacterial kingdom. Evident patterns emerged from our analysis, revealing popular taxa as prominent sources of both actual and potential biosynthetic diversity, and multiple yet uncommon taxa as promising producers.

## Results

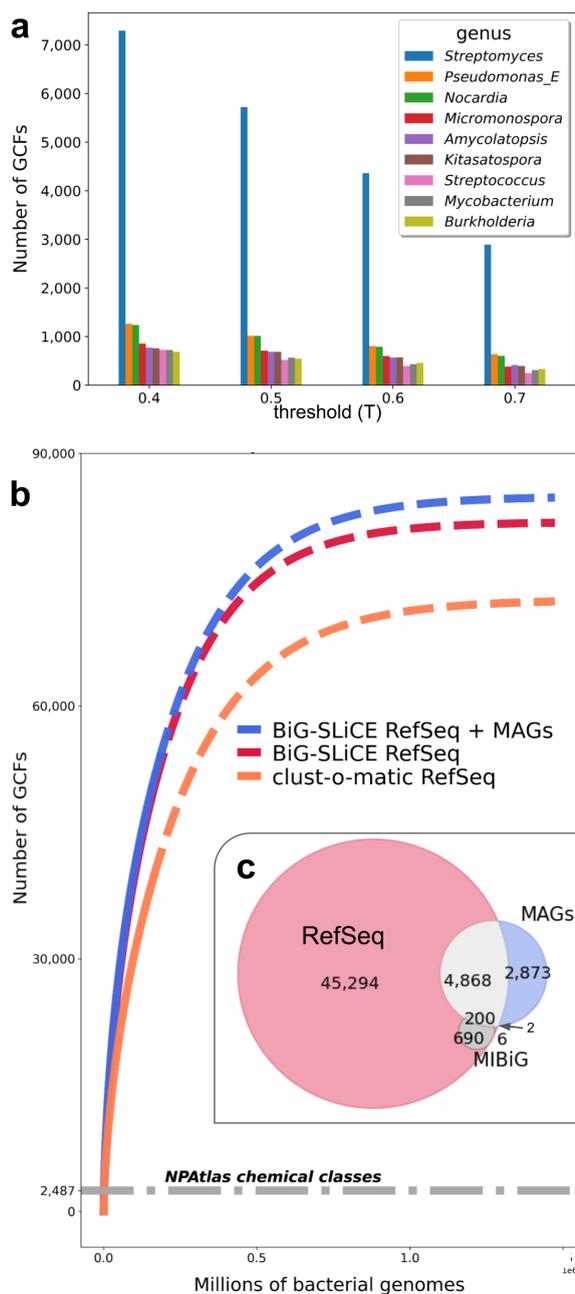
### *Biosynthetic diversity of the bacterial kingdom*

In order to assess the global number of Gene Cluster Families found in sequenced bacterial strains, we ran AntiSMASH on ~170,000 genomes from the NCBI RefSeq database<sup>20</sup> (Supplementary Table 1), spanning 48 bacterial phyla containing 464 families (according to the Genome Taxonomy DataBase classification - GTDB<sup>21</sup>). For the first step of the study, we also included more than 10,000 bacterial Metagenome Assembled Genomes (MAGs) from 5 metagenomic projects of various origins: bovine rumen<sup>22</sup>, chicken caecum<sup>23</sup>, human gut<sup>24</sup>, the ocean<sup>25</sup> and a dataset of uncultivated microbes<sup>26</sup> (Supplementary Table 1). To accurately group similar BGCs – which likely encode pathways towards the production of similar compounds – into Gene Cluster Families (GCFs) across such a large dataset, we used a slightly modified version of the BiG-SLiCE tool<sup>27</sup>, which has been calibrated to output GCFs that match the grouping of known compounds in the NP Atlas database<sup>19</sup> (see Methods). The resulting GCFs were then used to measure biosynthetic diversity across taxa.

**Table 1. Overview of input datasets and overall biosynthetic diversity under various BiG-SLiCE thresholds.** The dataset “Entire Bacterial Kingdom” was used for the computation of the actual and potential biosynthetic diversity found in all cultured (and some uncultured) bacteria. The dataset “RefSeq bacteria with known species taxonomy” was used for pinpointing the emergence of biosynthetic diversity, for which accurate taxonomic information was needed, and for identifying groups of promising producers. \*MAG sources: bovine rumen<sup>22</sup>, chicken caecum<sup>23</sup>, human gut<sup>24</sup>, ocean<sup>25</sup>, uncultivated bacteria<sup>26</sup>.

Dataset		Genomes	BGCs	Gene Cluster Families			
				T = 0.4	T = 0.5	T = 0.6	T = 0.7
Entire Bacterial kingdom	All RefSeq bacteria	170,549	1,060,592	<b>51,052</b>	37,785	28,057	19,152
	Bacterial MAGs*	13,620	34,285	<b>7,943</b>	6,516	5,014	3,519
	Total	184,169	1,094,877	<b>53,927</b>	40,497	31,998	24,446
RefSeq bacteria with known species taxonomy	Complete Genomes	16,004	94,904	<b>16,984</b>	13,546	10,399	7,151
	Draft Genomes	147,265	913,642	<b>37,123</b>	27,748	20,638	14,016
	Total	163,269	1,008,546	<b>41,870</b>	31,237	23,227	15,766

The number of GCFs in RefSeq ranged from 19,152 to 51,052 depending on the threshold used by BiG-SLiCE, while the diversity of GCFs from MAGs analyzed ranged from 3,519 to 7,943 (Table 1). While, as expected, the pure numbers of the analysis changed based on the threshold, the overall tendencies observed remained the same (Figure 1a, Supplementary Figure 1). The effect that the chosen threshold has on these results presented a challenge to our investigation, as previous estimations have also shown great heterogeneity when different thresholds were used<sup>14,16</sup>, precluding direct comparisons of their predictions. Consequently, we sought to relate the choice of our BGC clustering threshold to the structural relationships between their chemical products. After mapping the BiG-SLiCE groupings of 947 known BGCs from the MIBiG repository<sup>28</sup> at different thresholds against the compound-based clustering of their products as provided by the NPAtlas database<sup>19</sup> (Supplementary Figure 2), we chose a threshold of 0.4, as it provided the most congruent agreements between the two groupings, with v-score=0.94 (out of 1.00) and  $\Delta$ GCF=-17.



**Figure 1. Biosynthetic diversity of the bacterial kingdom.** Panel **a**: Bar plots of Gene Cluster Families (GCFs, as defined by BiG-SLiCE) of nine most biosynthetically diverse genera using different thresholds (T). The absolute number of GCFs changes from threshold to threshold, but the general tendencies (highest to lowest GCF count) are consistent between them. Panel **b**: Rarefaction curves of all RefSeq bacteria based on BiG-SLiCE (red) and based on clust-o-matic (orange), and rarefaction curve of the Entire Bacterial kingdom dataset, which includes bacterial MAGs (blue), based on BiG-SLiCE. BiG-SLiCE GCFs were calculated with T=0.4. Clust-o-matic GCFs were calculated with T=0.5. The number of chemical classes documented in NPAAtlas<sup>19</sup>, which come from bacterial producers (gray dotted line), corresponds to 2.9% - 3.3% of the predicted potential of the bacterial kingdom. Panel **c**: Venn Diagram of GCFs (as defined by BiG-SLiCE, T=0.4) of the bacterial RefSeq, Minimum Information about a Biosynthetic Gene cluster (MiBiG<sup>28</sup>) and bacterial MAGs datasets. More information on the MiBiG dataset can be found in Supplementary Table 6. About 36.2% of the GCFs of MAGs are unique (blue shape) to this dataset.

This calibration of thresholds of GCFs to families of chemical structures allowed us to perform a rarefaction analysis to assess how genomically encoded biochemical diversity (expressed as the number of distinct GCFs) increases with the number of sequenced and screened genomes (Figure 1b). The curve appears far from saturated, while the slope is steeper still if the bacterial MAGs are included in the analysis. When compared to the number of chemical classes documented in the NPAtlas<sup>19</sup> database (Figure 1b), it appears that to date only about 3% of the kingdom's biosynthetic diversity has been experimentally accessed.

In an attempt to evaluate the potential contribution of metagenomic data to Natural Product (NP) discovery, we studied how many of the GCFs found in the MAGs datasets were unique to this dataset (Figure 1c). Around 36.2% of GCFs in the MAGs were not found in the RefSeq strains or in the Minimum Information about a Biosynthetic Gene cluster database (MIBiG<sup>28</sup>). Paradoxically, in Figure 1b, the contribution of MAGs does not reflect this finding, but this is most likely because the metagenomic dataset is of limited size and does not cover the full microbial diversity of the biosphere. Considering that metagenomic data are a relatively new and rapidly growing source of genomic information<sup>14,29</sup> and that a high percentage of bacterial strains are still uncultivated<sup>30</sup>, this finding indicates metagenomes as a promising source of undiscovered GCFs.

