

ACCEPTED MANUSCRIPT



Global biogeographic sampling of bacterial secondary metabolism

Zachary Charlop-Powers, Jeremy G Owen, Boojala Vijay B Reddy, Melinda A Ternei, Denise O Guimarães, Ulysses A de Frias, Monica T Pupo, Prudy Seepe, Zhiyang Feng, Sean F Brady

DOI: <http://dx.doi.org/10.7554/eLife.05048>

Cite as: eLife 2015;10.7554/eLife.05048

Received: 5 October 2014
Accepted: 7 January 2015
Published: 19 January 2015

This PDF is the version of the article that was accepted for publication after peer review. Fully formatted HTML, PDF, and XML versions will be made available after technical processing, editing, and proofing.

Stay current on the latest in life science and biomedical research from eLife.
[Sign up for alerts](http://elife.elifesciences.org) at elife.elifesciences.org

1 **Title:** Global Biogeographic Sampling of Bacterial Secondary Metabolism

2
3 **Authors:** Zachary Charlop-Powers,^a Jeremy G. Owen,^a Boojala Vijay B. Reddy,^a Melinda Ternei,^a Denise
4 O. Guimarães,^b Ulysses A. de Frias,^c Monica T. Pupo,^c Prudy Seepe,^d Zhiyang Feng,^e and Sean F. Brady^{a*}

5
6 **Author affiliations:**

7 ^a Laboratory of Genetically Encoded Small Molecules, Howard Hughes Medical Institute, The Rockefeller
8 University, 1230 York Avenue, New York NY 10065.

9 ^b Laboratório de Produtos Naturais, Curso de Farmácia, Universidade Federal do Rio de Janeiro, Campus
10 Macaé, Pólo Novo Cavaleiro – IMMT, Rua Alcides da Conceição, 159, 27933-378, Macaé, RJ, Brazil

11 ^c School of Pharmaceutical Sciences of Ribeirão Preto, University of São Paulo (FCFRP-USP)

12 Av. do Cafe, s/n14040-903, Ribeirão Preto, SP, Brazil

13 ^d K-RITH: KwaZulu-Natal Research Institute for Tuberculosis & HIV, Nelson R. Mandela School of Medicine,
14 K-RITH Tower Building, 719 Umbilo Road, Durban

15 ^e College of Food Science and Technology, Nanjing Agricultural University, 1 Weigang, Nanjing 210095,
16 China

17
18 ***Corresponding Author:** Sean F. Brady

19 **Contact Information:** Laboratory of Genetically Encoded Small Molecules

20 The Rockefeller University

21 1230 York Avenue

22 New York, NY 10065

23 **Phone:** 212-327-8280

24 **Fax:** 212-327-8281

25 **Email:** sbrady@rockefeller.edu

26
27 **Acknowledgements:**

28 This work was supported by National Institutes of Health grant number GM077516 (S.F.B.), and F32
29 AI110029 (Z.C.P.). S.F.B. is a Howard Hughes Medical Institute Early Career Scientist. Brazilian research
30 was supported by São Paulo Research Foundation (FAPESP) grant #2011/50869-8. M. T. P. is a research
31 fellow of the Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq). D.O.G. was
32 supported by the Rio de Janeiro Research Foundation (FAPERJ) grant #E-26/110.281/2012 and CNPq
33 grant #477509/2013-4. The authors would also like to acknowledge the following people for assistance in
34 collecting samples: Rafael Bonante, Samyr Soares Viana, Vitor de Carli, Ronaldo de Carli, Erin Bishop, and
35 Vanessa Kowalski.

36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56

57 **Abstract:**

58 Recent bacterial (meta)genome sequencing efforts suggest the existence of an enormous untapped
59 reservoir of natural-product-encoding biosynthetic gene clusters in the environment. Here we use the pyro-
60 sequencing of PCR amplicons derived from both nonribosomal peptide adenylation domains and polyketide
61 ketosynthase domains to compare biosynthetic diversity in soil microbiomes from around the globe. We
62 see large differences in domain populations from all except the most proximal and biome-similar samples,
63 suggesting that most microbiomes will encode largely distinct collections of bacterial secondary metabolites.
64 Our data indicate a correlation between two factors, geographic distance and biome-type, and the
65 biosynthetic diversity found in soil environments. By assigning reads to known gene clusters we identify
66 hotspots of biomedically relevant biosynthetic diversity. These observations not only provide new insights
67 into the natural world, they also provide a road map for guiding future natural products discovery efforts.

