

1 Creating and leveraging bespoke large- 2 scale knowledge graphs for comparative 3 genomics and multi-omics drug 4 discovery with SocialGene

5

6 Chase M. Clark¹, Jason C. Kwan^{1*}

7 ¹Division of Pharmaceutical Sciences, School of Pharmacy, University of Wisconsin-Madison,

8 777 Highland Avenue, Madison, WI 53705, USA

9

10 * To whom correspondence should be addressed. Tel: +1 608-262-3829; Fax: +1 608-262-

11 5345; Email: jason.kwan@wisc.edu

12

13 Keywords: Genome mining, genomics, metagenomics, natural products, biosynthesis, software,

14 knowledge graphs, databases, drug discovery, neo4j, specialized metabolites, proteins, genes,

15 genomes, biosynthetic gene clusters, BGCs, multi-omics integration

16

17 Abstract

18 The rapid expansion of multi-omics data has transformed biological research, offering
19 unprecedented opportunities to explore complex genomic relationships across diverse
20 organisms. However, the vast volume and heterogeneity of these datasets presents significant
21 challenges for analyses. Here we introduce SocialGene, a comprehensive software suite
22 designed to collect, analyze, and organize multi-omics data into structured knowledge graphs,
23 with the ability to handle small projects to repository-scale analyses. Originally developed to
24 enhance genome mining for natural product drug discovery, SocialGene has been effective
25 across various applications, including functional genomics, evolutionary studies, and systems
26 biology. SocialGene's concerted Python and Nextflow libraries streamline data ingestion,
27 manipulation, aggregation, and analysis, culminating in a custom Neo4j database. The software
28 not only facilitates the exploration of genomic synteny but also provides a foundational
29 knowledge graph supporting the integration of additional diverse datasets and the development
30 of advanced search engines and analyses. This manuscript introduces some of SocialGene's
31 capabilities through brief case studies including targeted genome mining for drug discovery,
32 accelerated searches for similar and distantly related biosynthetic gene clusters in biobank-
33 available organisms, integration of chemical and analytical data, and more. SocialGene is free,
34 open-source, MIT-licensed, designed for adaptability and extension, and available from
35 github.com/socialgene.

36 Introduction

37 The advent of large-scale multi-omics datasets has ushered in a new era in biological research.
38 However, the volume and complexity of datasets present significant challenges for their
39 analysis. Here we present SocialGene, a suite of software for computing and organizing multi-
40 omics data—including genomics, metabolomics, and more—into structured knowledge graphs¹
41 ranging from small to repository-scale. While SocialGene is versatile and applicable across a
42 broad range of disciplines, its initial development was motivated by the need to enhance
43 genome mining for natural product drug discovery, the primary focus of this introductory
44 manuscript.

45
46 Searching genomes for proteins of similar function and synteny (defined hereafter as a set of
47 collinear, putatively-orthologous genes) is an essential task across multiple scientific disciplines.
48 In natural product drug discovery, the focus centers on biosynthetic gene clusters (BGCs) that
49 encode for the biosynthesis of specialized metabolites (SM). Developing new methods for
50 identifying and targeting orthologous BGCs across vast sets of public and private genomes can
51 enable a number of applications including finding additional sources of SM and their chemical
52 analogs.

53
54 This is especially important, and challenging, for BGCs that encode for the biosynthesis of
55 medicinally-important SM that have so far only been observed within the metagenome-
56 assembled genomes (MAGs) of microbial obligate symbionts.^{2,3} Obligate symbionts often have
57 reduced genomes^{4,5} and are recalcitrant to isolation and cultivation, necessitating the study of
58 their SMs through chemical extraction of the holobiont or genetic engineering of predicted BGCs
59 into heterologous hosts. These processes are resource-intensive but often necessary due to the
60 difficulty and economics of chemically synthesizing many SMs. Several examples of the

61 resources required to obtain SM from microbial endosymbionts include the antifungal SM
62 lagriamide (symbiont producer– *Burkholderia gladioli*⁶), which required collecting 28,000 *Lagria*
63 *villosa* beetle eggs to recover 600 µg of compound.⁶ The anticancer (and beetle-defense)
64 compound pederin (symbiont producer– *Pseudomonas* sp.⁷) whose structure determination
65 required the field collection of 25 million beetles (100 kg) over seven years, resulting in "many
66 hospitalizations".⁸ The anticancer SM bryostatin (symbiont *Candidatus* *Endobugula sertula*⁹)
67 which required 13,000 kg of *Bugula neritina* to produce 18 g of bryostatin 1 for clinical trials.^{10,11}
68 The anticancer compound halichondrin B (proposed to be symbiont produced), which required
69 600 kg of the marine sponge *Halichondria okadai* for structure elucidation¹² and an additional
70 1,000 kg collection of *Lissodendoryx* to recover approximately 300 mg of pure halichondrin for
71 clinical trials.¹¹ And the anticancer compound ET743 (symbiont producer–*Candidatus*
72 *Endoecteinascidia frumentensis*¹³), which required aquaculture of 100,000 kg of the tunicate
73 *Ecteinascidia turbinata* to recover the 100 g of compound needed up to Phase 2 of clinical
74 trials.¹⁴ When the BGC of a metagenomic SMs is known or suspected, finding similar BGCs in
75 cultured strains could be economically advantageous and influence the decision or speed of
76 pursuing clinical trials. Given the non-generalizable and lengthy process of finding suitable
77 vectors and hosts for genetic recombination, we sought a scalable framework for searching for
78 related proteins and BGCs across previously cultured organisms, independent of BGC
79 prediction frameworks, and where sequence identity and synteny might be low due to
80 evolutionary distance.

81
82 Existing tools for finding similar BGCs (e.g. clusterblast¹⁵, MultiGeneBlast¹⁶, cblaster¹⁷, FlaGs¹⁸,
83 CAGECAT¹⁹, BiG-SLiCE²⁰ etc.) have proven valuable but had limitations for our use cases.
84 Some require predicting BGCs within the target genomes and restricting searches to those BGC
85 regions. However, current ensemble and machine learning BGC predictors (e.g. antiSMASH¹⁵,
86 deepBGC²¹, GECCO²², etc.; review by Kim et al²³) have a limited ability to detect BGCs not

87 represented in their rule sets or training data, or when split across a sequence(s), necessitating
88 a full-genome search. Others don't scale well or rely on BLAST²⁴ sequence similarity, which has
89 limited ability to find distant homologs, particularly when search results are restrained and there
90 are many homologs in the target database. As a result, we needed the ability to compare
91 hundreds of millions of proteins and their positions across hundreds of thousands of genomes,
92 in near real-time, while retaining the ability to discover low sequence identity homologs.

93

94 Protein similarity is often determined through sequence-sequence alignment tools (e.g.
95 BLAST²⁴, DIAMOND^{25,26}, MMseqs2²⁷, etc.) which, though incredibly fast, have a limit of
96 detection for low sequence identity homologs, run considerably slower in high sensitivity modes,
97 and often produce a burdensome number of matches when searching large target databases.
98 This has forced a limitation in many current BGC search tools to only consider the top-n results
99 of searches. Another approach is sequence-model alignment (e.g. HMMER²⁸, HH-suite²⁹, etc.)
100 which provides detection of low sequence identity homologs but is too slow and compute-
101 intensive for just-in-time annotation at repository scale. As both approaches are often needed,
102 we aimed to develop a method that leverages each but performs the majority of the computation
103 upfront, conducting dynamic searches over the stored results.

104

105 To that end, we created SocialGene, a suite of software projects centered around three
106 elements: 1) a Python³⁰ library that defines and controls the majority of data transformations and
107 database interactions 2) a Nextflow³¹ workflow, written using nf-core³² templating and
108 standards, that allows users to reproducibly create Neo4j graph databases using input proteins
109 and/or genomes and 3) a Neo4j graph database, created in the final step of the Nextflow
110 workflow, that stores and organizes the data, facilitating complex queries and analyses. The
111 Nextflow workflow and Python library enable annotating proteins using profile hidden Markov
112 models (pHMMs), clustering proteins with MMseqs2, and/or creating all-vs-all protein similarity

113 networks with DIAMOND BLASTp. An option is available to run antiSMASH v7³³ on all input
114 genomes, allowing predicted BGC regions to be extracted and incorporated into the database.
115 Additionally, the software and database schema were written modularly to allow users with
116 programming experience the ability to extend the graph and perform custom analyses.

117
118 In this manuscript we present the software for the first time, as well as a limited set of potential
119 use cases. This includes searching thousands of known BGCs against more than 343,000
120 RefSeq³⁴ genomes, targeted genome mining for protein domains and functional co-occurrence,
121 developing query strategies for drug discovery, linking Minimum Information about a
122 Biosynthetic Gene cluster (MIBiG)³⁵ BGCs to chemicals in NPAtlas³⁶, linking genomes to LC-
123 MS/MS features, Global Natural Products Social Molecular Networking (GNPS)³⁷ clusters, and
124 reference libraries.

125
126 As a resource for the community, we have precomputed SocialGene databases of various
127 sizes, based on RefSeq as of November 14, 2023. This includes a SocialGene database
128 computed over all RefSeq genomes, specific subsets of Actinobacteria, *Streptomyces* spp., and
129 *Micromonospora* spp. genomes, and a database of 2,103,244 antiSMASH 7.0³³ predicted
130 BGCs— one of the largest public BGC compendia to date. Additionally, we provide an online and
131 interactive BGC atlas. The atlas contains the results of using SocialGene to search the full size
132 SocialGene RefSeq database for similar BGCS to each of the 2,502 BGCs in MIBiG³⁵, but
133 restricted to the >27,000 genomes associated with a strain available from a culture collection.