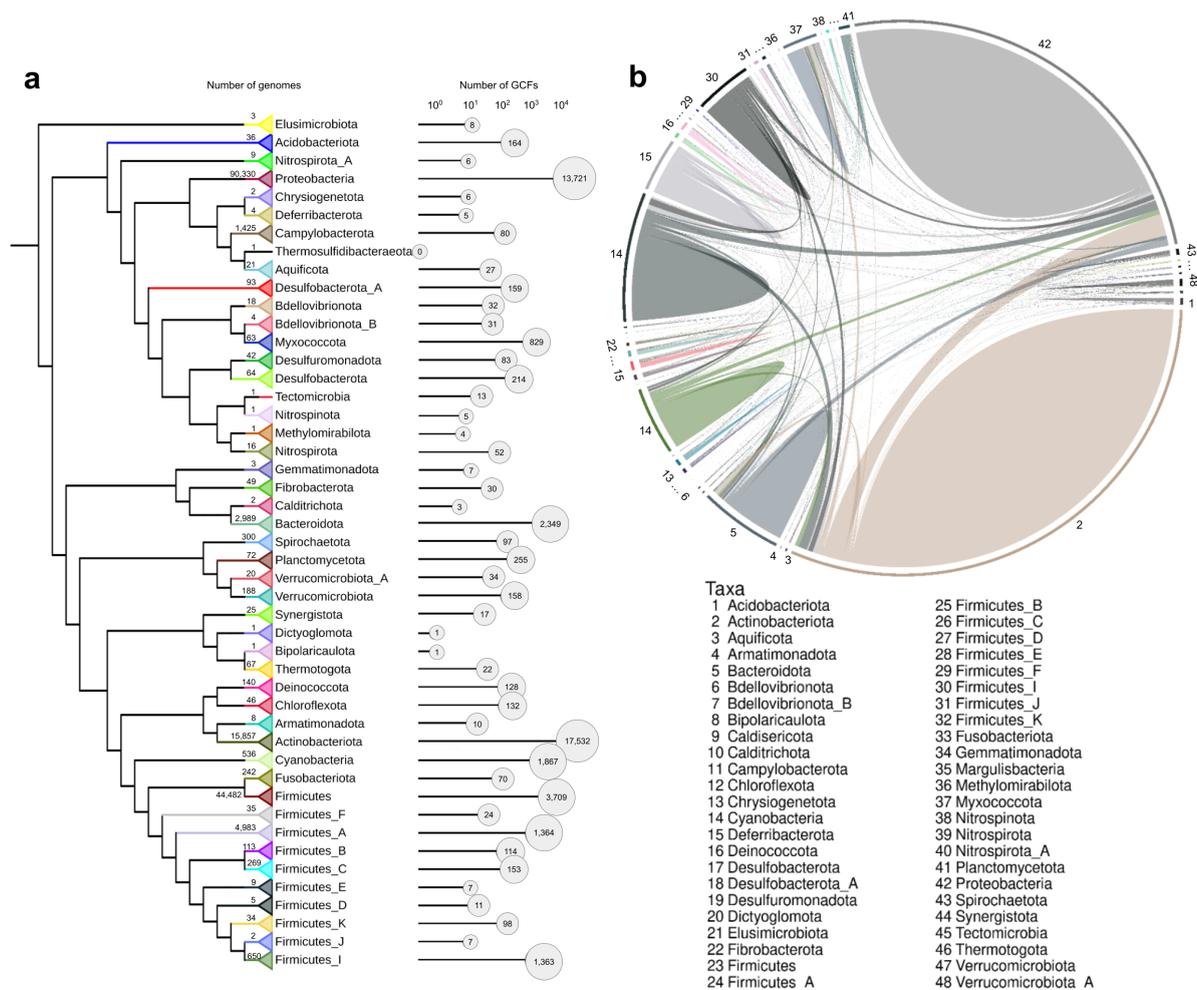
### *Genus as the most appropriate taxonomic rank to compare biosynthetic diversity*

In order to guide the choice of bacterial producers to target for NP discovery, it is important to compare them at a specific taxonomic level. Several studies indicate that there is significant discontinuity in how BGCs are distributed across taxonomy: while, over long periods of evolutionary time, closely related species have shared many BGCs through horizontal gene transfer, most GCFs are still confined to a specific taxonomic range<sup>14,31–33</sup>. Related to this and partly as a result, 'lower' taxonomic ranks like species within a genus carry very similar biosynthetic diversity, while 'higher' taxonomic ranks like phyla within a kingdom usually show large differences in this regard. To assess which taxonomic rank is the most appropriate to assess biosynthetic potential, we aimed to determine up to which taxonomic level biosynthetic diversity remains uniform. For this purpose we used, from our initial dataset, only the RefSeq bacterial strains, whose taxonomic assignment up to species rank (based on GTDB<sup>21</sup>) was known (Table 1), while the MAGs dataset was left out of the remaining analyses due to incomplete taxonomic information.

We first decorated the GTDB<sup>21</sup> bacterial tree with GCF values from the BiG-SLiCE analysis (Figure 2a), revealing the biosynthetic diversity found within currently sequenced genomes at the phylum rank. It immediately stood out that biosynthetic diversity was differently dispersed among the bacterial phyla, in accordance with published data<sup>14,34</sup>. The phyla Proteobacteria and Actinobacteria appeared

particularly diverse, as expected based on the current information regarding known NP producers<sup>16,35,36</sup>.

Next, we examined whether the diversity of each phylum contributed to the domain's total diversity, or if there was overlap among them. For this reason, we depicted the number of unique GCFs within each phylum, as well as the pairwise overlaps (Figure 2b). It appeared that in most phyla, the vast majority (on average  $73.81 \pm 20.35\%$ ) of their GCFs were unique to them and not found anywhere else.

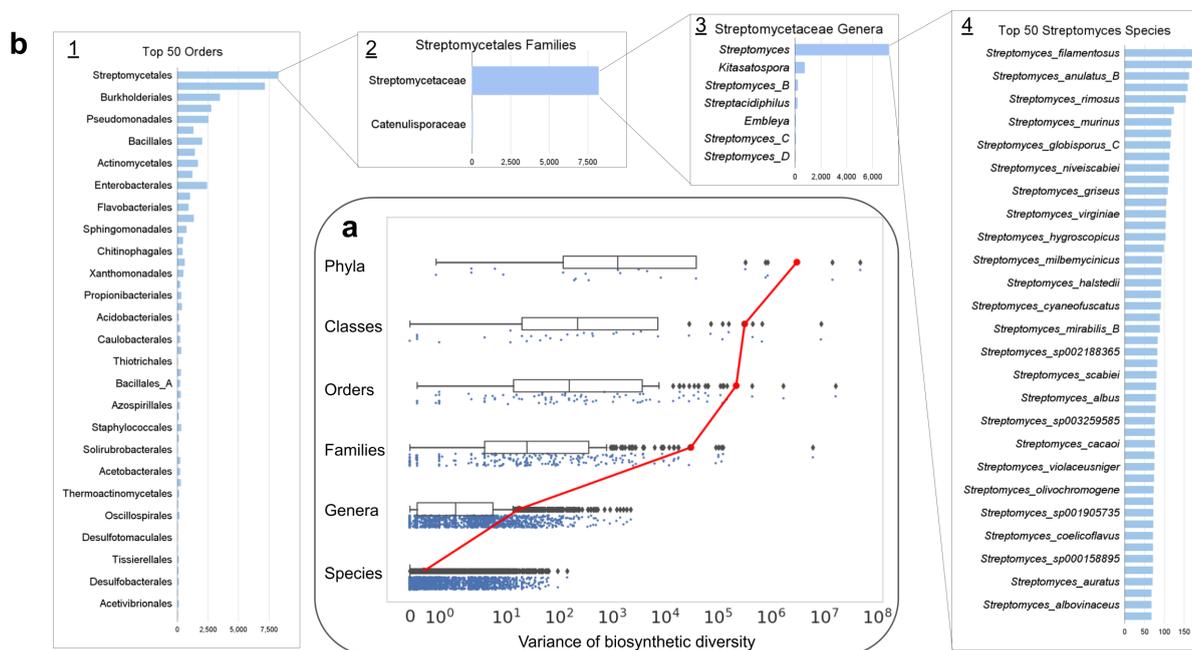


**Figure 2. Comparison of biosynthetic diversity among phyla.** Panel a: The Genome Taxonomy DataBase (GTDB<sup>21</sup>) bacterial tree was visualized with iTOL<sup>37</sup>, decorated with Gene Cluster Families (GCFs) values (as defined by BiG-SLiCE at T=0.4), collapsed at the phylum rank and accompanied by bar plot of GCFs in logarithmic scale ( $10^0$  to  $10^4$ ). The number of genomes belonging to each phylum is displayed next to the tree's leaf nodes. Panel b: GCFs, as defined by BiG-SLiCE (T=0.4), unique to phyla (solid shapes) and with pairwise overlaps between phyla (ribbons), visualized with circlize<sup>38</sup>. Each phylum has a distinct color. Actinobacteriota (2) and Proteobacteria (40) seem particularly rich in unique GCFs.

Once we obtained information on the diversity of different phyla, as well as the rest of the major taxonomic ranks (classes, orders, families, genera, species), we proceeded to determine from which taxonomic rank biosynthetic diversity levels no longer show high variability. After all, phyla are functionally and ecologically very

diverse and the number of GCFs in their underlying species are vastly different. We wanted to find the highest taxonomic rank in which the values of biosynthetic diversity can be considered uniform among the members. Therefore, we conducted a variance analysis that included each taxonomic rank, from phylum to species. For each rank, the variance value was computed based on the #GCFs values of immediately lower-ranked taxa (see Methods for more details). The distribution of these variance values for each rank is visualized in Figure 3a.

There is a noticeable drop in the range of variance values for each rank, while diversity becomes highly homogeneous at the species level (Figures 3a,b). The plunge is most striking from the family to the genus level (Figure 3a), with even the outliers all falling under the  $10^3$ - line in the genus rank. Additional statistical analysis confirmed the significance of this observation (Supplementary Figure 3). We therefore chose to investigate the genus rank in detail. Different genera have previously been studied as sources of unique biosynthetic diversity<sup>16,17</sup>. However, here we determined for the first time that the genus rank is the most appropriate for comparative analyses. Different species within a genus are likely to largely display uniform biosynthetic diversity, while much dissimilarity is observed between different genera belonging to the same family (Figure 3b); the latter point is examined in the following section.