68

69

70 **Introduction:**

71 Soil-dwelling bacteria produce many of the most important members of our pharmacy, including the majority
72 of our antibiotics as well as many of the cytotoxic compounds used in the treatment of cancers.(1) The
73 traditional approach for characterizing the biosynthetic potential of environmental bacteria has been to
74 examine metabolites produced by bacteria grown in monoculture in the lab. However, it is now clear that
75 this simple approach has provided access to only a small fraction of the global microbiome's biosynthetic
76 potential.(2-4) In most environments, uncultured bacteria outnumber their cultured counterparts by more
77 than two orders of magnitude, and among the small fraction of bacteria that has been cultured,(5, 6) only a
78 small subset of gene clusters found in these organisms is generally expressed in common fermentation
79 broths.(7, 8) The direct extraction and subsequent sequencing of DNA from environmental samples using
80 metagenomic methods provides a means of seeing this "biosynthetic dark matter" for the first time.
81 Unfortunately the genomic complexity of most metagenomes limits the use of the shotgun-sequencing and
82 assembly approaches (9, 10) that are now routinely used to study individual microbial genomes.(11, 12)
83 Although bacterial natural products represent an amazing diversity of chemical structures, the majority of
84 bacterial secondary metabolites, including most clinically useful microbial metabolites, arise from a very
85 small number of common biosynthetic themes (*e.g.* polyketides, ribosomal peptides, non-ribosomal
86 peptides, terpenes, etc.).(13) Because of the functional conservation of enzymes used by these common
87 systems, degenerate primers targeting the most common biosynthetic domains provide a means to broadly
88 study gene cluster diversity in the uncultured majority in a way similar to what is now regularly done for
89 bacterial species diversity using 16S rRNA gene sequences. Here we use this approach to conduct the first
90 global examination of non-ribosomal peptide synthetase (NRPS) adenylation domain (AD) and polyketide
91 synthase (PKS) ketosynthase (KS) domain biosynthetic diversity in soil environments. We chose to explore
92 NRPS and PKS biosynthesis because the highly modular nature of these biosynthetic systems has provided
93 a template for the production of a wide variety of gene clusters that give rise to a correspondingly diverse
94 chemical repertoire, including many of the most clinically useful microbial metabolites.(1)

95 **Results and discussion:**

96 With the help of a citizen science effort (www.drugsfromdirt.org), soil samples were collected from five
97 continents (North America, South America, Africa, Asia, Australia) and several oceanic islands (Hawaii,
98 Dominican Republic), covering biomes that include multiple rainforests, temperate forests, deserts and
99 coastal sediments (Supp. Mat. Table1, Map 1). DNA was extracted directly from these soils as previously
100 described (14) and 96 samples were chosen for analysis of NRPS/PKS diversity using 454 pyro-sequencing
101 of AD and KS domain PCR amplicons. Samples were chosen on the basis of DNA quality and biome
102 diversity; raw sequence reads from these samples were combined with existing amplicon datasets derived
103 from other biomes using the same DNA isolation, PCR and sequencing protocols. (15) The entire dataset
104 representing 185 biomes was clustered into operational taxonomic units (OTUs) at a sequence distance of
105 five percent. Despite millions of unique sequencing reads yielding a predicted Chao1 OTU estimate of
106 greater than 350,000 for each domain, rarefaction analysis suggests that we have not yet saturated the
107 sequence space of either domain (Figure 1A, 1C).

108

109 The first question we sought to address with this data was how biosynthetic sequence composition varies
110 by geographic distance. To do this we calculated the pairwise Jaccard distances between AD/KS sequence
111 sets derived from each sampling site and used these metrics to compare samples. The Jaccard distance, a
112 widely used metric for comparing the fraction of shared OTUs between samples, was chosen over
113 alternative metrics due to its simplicity and to the lack of a comprehensive reference phylogenetic tree for
114 AD and KS domains as exists for 16S analyses. Most Jaccard distances were found to be quite small (<
115 3%), indicating large differences in secondary metabolite gene sequence composition between almost all
116 sample collection sites (Figure 1B, 1D). Although the OTU overlap between our individual experimental
117 samples is generally small, these relationships allow us to begin to develop a picture of how biosynthetic
118 diversity varies globally. On a global level, the strongest biosynthetic sequence composition relationships
119 are seen between samples collected in close physical proximity to one another (Figure 1: B, D, E, F) as
120 opposed to between samples from similar biomes in different geographic locations. For example, at a cutoff
121 of even as low as 3% shared KS or AD OTUs, essentially all inter-sample relationships are observed
122 between immediate geographic neighbors and not similar biomes in different global locations (Figure 1E,
123 1F). This likely explains the limited inter-sample relationships we observe between samples from the