134 Methods

135 Nextflow workflow

136 Nextflow³¹ is a domain-specific language for producing reproducible scientific workflows.
137 Nextflow was chosen for the promise of creating a single SocialGene workflow that would
138 provide reproducibility, parallelism, checkpointing, and ability to run on local and cloud
139 computing platforms. To provide standardization, SocialGene's database-building workflow was
140 designed to nf-core³² standards. nf-core is a "framework for community-curated bioinformatics
141 pipelines"³² and, while SocialGene was not submitted as an official nf-core workflow, it was built
142 using the framework and therefore benefits from the surrounding tooling. This includes a GUI for
143 launching the workflow and the ability to interface with nf-core "tools", "modules",
144 "subworkflows", etc. Additionally, the workflow can be run using Seqera's Nextflow Tower, an
145 online Nextflow workflow orchestrator. SocialGene's Nextflow workflow
146 (github.com/socialgene/sgnf) handles downloading data from a number of sources (e.g. NCBI
147 genomes, MIBiG BGCs, multiple public pHMM databases, etc.), the extraction, transformation,
148 and loading (ETL) of input and computed data, and culminates in the building of a custom Neo4j
149 graph database (Fig. 1 and [Supplementary Fig. 1](#)). The SocialGene Nextflow workflow and
150 Python library were designed modularly so that users can choose to run any or all analyses.
151 The configuration files used for creating the Neo4j databases in this manuscript are available
152 within the archived codebase. As of writing, SocialGene was built against Nextflow version 24;
153 and nf-core tools template version 2.10.

154 Python library

155 The Nextflow workflow contains several independent Python scripts but also makes use of
156 command line entry points defined within the SocialGene python library
157 (github.com/socialgene/sgpy; pypi.org/project/socialgene). The library was written with entry
158 points for limited use as a command line tool, as utilized in the Nextflow workflow, and as a
159 Python library directly for a number of bio- and cheminformatic tasks. All code changes are
160 checked for breaking changes (pytest) and code style (Flake8 and Black) through continuous
161 integration and continuous delivery (CI/CD) workflows using GitHub Actions and "Release
162 Please"³⁸ which automate new releases and deployments to PyPI. Test coverage is monitored
163 with Codecov. With the SocialGene software split across several git repositories and software
164 languages it was important to coordinate a consistent set of parameters when using each (e.g.
165 parameters passed to HMMER's `hmmsearch` when creating a SocialGene database and when
166 annotating a query protein later). To maintain consistent settings across database creation (via
167 Nextflow), notebook analysis (Python), and future interfaces (Django), the Python library
168 contains a file of environment variables "common_parameters.env" which are read and modified
169 at runtime from within Nextflow, Django, etc. These parameters are also saved within the Neo4j
170 database at the time of creation.

171 Neo4j graph database

172 Neo4j is a company that maintains graph database software of the same name. Neo4j
173 maintains support of docker images of community and enterprise editions, drivers in popular
174 languages, and in-database graph data science and machine learning libraries. The SocialGene
175 Nextflow workflow and Python library automate the creation of bespoke Neo4j databases, which
176 can then be interrogated directly in a web browser, with the SocialGene Python library,
177 Cytoscape, or other third party tools. While SocialGene makes use of Neo4j graph databases,

178 the Nextflow workflow gathers all intermediate files as tab-separated flat files that could be
179 imported into an alternate database system.

180

181 Input genomes

182 The Nextflow workflow can download genomes from NCBI or use local GenBank files and/or
183 protein FASTA files. To identify redundant proteins, as well as provide a consistent, cross-
184 source sequence-based identifier, SocialGene uses sequence hashes as universal identifiers.
185 When genomes are provided in GenBank format SocialGene uses BioPython⁴⁰ and custom
186 scripts to parse genome and sequence data. Additionally, we have found highly-relevant
187 pseudogenes within BGCs (and elsewhere) which lack translated sequence data in GenBank
188 files. With this observation, and recent studies showing some portion of PGAP-labeled⁴¹
189 pseudogenes are misassembled coding genes,⁴² we decided to attempt to include annotated
190 pseudogenes. Therefore, SocialGene attempts to include pseudogene content via extracting the
191 relevant nucleotide sequence and employing BioPython's Bio.Seq.translate. As some
192 pseudogenes simply contain a potential early stop codon(s) these may also be physiologically
193 relevant via translational read-through⁴³ and other mechanisms. However, the correct
194 translation of pseudogenes that aren't transcribed or the incorrect translation of pseudogenes is
195 a data inclusion bias users should be aware of. SocialGene tracks which sequences were
196 derived from pseudogenes by prepending "pseudo_" to the locus description, which can be
197 used to filter results in the SocialGene database. Additionally, if available in the GenBank file,
198 SocialGene will attempt to include the reason the gene was marked as pseudo, (e.g internal
199 stop, frameshift, etc.).

200 Representing proteins as hashes

201 SocialGene Neo4j database protein entries use sha512t24u⁴⁴ hashes as universal identifiers but
202 are also assigned a CRC64 hash for fast cross-referencing with UniProt.^{44,45} Hashing is a
203 process that takes an input string of characters and transforms it into a uniquely identifiable
204 hash. This is often used to assign a short, unique identifier to a large quantity of information. For
205 example, the human protein titin (UniprotKB⁴⁶ Q8WZ42) contains 34,350 amino acids but can
206 be represented by its CRC64 hash: DEB216410AD560D9. Different hashing algorithms have
207 different probabilities for the scenario that two different inputs produce an identical hash (hash
208 collision). While preparing this manuscript we switched to using CRC64 due to its use in UniProt
209 which would provide the ability to crosslink, and link out to, UniProt information/resources⁴⁵.
210 However, it was discovered there were 1,704 hash collisions across UniParc (out of
211 517,621,195 total proteins) and, more concerning, tens of thousands of collisions when hashing
212 SocialGene's internally used string "{genbank_accession}__{genbank_locus}" across all RefSeq
213 nucleotide sequences. For this reason we switched back to using sha512t24u due to its
214 predicted collision probability (no collisions were detected in UniParc or internal SocialGene
215 identifiers), speed (~2x SEGUID), and being url-safe.⁴⁴

216 Hidden Markov models and annotation

217 SocialGene's Nextflow workflow can download and format pHMMs from any or all of
218 antiSMASH⁴⁷, AMRFinder⁴⁸, BiG-SLiCE²⁰, ClassiPhage⁴⁹, Pfam⁵⁰, PRISM⁵¹, Resfams⁵², and
219 TIGRFAMs⁵³; as well as user-provided pHMMs in HMMER3 format. The Python library reduces
220 input pHMMs to a less-redundant set of models by hashing the models' emissions and
221 transitions and, for Pfam, uses only the latest version of a model. For example, if the user or
222 combination of above reference databases try to include Pfam models PF00001.23 and
223 PF00001.24, SocialGene will only annotate proteins with PF00001.24, and in the resulting

224 Neo4j database will note that the PF00001.23 model was specified by the input source but
225 PF00001.24 was used for domain prediction. For compatibility with HMMER's hmmsearch, less-
226 redundant models are output in two files, one for models with gathering cutoffs and one for
227 models without. Within Nextflow, less-redundant fasta files are split into n-files and run against
228 the two pHMM files in parallel using HMMER's hmmsearch. Through extensive testing we found
229 that, with fast hard drives, splitting an input fasta into multiple files, assigning 1 logical cpu to
230 hmmsearch, and running in a highly parallel fashion provides the fastest results for this step.

231 antiSMASH

232 SocialGene's Nextflow workflow has the ability to annotate input genomes with antiSMASH
233 version 7³³. A custom Python script reduces resulting antiSMASH json files into a minimal
234 JSONL file (newline-delimited JSON) that describes the assembly, locus, coordinates and
235 minimal metadata for all predicted BGCs. While at first this may seem unnecessary, the gzipped
236 tar archive of unmodified antiSMASH output for all successfully annotated RefSeq genomes
237 was >1.5 TB. The gzipped, summarizing, gzipped minimal JSONL for the same was 86 MB, a
238 >16,000x reduction in storage size.

239 MMSeqs2

240 SocialGene's Nextflow workflow performs cascaded clustering using MMseqs2. For example,
241 clustering non-redundant proteins to 90% and 50% sequence identity first clusters proteins to
242 90%, followed by taking the 90% cluster representatives and clustering them to 50%. This is
243 important because it means to find proteins in the database with less than 90%, but greater than
244 50%, sequence similarity will require a two-hop traversal, first traversing "MMSEQS_90"
245 relationships then "MMSEQS_50". To allow users to cluster input proteins to multiple, custom
246 identity levels the Nextflow module was written to take a delimited string of identity levels (e.g.

247 '90,70,50', representing 90%, 70%, and 50% sequence identity). Depending on the number of
248 proteins and clustering levels this process can require a significant amount of disk space and
249 RAM (100s of GBs). This Nextflow process outputs a single flat file edge list representing the
250 protein clusters, as well as MMseqs databases for each level.

251 Hardware

252 Data was created/analyzed on either: "Desktop 1": single AMD® Ryzen 9 3900xt 12-core
253 processor with 62 GB of RAM; or "Server 1": dual AMD® EPYC 7352 24-Core processors, with
254 1 TB RAM. Both machines used SABRENT 4 TB Rocket NVMe PCIe M.2 2280 as working
255 drives. NCBI RefSeq genomes were stored on, and processed from, a Western Digital 18TB
256 WDC_WD181KRYZ disk drive. Large scale pHMM annotations were computed using the
257 University of Wisconsin-Madison's Center for High Throughput Computing (CHTC) and the
258 Open Science Grid (OSG). While not required, Neo4j database creation, initialization, and large-
259 scale read/write benefit from fast hard drive storage.