**Figure 3. Relations of taxonomic levels to variability in biosynthetic diversity.** Panel **a**: Modified “raincloud plots”<sup>39</sup> of major taxonomic ranks (X axis in logarithmic scale). Each boxplot represents the dispersion of variance values of a certain taxonomic rank, computed from the number of Gene Cluster Families (GCFs as defined by BiG-SLiCE at T=0.4) of the immediately lower rank. The boxplots’ center line represents the median value; the box limits represent the upper and lower quartiles. Whiskers represent a 1.5x interquartile range. Points outside of the whiskers are outliers. Jittered raw data points are plotted under the boxplots for better visualization of the values’ distribution. The red line connects the mean variance values of each rank. There is a noticeable drop in dispersion of variance values from the family rank to the genus rank (see also Supplementary Figure 3), indicating that the genera are suitable taxonomic groups to be characterised as diverse and be compared to each other. Panel **b**: Biosynthetic diversity of various taxa, measured in absolute numbers of distinct GCFs as defined by BiG-SLiCE (T=0.4) from currently sequenced genomes. Top 50 most diverse orders (1), Streptomycetales families (2), Streptomycetaceae genera (3), top 50 most diverse *Streptomyces* species (4). The difference in variance is visible in the graphs 1,2,3, but becomes homogeneous at the species level as is shown in graph 4.

### *Well-known as well as overlooked taxa as sources of biosynthetic diversity*

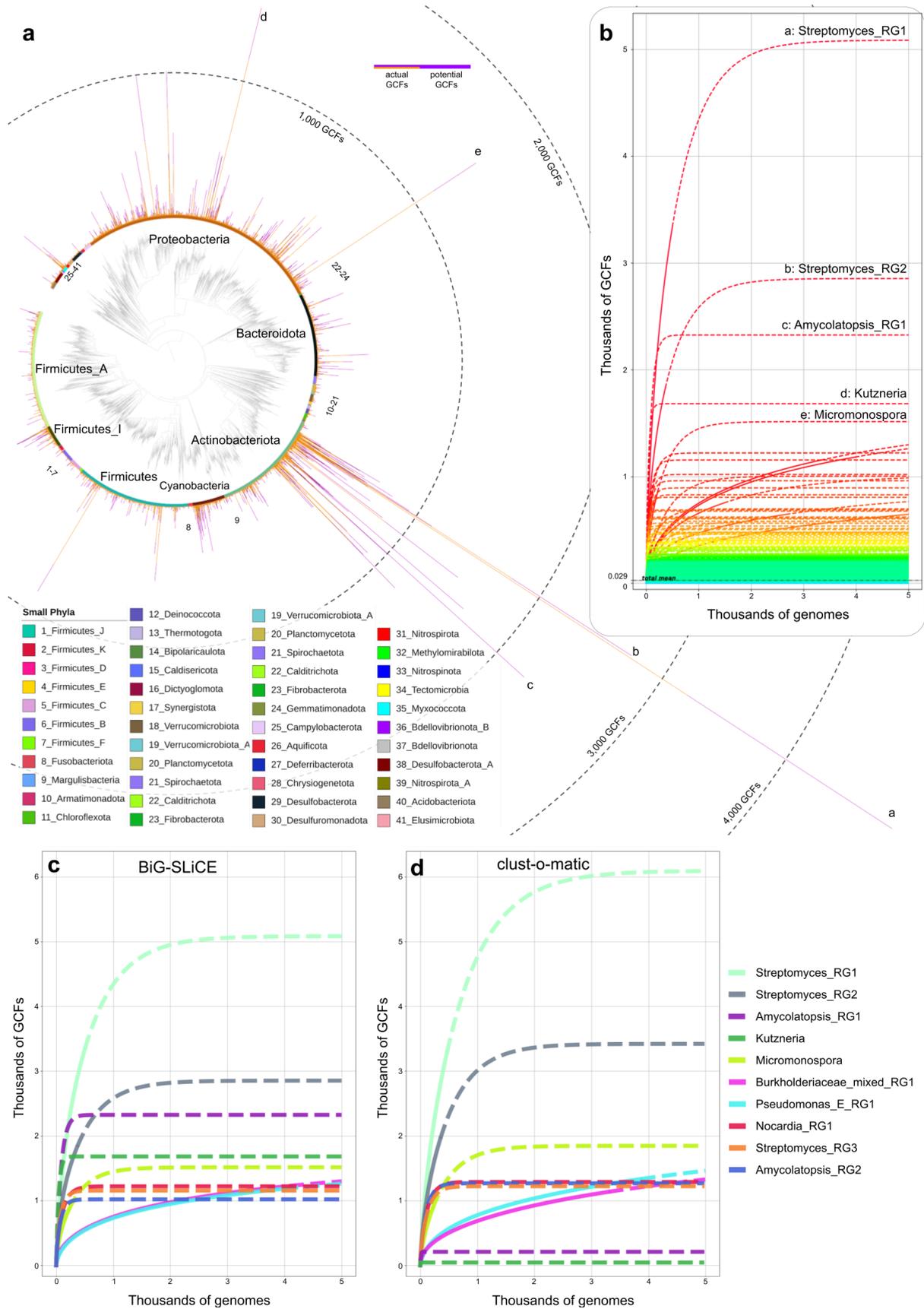
The identification of the genus level as the most informative rank to measure biosynthetic diversity across taxonomy paved the way for a comprehensive comparative analysis of biosynthetic potential across the bacterial tree of life; however, to be able to systematically compare diversity values among groups, said groups need to be uniform. In this case, a common phylogenetic metric was necessary. We chose Relative Evolutionary Divergence (RED) and a specific threshold that was based on the GTDB’s range of RED values for the genus rank<sup>21</sup> to define REDgroups: groups of bacteria analogous to genera but characterized by equal evolutionary distance (see Methods: Definition of REDgroups). Our classification revealed the inequalities in within-taxon phylogenetic similarities among the genera, with some being divided into multiple REDgroups (for example the *Streptomyces* genus was split into 21 REDgroups: *Streptomyces*\_RG1, *Streptomyces*\_RG2 etc.) and some being joined together with other genera to form mixed REDgroups (for example *Burkholderiaceae*\_mixed\_RG1 includes the genera *Paraburkholderia*, *Paraburkholderia*\_A, *Paraburkholderia*\_B, *Burkholderia*, *Paraburkholderia*\_E and *Caballeronia*). This disparity among the genera reaffirmed the importance of defining the REDgroups as a technique that allowed for fair comparisons among bacterial producers.

The resulting 3,779 REDgroups showed huge differences in biosynthetic diversity as measured by the numbers of GCFs found in genomes sequenced from these groups so far, with the maximum diversity at 3,339 GCFs, average at 17 GCFs and minimum at 1 GCF. Nevertheless, the variance of diversity within the REDgroups was even more uniform than in the genera (Supplementary Figure 4). Some of the top groups (Supplementary Table 7) included known rich NP producers, such as *Streptomyces*, *Pseudomonas*\_E and *Nocardia*<sup>28,35,36,40</sup>.

Although very informative, this analysis is biased because of large differences in the number of sequenced strains among the groups, with the economically or medically important strains having been sequenced more systematically than others, e.g., the

biggest Enterobacteriaceae REDgroup includes 45,022 genomes from 33 genera, while most Desulfovibrionaceae REDgroups have a single member. To overcome this bias and allow comparisons across taxa, rarefaction analyses were conducted for each REDgroup, as performed in previous studies<sup>41,42</sup> (Figure 4b). With this information, and in order to provide a global overview of the actual biosynthetic diversity and the potential number of GCFs, we modified and complemented the bacterial tree from Parks *et. al.*<sup>21</sup>, as shown in Figure 4a. The dispersion of these values across the various phyla can also be seen, with the exceptional outliers standing out: Streptomyces\_RG1, Streptomyces\_RG2, Amycolatopsis\_RG1, Kutzneria, and Micromonospora. All these are groups known for their NP producers<sup>16,35,36,43</sup> and they remain in the top (Extended Data Table 1, Supplementary Table 1), seemingly having much unexplored biosynthetic potential.