124 Eastern hemisphere as most samples from this part of the world were collected from sites at a significant
125 geographic distance from one another. The only exception is the set of soil samples from South Africa, of
126 which a number were collected in relatively close geographic proximity. These samples exhibit similar
127 pairwise Jaccard metrics to those observed between geographically proximal samples collected in the
128 Western hemisphere (Figure 1E, 1F).

129

130 Although differences in biosynthetic composition of microbiomes appear to depend at least in part on the
131 geographic distance between samples, our data suggests that change in the biome type is an important
132 additional factor for the differentiation of biosynthetic diversity on a more local level (Figure 1G, 1H). For
133 example, at a cutoff of 3% shared OTUs, essentially all inter-sample relationships are observed between
134 immediate geographic neighbors when this is raised to 10% shared OTUs (Figure 1E, 1F), relationships are
135 only seen between nearby samples belonging to the same biome. This phenomenon is highlighted by the
136 two examples shown in Figure 1G and 1H. In the first example, Brazilian soils were collected from Atlantic
137 rainforest, saline or cerrado (savanna-like) sites located only a few miles from one another. Our AD and KS
138 data show these sample are i) distinct from other globally distributed samples, ii) most strongly related to
139 the samples from the same Brazilian biome and iii) only distantly related to the samples from other Brazilian
140 biomes. In the second example, a sample collected from a New Mexican hot spring where the soil is
141 heated continuously by subterranean water is compared with samples derived from the dry soils of the
142 surrounding environment. Once again our amplicon data show that these samples are i) distinct from other
143 globally distributed samples, ii) most strongly related to other samples from the same biome and iii) only
144 distantly related to samples from other nearby biomes. Although it is possible that at a much greater
145 sampling depth all AD and KS domains will be found at all sites as predicted by Baas-Becking's "everything
146 is everywhere but the environment selects" hypothesis of global microbial distribution (16, 17), our PCR-
147 based data suggest that both geography and ecology play a role in determining the major biosynthetic
148 components of a microbiome.

149

150

151 The vast majority of AD and KS domain sequences coming from environmental DNA (eDNA) are only
152 distantly related to functionally characterized NRP/PK gene clusters, precluding precise predictions about
153 the specific natural products encoded by the gene clusters from which most amplicons arise. However, in
154 cases where eDNA sequence tags show high sequence similarity to domains found in functionally
155 characterized gene clusters, this information can be used to predict the presence of specific gene cluster
156 families within a specific microbiome. This type of phylogenetic analysis is the basis of the recently
157 developed eSNaPD program, a BLAST-based algorithm for classifying the gene cluster families that are
158 associated with eDNA-derived sequence tags. (18, 19) When an eDNA sequence tag clades with, but is
159 not identical to, a reference sequence in an eSNaPD-type analysis, it is considered to be indicative of the
160 presence of a gene cluster that encodes a congener (*i.e.*, a derivative) of the metabolite encoded by the
161 reference cluster.

162

163 Interestingly, eSNaPD analysis of the data from all sites reveals there are two distinct types of biomedically
164 relevant natural product gene cluster “hot spots” within our data (Figure 2A, 2B, 2D). These include
165 “specific gene cluster hotspots” and “gene cluster family hotspots”. Metagenomes from “specific gene
166 cluster hotspots” are predicted to be enriched for a gene cluster that encodes a congener of the target
167 natural product, while metagenomes from “gene cluster family hotspots” are predicted to encode multiple
168 congeners related to the target natural product. Figure 2A shows several of the strongest examples of
169 “specific gene cluster hotspots” where reads falling into an OTU related to a specific biomedically relevant
170 gene cluster or gene cluster family are disproportionately represented in the sequence data from individual
171 microbiomes. These examples highlight the different enrichment patterns that we observe in the
172 environment – hotspots are either local in nature, consisting of only one or two samples containing
173 sequence reads mapping to the target (epoxomycin, oocydin); regional (tiacumicinB); or global with
174 punctuated increases in diversity (glycopeptides). We would predict “specific gene cluster hotspots” (Figure
175 2D) are naturally enriched for bacteria that encode congeners of the biomedically relevant target
176 metabolites, thereby potentially simplifying the discovery of new congeners. Figure 2B shows examples of
177 “gene cluster family hotspots,” where metagenomes having a disproportionately high number of OTUs
178 mapping to a specific biomedically relevant target molecule family (*e.g.*, nocardicin, rifamycin, bleomycin,