260 Scaling

261 Databases used in this manuscript were computed using a combination of Desktop 1, Server 1
262 and CHTC resources. For inputs over 1,000 genomes, data aggregation steps can be computed
263 on a mid-tier laptop or desktop computer, but the non-distributed DIAMOND and MMseqs2
264 protein comparisons begin to require a high amount of RAM and it becomes best to shard the
265 input FASTA and run the pHMM annotation step on a high throughput computing cluster, if
266 available. Since Nextflow is not currently supported on the University of Wisconsin-Madison's
267 CHTC submit server, a flag "--htcondor" was created in SocialGene's Nextflow workflow which
268 signals for the organization and output of the bundled set of processed non-redundant fasta
269 files, pHMMs, generated scripts, and instructions for submitting the jobs via HTCondor (but

270 generalizable to other computing environments). The Nextflow workflow can then be run a
271 second time with the "--htcondor" flag removed and the path to HMMER results provided to the
272 command line flag "--domtblout_path". Adding the "--resume" flag allows this second run to
273 continue where the first left off, reusing already completed computations. Utilizing this
274 technique, combined with resources available through CHTC and OSG, has allowed us to
275 create SocialGene databases with >340,000 genomes, requiring tens of thousands of CPU
276 hours, in under two days, instead of months to years.

277 Precomputed databases

278 To test SocialGene's ability to scale to large collections of genomes, we ran the workflow on all
279 genomes available in RefSeq (including non-bacterial). While it's possible to use SocialGene to
280 download all RefSeq genomes, doing so requires a substantial amount of disk space (>1.5 Tb)
281 and thus we used an existing local copy of the 343,381 genomes, updated on November 14,
282 2023. The SocialGene Nextflow workflow was run on Server 1 with settings to annotate all
283 genomes with antiSMASH ⁷³³; annotate all non-redundant proteins with pHMMs from
284 antiSMASH³³, AMRFinder⁴⁸, Pfam⁵⁰, Resfams⁵², and TIGRFAM⁵³; and cluster non-redundant
285 proteins to 90%, 70%, 50%, 30% with MMSeqs2. To run hmmsearch on CHTC/OSG we
286 instructed the Nextflow workflow to split the non-redundant protein FASTA into 3000 files (using
287 SeqKit⁵⁴ split). SocialGene's Nextflow flag "--htcondor" then instructed the workflow to package
288 the resulting FASTA files, two non-redundant pHMM model files (those with and without
289 gathering cutoffs) and customized scripts, for submission with HTCondor. The resulting 6,000
290 hmmsearch jobs required 14,726 cpu hours to complete and the total workflow required
291 approximately 17,000 CPU hours. This does not account for the more than 10,000 CPU hours
292 to compute antiSMASH BGCs across all 343k genomes on Server 1. Apart from downloading
293 input genomes and antiSMASH predictions, due to its parallel design, the workflow completed

294 start-to-finish in less than 48 hours. However, it should be noted this is highly variable and
295 dependent on the compute resources, especially the number of CPUs. Supplementary Table 1
296 shows the number of nodes and relationships in the resulting database. The full graph database
297 occupies 650 GB of disk space and is available for download as a 220 GB Neo4j database
298 dump (see Data and Code Availability). We are also making available a separate 30 GB
299 SocialGene database of the more than two million antiSMASH predicted BGCs.

300
301 Three additional RefSeq databases (named "actinomycetota", "streptomyces", and
302 "micromonospora") were created with the intention of providing smaller precomputed databases
303 for those without access to adequate computational resources. Each was built independently
304 with the SocialGene Nextflow workflow, making use of the NCBI Datasets module (e.g. --
305 ncbi_datasets_command 'genome taxon "actinomycetota" --assembly-source refseq --exclude-
306 atypical'). Databases in this manuscript have been labeled as version 2023_v0.4.1.

307 Representing and linking chemistry

308 SocialGene has the ability to incorporate and crosslink non-redundant chemical compounds
309 from a variety of sources, using RDKit⁵⁵ and custom scripts. As of writing, redundancy is based
310 on unique InChI⁵⁶ strings, as most NP databases don't contain more detailed structural
311 information than InChI⁵⁶ or SMILES⁵⁷. Additionally, SocialGene links similar compounds within
312 the database using an all-vs-all comparison of Morgan fingerprints⁵⁸ (radius 2, 2048 bits) and
313 Tanimoto similarity.

314 Results and Discussion

315 HMM outdegree accelerated BGC search

316 The SocialGene RefSeq database (version 2023_v0.4.1) contains data on >340 thousand
317 genomes, >300 million non-redundant proteins, >25 thousand less-redundant pHMMs (from
318 antiSMASH⁴⁷, AMRFinder⁴⁸, Pfam⁵⁰, Resfams⁵², and TIGRFAMs⁵³), and >840 million pHMM-to-
319 protein annotations. Evenly distributed, this would result in 35,919 annotations per pHMM. But,
320 as shown in Supplementary [Fig. 2](#), the actual distribution of annotations per model is right-
321 skewed and log-normal distributed, with a mean of 33,163 and median of 2,948.

322

323 SocialGene's BGC search first annotates a query BGC's proteins using the same pHMM models
324 used to annotate proteins in the database, either pulling annotations from the database when
325 the protein is present, or using HMMER²⁸. To compare proteins by their pHMM annotations
326 reduces the initial search space to the 25,566 pHMM nodes, but the number of outgoing
327 relationships from pHMM nodes is over 847 million. Consequently, searches can quickly begin
328 traversing an excessive percentage of the database. To alleviate this SocialGene's BGC search
329 algorithm first calculates, and sets as a node property, the outdegree of pHMM nodes. The input
330 proteins are then prioritized by the lowest to highest summed outdegree of their pHMM
331 annotation nodes.

332

333 The database is then searched for all proteins with similar domains (pHMM annotations). These
334 similar proteins and their gene coordinates within all genomes are clustered and filtered in
335 Python based on a threshold number of hits to the input BGC's proteins. After filtering, the
336 remaining nucleotide sequences are divided into multiple regions based on the user-specified

337 'break_bgc_on_gap_of' parameter, which splits a nucleotide sequence where any region of the
338 specified length has no hits to an input BGC protein. Regions are filtered again by a threshold
339 number of hits to the input BGC proteins. Remaining regions are evaluated by reciprocal best hit
340 (RBH) analysis using either DIAMOND BLASTp or pHMM annotation similarity (user-selected).
341 The resulting putatively similar BGCs are then evaluated and ranked based on the similarity of
342 RBH content (Jaccard) and order (Levenshtein) compared to the input BGC. The search
343 can be done either within an interactive Python terminal or Jupyter notebook, enabling further
344 computation, or as a standalone command line function which outputs a JSON file for
345 visualization with clustermap.js⁵⁹.

346

347 This outdegree prioritization can dramatically speed up a search and essentially prioritizes less
348 common pHMM annotations (and, thereby, domains and proteins). However, this strategy can
349 miss target clusters if the only related proteins between query and target BGCs are those
350 excluded in the prioritized search. [Fig. 2](#) provides a visual guide of this relationship of pHMMs
351 (and their outdegree) to protein and nucleotide sequences (labeled as BGC in the figure).
352 Further explanation is available in [Supplementary Text 1](#).

353

354 Multiple methods are required for measuring protein similarity

355 To justify the protein similarity search strategy for biosynthetic gene clusters (BGCs) we
356 explored the correlation between DIAMOND's²⁵ BLASTp protein-protein sequence identity
357 scores, MMseqs2 clustering, and the Jaccard and Levenshtein similarity of HMMER²⁸ pHMM
358 annotations. While MMseqs2²⁷ and DIAMOND were comparable (See Supplementary Figs.
359 [3,4](#)), there was little, if any, global correlation between pHMM annotations and DIAMOND
360 BLASTp identities (Supplementary Fig. [5](#)). This lack of correlation is due to the algorithms used

361 for pHMM annotation similarity which don't account for model or sequence coverage, or detailed
362 domain position. For single domain proteins, perfectly-similar pHMM annotation often consists
363 of only a single pHMM model annotation. Thus, while single domain proteins can have a range
364 of sequence alignment identities they usually only have binary pHMM Jaccard and Levenshtein
365 similarity scores.

366 For example, UniProtKB⁴⁶ proteins Q8X5K5 and Q8XCP8 are encoded by *Escherichia coli*
367 O157:H7 genes *lpfA* and *yfcQ*. Both proteins are potentially highly relevant to human health, as
368 *lpfA* is part of the *lpfABCC'DE* fimbrial operon and has been shown to promote
369 enterohemorrhagic *E. coli* cells' interaction and adherence to eukaryotic cells.⁶¹⁻⁶⁶ However,
370 while *yfcQ* from the laboratory-cryptic *yfcOPQRSTUVWXYZ* operon has been computationally inferred
371 to also be a fimbrial-like adhesin protein, there have been limited studies on its role in
372 pathogenesis or adhesion.⁶⁷⁻⁶⁹ While NCBI's BLASTp^{24,70} was unable to align these protein
373 sequences due to their low sequence similarity, their predicted AlphaFold⁷¹ 3D protein
374 structures did align (see Supplementary Fig. [6](#)). Additionally, when looking at pHMM
375 annotations, over 80% of the AAs in both proteins were annotated by the PF00419.23 (Fimbrial)
376 Pfam model. Therefore, a search strategy starting with one of these proteins would fail to find
377 the other when using BLASTp, while a strategy employing pHMM annotations would succeed.