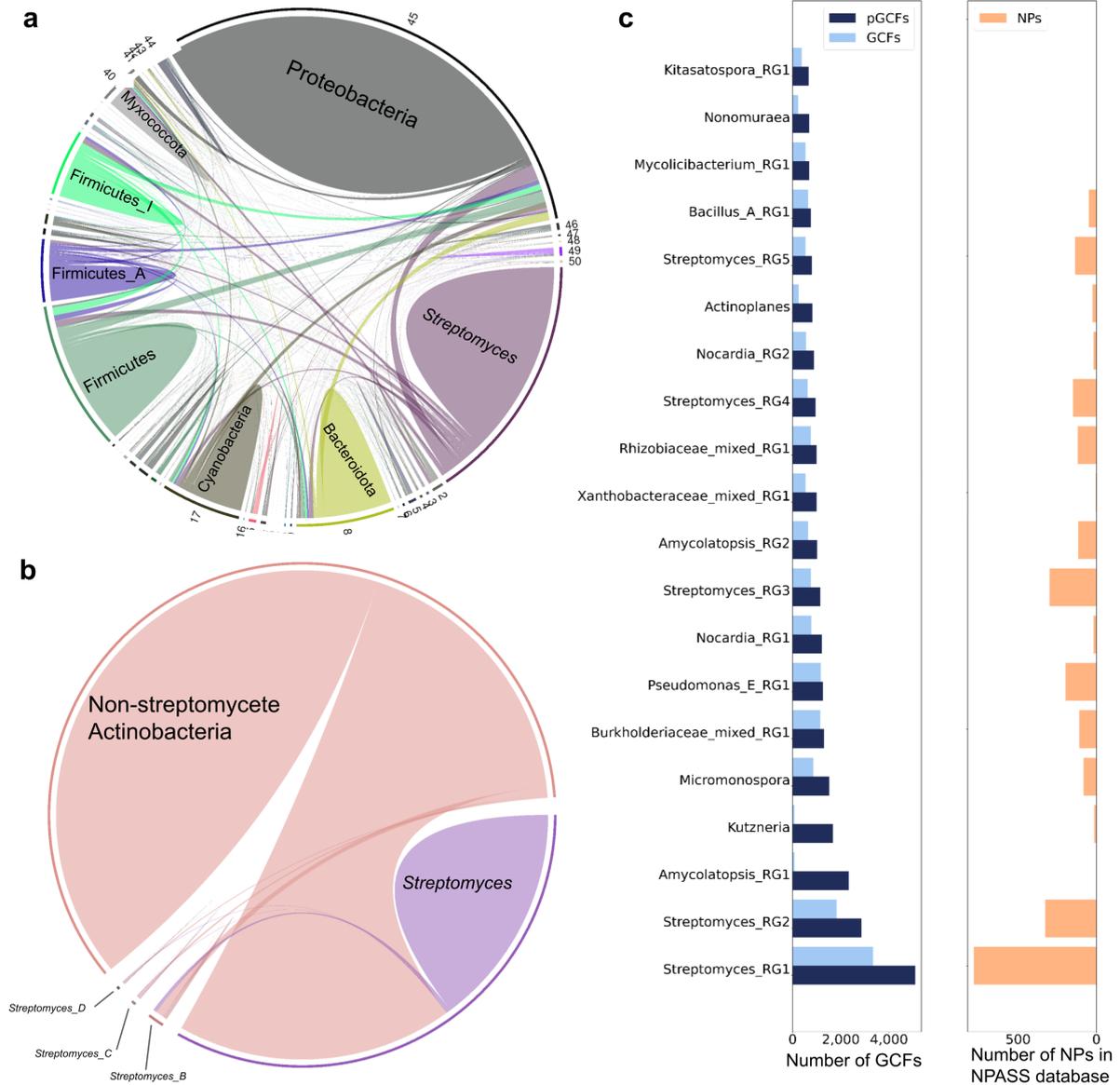
To ensure that our conclusions are not the product of algorithmic artifacts, we decided to rerun the analysis using an alternative method of quantifying biosynthetic diversity, which was developed independently, yet for the same purpose. This alternative approach, called clust-o-matic, is based on a sequence similarity all-versus-all distance matrix of BGCs and subsequent agglomerative hierarchical clustering in order to form GCFs (see Methods). Like for BiG-SLiCE, we calibrated the threshold used for clust-o-matic based on NP Atlas clusters. When comparing the results of the two algorithms (Figure 4c,d, Supplementary Table 7) it is apparent that, although the absolute numbers differ slightly, the trends identified by both methods are very much alike.



**Figure 4. Overview of actual and potential biosynthetic diversity of bacterial kingdom, compared at REDgroup level.** Panel **a**: GTDB<sup>21</sup> bacterial tree up to REDgroup level, visualized with iTOL<sup>37</sup>, colour coded by phylum, decorated with barplots of actual (orange) and potential (purple) Gene Cluster Families (GCFs), as defined by BiG-SLiCE (T=0.4). Top REDgroups with most potential GCFs include the following: A: *Streptomyces*\_RG1, B: *Streptomyces*\_RG2, C: *Amycolatopsis*\_RG1, D: *Kutzneria*, E: *Pseudomonas*\_E. Phyla known to be enriched in NP producers are immediately visible (Actinobacteriota, Protobacteriota), with the most promising groups coming from the Actinobacteriota phylum (the highest peak belongs to a REDgroup containing *Streptomyces* strains). Simultaneously, within the underexplored phyla, there seems to be significant biosynthetic diversity and potential. An interactive version of Figure 4a can be accessed online (Supplementary Figure 5). Panel **b**: Rarefaction curves of REDgroups (BiG-SLiCE T=0.4). The genera included in the most promising REDgroups are indicated (the letters a-e correspond to the peaks in A). *Streptomyces* strains are included in several of the top most promising REDgroups. Panel **c**: Rarefaction curves of the most promising REDgroups (BiG-SLiCE T=0.4). Panel **d**: Rarefaction curves of the most promising REDgroups (clust-o-matic T=0.5). Though the exact numbers differ, the similarities between the two methods are apparent.

*Streptomyces*, even when split into multiple REDgroups, appears in the top groups both based on the known biosynthetic diversity and based on the estimated potential values. It appears to have 5,908 (+103 *Streptomyces*\_B, +39 *Streptomyces*\_C, +16 *Streptomyces*\_D) GCFs that are unique to the group, even among other phyla (Figure 5a). This is in agreement with previous studies investigating how much overlap there is among the main groups of producers<sup>44</sup>. What is more, streptomycetes appear to be the source of a good percentage of the biosynthetic diversity attributed to the Actinobacteria phylum, as seen in Figure 5b.

However, taxa less popular for NP discovery also show promise, as was evident by a comparison of our results with data from the NPASS database of Natural Products<sup>45</sup> (Figure 5c). Among the 20 overall most promising REDgroups we found at least 6 groups that show promise but whose members are either not catalogued in the database as NP sources or are connected too few (<15) known compounds: *Amycolatopsis*\_RG1, *Kutzneria*, *Xanthobacteriaceae*\_mixed\_RG1, *Mycolicibacterium*\_RG1, *Nonomuraea*, *Kitasatospora*\_RG1. The *Amycolatopsis*\_RG1 group only includes three rare species: *Amycolatopsis antarctica*, *marina* and *nigrescens*. Other promising REDgroups with very few known producers include *Cupriavidus* (from Proteobacteria phylum), *Weeksellaceae*\_mixed\_RG1 (from Bacteroidota phylum) and *Pleurocapsa* (from Cyanobacteria phylum). More information about the promising underexplored taxa can be found in (Supplementary Table 7).



**Figure 5. Unique diversity in the known producer *Streptomyces* and promising potential of less popular taxa.** Panel **a**: Unique Gene Cluster Families (GCFs) as defined by BiG-SLiCE (T=0.4), of phyla and *Streptomyces* (solid shapes) and pairwise overlaps of phyla - phyla and phyla - *Streptomyces* (ribbons), visualized with circlize<sup>38</sup>. Each taxon has a distinct color. The smaller shapes and ribbons represent smaller phyla that can be seen in Supplementary Figure 6. The genus *Streptomyces* appears to have a very high amount of unique GCFs comparable to entire phyla, such as Proteobacteria. Panel **b**: Unique GCFs as defined by BiG-SLiCE (T=0.4), of non-streptomycete Actinobacteriota and all *Streptomyces* genera (solid shapes) and pairwise overlaps between Actinobacteriota and *Streptomyces* (ribbons), visualized with circlize<sup>38</sup>. The *Streptomyces* genus, only one of many belonging to the Actinobacteriota phylum, appears to be responsible for a big percentage of the phylum's unique diversity (comparison of solid shapes of group 1 and 5). Panel **c**: Left: Potential (pGCFs) and actual (GCFs) number of Gene Cluster Families as defined by BiG-SLiCE (T=0.4), of top 20 most promising REDgroups. Right: number of Natural Products (NPs) found in the NPASS database<sup>45</sup>, that originate from species included in each REDgroup. Several of these REDgroups appear to have few (< 15) to no known NPs associated with them: Amycolatopsis\_RG1, Kutzneria, Xanthobacteraceae\_mixed\_RG1 (containing the genera *Bradyrhizobium*, *Rhodopseudomonas*, *Tardiphaga* and *Nitrobacter*), Mycolicibacterium\_RG1, Nonomuraea, Kitasatospora\_RG1.

## Discussion

We made use of the large amounts of sequencing data that have become publicly available to identify microbial Biosynthetic Gene Clusters (BGCs) and group them into gene cluster families (GCFs) using two different independent algorithms. The GCF identification was standardized by calibrating the thresholds of both algorithms used to NPAtlas chemical classes, making sure that GCFs mirror chemical families of encoded compounds. We studied the biosynthetic diversity of bacteria and located the emergence of the highest diversity close to the genus rank, hence choosing to further investigate analogous taxonomic groups (REDgroups). Rarefaction analysis was conducted to infer the full biosynthetic potential of the bacterial kingdom as well as to determine the most promising bacterial taxa. Our analysis led to the identification of many known diverse groups as well as multiple promising understudied producers. To the best of our knowledge, this is not only the largest investigation of this sort including all available genomes so far, as well as published Metagenome Assembled Genomes (MAGs) from diverse environments, but it also provides a reproducible pipeline to determine the metabolic diversity of comparable groups of bacteria and propounds a rationale for drug discovery efforts.

The full biosynthetic capacity of the bacterial kingdom has been assessed in a previous study by Cimermanic et al.<sup>14</sup>, already providing a glimpse into the huge untapped biosynthetic potential within bacterial genomes. However, besides their much smaller dataset (about 33,000 BGCs vs 1,094,877 BGCs here), they used ClusterFinder as an identification tool, a more exploratory software algorithm which can include false positive gene clusters. Projects that exploit publicly available genomic data are reliant on the quality of genomes sequenced as well as the efficiency of available genome mining methods, which - though potent as ever - still have some limitations<sup>46</sup>. For instance, the study of GCF uniqueness among taxa may

be affected by antiSMASH's imperfect BGC boundary prediction<sup>47</sup>. Even though BiG-SLICE converts BGCs into features based only on domains related to biosynthesis<sup>27</sup>, genomic context unrelated to the biosynthetic pathway of a BGC could still play a role in the GCF assignment; this issue cannot be fully addressed yet with currently available tools. However, antiSMASH's ability to discern cluster limits and detect BGCs from cultured strains and MAGs is comparable to alternative tools, while its ability to predict so many different BGC types is unparalleled<sup>48</sup>, as is apparent from its common use in Natural Product (NP) research<sup>14,17,34,40,42,49,50</sup>. What is more, the fact that it is rule-based<sup>47</sup> implies the possibility of undetected types of clusters and increases the likelihood that our calculations have underestimated the true biosynthetic potential of bacterial organisms, which could be even more impressive.