179 and daptomycin families are shown) are highlighted. This analysis identifies specific sample sites, from
180 among those surveyed, that are predicted to contain the most diverse collection of gene clusters associated
181 with a target molecule of interest (Figure 2B). Both types of hotspots should represent productive starting
182 points for future natural product discovery efforts aimed at expanding the structural diversity and potential
183 utility of specific biomedically relevant natural product families.

184

185 Biosynthetic domain sequence tag data are not only useful for pinpointing environments that are rich in
186 specific biosynthetic targets of interest but also as a metric for natural product biosynthetic diversity in
187 general. As only a small fraction (5-10%) of total AD and KS sequences can be confidently assigned by the
188 eSNaPD algorithm, samples showing the largest collection of unique OTUs (at a common sequencing
189 depth) might be expected to contain the most diverse collection of novel biosynthetic gene clusters (Figure
190 2C) and therefore be the most productive sites to target for future novel molecule discovery efforts. Once
191 normalized for sequencing depth, the number of unique KS and AD sequence tags observed per collection
192 site differs by almost an order of magnitude between environments (Figure 2C), with the most diverse
193 samples mapping to Atlantic forest and Desert environments (Figure 2C, 2D teal spots, Supplementary File
194 7).

195

196 The development of cost effective high-throughput DNA sequencing methodologies and powerful
197 biosynthesis focused bioinformatics algorithms allow for the direct interrogation and systematic mapping of
198 global microbial biosynthetic diversity. Our analyses of hundreds of distinct soil microbiomes suggests that
199 geographic distance and local environment play important roles in the sample-to-sample differences we
200 detected in biosynthetic gene populations. As variations in biosynthetic gene content are expected to
201 correlate with variations in the small-molecule producing capabilities of a microbiome, the broader
202 implication of these observations from a drug discovery perspective is that the dominant biosynthetic
203 systems of geographically distinct soil microbiomes are expected to encode orthogonal, largely unexplored
204 collections of natural products. Taken together, our biosynthetic domain hotspot and OTU diversity analyses
205 represent a starting point in the creation of a global natural products atlas that will use sequence data to
206 guide natural product discovery in the future. Based on the historical success of natural products as

207 therapeutics, microbial “biosynthetic dark matter” is likely to hold enormous biomedical potential. The key
208 will be learning how to harvest molecules encoded by the biosynthetic diversity we are now able to find
209 through sequencing.
210

211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229

Figure 1. Global Abundance and Comparative Distribution of AD/KS Sequences

The global abundance (A, C), sample-to-sample variation (B, D), and geographic distribution (E, F, G, and H) of adenylation domains (AD) and ketosynthase domains (KS) were assessed by pyro-sequencing of amplicons generated using degenerate primers targeting AD and KS domains found in 185 soils/sediments from around the world.

(A, C) Global AD (A) or KS (C) domain diversity estimates were obtained by rarefying the global OTU table (*de novo* clustering at 95%) for AD and KS sequences and calculating the average Chao1 diversity metric at each sampling depth.

(B, D) The ecological distance (*i.e.*, Jaccard dissimilarity) between AD (B) or KS (D) domain populations sequenced from each metagenome was determined as a function of the great circle distance between sample collection sites (km). Insets show local relationships (< 500 km) in more detail.

(E, F) All sample collection sites are shown on each world map and lines are used to connect sample sites that share at least the indicated fraction (3%, 10%) of AD (E) or KS (F) OTUs.

(G, H) Biome-specific relationships within domain OTU populations sequenced from geographically proximal samples assessed by Jaccard similarity. Samples were collected from (G) Atlantic forest, saline or cerrado environments or from the (H) New Mexican desert topsoils or hot springs sediments.