378

379 Conversely, a nearly-perfect BLASTp alignment doesn't necessitate similar pHMM annotation.
380 For example, UniProtKB⁴⁶ A0A0H3JI96 and A0A0H3JGM8 are phage tail proteins encoded in
381 the *Escherichia coli* O157:H7 genome. While BLASTp alignment revealed matches in 233 of
382 238 positions (97.9% identity), only a third of their pHMM annotations overlap in SocialGene's
383 RefSeq database (see [Supplementary Text 2](#)).

384

385 Therefore, it is important to consider using both sequence and model approaches to protein
386 similarity in SocialGene. For large databases, we recommend utilizing SocialGene's MMseqs2
387 cascaded clustering method rather than all-vs-all BLASTp, as the latter can result in an
388 excessive number of relationships.

389 Finding metagenomic, fragmented and multiple copy BGCs

390 Searching for metagenome-assembled genome (MAG) derived BGCs in the sequences of
391 cultivated organisms is challenging for a variety of reasons including the sheer number of public
392 genomes and the low quality of many MAG BGCs and public genomes. To examine the ability
393 of SocialGene's BGC search algorithm to look through hundreds of thousands of public
394 genomes for metagenomic BGC homologs we looked at a recently verified example: lagriamide.
395 Lagriamide is an antifungal SM whose BGC is encoded in the reduced genome of the *Lagria*
396 *villosa* beetle endosymbiont, *Burkholderia gladioli* Lv-StB.^{6,72} Through a combination of
397 individual BLASTp searches against NCBI's nr database and manual bioinformatic analysis we
398 and several laboratories recently collaborated to find two free-living strains of *Paraburkholderia*
399 *acidicola* that contained a partial match to the metagenomic-derived BGC encoding for
400 lagriamide.^{6,72}

401
402 To evaluate SocialGene's BGC search algorithm we took the MAG-derived lagriamide BGC
403 (MIBiG BGC0001646) and ran the search against the SocialGene RefSeq database (343,381
404 public genomes) and were able to recover the aforementioned *P. acidicola* BGC
405 (Supplementary Fig. [7](#)).

406
407 SocialGene was also able to recover the BGC of the immunosuppressant SM rapamycin, even
408 when fragmented and/or containing corrupted-genes, as shown in Supplementary Fig. [8](#). While

409 SocialGene is able to find fragmented BGCs, the default search returns only the highest-scoring
410 fragment due to limitations of plotting in clustermap.js⁵⁹. Lastly, the BGC search function was
411 able to recover the multiple integrations of a nybomycin-encoding plasmid in a genome
412 engineered strain (Supplementary Fig. 9).

413 Finding syntenic but distantly related BGCs

414 There are few examples of finding endosymbiont-derived, metagenomic BGCs in free-living
415 relatives and few references of the extent of sequence divergence of orthologous BGCs over
416 large evolutionary distances. While the lagriamide example above was reported to have 93.7%
417 pairwise identity,⁷² the individual proteins in the public genome assembly have amino acid
418 identities around 70 to 80%. And though we hypothesized the need to find syntenic BGCs
419 where individual ortholog sequence similarities were low we were unsure if this existed in
420 nature. Additionally, while finding collinear, putatively-orthologous genes is suggestive of
421 common ancestry and conserved function it is important to consider the likelihood of convergent
422 evolution, though the probability of the later assumedly decreases as the sequence similarity,
423 count, and synteny of shared genes increases.

424

425 To test our hypothesis, SocialGene's automated BGC search algorithm was used to search
426 each of 2,502 MIBiG BGCs as queries against the entire SocialGene RefSeq database. While
427 the algorithm uses pHMM annotations for the primary search, the lower bound of sequence
428 identity was limited by the final step, where putative target BGCs were compared with reciprocal
429 best hit (RBH) analysis using DIAMOND BLASTp in "ultra-sensitive" mode. Though confounded
430 due to biases in the RefSeq and MIBiG databases,⁷³ the majority of query BGCs that had
431 targets with high synteny and low sequence similarity also had a large number of total hits (i.e.
432 BGCs that are highly prevalent across RefSeq). These abundant BGC classes were similar to

433 Cimermancic, Medema, Claesen et al's observation of the widespread occurrence of "O-
434 antigens, capsular polysaccharides, carotenoids and NRPS-independent siderophores".⁷⁴

435
436 To that end we looked for the longest MIBiG BGC with the highest synteny and lowest median
437 RBH identities. MIBiG BGC0000182 is a BGC from a *Pseudomonas fluorescens* bacterium with
438 36 protein-coding genes that encode the biosynthesis of the polyketide antibiotic pseudomonic
439 acid A (mupirocin). Mupirocin is a clinically-important antibiotic that continues to be included on
440 the World Health Organization's List of Essential Medicines.⁷⁵ Using SocialGene's BGC search
441 function, we searched the SocialGene RefSeq database for BGC0000182. While most of the
442 resulting 17 target BGCs were highly similar, two had median RBH identity values of 73.8% and
443 58.5%, while still containing a RBH to every BGC0000182 protein ([Fig. 3](#)). While the strain with
444 a median of 73.8% protein sequence identity was also a *Pseudomonas* sp., the strain with
445 58.5% median identity was *Chromobacterium* IIBBL 290-4, which belongs to a different
446 taxonomic Class. The *Chromobacterium* sp. BGC was flanked by transposases
447 (NKT35_RS10105/NKT35_RS10110 and NKT35_RS10295) suggesting potential mobility of the
448 BGC. Interestingly, the region between transposons only contains 34 of the 36 proteins, with
449 MupR and MupX homologs occurring directly adjacent to, but outside of, NKT35_RS10295 (a
450 pseudo IS1380 family transposase). While *Pseudomonas* sp. QS1027 is a known producer of
451 mupirocin.⁷⁶ it is currently unknown whether *Chromobacterium* IIBBL 290-4 produces mupirocin
452 or a mupirocin chemical analog.

453

454 Following the evolution of BGCs

455 While some BGCs have few matches (e.g. mupirocin, mentioned above), others are
456 overrepresented due to organism bias in RefSeq, wider phylogenetic distribution, or both. Using

457 the same search strategy as with mupirocin, but with BGC0000946 (*Vibrio parahaemolyticus*
458 BGC encoding for vibrioferrin) as the query BGC, resulted in 6,577 complete and syntenic target
459 BGCs across 6,571 genome assemblies, along with 4 MIBiG BGCs. These BGCs were
460 distributed across 968 species, 81 genera, 46 families, 6 classes, and 3 phyla; with the median
461 percent identities of RBHs occurring in a stepped gradient (Fig. 4).

462

463 Though tempting to believe the gradients would represent functional evolution and
464 diversification of end-product SMs, one of the lowest median RBHs (46.6%) belonged to MIBiG
465 BGC0002527, a vibrioferrin-producing BGC from *Azotobacter vinelandii* CA. The actual lowest
466 median RBH of 32.5% was found in a *Facilibium subflavum* assembly. While there's no
467 evidence this *F. subflavum* strain produces the vibrioferrin siderophore, the flanking genes
468 suggest the region is involved in metal acquisition and homeostasis (Supplementary Fig. 11).

469

470 While Fig. 4 shows 6,581 intact and syntenic BGCs, it is also possible these are situated within
471 broader genomic contexts that catalyze the modification of vibrioferrin or an alike molecule.
472 However, it is unclear what proportion of the BGCs this is, if any. Further studies are needed to
473 determine the cause of the stepped gradients and whether they are due to speciation, horizontal
474 transfer (see Supplementary Fig. 12), or other mechanisms. While outside the scope of the
475 current study it is possible to create comprehensive in-database similarity links between BGCs
476 for studying phylogenetic histories, especially those that are difficult to express as a phylogram.

477

478 Extending capability by computing new nodes and relationships, 479 in-database

480 SocialGene's default schema is useful on its own but also designed to serve as a foundation
481 from which new nodes and relationships can be created, a strength of Neo4j. For example, the
482 SocialGene RefSeq database contains all MIBiG BGCs, antiSMASH predictions across 343,381
483 genomes, and MMseqs2 protein clustering to 90%, 70%, 50%, and 30% sequence identities. By
484 traversing the existing nodes and relationships in the graph database a new type of relationship
485 can be calculated that directly connects MIBiG BGCs to any genome assembly containing a
486 similar BGC (Fig. 5). As shown in Fig. 5, the ability to filter subgraphs by additional metadata
487 (e.g. taxon, host, etc.) enables researchers to create hypotheses about patterns of BGC
488 distribution. Further, filtering or coloring availability in a culture collection provides a fast route to
489 procuring strains for further experiments (Supplementary Fig. [14](#)).

490

491 Atlas of BGCs available in culture collections

492 One goal of creating the BGC search function in SocialGene was to enable repository-scale
493 searches for BGCs across public and private culture collections. This aimed to uncover new
494 sources of previously inaccessible BGCs, identify higher-yield strains, and provide a tool for
495 hypothesis testing. However, we recognize that not everyone will have the necessary
496 computational resources or expertise to install the SocialGene RefSeq database.

497

498 To address this, an online interactive atlas of MIBiG BGCs was developed, where similar BGCs
499 can be found in various strain collections (e.g., NRRL, ATCC, DSMZ, etc.). This was achieved

500 by searching each of the 2,502 MIBiG BGCs against the SocialGene RefSeq database,
501 focusing specifically on the 27,406 genomes with metadata indicating availability in a culture
502 collection. The resulting clustermap.js⁵⁹ plots, restricted to 100 target BGCs per query MIBiG
503 BGC (limited by visualization), include a total of 92,936 target BGCs spread across 2,112 MIBiG
504 BGCs. For access to the atlas, refer to the Data and Code Availability section.