Furthermore, our pipeline was the first that used GTDB<sup>21</sup> taxonomy for studying bacterial biosynthetic diversity globally. This enabled us to avoid the established misclassifications of the NCBI taxonomic placement<sup>51–54</sup>. The use of rarefaction curves allowed us to infer the biosynthetic potential of bacterial groups as done in some smaller-scaled projects<sup>14,16,41,42</sup>. This method aims to enable fair comparisons among incomplete samples<sup>55</sup>. However, for those groups that contain very few genomes, there is a tendency to underestimate their potential capacity<sup>55</sup>, so sequencing bias of popular taxa might still affect our results. We tried to minimize the bias within the pipeline as much as possible while retaining high diversity of bacterial taxa; therefore, we decided not to exclude REDgroups with very few members from the dataset. Nonetheless, the remaining bias will only be eliminated when sequencing projects will further aim for inclusion of increased biodiversity<sup>24,56</sup>. Until then, undersequenced taxa with high biosynthetic diversity will inevitably be underestimated (while overestimation is not expected to happen).

Our analysis provided an answer to the main question posed at the beginning of the project: there are a plethora of taxonomic groups that have not yet received adequate attention from researchers, though they seem to be hiding uninvestigated biosynthetic potential, rendering their study essential<sup>10,17,18,30,57–60</sup>. At the same time, the popular targets for NP discovery have still undiscovered biosynthetic diversity. Indicatively, multiple Proteobacteria taxa were identified among the top producers: *Pseudomonas*, *Pseudoalteromonas*, *Paracoccus*, *Serratia* among others. This is in accordance with the known biosynthetic potential of the Proteobacteria phylum<sup>43</sup>. Furthermore, there were phyla that did not have a high amount of genomic data publicly available and still included some promising groups, like the myxobacterial genera *Cystobacter*, *Melittangium*, *Archangium*, *Vitiosangium*, *Sorangium* and *Myxococcus*<sup>17,40,61</sup>, or the genera *Chryseobacterium* and *Chryseobacterium\_A*<sup>62–64</sup> from the Bacteroidota phylum. However, the majority of most diverse groups comprise actinobacterial strains of well-known and well-studied NP producers from genera such as *Actinoplanes*, *Amycolatopsis*, *Micromonospora*, *Mycobacterium*, *Nocardia* and *Streptomyces*<sup>16,35,36,58,65</sup>. These bacteria produce the majority of Natural

Products used as antibiotics<sup>35,58</sup> and our analysis confirms the notion that there is still much more natural product diversity to be discovered within this group<sup>35,36</sup>. One of the first comprehensive studies of the biosynthetic diversity of the Actinobacteria phylum<sup>16</sup> provided a roadmap to drug discovery for this taxon. However, at the time their analysis indicated that all biosynthetic diversity encoded in Actinobacteria would be identified when 15,000 varied genomes had sequenced<sup>16</sup>. Our data indicate that although the raw number has already been reached (we analyzed 15,857 Actinobacteria genomes from 2,505 different species), more biosynthetic diversity is expected from this taxon as more diversified strains get sequenced. At the same time, for specific purposes like antibiotic discovery, for which many large-scale bioactivity-based screens have already been performed on actinobacterial strain collections<sup>18</sup>, it may make sense to aim at a diversified portfolio of taxa that also include other high-potential bacterial clades.

*Streptomyces*, which has been showcased in our investigation, is a genus of the Actinobacteria phylum that contains some of the most complex bacteria that we know of, though by far not the most sequenced in our dataset (Supplementary Figure 8). These bacteria have been thoroughly studied for multiple reasons, such as their abundance and important role in soil environments and of course their reputation for being competent NP producers<sup>35,44</sup>. The biosynthetic diversity displayed by members of this genus and their hidden potential, which was confirmed in the present study, have been noted for decades now; for some, the genus seems to be an inexhaustible source<sup>65</sup>. The factors that cause this taxonomic group to stand out are not completely clear but probably related to their sophisticated lifestyle. Many observations suggest that NP biosynthesis drives speciation within the *Streptomyces* genus<sup>16</sup>. The exploration of factors that led to the rise of biosynthetic diversity in *Streptomyces* to such an impressive degree will be the subject of further investigations in the future.

Having the genomic capacity for the biosynthesis of secondary metabolites does not always herald the discovery of a new compound though<sup>66,67</sup>. Sometimes, the bacterium in question cannot be grown in laboratory conditions, a problem that is usually approached with attempts to heterologously express BGCs in another host<sup>35,59,66,67</sup>. But very often the conditions in which BGC hosts, native or heterologous, are grown do not resemble the natural circumstances of NP production and therefore its biosynthesis does not take place<sup>60,66</sup>. This issue is related to the complexity of BGCs, whose intricate regulation and connection to the primary metabolism<sup>12,59,66</sup> are not yet fully understood to overcome these obstacles<sup>68</sup>. However, efforts to decode biosynthetic mechanisms for the activation of silent clusters need to be tailored to specific producer groups<sup>35,36,67,69</sup> and the present study can help with this selection.

Original approaches to the prioritization issue of NP research continue to emerge, as indicated by novel strategies, like the study of synergistic interactions of bioactive compounds and transcriptomics analyses<sup>70</sup>, the search for producers among

archaea and fungi<sup>71</sup> and by advances in metagenomics. New computational tools are constantly being developed to make use of the biosynthetic potential of unculturable bacteria from environmental samples. Furthermore, apart from the few metagenomic projects whose MAGs we incorporated in the first part of our analysis, there are multiple such projects publicly available, some of which have been the focus of NP studies<sup>72,73</sup>. Various different habitats have been studied and new compounds have been revealed. Metagenomics is proving a promising source of information on NPs and their producers<sup>14,29,44,59,70,72,74</sup>, as made apparent in the present investigation, and we expect the effect of this field on NP research to become more evident in the following years.

The collection of microbial data from increasingly new habitats points to another interesting aspect, namely the relation between the environmental niche of the producers and the uniqueness of their biosynthetic diversity. Past studies have examined this connection to a limited extent<sup>34,42,43,60</sup> but it would be even more interesting to deduct conclusions based on a wider-scale dataset. This type of analysis will only be possible when much more detailed and standardized annotations of metadata of producers' genomes will become available.

Our analysis provides a global overview of diverse known and promising understudied NP-producing taxa. We expect this to greatly help overcome one of the main bottlenecks of antibiotic discovery: the prioritization of producers for research<sup>70</sup>.

## Methods

### *BGC data set*

We obtained 170,585 complete and draft bacterial genomes (Table 1) from RefSeq<sup>20</sup> on 27 March 2020. Furthermore, a dataset of 11,143 MAGs was included in the first part of the analysis (see Results: Biosynthetic diversity of the bacterial kingdom). For the rest of the study, we used only 161,290 RefSeq bacterial genomes whose taxonomic classification up to the species level was known (Table 1). All genomes were analyzed with antiSMASH (version 5)<sup>47</sup>, which identified their BGCs (Extended Data Table 1, Supplementary Table 1).

### *Taxonomic classification*

Due to multiple indications regarding a lack of accuracy of NCBI's taxonomic classification of bacterial genomes<sup>51-54</sup>, we chose to use the Genome Taxonomy Database (GTDB<sup>21</sup>) instead. The bacterial tree of 120 concatenated proteins (GTDB release 89), as well as the classifications of organisms up to the species level, were included in the analysis.

### *Quantification of biosynthetic diversity with BiG-SLiCE*

For a bacterium to be considered biosynthetically diverse, we considered not the number of BGCs important, but rather how different these BGCs are to each other. In order to quantify this diversity, we analyzed all BGCs with the new BiG-SLiCE tool<sup>27</sup>, which groups similar clusters into Gene Cluster Families (GCFs). However, the first version of this tool have an inherent bias towards multi-protein families BGCs, producing uneven coverage between BGCs of different classes (i.e., due to their lack of biosynthetic domain diversity, all lanthipeptide BGCs may be grouped together using the Euclidean threshold of  $T=900$ , which in contrast is ideal for clustering Type-I Polyketide BGCs). To alleviate this issue and provide a fair measurement of biosynthetic diversity between the taxa, we modified the original distance measurement by normalizing the BGC features under  $L^2$ -norm, which will produce a cosine-like distance when processed by the Euclidean-based BIRCH algorithm. This usage of cosine-like distance will virtually balance the measured distance between BGCs with “high” and “low” feature counts (Supplementary Figure 7A), in the end providing an improved clustering performance when measured using the reference data of manually-curated MIBiG GCFs (Supplementary Figure 7B).

The GTDB<sup>21</sup> (release 89) bacterial tree was pruned so that it included only the organisms that are part of our dataset. Then, having both the taxonomic classification of all bacteria, as well as how many GCFs their BGCs group into, the pruned GTDB tree was decorated with #GCFs values at each node. This allowed for the evaluation of the biosynthetic diversity of any clade, including the main taxonomic ranks. To pick a single threshold for subsequent taxonomy richness analysis, we compared BiG-SLiCE results on 947 MIBiG BGCs versus the compound-based clustering provided by the NPAtlas database<sup>19</sup> (Supplementary Figure 2). A final threshold of  $T=0.4$  was chosen based on its similarity to NPAtlas's compound clusters (V-score=0.9X, GCF counts difference=+XX).