230 **Figure 2. Biomedically Relevant Natural Product Hotspots and Diversity**

231 Hotspot analysis of natural product biosynthetic diversity to identify samples with a high total proportion of
232 reads corresponding to a natural product family of interest (A, D), the maximum unique OTUs
233 corresponding to a natural product family of interest (B, D), or the estimated sample biodiversity (C, D). In A
234 and B samples are arranged by longitude and hemisphere as is shown in the Sample Key.

235 (A) For each sample, sequence reads assigned by eSNaPD are expressed as a percentage of total
236 reads obtained for that sample. A sample is designated a hotspot if more than one percent (0.01; horizontal
237 line) of its reads map to a specific gene cluster. Fractional observance data for five representative gene
238 clusters or gene cluster families (zorbamycin, oocydin, tiacumicinB, epoxomicin, glycopeptides) that show
239 significant sample dependent difference in read frequency are shown.

240 (B) Hotspots of elevated gene cluster family diversity can be identified by determining the number of
241 unique OTUs occurring in each sample that, by eSNaPD, map to a natural product gene cluster of interest.
242 Sample specific OTU counts for nocardicin, rifamycin, bleomycin, and daptomycin clusters are shown.
243 Samples containing greater than 50% of the maximum observed OTU value are colored and mapped in (C).
244 OTU diversity measurements do not predict the abundance of a specific cluster in a metagenome [as
245 predicted in (A)], but instead are used to identify locations where the largest number of congener-encoding
246 clusters may be found. These sites are predicted to be most useful for increasing the structural diversity
247 and therefore potential clinical utility of these medically important families of natural products.

248 (C) Estimated diversity of AD/KS reads by sample. AD and KS OTU tables were combined and for
249 each sample the Chao1 diversity metric was calculated at 5,000 reads, providing a baseline metric for
250 comparing sample biosynthetic diversity. The average number of unique OTUs observed over 10
251 rarefactions analyses is shown (also see Supplementary File 7).

252 (D) Hotspot map of samples identified in A, B and C.

253 (E) Representative structures of target molecule families highlighted in A and B.

254
255

256 **Materials and Methods:**

257

258 **Soil Collection.** Soil from the top 6 inches of earth was collected at unique locations in the continental
259 United States, China, Brazil, Alaska, Hawaii, Costa Rica, Ecuador, the Dominican Republic, Australia and
260 South Africa. The full sample table is available in Supplementary File1.

261

262 **Soil DNA extraction.** To reduce the potential for cross contamination, DNA was extracted from soil using a
263 simplified version of our previously published DNA isolation protocol (14, 20). The modified protocol was as
264 follows: 250 grams of each soil sample was incubated at 70°C in 150 ml of lysis buffer (2% sodium dodecyl
265 sulfate [wt/vol], 100 mM Tris-HCl, 100 mM EDTA, 1.5 M NaCl, 1% cetyl trimethyl-ammonium bromide
266 [wt/vol]) for 2 h. Large particulates were then removed by centrifugation (4,000 x g, 30 min), and crude
267 eDNA was precipitated from the resulting supernatant with the addition of 0.6 volumes of isopropyl alcohol.
268 Precipitated DNA was collected by centrifugation (4,000 x g, 30 min), washed with 70% ethanol and
269 resuspended in a minimum volume of TE (10 mM Tris, 1 mM EDTA [pH 8]). Crude environmental DNA was
270 passed through two rounds of column purification using the PowerClean system (MO BIO, Carlsbad,
271 California). Purified environmental DNA was then diluted to 30 ng/μl and archived for use in PCR reactions.

272

273 **PCR amplification.** Degenerate primers targeting conserved regions of AD [A3F (5'-
274 GCSTACSYSATSTACACSTCSGG) and A7R (5'-SASGTCVCCSGTSCGGTA) (21)] and KS [degKS2F.i (5'-
275 GCIATGGAYCCICARCARMGIVT) and degKS2R.i (5'-GTICCGTICCRTGISCYTCIAC) (22)] domains were
276 used to amplify gene fragments from crude eDNA. Forward primers were designed to contain a 454
277 sequencing primer (CGTATCGCCTCCCTCGCGCCATCAG) followed by a unique 8 bp barcode that
278 allowed simultaneous sequencing of up to 96 different AD- or KS- samples in a single GS-FLX Titanium
279 region. PCR reaction consisted of 25 μl of FailSafe PCR Buffer G (Epicentre, Madison, Wisconsin), 1 μl
280 recombinant *Taq* Polymerase (Bulldog Bio, Portsmouth, New Hampshire), 1.25 μl of each primer (100 mM),
281 14.5 μl of water and 6.5 μl of purified eDNA. PCR conditions for AD domain primers were as follows: 95 °C
282 for 4 min followed by 40 cycles of 94 °C for 0.5 min., 67.5 °C for 0.5 min, 72 °C for 1 min and finally 72°C for
283 5 min. PCR conditions for KS domain primers were as follows: 95 °C for 4 min followed by 40 cycles of 54