505 Supporting narrow and broad meta-analyses

506 Version 3.0 of the MIBiG repository contains 2,502 BGCs and, like many natural product
507 databases (e.g. npatlas^{36,77}), entries are skewed towards well-studied taxa (e.g. Actinobacteria,
508 especially *Streptomyces* spp.) and biosynthetic classes (e.g. PKS, NRPS, etc.). Despite this,
509 early versions of MIBiG have been invaluable for building software and evaluating how
510 computational methods and models behave with validated BGCs. SocialGene's Nextflow
511 workflow contains an optional flag (“--mibig”) which signals for the incorporation of all MIBiG
512 BGCs into a SocialGene database, with or without additional input genomes.

513

514 For development and proof of concept work a SocialGene database was created containing all
515 MIBiG BGCs. This resulted in a modest-sized graph database with 2.7 million nodes and 4.9
516 million edges, including more than 40,000 non-redundant protein nodes and more than 500,000
517 pHMM annotation relationships. Additionally, as many MIBiG BGCs contain NCBI taxonomy
518 identifiers, SocialGene's Nextflow “--ncbi_taxonomy” flag was used, which downloads and
519 parses the entire NCBI taxonomy database, and links input BGCs/genomes to the source
520 organism in the taxonomy graph. Supplementary Fig. [16](#) visualizes this placement of all MIBiG
521 BGCs onto the taxonomic graph within a SocialGene database and highlights the taxonomic
522 bias. We also exported a subgraph of all non-redundant proteins, less-redundant pHMMs, and
523 the annotation links connecting the two, for import and layout in Gephi (Supplementary Fig. [17](#)).

524 As expected, proteins were primarily clustered by function, but the graph also excelled in
525 displaying the complicated evolutionary relationships between both large multidomain and
526 smaller accessory proteins. Similar analysis allows for putative functional transfer to
527 hypothetical proteins.

528 Targeted antibiotic drug discovery

529 Intentional query engineering leveraging *in silico*, *in vitro*, and *in vivo* domain knowledge
530 enables targeted large-scale searches of SocialGene databases across biochemistry,
531 chemistry, and modes of action. These targeted searches and analyses can be designed and
532 used to inform wet-lab experiments pre-, ad-, and post hoc. For example, a customized search
533 for peptidic and halogenated antibiotics can guide the choice of isolation and bioassay
534 techniques.

535
536 To engineer such a search we can exploit the fact that bacteria often encode resistance
537 mechanisms within a BGC to counteract the toxicity of the produced specialized metabolite(s).
538 Searching for self resistance proteins is a strategy incorporated into other genome mining
539 software such as ARTS.⁷⁸ As protein functional information is present in SocialGene, in the form
540 of pHMM annotations, similar strategies can be employed in-database without prior BGC
541 detection. Using a similar strategy we can also detect putative halogenase and NRPS enzymes
542 and their co-occurrence.

543
544 For example, Supplementary Fig. [18](#) displays a query for any nucleotide sequence containing a
545 protein annotated by a tryptophan halogenase pHMM, within 10 kb of a nonribosomal peptide
546 synthetase (NRPS) protein (detected with antiSMASH's pHMM detection rule, performed in-
547 database), and within 50 kb of a protein annotated by an AMRfinder^{48,79,80} pHMM (antimicrobial

548 resistance gene detection). When limited to MIBiG sequences, the query only returned
549 halogenated NRP antibiotics such as vancomycin (Supplementary Fig. [19](#)). Because chlorinated
550 and brominated natural products provide characteristic isotopologue distributions, peptidic SM
551 often fragment well in ESI LC-MS/MS, and antibiotic activity is suggested, subsequent lab work
552 can be intentional and directed.

553 Targeted drug discovery

554 SocialGene facilitates targeted and untargeted drug discovery beyond microbial antibiotics. For
555 instance, Pfam PF00227 is a multi-kingdom pHMM model of the "proteasome subunit".
556 Proteasomes are ancient multi-subunit proteases^{81,82} involved in controlled protein degradation
557 and recycling, and small molecule inhibitors targeting proteasomes have provided promising
558 candidates for cancer therapeutics.^{83,84} As proof of concept, the SocialGene RefSeq database
559 was searched for MIBiG BGCs containing a protein annotated by the Pfam pHMM model
560 PF00227, of which there were eight. All eight of the BGCS produce proteasome inhibitors:
561 fellutamide B⁸⁵, cinnabaramide A⁸⁶, landepoxcin⁸⁷, salinosporamide A⁸⁸⁻⁹⁰ (two BGCs),
562 clarepoxcin⁸⁷, eponemycin⁹¹⁻⁹³, and TMC-86A⁹⁴. Next, an identical search was performed but
563 restricted instead to the over two million antiSMASH-predicted BGC regions in the RefSeq
564 genomes. While the MIBiG BGCs consisted of PKS, NRPS, and PKS/NRPS hybrids, this larger
565 search revealed 1,595 diverse BGCs from 25 phyla across Eukaryota, Archaea, and Bacteria.
566 BGC type counts are displayed in [Supplementary Table 1](#). Enabling fast searches (the above
567 search takes milliseconds) allows users to quickly iterate over potential targets and hypotheses.

568 Restricting pHMM annotations to specific motifs

569 While pHMM annotations provide a rapid means for discovering proteins with specific functions
570 they may prove too general for some tasks. However, highly targeted searches can be

571 performed using in-database regex filtering of the non-redundant protein amino acid sequences,
572 either alone, or in addition to pHMM annotation queries.

573
574 For example, proteins in the CAP protein superfamily⁹⁵ family are implicated in various
575 biological roles, including virulence and facilitating host-symbiont/pathogen relationships⁹⁶.
576 Based on this, one could hypothesize that BGCs containing a CAP might also be involved in
577 virulence and/or cross-kingdom interactions. However, across the 343,381 genomes in the
578 SocialGene RefSeq database, 135,595 genomes contain a protein(s) annotated by Pfam pHMM
579 PF00188 (CAP superfamily), and in 5,593 of those a CAP is in an antiSMASH 7-predicted BGC.
580 As this is still an unmanageable number, we could further restrict our hypothesis. For example,
581 Hirano et al. recently suggested that insect cysteine-rich secretory proteins (CAPs) may induce
582 the formation of plant galls⁹⁷ and, specifically, CAPs with the core sequence (F/Y-T-Q-I/V-V-W),
583 which can be expressed with the regex pattern "[FY]TQ[IV]VW.*". Filtering on PF00188 and
584 this pattern reduced the number of genomes to 2,337; and limiting the results to antiSMASH 7-
585 annotated regions returned a reasonable 43 genomes (Supplementary Fig. [20](#)). Each of these
586 searches completed in less than one second, which allows fast iteration over hypotheses.

587 Extending SocialGene

588 It is easy to extend the current graph schema both at database creation and within an active
589 database. Neo4j has a number of import and cross-database integration tools, including being
590 able to read data directly from web-hosted SPARQL endpoints. This provides the opportunity for
591 future integration of additional data from sources such as ChEMBL⁹⁸, UniProt⁴⁵, etc.
592 Additionally, the Python library and Nextflow workflows were written modularly to allow add-ons
593 of new node and relationship types.

594 Connecting chemistry to biology

595 Socialgene's Python library has an add-on that parses and integrates nodes and relationships
596 representing the Natural Products Atlas (NP Atlas) into an active SocialGene database
597 (Supplementary Fig. [21](#)). Many NP Atlas chemical structures are linked to a species-level NCBI
598 taxonomic identifier of the organism from which the compound was first reported. However,
599 species level taxonomic IDs are often not specific enough to correlate a chemical compound
600 with the genome assembly of a specific producer. To alleviate this, the add-on creates an
601 additional link using simple text-similarity measures on taxa names between NP Atlas and NCBI
602 taxonomy. In addition to including NP Atlas metadata, SocialGene creates non-redundant
603 chemical nodes within the graph database which are linked by Tanimoto similarity
604 (Supplementary Figs. [22](#), [23](#)).

605
606 For users with paired genome and mass spectrometry data, GNPS networking³⁷ results can be
607 downloaded and integrated directly. For example, Supplementary Fig. [24](#) shows the ingestion
608 and incorporation of the GNPS molecular networking and library search results from 122 LC-
609 MS/MS runs across 120 bacterial isolates (data previously published⁹⁹). Additionally,
610 Supplementary Figs. [25](#), [26](#) show the resulting interconnections between genomes; mass
611 spectra; GNPS clusters, networks and library hits; and NP Atlas entries. Future research will
612 aim to build models correlating MS features, clusters, and chemical moieties to BGCs, sub-
613 BGCs, proteins, and protein domains.

614 Limitations and improvements

615 Machine learning/artificial intelligence for protein function inference has significantly advanced
616 during the years-long creation of SocialGene, including AlphaFold protein structure prediction⁷¹
617 and the subsequent dramatic increase in available predicted structures.¹⁰⁰ While future additions

618 to SocialGene should include methods such as Foldseek¹⁰¹ (combination of protein sequence
619 and structure alignment) it would rely on input sequences being present in the DeepMind/EMBL-
620 EBI's AlphaFold Protein Structure Database or running AlphaFold. The latter would require an
621 increase in SocialGene's complexity and compute requirements and including full models would
622 require hundreds of additional gigabytes of storage, even with compression.¹⁰² Additionally, the
623 AlphaFold Protein Structure Database is currently limited to proteins less than 2,700 AA for
624 proteomes/Swiss-Prot and 1,280 AA for the rest of UniProt⁴⁶. Another promising avenue, with a
625 preliminary analysis by Schütze et al.¹⁰³ is to utilize the vector similarities of protein model
626 embeddings directly. This is possible in the latest versions of Neo4j, which incorporate vector
627 indexing, cosine similarity, and K-Nearest Neighbors searches. However, as discussed by
628 Schütze et al,¹⁰³ and our own experimentation, the speed of generating embeddings can be
629 slow for large scale experiments, especially on commodity hardware. UniProt currently
630 distributes embeddings generated from prottrans_t5_xl_u50, though this only covers
631 UniProtKB/Swiss-Prot⁴⁶; however, it may be feasible to include publicly released ESM-2
632 embeddings¹⁰⁴ in the future.