### *Quantification of biosynthetic diversity with clust-o-matic*

We aimed to repeat and evaluate the reproducibility of the BGC-to-GCF quantification step of BiG-SLiCE with an alternative, independently derived algorithm. For that instead of grouping BGCs into GCFs based on biosynthetic domain diversity, we developed an algorithm that considers full core biosynthetic genes. Biosynthetic gene clusters that were detected in the input data by antiSMASH 5.1 were parsed to deliver core biosynthetic protein sequences. Those protein sequences were subjected to all-against-all multi-gene sequence similarity search with DIAMOND 2.0 using default settings. Only one best hit per query core gene per BGC was allowed divided by a total core protein length, resulting in the final pairwise BGC score always being within range of 0 to 1. Pairwise BGC similarity scores were used to build a distance matrix that was later subjected to agglomerative hierarchical clustering in python programming language (package `scipy.cluster.hierarchy`). The determined optimal threshold (see paragraph above) of 0.5 was then used to

generate GCFs, which were then fed into the next steps in parallel to the original set of GCFs obtained from BiG-SLiCE.

### *Variance Analysis*

In order to pinpoint the emergence of biosynthetic diversity, the within-taxon homogeneity was compared among the main taxonomic ranks. For each rank, the variance value was computed (with NumPy<sup>75</sup>) based on the #GCFs values of immediately lower-ranked taxa, as long as there were at least two such taxa. For example, a phylum that includes only one class in our dataset was omitted from this computation. But a phylum with two or more classes would be assigned a variance value computed from its classes' #GCFs values. The distribution of these variance values was plotted for each rank in Figure 3a. We noticed a significant reduction in variance from the family to the genus rank, which was confirmed with an additional statistical test (Supplementary Figure 3). A similar variance analysis was performed to compare genera and REDgroups (Supplementary Figure 4) but in this case variance was calculated based on the strains' biosynthetic diversity.

### *Definition of REDgroups*

To study the biosynthetic diversity of genera, we attempted to achieve uniform taxa. The creators of GTDB used Relative Evolutionary Divergence (RED) for taxonomic rank normalization<sup>21</sup>; it is a metric that relies heavily on the branch length of a phylogenetic tree and is consequently dependent on the rooting. The GTDB developers provided us with a bacterial tree decorated with the average RED values of all plausible rootings at each node. Since GTDB accepts a range of RED values for each taxonomic rank placement<sup>21</sup>, we chose the median of GTDB genus RED values, namely 0.934, as a cutoff threshold. Any clade in the GTDB bacterial tree with an assigned RED value higher than the threshold was considered one group (Supplementary Figure 9) that we named "REDgroup". For REDgroup naming conventions, see Supplementary Figure 9.

### *Rarefaction analysis*

The extrapolation of potential #GCFs values was achieved by conducting rarefaction analyses, by use of the iNEXT R package<sup>76</sup>. A GCF presence/absence table (GCF-by-strain matrix) was constructed for each group considered and was then used as "incidence-raw" data in the iNEXT main function, where 500 points were inter- or extrapolated with an endpoint of 5000 for the REDgroups, and of 8 times the number of strains in each group for the RefSeq analyses (where 2000 points were inter- or extrapolated). By default, the number of bootstrap replications is 50.

### *Identification of unknown producers*

We investigated the genera included in the most promising REDgroups, to find out whether they include species that are producers of known compounds. Hence, the species names were cross-referenced with the species named as producers in the

NPASS depository<sup>45</sup> (accessed on 15 October 2020), taking care to match the GTDB-given names to the NCBI-given names that the database uses.

## References

1. O'Connor, S. E. Engineering of Secondary Metabolism. *Annu. Rev. Genet.* **49**, 71–94 (2015).
2. Newman, D. J. & Cragg, G. M. Natural Products as Sources of New Drugs over the Nearly Four Decades from 01/1981 to 09/2019. *J. Nat. Prod.* **83**, 770–803 (2020).
3. Atanasov, A. G., Zotchev, S. B., Dirsch, V. M., International Natural Product Sciences Taskforce & Supuran, C. T. Natural products in drug discovery: advances and opportunities. *Nat. Rev. Drug Discov.* **20**, 200–216 (2021).
4. Hutchings, M. I., Truman, A. W. & Wilkinson, B. Antibiotics: past, present and future. *Current Opinion in Microbiology* vol. 51 72–80 (2019).
5. Brown, E. D. & Wright, G. D. Antibacterial drug discovery in the resistance era. *Nature* **529**, 336–343 (2016).
6. Wright, G. D. Unlocking the potential of natural products in drug discovery. *Microb. Biotechnol.* **12**, 55–57 (2019).
7. Fedorenko, V. *et al.* Antibacterial Discovery and Development: From Gene to Product and Back. *Biomed Res. Int.* **2015**, 591349 (2015).
8. Michael, C. A., Dominey-Howes, D. & Labbate, M. The antimicrobial resistance crisis: causes, consequences, and management. *Front Public Health* **2**, 145 (2014).
9. Lyddiard, D., Jones, G. L. & Greatrex, B. W. Keeping it simple: lessons from the golden era of antibiotic discovery. *FEMS Microbiol. Lett.* **363**, (2016).
10. Pye, C. R., Bertin, M. J., Lokey, R. S., Gerwick, W. H. & Linington, R. G. Retrospective analysis of natural products provides insights for future discovery trends. *Proc. Natl. Acad. Sci. U. S. A.* **114**, 5601–5606 (2017).
11. Katz, L. & Baltz, R. H. Natural product discovery: past, present, and future. *J. Ind. Microbiol. Biotechnol.* **43**, 155–176 (2016).

12. Ziemert, N., Alanjary, M. & Weber, T. The evolution of genome mining in microbes - a review. *Nat. Prod. Rep.* **33**, 988–1005 (2016).
13. Medema, M. H., de Rond, T. & Moore, B. S. Mining genomes to illuminate the specialized chemistry of life. *Nat. Rev. Genet.* (2021) doi:10.1038/s41576-021-00363-7.
14. Cimermancic, P. *et al.* Insights into secondary metabolism from a global analysis of prokaryotic biosynthetic gene clusters. *Cell* **158**, 412–421 (2014).
15. Chevrette, M. G. *et al.* Evolutionary dynamics of natural product biosynthesis in bacteria. *Nat. Prod. Rep.* **37**, 566–599 (2020).
16. Doroghazi, J. R. *et al.* A roadmap for natural product discovery based on large-scale genomics and metabolomics. *Nat. Chem. Biol.* **10**, 963–968 (2014).
17. Hoffmann, T. *et al.* Correlating chemical diversity with taxonomic distance for discovery of natural products in myxobacteria. *Nat. Commun.* **9**, 803 (2018).
18. Lewis, K. The Science of Antibiotic Discovery. *Cell* vol. 181 29–45 (2020).
19. van Santen, J. A. *et al.* The Natural Products Atlas: An Open Access Knowledge Base for Microbial Natural Products Discovery. *ACS Cent Sci* **5**, 1824–1833 (2019).
20. Haft, D. H. *et al.* RefSeq: an update on prokaryotic genome annotation and curation. *Nucleic Acids Res.* **46**, D851–D860 (2018).
21. Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of life. *Nat. Biotechnol.* **36**, 996–1004 (2018).
22. Stewart, R. D. *et al.* Compendium of 4,941 rumen metagenome-assembled genomes for rumen microbiome biology and enzyme discovery. *Nat. Biotechnol.* **37**, 953–961 (2019).
23. Glendinning, L., Stewart, R. D., Pallen, M. J., Watson, K. A. & Watson, M. Assembly of hundreds of novel bacterial genomes from the chicken caecum. *Genome Biol.* **21**, 34 (2020).
24. Almeida, A. *et al.* A unified sequence catalogue of over 280,000 genomes obtained from the human gut microbiome. doi:10.1101/762682.
25. Tully, B. J., Graham, E. D. & Heidelberg, J. F. The reconstruction of 2,631 draft metagenome-assembled genomes from the global oceans. *Sci Data* **5**, 170203 (2018).

26. Parks, D. H. *et al.* Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life. *Nat Microbiol* **2**, 1533–1542 (2017).
27. Kautsar, S. A., van der Hooff, J. J. J., de Ridder, D. & Medema, M. H. BiG-SLiCE: A Highly Scalable Tool Maps the Diversity of 1.2 Million Biosynthetic Gene Clusters. *Bioinformatics* (2020).
28. Kautsar, S. A. *et al.* MIBiG 2.0: a repository for biosynthetic gene clusters of known function. *Nucleic Acids Res.* **48**, D454–D458 (2020).
29. Land, M. *et al.* Insights from 20 years of bacterial genome sequencing. *Funct. Integr. Genomics* **15**, 141–161 (2015).
30. Ling, L. L. *et al.* A new antibiotic kills pathogens without detectable resistance. *Nature* **517**, 455–459 (2015).
31. Chase, A. B., Sweeney, D., Muskat, M. N., Guillén-Matus, D. & Jensen, P. R. Vertical inheritance governs biosynthetic gene cluster evolution and chemical diversification. doi:10.1101/2020.12.19.423547.
32. Medema, M. H., Cimermancic, P., Sali, A., Takano, E. & Fischbach, M. A. A systematic computational analysis of biosynthetic gene cluster evolution: lessons for engineering biosynthesis. *PLoS Comput. Biol.* **10**, e1004016 (2014).
33. Undabarrena, A. *et al.* Rhodococcus comparative genomics reveals a phylogenomic-dependent non-ribosomal peptide synthetase distribution: insights into biosynthetic gene cluster connection to an orphan metabolite. *Microbial Genomics* vol. 7 (2021).
34. Sharrar, A. M. *et al.* Bacterial Secondary Metabolite Biosynthetic Potential in Soil Varies with Phylum, Depth, and Vegetation Type. *MBio* **11**, (2020).
35. Barka, E. A. *et al.* Taxonomy, Physiology, and Natural Products of Actinobacteria. *Microbiol. Mol. Biol. Rev.* **80**, 1–43 (2016).
36. Genilloud, O. Actinomycetes: still a source of novel antibiotics. *Nat. Prod. Rep.* **34**, 1203–1232 (2017).
37. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic

- tree display and annotation. *Nucleic Acids Res.* **49**, W293–W296 (2021).
38. Gu, Z., Gu, L., Eils, R., Schlesner, M. & Brors, B. circlize implements and enhances circular visualization in R. *Bioinformatics* vol. 30 2811–2812 (2014).
  39. Allen, M., Poggiali, D., Whitaker, K., Marshall, T. R. & Kievit, R. A. Raincloud plots: a multi-platform tool for robust data visualization. *Wellcome Open Res* **4**, 63 (2019).
  40. Männle, D. *et al.* Comparative Genomics and Metabolomics in the Genus *Nocardia*. *mSystems* **5**, (2020).
  41. Ziemert, N. *et al.* Diversity and evolution of secondary metabolism in the marine actinomycete genus *Salinispora*. *Proc. Natl. Acad. Sci. U. S. A.* **111**, E1130–9 (2014).
  42. Adamek, M. *et al.* Comparative genomics reveals phylogenetic distribution patterns of secondary metabolites in *Amycolatopsis* species. *BMC Genomics* **19**, 426 (2018).
  43. Buijs, Y. *et al.* Marine Proteobacteria as a source of natural products: advances in molecular tools and strategies. *Nat. Prod. Rep.* **36**, 1333–1350 (2019).
  44. Bérdy, J. Bioactive microbial metabolites. *J. Antibiot.* **58**, 1–26 (2005).
  45. Zeng, X. *et al.* NPASS: natural product activity and species source database for natural product research, discovery and tool development. *Nucleic Acids Res.* **46**, D1217–D1222 (2018).
  46. Miller, M. E. *et al.* Increased virulence of *Puccinia coronata* f. sp. *avenae* populations through allele frequency changes at multiple putative Avr loci. *PLoS Genet.* **16**, e1009291 (2020).
  47. Blin, K. *et al.* antiSMASH 5.0: updates to the secondary metabolite genome mining pipeline. *Nucleic Acids Res.* **47**, W81–W87 (2019).
  48. Chavali, A. K. & Rhee, S. Y. Bioinformatics tools for the identification of gene clusters that biosynthesize specialized metabolites. *Brief. Bioinform.* **19**, 1022–1034 (2018).
  49. Ueoka, R. *et al.* Genome-Based Identification of a Plant-Associated Marine Bacterium as a Rich Natural Product Source. *Angew. Chem. Int. Ed Engl.* **57**, 14519–14523 (2018).
  50. Adamek, M., Alanjary, M. & Ziemert, N. Applied evolution: Phylogeny-based approaches

- in natural products research. *Natural Product Reports* vol. 36 1295–1312 (2019).
51. Ciufu, S. *et al.* Using average nucleotide identity to improve taxonomic assignments in prokaryotic genomes at the NCBI. *Int. J. Syst. Evol. Microbiol.* **68**, 2386–2392 (2018).
  52. Martínez-Romero, E. *et al.* Genome misclassification of *Klebsiella variicola* and *Klebsiella quasipneumoniae* isolated from plants, animals and humans. *Salud Publica Mex.* **60**, 56–62 (2018).
  53. Mateo-Estrada, V., Graña-Miraglia, L., López-Leal, G. & Castillo-Ramírez, S. Phylogenomics Reveals Clear Cases of Misclassification and Genus-Wide Phylogenetic Markers for *Acinetobacter*. *Genome Biol. Evol.* **11**, 2531–2541 (2019).
  54. Rekadwad, B. N. & Gonzalez, J. M. Correcting names of bacteria deposited in National Microbial Repositories: an analysed sequence data necessary for taxonomic re-categorization of misclassified bacteria-ONE example, genus *Lysinibacillus*. *Data Brief* **13**, 761–778 (2017).
  55. Chao, A. *et al.* *Rarefaction and extrapolation with Hill numbers: a framework for sampling and estimation in species diversity studies.* vol. 84 45–67 <http://purl.oclc.org/estimates> (2014).
  56. Nayfach, S. *et al.* Author Correction: A genomic catalog of Earth's microbiomes. *Nat. Biotechnol.* **39**, 521 (2021).
  57. Zheng, Y. *et al.* Genome Features and Secondary Metabolites Biosynthetic Potential of the Class Ktedonobacteria. *Front. Microbiol.* **10**, 893 (2019).
  58. Amin, D. H., Abdallah, N. A., Abolmaaty, A., Tolba, S. & Wellington, E. M. H. Microbiological and molecular insights on rare Actinobacteria harboring bioactive prospective. *Bulletin of the National Research Centre* **44**, 1–12 (2020).
  59. Hug, J. J., Bader, C. D., Remškar, M., Cirnski, K. & Müller, R. Concepts and Methods to Access Novel Antibiotics from Actinomycetes. *Antibiotics (Basel)* **7**, (2018).
  60. Subramani, R. & Sipkema, D. Marine rare actinomycetes: A promising source of structurally diverse and unique novel natural products. *Marine Drugs* vol. 17 (2019).
  61. Weissman, K. J. & Müller, R. Myxobacterial secondary metabolites: bioactivities and

- modes-of-action. *Nat. Prod. Rep.* **27**, 1276–1295 (2010).
62. Cimmino, T. & Rolain, J.-M. Whole genome sequencing for deciphering the resistome of *Chryseobacterium indologenes*, an emerging multidrug-resistant bacterium isolated from a cystic fibrosis patient in Marseille, France. *New Microbes New Infect* **12**, 35–42 (2016).
63. Kang, D., Shoaie, S., Jacquioid, S., Sørensen, S. J. & Ledesma-Amaro, R. Comparative genomics analysis of *Chryseobacterium* sp. KMC2 reveals metabolic pathways involved in keratinous utilization and natural product biosynthesis. *bioRxiv* (2021)  
doi:10.1101/2021.02.23.432615.
64. Dahal, R. H., Chaudhary, D. K., Kim, D.-U., Pandey, R. P. & Kim, J. *Chryseobacterium antibioticum* sp. nov. with antimicrobial activity against Gram-negative bacteria, isolated from Arctic soil. *J. Antibiot.* **74**, 115–123 (2021).
65. Berdi, J. Bioactive Microbial Metabolites: A Personal View. *Journal of Antibiotics*. *Antibiotics* **58**, 1–26 (2005).
66. Seyedsayamdost, M. R. Toward a global picture of bacterial secondary metabolism. *Journal of Industrial Microbiology and Biotechnology* vol. 46 301–311 (2019).
67. Wohlleben, W., Mast, Y., Stegmann, E. & Ziemert, N. Antibiotic drug discovery. *Microb. Biotechnol.* **9**, 541–548 (2016).
68. van Bergeijk, D. A., Terlouw, B. R., Medema, M. H. & van Wezel, G. P. Ecology and genomics of Actinobacteria: new concepts for natural product discovery. *Nat. Rev. Microbiol.* **18**, 546–558 (2020).
69. Kalkreuter, E., Pan, G., Cepeda, A. J. & Shen, B. Targeting Bacterial Genomes for Natural Product Discovery. *Trends Pharmacol. Sci.* **41**, 13–26 (2020).
70. Tracanna, V., de Jong, A., Medema, M. H. & Kuipers, O. P. Mining prokaryotes for antimicrobial compounds: From diversity to function. *FEMS Microbiology Reviews* vol. 41 417–429 (2017).
71. Gluck-Thaler, E. *et al.* The Architecture of Metabolism Maximizes Biosynthetic Diversity in the Largest Class of Fungi. *Mol. Biol. Evol.* **37**, 2838–2856 (2020).

72. Chen, R. *et al.* Discovery of an Abundance of Biosynthetic Gene Clusters in Shark Bay Microbial Mats. *Front. Microbiol.* **11**, 1950 (2020).
73. Milshteyn, A., Colosimo, D. A. & Brady, S. F. Accessing Bioactive Natural Products from the Human Microbiome. *Cell Host Microbe* **23**, 725–736 (2018).
74. Shen, B. A New Golden Age of Natural Products Drug Discovery. *Cell* **163**, 1297–1300 (2015).
75. Harris, C. R. *et al.* Array programming with NumPy. *Nature* **585**, 357–362 (2020).
76. Hsieh, T. C., Ma, K. H. & Chao, A. iNEXT: an R package for rarefaction and extrapolation of species diversity ( Hill numbers). *Methods Ecol. Evol.* **7**, 1451–1456 (2016).

## Fundings

A.G. is grateful for the support of the Deutsche Forschungsgemeinschaft (DFG; Project ID # 398967434-TRR 261). N.Zi. is supported by the German Center for Infection Research (DZIF) (TTU 09.716). M.H.M. is supported by an European Research Council Starting Grant 948770-DECIPHER. Work in the lab of R.M. is supported by BMBF, DFG and DZIF.

## Acknowledgements

A.G. and N.Zi. thank the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy – EXC 2124 – 390838134 for the infrastructural support. A.G. thanks Direnc M. Mungan for valued discussions on optimizing the analysis. We also thank Dr. Libera do Presti for helpful comments on the manuscript.

## Author contributions

A.G., S.A.K., N.Za. and D.K. have performed the analysis. S.A.K. and N.Za. have contributed analysis tools. A.G., D.K., R.M., M.H.M. and N.Zi. have written the paper. All authors have contributed to the conception and design of the analysis. All authors have read and agreed to the published version of the manuscript.

## Competing interests

M.H.M. is a co-founder of Design Pharmaceuticals and a member of the scientific advisory board of Hexagon Bio. The other authors declare no competing interests.