284 °C for 40 seconds, 56.3 °C for 40 seconds, 72 °C for 75 seconds and finally 72 °C for 5 min. PCR reactions
285 were examined by 2% agarose gel electrophoresis to determine the concentration and purity of each
286 amplicon. Amplicons were pooled in equal molar ratios, gel purified using the Invitrogen eGel system and
287 DNA of the appropriate size was recovered using Agencourt Ampure XP beads (Beckman Coulter, Brea,
288 California). Amplicons were sequenced using the 454 GS-FLX Titanium platform. Raw flowgram files from
289 454's shotgun processing routine were used for downstream analysis.

290

291 **Processing 454 data.** Raw reads were assigned to samples using the unique primer barcodes and filtered
292 by quality (50 bp rolling window PHRED cutoff of 20) using Qiime (version 1.6).(23) USEARCH (version 7),
293 which implements the improved UPARSE clustering algorithm (24), was used to remove Chimeric
294 sequences with the default 1.9 value of the *de novo* chimera detection tool. UPARSE clustering requires all
295 sequences to be of the same length. In an effort to balance read quality and abundance with the ability to
296 phylogenetically discriminate gene clusters we used 419 bp as our read length cutoff. The trimmed fasta file
297 was then clustered to 5% to compensate for sequencing error and natural polymorphism that is often
298 observed in gene clusters found in natural bacterial populations. Clustering proceeded as per the
299 USEARCH manual by clustering at a distance of 3% and using representative sequences from each cluster
300 to cluster again at 5%. The resulting "5%" AD and KS OTU tables were used for all subsequent rarefaction
301 and diversity analyses.

302

303 **Rarefaction and Diversity Analyses.** To assess global AD and KS diversity in our sample set we sought
304 to assess the global number of AD and KS domains we might expect to see if all of our data had been
305 generated from a single sample. To do this, all reads assigned to an OTU were consolidated to generate a
306 single-column OTU table where each row contains the sum of all sequences assigned to that OTU from any
307 of the 185 samples. To assess the global diversity we subsampled this table at multiple depths using Qiime
308 (23) and used the Chao1 formula to estimate the expected number of OTUs at this depth. This rarefaction
309 analysis was performed ten times at each subsampling depth (Figure1A, 1C; Supplementary Files 3 and 4)
310 and the curves were fit to the data using the following equation: $y = 1 + \log(x) + \log(x^2) + \log(x^3)$ where x
311 is the read value and y is the Chao1 diversity.

312 Ecological distances are calculated using the Jaccard $[1 - (OTU_{A\&B}) / (OTU_A + OTU_B - OTU_{A\&B})]$ or
313 inverse Jaccard metric (25) and geographic distances were calculated using great circle (spherical) distance
314 derived from the latitude/longitude values of each set of points (26)(Supplementary File 5). Pairwise
315 ecological and geographic distances were used to create Figure 1B, 1D. Network plots of subsamples
316 (Figure 1: G, H) were generated using Phyloseq (27) to calculate the intersample Jaccard distance. As
317 expected, the strongest relationships are observed between sample proximity controls where soils were
318 collected approximately 10 meters from one another and processed independently, demonstrating that
319 closely related samples do in fact group together in our analysis pipeline.

320

321 **Assignment of AD and KS domains to known gene clusters.** AD and KS amplicon reads were assigned
322 to known biosynthetic gene clusters using the eSNaPD algorithm at an e-value cutoff of 10^{-45} .(18) At this
323 threshold eSNaPD has been used to successfully assign-and-recover gene clusters that encode congeners
324 of multiple natural product families using only the sequence from a single domain amplicon.(19, 28, 29)
325 NRPS/PKS clusters typically have multiple KS or AD domains. Hits to all domains in a cluster were
326 aggregated in our analyses. Data for eSNaPD hits broken down by sample and molecule are included as
327 Supplementary File 6.