633

634 In the SocialGene RefSeq database 81% of the more than 304 million proteins are annotated by
635 at least one pHMM model, and 79.4% by Pfam models alone; this is similar to a report in 2019
636 that 77.2% of UniprotKB proteins had at least one Pfam annotation.¹⁰⁵ Along with greater pHMM
637 coverage, SocialGene could be improved by community analysis and clustering of proteins by
638 domain architecture in-database, or calculated externally using tools such as DAMA.¹⁰⁶

639

640 Apart from MIBiG, the datasets in this manuscript were created from all, or subsets, of RefSeq.
641 RefSeq was chosen for proof-of-concept as it had the highest number of publicly available
642 genomes with mostly-consistent genome annotation.⁴¹ However, there are a number of
643 limitations with applying at-scale analyses on current public databases, including RefSeq. We

644 especially caution anyone using our premade RefSeq databases, or other non-curated
645 genomes, for ecology and evolution studies. For such studies SocialGene has the ability to use
646 nucleotide FASTA files as input to the Nextflow workflow, which are then analyzed with
647 Prokka¹⁰⁷ for consistent gene annotation; or users can provide their own gene-called genomes.

648 Cost to create a database

649 Nextflow Tower was used to run the SocialGene Nextflow workflow on 314 *Micromonospora*
650 genomes. This included the compute-intensive steps of: annotating proteins with antiSMASH³³,
651 AMRFinder⁴⁸, Pfam⁵⁰, Resfams⁵², and TIGRFAM⁵³ pHMMs; clustering proteins to three levels
652 with MMseqs2; and running antiSMASH v7³³ on all 314 genomes. Running the full workflow on
653 AWS Batch cost less than \$5.00 USD, unoptimized. However, with the ability to run multiple
654 analyses on large amounts of data, we suggest users using fee-based computing resources
655 conduct their own cost estimate experiments before scaling.

656 Data and Code Availability

657 SocialGene's source code is hosted at github.com/SocialGene, with documentation at
658 socialgene.github.io. The Python library is available at github.com/socialgene/sgpy and
659 doi.org/10.5281/zenodo.12207092; the Nextflow workflow is available at
660 github.com/socialgene/sgnf and doi.org/10.5281/zenodo.12207039.
661 Code and notebooks used to generate manuscript figures and analyses are available at
662 github.com/socialgene/manuscript_1 and doi.org/10.5281/zenodo.13333842. All SocialGene
663 databases computed for this manuscript are archived as Neo4j dumps at:
664 doi.org/10.5061/dryad.ns1rn8q2k. The BGC atlas is available interactively at
665 bgcatlas.pages.dev and archived at doi.org/10.5281/zenodo.12775149.

666 Acknowledgements

667 This research was supported by an NIGMS grant R35GM133776. C.M.C was additionally
668 supported by an NLM training grant to the Computation and Informatics in Biology and Medicine
669 Training Program (NLM 5T15LM007359). Portions of this research was performed using the
670 compute resources and assistance of the UW-Madison Center For High Throughput Computing
671 (CHTC) in the Department of Computer Sciences. The CHTC is supported by UW-Madison, the
672 Advanced Computing Initiative, the Wisconsin Alumni Research Foundation, the Wisconsin
673 Institutes for Discovery, and the National Science Foundation, and is an active member of the
674 OSG Consortium, which is supported by the National Science Foundation and the U.S.
675 Department of Energy's Office of Science. Portions of this research was done using services
676 provided by the OSG Consortium,^{108,109} which is supported by the National Science Foundation
677 awards #2030508 and #1836650.

678 Author Contributions

679 C.M.C. Conceptualization, Data curation, Formal Analysis, Investigation, Methodology,
680 Visualization, Writing – original draft, Writing – review & editing, Funding acquisition.
681 JCK: Conceptualization, Editing, Supervision, Funding acquisition

682

683 References

- 684 1. Meijer, D. *et al.* Empowering natural product science with AI: leveraging multimodal data
685 and knowledge graphs. *Nat. Prod. Rep.* (2024) doi:10.1039/d4np00008k.
- 686 2. Piel, J. Bacterial symbionts: prospects for the sustainable production of invertebrate-derived
687 pharmaceuticals. *Curr. Med. Chem.* **13**, 39–50 (2006).
- 688 3. Newman, D. J. Predominately uncultured microbes as sources of bioactive agents. *Front.*
689 *Microbiol.* **7**, 1832 (2016).
- 690 4. Sloan, D. B. & Moran, N. A. Genome reduction and co-evolution between the primary and
691 secondary bacterial symbionts of psyllids. *Mol. Biol. Evol.* **29**, 3781–3792 (2012).
- 692 5. Morris, J. J., Lenski, R. E. & Zinser, E. R. The Black Queen Hypothesis: evolution of
693 dependencies through adaptive gene loss. *MBio* **3**, (2012).
- 694 6. Flórez, L. V. *et al.* An antifungal polyketide associated with horizontally acquired genes
695 supports symbiont-mediated defense in *Lagria villosa* beetles. *Nat. Commun.* **9**, 2478
696 (2018).
- 697 7. Piel, J. A polyketide synthase-peptide synthetase gene cluster from an uncultured bacterial
698 symbiont of *Paederus* beetles. *Proc. Natl. Acad. Sci. U. S. A.* **99**, 14002–14007 (2002).
- 699 8. Pavan, M. *Defensive Secretions of Arthropoda*. <https://apps.dtic.mil/sti/pdfs/AD0832802.pdf>
700 (1968).
- 701 9. Davidson, S. K., Allen, S. W., Lim, G. E., Anderson, C. M. & Haygood, M. G. Evidence for
702 the biosynthesis of bryostatins by the bacterial symbiont ‘*Candidatus Endobugula sertula*’ of
703 the bryozoan *Bugula neritina*. *Appl. Environ. Microbiol.* **67**, 4531–4537 (2001).
- 704 10. Mendola, D. Aquaculture of three phyla of marine invertebrates to yield bioactive
705 metabolites: process developments and economics. *Biomol. Eng.* **20**, 441–458 (2003).
- 706 11. Newman, D. J. & Cragg, G. M. Marine natural products and related compounds in clinical

- 707 and advanced preclinical trials. *J. Nat. Prod.* **67**, 1216–1238 (2004).
- 708 12. Hirata, Y. & Uemura, D. Halichondrins - antitumor polyether macrolides from a marine
709 sponge. *J. Macromol. Sci. Part A Pure Appl. Chem.* **58**, 701–710 (1986).
- 710 13. Schofield, M. M., Jain, S., Porat, D., Dick, G. J. & Sherman, D. H. Identification and
711 analysis of the bacterial endosymbiont specialized for production of the chemotherapeutic
712 natural product ET-743. *Environ. Microbiol.* **17**, 3964–3975 (2015).
- 713 14. Newman, D. J. From large-scale collections to the potential use of genomic techniques for
714 supply of drug candidates. *Frontiers in Marine Science* **5**, (2018).
- 715 15. Medema, M. H. *et al.* antiSMASH: Rapid identification, annotation and analysis of
716 secondary metabolite biosynthesis gene clusters in bacterial and fungal genome
717 sequences. *Nucleic Acids Res.* **39**, W339–46 (2011).
- 718 16. Medema, M. H., Takano, E. & Breitling, R. Detecting sequence homology at the gene
719 cluster level with MultiGeneBlast. *Mol. Biol. Evol.* **30**, 1218–1223 (2013).
- 720 17. Gilchrist, C. L. M. *et al.* cblaster: a remote search tool for rapid identification and
721 visualization of homologous gene clusters. *Bioinform Adv* **1**, vbab016 (2021).
- 722 18. Saha, C. K., Sanches Pires, R., Brolin, H., Delannoy, M. & Atkinson, G. C. FlaGs and
723 webFlaGs: Discovering novel biology through the analysis of gene neighbourhood
724 conservation. *Bioinformatics* **37**, 1312–1314 (2021).
- 725 19. van den Belt, M. *et al.* CAGECAT: The CompARative GEne Cluster Analysis Toolbox for
726 rapid search and visualisation of homologous gene clusters. *BMC Bioinformatics* **24**, 1–8
727 (2023).
- 728 20. Kautsar, S. A., van der Hoof, J. J. J., de Ridder, D. & Medema, M. H. BiG-SLiCE: A highly
729 scalable tool maps the diversity of 1.2 million biosynthetic gene clusters. *Gigascience* **10**,
730 (2021).
- 731 21. Hannigan, G. D. *et al.* A deep learning genome-mining strategy for biosynthetic gene
732 cluster prediction. *Nucleic Acids Res.* **47**, e110 (2019).