## Extended Data

### Extended Data Table 1

Information on the most biosynthetically promising REDgroups (BiG-SLiCE T=0.4)

## Supplementary information

### Supplementary Figures

#### Supplementary Table 1

Accession numbers, GTDB-based taxonomic information and BGC IDs of all genomes from all datasets used in the analysis

#### Supplementary Table 2

BGC to BiG-SLiCE GCF assignment and centroid distance for T=0.4 (which proved to be the most suitable threshold)

#### Supplementary Table 3

BGC to BiG-SLiCE GCF assignment and centroid distance for T=0.5

#### Supplementary Table 4

BGC to BiG-SLiCE GCF assignment and centroid distance for T=0.6

#### Supplementary Table 5

BGC to BiG-SLiCE GCF assignment and centroid distance for T=0.7

#### Supplementary Table 6

BGC IDs, MiBIG IDs and producer GTDB-based taxonomic information for all BGCs included in the creation of Figure 1C.

#### Supplementary Table 7

REDgroup full metadata: Node IDs (can be used in the exploration of the tree in Supplementary Figure 5), labels, number of members, number of BGCs, number of GCFs and potential GCFs (pGCFs) as defined by BiG-SLiCE (T=0.4) and clust-o-matic (T=0.5), GTDB taxonomic information and number of products in the NPASS database whose producer is a member of the REDgroup (NPASS\_hits).

(Note to reviewers: These data in Supplementary Tables 1 to 5, which are especially large files, would be publicly released upon completion of peer review of this manuscript under the following link: <https://doi.org/10.5281/zenodo.5159211>)

## Code availability

All custom code will be made available once the paper is accepted in principle for publication in Nature Microbiology.

## Figure Legends

**Figure 1. Biosynthetic diversity of the bacterial kingdom.** Panel **a**: Bar plots of Gene Cluster Families (GCFs, as defined by BiG-SLiCE) of nine most biosynthetically diverse genera using different thresholds (T). The absolute number of GCFs changes from threshold to threshold, but the general tendencies (highest to lowest GCF count) are consistent between them. Panel **b**: Rarefaction curves of all RefSeq bacteria based on BiG-SLiCE (red) and based on clust-o-matic (orange), and rarefaction curve of the Entire Bacterial kingdom dataset, which includes bacterial MAGs (blue), based on BiG-SLiCE. BiG-SLiCE GCFs were calculated with T=0.4. Clust-o-matic GCFs were calculated with T=0.5. The number of chemical classes documented in NPAtlas<sup>19</sup>, which come from bacterial producers (gray dotted line), corresponds to 2.9% - 3.3% of the predicted potential of the bacterial kingdom. Panel **c**: Venn Diagram of GCFs (as defined by BiG-SLiCE, T=0.4) of the bacterial RefSeq, Minimum Information about a Biosynthetic Gene cluster (MiBiG<sup>28</sup>) and bacterial MAGs datasets. More information on the MiBiG dataset can be found in Supplementary Table 6. About 36.2% of the GCFs of MAGs are unique (blue shape) to this dataset.

**Figure 2. Comparison of biosynthetic diversity among phyla.** Panel **a**: The Genome Taxonomy DataBase (GTDB<sup>21</sup>) bacterial tree was visualized with iTOL<sup>37</sup>, decorated with Gene Cluster Families (GCFs) values (as defined by BiG-SLiCE at T=0.4), collapsed at the phylum rank and accompanied by bar plot of GCFs in logarithmic scale ( $10^0$  to  $10^4$ ). The number of genomes belonging to each phylum is displayed next to the tree's leaf nodes. Panel **b**: GCFs, as defined by BiG-SLiCE (T=0.4), unique to phyla (solid shapes) and with pairwise overlaps between phyla (ribbons), visualized with circlize<sup>38</sup>. Each phylum has a distinct color. Actinobacteriota (2) and Proteobacteria (40) seem particularly rich in unique GCFs.

**Figure 3. Relations of taxonomic levels to variability in biosynthetic diversity.** Panel **a**: Modified "raincloud plots"<sup>39</sup> of major taxonomic ranks (X axis in logarithmic scale). Each boxplot represents the dispersion of variance values of a certain taxonomic rank, computed from the number of Gene Cluster Families (GCFs as defined by BiG-SLiCE at T=0.4) of the immediately lower rank. The boxplots' center line represents the median value; the box limits represent the upper and lower quartiles. Whiskers represent a 1.5x interquartile range. Points outside of the whiskers are outliers. Jittered raw data points are plotted under the boxplots for better visualization of the values' distribution. The red line connects the mean variance values of each rank. There is a noticeable drop in dispersion of variance values from the family rank to the genus rank (see also Supplementary Figure 3), indicating that the genera are suitable taxonomic groups to be characterised as diverse and be compared to each other. Panel **b**: Biosynthetic diversity of various taxa, measured in absolute numbers of distinct GCFs as defined by BiG-SLiCE (T=0.4) from currently sequenced genomes. Top 50 most diverse orders (1), Streptomycetales families (2), Streptomycetaceae genera (3), top 50 most diverse

*Streptomyces* species (4). The difference in variance is visible in the graphs 1,2,3, but becomes homogeneous at the species level as is shown in graph 4.

**Figure 4. Overview of actual and potential biosynthetic diversity of bacterial kingdom, compared at REDgroup level.** Panel a: GTDB<sup>21</sup> bacterial tree up to REDgroup level, visualized with iTOL<sup>37</sup>, colour coded by phylum, decorated with barplots of actual (orange) and potential (purple) Gene Cluster Families (GCFs), as defined by BiG-SLiCE (T=0.4). Top REDgroups with most potential GCFs include the following: A: *Streptomyces*\_RG1, B: *Streptomyces*\_RG2, C: *Amycolatopsis*\_RG1, D: *Kutzneria*, E: *Pseudomonas*\_E. Phyla known to be enriched in NP producers are immediately visible (Actinobacteriota, Protobacteriota), with the most promising groups coming from the Actinobacteriota phylum (the highest peak belongs to a REDgroup containing *Streptomyces* strains). Simultaneously, within the underexplored phyla, there seems to be significant biosynthetic diversity and potential. An interactive version of Figure 4a can be accessed online (Supplementary Figure 5). Panel b: Rarefaction curves of REDgroups (BiG-SLiCE T=0.4). The genera included in the most promising REDgroups are indicated (the letters a-e correspond to the peaks in A). *Streptomyces* strains are included in several of the top most promising REDgroups. Panel c: Rarefaction curves of the most promising REDgroups (BiG-SLiCE T=0.4). Panel d: Rarefaction curves of the most promising REDgroups (clust-o-matic T=0.5). Though the exact numbers differ, the similarities between the two methods are apparent.

**Figure 5. Unique diversity in the known producer *Streptomyces* and promising potential of less popular taxa.** Panel a: Unique Gene Cluster Families (GCFs) as defined by BiG-SLiCE (T=0.4), of phyla and *Streptomyces* (solid shapes) and pairwise overlaps of phyla - phyla and phyla - *Streptomyces* (ribbons), visualized with circlize<sup>38</sup>. Each taxon has a distinct color. The smaller shapes and ribbons represent smaller phyla that can be seen in Supplementary Figure 6. The genus *Streptomyces* appears to have a very high amount of unique GCFs comparable to entire phyla, such as Proteobacteria. Panel b: Unique GCFs as defined by BiG-SLiCE (T=0.4), of non-streptomycete Actinobacteriota and all *Streptomyces* genera (solid shapes) and pairwise overlaps between Actinobacteriota and *Streptomyces* (ribbons), visualized with circlize<sup>38</sup>. The *Streptomyces* genus, only one of many belonging to the Actinobacteriota phylum, appears to be responsible for a big percentage of the phylum's unique diversity (comparison of solid shapes of group 1 and 5). Panel c: Left: Potential (pGCFs) and actual (GCFs) number of Gene Cluster Families as defined by BiG-SLiCE (T=0.4), of top 20 most promising REDgroups. Right: number of Natural Products (NPs) found in the NPASS database<sup>45</sup>, that originate from species included in each REDgroup. Several of these REDgroups appear to have few (< 15) to no known NPs associated with them: *Amycolatopsis*\_RG1, *Kutzneria*, *Xanthobacteraceae*\_mixed\_RG1 (containing the genera *Bradyrhizobium*, *Rhodopseudomonas*, *Tardiphaga* and *Nitrobacter*), *Mycolicibacterium*\_RG1, *Nonomuraea*, *Kitasatospora*\_RG1.

## Tables

**Table 1. Overview of input datasets and overall biosynthetic diversity under various BiG-SLiCE thresholds.** The dataset “Entire Bacterial Kingdom” was used for the computation of the actual and potential biosynthetic diversity found in all cultured (and some uncultured) bacteria. The dataset “RefSeq bacteria with known species taxonomy” was used for pinpointing the emergence of biosynthetic diversity, for which accurate taxonomic information was needed, and for identifying groups of promising producers. \*MAG sources: bovine rumen<sup>22</sup>, chicken caecum<sup>23</sup>, human gut<sup>24</sup>, ocean<sup>25</sup>, uncultivated bacteria<sup>26</sup>.

Dataset		Genomes	BGCs	Gene Cluster Families			
				T = 0.4	T = 0.5	T = 0.6	T = 0.7
Entire Bacterial kingdom	All RefSeq bacteria	170,549	1,060,592	<b>51,052</b>	37,785	28,057	19,152
	Bacterial MAGs*	13,620	34,285	<b>7,943</b>	6,516	5,014	3,519
	Total	184,169	1,094,877	<b>53,927</b>	40,497	31,998	24,446
RefSeq bacteria with known species taxonomy	Complete Genomes	16,004	94,904	<b>16,984</b>	13,546	10,399	7,151
	Draft Genomes	147,265	913,642	<b>37,123</b>	27,748	20,638	14,016
	Total	163,269	1,008,546	<b>41,870</b>	31,237	23,227	15,766