328

329 **Hotspot Analysis.** AD and KS OTU tables were analyzed for the presence of eSNaPD hits. For each
330 sample the abundance of each eSNaPD hit (i.e. a particular molecule) was calculated as either a
331 percentage of total reads (Figure 2A, C) or as the total number of unique OTUs assigned to the molecule
332 that were found in that sample (Figure 2B, C), or as the total number of OTUs mapped to a molecule in
333 each sample. In the read-based hotspot analysis, the number of reads assigned by eSNaPD to a specific
334 gene cluster is expressed as a fraction of total per sample reads: (reads-to-cluster-of-interest)/total sample
335 reads). In the OTU-based hotspot analysis we calculated the number of unique eSNaPD assigned OTUs
336 found in each sample that map to a specific gene cluster. The full eSNaPD dataset is available in
337 Supplementary File 6. To compare global biosynthetic diversity of each sample, the AD and KS OTU tables
338 were combined and for each sample they were subsampled ten times to a depth of 5000 reads. The Chao1

339 diversity metric was calculated for each sample and the average was used to compare the expected
340 biodiversity in different samples at the same sampling depth (Figure 1C, Supplementary File 7).

341

342 **Supplementary Files:**

343 Supplementary File 1: Sample Location Data

344 Supplementary File 2: Sample Read and 95% OTU Count

345 Supplementary File 3: Adenylation Domain Rarefaction Data (Figure 1A)

346 Supplementary File 4: Ketosynthase Domain Rarefaction Data (Figure 1C)

347 Supplementary File 5: Pairwise Sample Distances. Great Circle Distance and Jaccard Distance for AD and
348 KS Amplicons

349 Supplementary File 6: eSNaPD Hits Broken Down by Sample and Molecule

350 Supplementary File 7: Per Sample Chao1 Biodiversity Estimates at a Rarefaction Depth of 5,000 Reads

351

- 355 1. G. M. Cragg, D. J. Newman, Natural products: a continuing source of novel drug leads. *Biochimica*
356 *et biophysica acta* **1830**, 3670 (Jun, 2013).
- 357 2. M. S. Rappe, S. J. Giovannoni, The uncultured microbial majority. *Annual review of microbiology*
358 **57**, 369 (2003).
- 359 3. J. Rajendhran, P. Gunasekaran, Microbial phylogeny and diversity: small subunit ribosomal RNA
360 sequence analysis and beyond. *Microbiological research* **166**, 99 (Feb 20, 2011).
- 361 4. J. A. Gilbert, C. L. Dupont, Microbial metagenomics: beyond the genome. *Annual review of marine*
362 *science* **3**, 347 (2011).
- 363 5. V. Torsvik, F. L. Daae, R. A. Sandaa, L. Ovreas, Novel techniques for analysing microbial diversity
364 in natural and perturbed environments. *Journal of biotechnology* **64**, 53 (Sep 17, 1998).
- 365 6. V. Torsvik, J. Goksoyr, F. L. Daae, High diversity in DNA of soil bacteria. *Applied and*
366 *environmental microbiology* **56**, 782 (Mar, 1990).
- 367 7. S. D. Bentley *et al.*, Complete genome sequence of the model actinomycete *Streptomyces*
368 *coelicolor* A3(2). *Nature* **417**, 141 (May 9, 2002).
- 369 8. H. Ikeda *et al.*, Complete genome sequence and comparative analysis of the industrial
370 microorganism *Streptomyces avermitilis*. *Nat Biotechnol* **21**, 526 (May, 2003).
- 371 9. A. C. Howe *et al.*, Tackling soil diversity with the assembly of large, complex metagenomes.
372 *Proceedings of the National Academy of Sciences of the United States of America* **111**, 4904 (Apr
373 1, 2014).
- 374 10. V. Iverson *et al.*, Untangling genomes from metagenomes: revealing an uncultured class of
375 marine Euryarchaeota. *Science* **335**, 587 (Feb 3, 2012).
- 376 11. P. Cimermancic *et al.*, Insights into secondary metabolism from a global analysis of prokaryotic
377 biosynthetic gene clusters. *Cell* **158**, 412 (Jul 17, 2014).
- 378 12. S. Donadio, P. Monciardini, M. Sosio, Polyketide synthases and nonribosomal peptide
379 synthetases: the emerging view from bacterial genomics. *Natural product reports* **24**, 1073
380 (2007).
- 381 13. P. M. Dewick, *Medicinal Natural Products: A Biosynthetic Approach*. (John Wiley & Sons, 2002).
- 382 14. S. F. Brady, Construction of soil environmental DNA cosmid libraries and screening for clones
383 that produce biologically active small molecules. *Nature protocols* **2**, 1297 (2007).
- 384 15. Z. Charlop-Powers, J. G. Owen, B. V. Reddy, M. A. Ternei, S. F. Brady, Chemical-biogeographic
385 survey of secondary metabolism in soil. *Proceedings of the National Academy of Sciences of the*
386 *United States of America* **111**, 3757 (Mar 11, 2014).
- 387 16. M. A. O'Malley, The nineteenth century roots of 'everything is everywhere'. *Nature reviews.*
388 *Microbiology* **5**, 647 (Aug, 2007).
- 389 17. R. de Wit, T. Bouvier, 'Everything is everywhere, but, the environment selects'; what did Baas
390 Becking and Beijerinck really say? *Environ Microbiol* **8**, 755 (Apr, 2006).
- 391 18. B. V. Reddy, A. Milshteyn, Z. Charlop-Powers, S. F. Brady, eSNaPD: A Versatile, Web-Based
392 Bioinformatics Platform for Surveying and Mining Natural Product Biosynthetic Diversity from
393 Metagenomes. *Chemistry & biology*, (Jul 23, 2014).
- 394 19. J. G. Owen *et al.*, Mapping gene clusters within arrayed metagenomic libraries to expand the
395 structural diversity of biomedically relevant natural products. *Proceedings of the National*
396 *Academy of Sciences of the United States of America*, (Jul 3, 2013).
- 397 20. B. V. B. Reddy *et al.*, Natural product biosynthetic gene diversity in geographically distinct soil
398 microbiomes. *Applied and environmental microbiology* **78**, 3744 (2012).