- 733 22. Carroll, L. M. *et al.* Accurate de novo identification of biosynthetic gene clusters with
734 GECCO. *bioRxiv* 2021.05.03.442509 (2021) doi:10.1101/2021.05.03.442509.
- 735 23. Kim, H. U., Blin, K., Lee, S. Y. & Weber, T. Recent development of computational
736 resources for new antibiotics discovery. *Curr. Opin. Microbiol.* **39**, 113–120 (2017).
- 737 24. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment
738 search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
- 739 25. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND.
740 *Nat. Methods* **12**, 59–60 (2015).
- 741 26. Buchfink, B., Reuter, K. & Drost, H.-G. Sensitive protein alignments at tree-of-life scale
742 using DIAMOND. *Nat. Methods* **18**, 366–368 (2021).
- 743 27. Steinegger, M. & Söding, J. MMseqs2 enables sensitive protein sequence searching for the
744 analysis of massive data sets. *Nat. Biotechnol.* **35**, 1026–1028 (2017).
- 745 28. Eddy, S. R. Accelerated profile HMM searches. *PLoS Comput. Biol.* **7**, e1002195 (2011).
- 746 29. Steinegger, M. *et al.* HH-suite3 for fast remote homology detection and deep protein
747 annotation. *BMC Bioinformatics* **20**, 473 (2019).
- 748 30. Guido, V. R. & Drake, F. L. Python 3 reference manual. *CreateSpace: Scotts Valley, CA,*
749 *USA.*
- 750 31. Di Tommaso, P. *et al.* Nextflow enables reproducible computational workflows. *Nat.*
751 *Biotechnol.* **35**, 316–319 (2017).
- 752 32. Ewels, P. A. *et al.* The nf-core framework for community-curated bioinformatics pipelines.
753 *Nat. Biotechnol.* **38**, 276–278 (2020).
- 754 33. Blin, K. *et al.* antiSMASH 7.0: new and improved predictions for detection, regulation,
755 chemical structures and visualisation. *Nucleic Acids Res.* **51**, W46–W50 (2023).
- 756 34. O’Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: Current status,
757 taxonomic expansion, and functional annotation. *Nucleic Acids Res.* **44**, D733–45 (2016).
- 758 35. Terlouw, B. R. *et al.* MIBiG 3.0: A community-driven effort to annotate experimentally

- 759 validated biosynthetic gene clusters. *Nucleic Acids Res.* **51**, D603–D610 (2023).
- 760 36. van Santen, J. A. *et al.* The Natural Products Atlas 2.0: A database of microbially-derived
761 natural products. *Nucleic Acids Res.* **50**, D1317–D1323 (2022).
- 762 37. Wang, M. *et al.* Sharing and community curation of mass spectrometry data with Global
763 Natural Products Social Molecular Networking. *Nat. Biotechnol.* **34**, 828–837 (2016).
- 764 38. Release Please. <https://github.com/googleapis/release-please>.
- 765 39. Wiese, R., Eiglsperger, M. & Kaufmann, M. yFiles — Visualization and automatic layout of
766 graphs. in *Graph Drawing Software* (eds. Jünger, M. & Mutzel, P.) 173–191 (Springer Berlin
767 Heidelberg, Berlin, Heidelberg, 2004).
- 768 40. Cock, P. J. A. *et al.* Biopython: freely available Python tools for computational molecular
769 biology and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).
- 770 41. Tatusova, T. *et al.* NCBI prokaryotic genome annotation pipeline. *Nucleic Acids Res.* **44**,
771 6614–6624 (2016).
- 772 42. Cooley, N. P. & Wright, E. S. Many purported pseudogenes in bacterial genomes are bona
773 fide genes. *BMC Genomics* **25**, 1–12 (2024).
- 774 43. Belinky, F., Ganguly, I., Poliakov, E., Yurchenko, V. & Rogozin, I. B. Analysis of stop
775 codons within prokaryotic protein-coding genes suggests frequent readthrough events. *Int.*
776 *J. Mol. Sci.* **22**, (2021).
- 777 44. Hart, R. K. & Prlić, A. SeqRepo: A system for managing local collections of biological
778 sequences. *PLoS One* **15**, (2020).
- 779 45. Bairoch, A. *et al.* The universal protein resource (UniProt). *Nucleic Acids Res.* **33**, D154–9
780 (2005).
- 781 46. UniProt Consortium. UniProt: The universal protein knowledgebase in 2023. *Nucleic Acids*
782 *Res.* **51**, D523–D531 (2023).
- 783 47. Blin, K. *et al.* antiSMASH 6.0: Improving cluster detection and comparison capabilities.
784 *Nucleic Acids Res.* **49**, W29–W35 (2021).

- 785 48. Feldgarden, M. *et al.* AMRFinderPlus and the Reference Gene Catalog facilitate
786 examination of the genomic links among antimicrobial resistance, stress response, and
787 virulence. *Sci. Rep.* **11**, 12728 (2021).
- 788 49. Chibani, C. M., Farr, A., Klama, S., Dietrich, S. & Liesegang, H. Classifying the unclassified:
789 A phage classification method. *Viruses* **11**, (2019).
- 790 50. Finn, R. D. *et al.* The Pfam protein families database. *Nucleic Acids Res.* **36**, D281–8
791 (2008).
- 792 51. Skinnider, M. A. *et al.* Genomes to natural products PRediction Informatics for Secondary
793 Metabolomes (PRISM). *Nucleic Acids Res.* **43**, 9645–9662 (2015).
- 794 52. Gibson, M. K., Forsberg, K. J. & Dantas, G. Improved annotation of antibiotic resistance
795 determinants reveals microbial resistomes cluster by ecology. *ISME J.* **9**, 207–216 (2015).
- 796 53. Haft, D. H. *et al.* TIGRFAMs: a protein family resource for the functional identification of
797 proteins. *Nucleic Acids Res.* **29**, 41–43 (2001).
- 798 54. Shen, W., Le, S., Li, Y. & Hu, F. SeqKit: A cross-platform and ultrafast toolkit for FASTA/Q
799 file manipulation. *PLoS One* **11**, e0163962 (2016).
- 800 55. Landrum, G. *RDKit: Open-Source Cheminformatics*. doi:10.5281/zenodo.10460537.
- 801 56. Heller, S. R., McNaught, A., Pletnev, I., Stein, S. & Tchekhovskoi, D. InChI, the IUPAC
802 International Chemical Identifier. *J. Cheminform.* **7**, 23 (2015).
- 803 57. Weininger, D. SMILES, a chemical language and information system. 1. Introduction to
804 methodology and encoding rules. *J. Chem. Inf. Comput. Sci.* **28**, 31–36 (1988).
- 805 58. Morgan, H. L. The Generation of a Unique Machine Description for Chemical Structures-A
806 Technique Developed at Chemical Abstracts Service. *J. Chem. Doc.* **5**, 107–113 (1965).
- 807 59. Gilchrist, C. L. M. & Chooi, Y.-H. Clinker & clustermap.js: Automatic generation of gene
808 cluster comparison figures. *Bioinformatics* (2021) doi:10.1093/bioinformatics/btab007.
- 809 60. Braesel, J., Lee, J.-H., Arnould, B., Murphy, B. T. & Eustáquio, A. S. Diazaquinomycin
810 biosynthetic gene clusters from marine and freshwater Actinomycetes. *J. Nat. Prod.* **82**,

- 811 937–946 (2019).
- 812 61. Torres, A. G. *et al.* Identification and characterization of IpfABCC'DE, a fimbrial operon of
813 enterohemorrhagic *Escherichia coli* O157:H7. *Infect. Immun.* **70**, 5416–5427 (2002).
- 814 62. Torres, A. G., Kanack, K. J., Tutt, C. B., Popov, V. & Kaper, J. B. Characterization of the
815 second long polar (LP) fimbriae of *Escherichia coli* O157:H7 and distribution of LP fimbriae
816 in other pathogenic *E. coli* strains. *FEMS Microbiol. Lett.* **238**, 333–344 (2004).
- 817 63. Torres, A. G., Zhou, X. & Kaper, J. B. Adherence of diarrheagenic *Escherichia coli* strains
818 to epithelial cells. *Infect. Immun.* **73**, 18–29 (2005).
- 819 64. Bäumler, A. J., Tsolis, R. M. & Heffron, F. The Ipf fimbrial operon mediates adhesion of
820 *Salmonella typhimurium* to murine Peyer's patches. *Proc. Natl. Acad. Sci. U. S. A.* **93**, 279–
821 283 (1996).
- 822 65. Osek, J., Weiner, M. & Hartland, E. L. Prevalence of the IpfO113 gene cluster among
823 *Escherichia coli* O157 isolates from different sources. *Vet. Microbiol.* **96**, 259–266 (2003).
- 824 66. Newton, H. J. *et al.* Contribution of long polar fimbriae to the virulence of rabbit-specific
825 enteropathogenic *Escherichia coli*. *Infect. Immun.* **72**, 1230–1239 (2004).
- 826 67. Korea, C.-G., Badouraly, R., Prevost, M.-C., Ghigo, J.-M. & Beloin, C. *Escherichia coli* K-12
827 possesses multiple cryptic but functional chaperone-usher fimbriae with distinct surface
828 specificities. *Environ. Microbiol.* **12**, 1957–1977 (2010).
- 829 68. Wurpel, D. J., Beatson, S. A., Totsika, M., Petty, N. K. & Schembri, M. A. Chaperone-usher
830 fimbriae of *Escherichia coli*. *PLoS One* **8**, e52835 (2013).
- 831 69. Qiao, J. *et al.* Construction of an *Escherichia coli* strain lacking fimbriae by deleting 64
832 genes and its application for efficient production of poly(3-hydroxybutyrate) and L-threonine.
833 *Appl. Environ. Microbiol.* **87**, e0038121 (2021).
- 834 70. Sayers, E. W. *et al.* Database resources of the national center for biotechnology
835 information. *Nucleic Acids Res.* **50**, D20–D26 (2022).
- 836 71. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**,