- 399 21. A. Ayuso-Sacido, O. Genilloud, New PCR primers for the screening of NRPS and PKS-I systems in
400 actinomycetes: detection and distribution of these biosynthetic gene sequences in major
401 taxonomic groups. *Microbial ecology* **49**, 10 (2005).
- 402 22. A. Schirmer *et al.*, Metagenomic Analysis Reveals Diverse Polyketide Synthase Gene Clusters in
403 Microorganisms Associated with the Marine Sponge *Discodermia dissoluta*. *Applied and*
404 *environmental microbiology* **71**, 4840 (2005).
- 405 23. J. G. Caporaso *et al.*, QIIME allows analysis of high-throughput community sequencing data.
406 *Nature methods* **7**, 335 (2010).
- 407 24. R. C. Edgar, UPARSE: highly accurate OTU sequences from microbial amplicon reads. *Nature*
408 *methods* **10**, 996 (Oct, 2013).
- 409 25. J. B. Oksanen, F. Guillaume ; Kindt, Roeland; Legendre, Pierre; Minchin, Peter R.; O'Hara, R. B.;
410 Simpson, Gavin L.; Solymos, Peter; Henry, M. ; Stevens, H. ; Wagner, Helene. (2013).
- 411 26. R. S. B. Edzer J. Pebesma, Classes and methods for spatial data in R. *R News* **5**, 9 (2005).
- 412 27. P. J. McMurdie, S. Holmes, phyloseq: An R Package for Reproducible Interactive Analysis and
413 Graphics of Microbiome Census Data. *PloS one* **8**, (2013).
- 414 28. H. S. Kang, S. F. Brady, Arixanthomycins A-C: Phylogeny-guided discovery of biologically active
415 eDNA-derived pentangular polyphenols. *ACS chemical biology* **9**, 1267 (Jun 20, 2014).
- 416 29. F. Y. Chang, S. F. Brady, Characterization of an environmental DNA-derived gene cluster that
417 encodes the bisindolylmaleimide methylarcyriarubin. *Chembiochem : a European journal of*
418 *chemical biology* **15**, 815 (Apr 14, 2014).
- 419
420

421

422

Adenylation Domains

Ketosynthase Domains