- 837 583–589 (2021).
- 838 72. Fergusson, C. H. *et al.* Discovery of a lagriamide polyketide by integrated genome mining,
839 isotopic labeling, and untargeted metabolomics. *Chem. Sci.* (2024)
840 doi:10.1039/D4SC00825A.
- 841 73. Albright, S. & Louca, S. Trait biases in microbial reference genomes. *Sci Data* **10**, 84
842 (2023).
- 843 74. Cimermancic, P. *et al.* Insights into secondary metabolism from a global analysis of
844 prokaryotic biosynthetic gene clusters. *Cell* **158**, 412–421 (2014).
- 845 75. *The Selection and Use of Essential Medicines: Report of the WHO Expert Committee on*
846 *Selection and Use of Essential Medicines, 2023.*
847 <https://www.who.int/publications/i/item/WHO-MHP-HPS-EML-2023.01> (2024).
- 848 76. Arp, J. *et al.* Synergistic activity of cosecreted natural products from amoebae-associated
849 bacteria. *Proc. Natl. Acad. Sci. U. S. A.* **115**, 3758–3763 (2018).
- 850 77. van Santen, J. A. *et al.* The Natural Products Atlas: An open access knowledge base for
851 microbial natural products discovery. *ACS Cent Sci* **5**, 1824–1833 (2019).
- 852 78. Mungan, M. D. *et al.* ARTS 2.0: feature updates and expansion of the Antibiotic Resistant
853 Target Seeker for comparative genome mining. *Nucleic Acids Res.* **48**, W546–W552
854 (2020).
- 855 79. Haft, D. H. *et al.* RefSeq: an update on prokaryotic genome annotation and curation.
856 *Nucleic Acids Res.* **46**, D851–D860 (2018).
- 857 80. Feldgarden, M. *et al.* Validating the AMRFinder tool and Resistance Gene Database by
858 using antimicrobial resistance genotype-phenotype correlations in a collection of isolates.
859 *Antimicrob. Agents Chemother.* **63**, (2019).
- 860 81. Humbard, M. A. & Maupin-Furlow, J. A. Prokaryotic proteasomes: nanocompartments of
861 degradation. *J. Mol. Microbiol. Biotechnol.* **23**, 321–334 (2013).
- 862 82. Valas, R. E. & Bourne, P. E. Rethinking proteasome evolution: two novel bacterial

- 863 proteasomes. *J. Mol. Evol.* **66**, 494–504 (2008).
- 864 83. Cromm, P. M. & Crews, C. M. The proteasome in modern drug discovery: Second life of a
865 highly valuable drug target. *ACS Cent Sci* **3**, 830–838 (2017).
- 866 84. Adams, J. The proteasome: a suitable antineoplastic target. *Nat. Rev. Cancer* **4**, 349–360
867 (2004).
- 868 85. Yeh, H.-H. *et al.* Resistance gene-guided genome mining: Serial promoter exchanges in
869 *Aspergillus nidulans* reveal the biosynthetic pathway for fellutamide B, a proteasome
870 inhibitor. *ACS Chem. Biol.* **11**, 2275–2284 (2016).
- 871 86. Rachid, S. *et al.* Mining the cinnabaramide biosynthetic pathway to generate novel
872 proteasome inhibitors. *Chembiochem* **12**, 922–931 (2011).
- 873 87. Owen, J. G. *et al.* Multiplexed metagenome mining using short DNA sequence tags
874 facilitates targeted discovery of epoxyketone proteasome inhibitors. *Proc. Natl. Acad. Sci.*
875 *U. S. A.* **112**, 4221–4226 (2015).
- 876 88. Gulder, T. A. M. & Moore, B. S. Salinosporamide natural products: Potent 20 S proteasome
877 inhibitors as promising cancer chemotherapeutics. *Angew. Chem. Int. Ed Engl.* **49**, 9346–
878 9367 (2010).
- 879 89. Eustáquio, A. S. *et al.* Biosynthesis of the salinosporamide A polyketide synthase substrate
880 chloroethylmalonyl-coenzyme A from S-adenosyl-L-methionine. *Proc. Natl. Acad. Sci. U. S.*
881 *A.* **106**, 12295–12300 (2009).
- 882 90. Fenical, W. *et al.* Discovery and development of the anticancer agent salinosporamide A
883 (NPI-0052). *Bioorg. Med. Chem.* **17**, 2175–2180 (2009).
- 884 91. Schorn, M. *et al.* Genetic basis for the biosynthesis of the pharmaceutically important class
885 of epoxyketone proteasome inhibitors. *ACS Chem. Biol.* **9**, 301–309 (2014).
- 886 92. Sugawara, K. *et al.* Eponemycin, a new antibiotic active against B16 melanoma. I.
887 Production, isolation, structure and biological activity. *J. Antibiot.* **43**, 8–18 (1990).
- 888 93. Meng, L., Kwok, B. H., Sin, N. & Crews, C. M. Eponemycin exerts its antitumor effect

- 889 through the inhibition of proteasome function. *Cancer Res.* **59**, 2798–2801 (1999).
- 890 94. Zabala, D. *et al.* A flavin-dependent decarboxylase-dehydrogenase-monooxygenase
891 assembles the warhead of α,β -epoxyketone proteasome inhibitors. *J. Am. Chem. Soc.* **138**,
892 4342–4345 (2016).
- 893 95. Tadokoro, T., Modahl, C. M., Maenaka, K. & Aoki-Shioi, N. Cysteine-rich secretory proteins
894 (CRISPs) from venomous snakes: An overview of the functional diversity in a large and
895 underappreciated superfamily. *Toxins* **12**, (2020).
- 896 96. Schneiter, R. & Di Pietro, A. The CAP protein superfamily: function in sterol export and
897 fungal virulence. *Biomol. Concepts* **4**, 519–525 (2013).
- 898 97. Hirano, T. *et al.* CAP peptide artificially induces insect gall. *bioRxiv* 2024.01.06.574462
899 (2024) doi:10.1101/2024.01.06.574462.
- 900 98. Jupp, S. *et al.* The EBI RDF platform: linked open data for the life sciences. *Bioinformatics*
901 **30**, 1338–1339 (2014).
- 902 99. Chevrette, M. G. *et al.* The antimicrobial potential of *Streptomyces* from insect
903 microbiomes. *Nat. Commun.* **10**, 516 (2019).
- 904 100. Varadi, M. *et al.* AlphaFold Protein Structure Database: massively expanding the structural
905 coverage of protein-sequence space with high-accuracy models. *Nucleic Acids Res.* **50**,
906 D439–D444 (2022).
- 907 101. van Kempen, M. *et al.* Fast and accurate protein structure search with Foldseek. *Nat.*
908 *Biotechnol.* **42**, 243–246 (2024).
- 909 102. Kim, H., Mirdita, M. & Steinegger, M. Foldcomp: a library and format for compressing and
910 indexing large protein structure sets. *Bioinformatics* **39**, (2023).
- 911 103. Schütze, K., Heinzinger, M., Steinegger, M. & Rost, B. Nearest neighbor search on
912 embeddings rapidly identifies distant protein relations. *Front Bioinform* **2**, 1033775 (2022).
- 913 104. Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a language
914 model. *Science* **379**, 1123–1130 (2023).

- 915 105. El-Gebali, S. *et al.* The Pfam protein families database in 2019. *Nucleic Acids Res.* **47**,
916 D427–D432 (2019).
- 917 106. Bernardes, J. S., Vieira, F. R. J., Zaverucha, G. & Carbone, A. A multi-objective
918 optimization approach accurately resolves protein domain architectures. *Bioinformatics* **32**,
919 345–353 (2016).
- 920 107. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069
921 (2014).
- 922 108. Sfiligoi, I. *et al.* The pilot way to grid resources using glideinWMS. in *2009 WRI World*
923 *Congress on Computer Science and Information Engineering* vol. 2 428–432 (IEEE, 2009).
- 924 109. The Open Science Grid Executive Board on behalf of the Osg Consortium: Ruth Pordes *et*
925 *al.* The open science grid. *J. Phys. Conf. Ser.* **78**, 012057 (2007).

926

927

928 **Figure 1.** A partial schema illustration of a SocialGene database, showing nodes (circles) and
929 their relationships (lines between circles). The visualization was auto-generated and formatted
930 by connecting the RefSeq-based Neo4j database to yFiles¹³⁹ Neo4j Explorer. A high-resolution
931 version is available (see Data and Code Availability), and an up-to-date, auto-generated, full
932 node-relationship schema is available in SocialGene's online documentation.

933

934 **Figure 2.** A simplified illustration of two BGCs from MIBiG, their encoded proteins, and shared
935 pHMM annotations (gray lines). pHMMs models are labeled with numbers representing the log
936 of their outdegree (e.g. 4 is approximately 1,000 relationships, 6 is approximately 100,000).
937 Searching SocialGene databases for similar BGCs is accelerated by a first-pass search of
938 annotations by low-outdegree pHMMs. The illustration also highlights the complexity of shared
939 pHMM annotations between proteins within a single BGC. A comprehensive comparison of the
940 two displayed BGCs was previously published by Braesel et. al.⁶⁰

941

942 **Figure 3.** SocialGene's BGC search function outputs a clustermap.js⁵⁹ plot, as shown. The
943 abridged plot here displays two target BGCs, obtained from searching >343,000 genomes for
944 BGCs similar to the mupirocin-producing BGC (BGC0000182, middle row). Two result hits are
945 displayed in the top and bottom rows, with median RBHs of 73.8% and 58.5% to BGC0000182,
946 respectively, as calculated by DIAMOND BLASTp. Individual alignment identities are shown in
947 Supplementary Fig. [10](#).

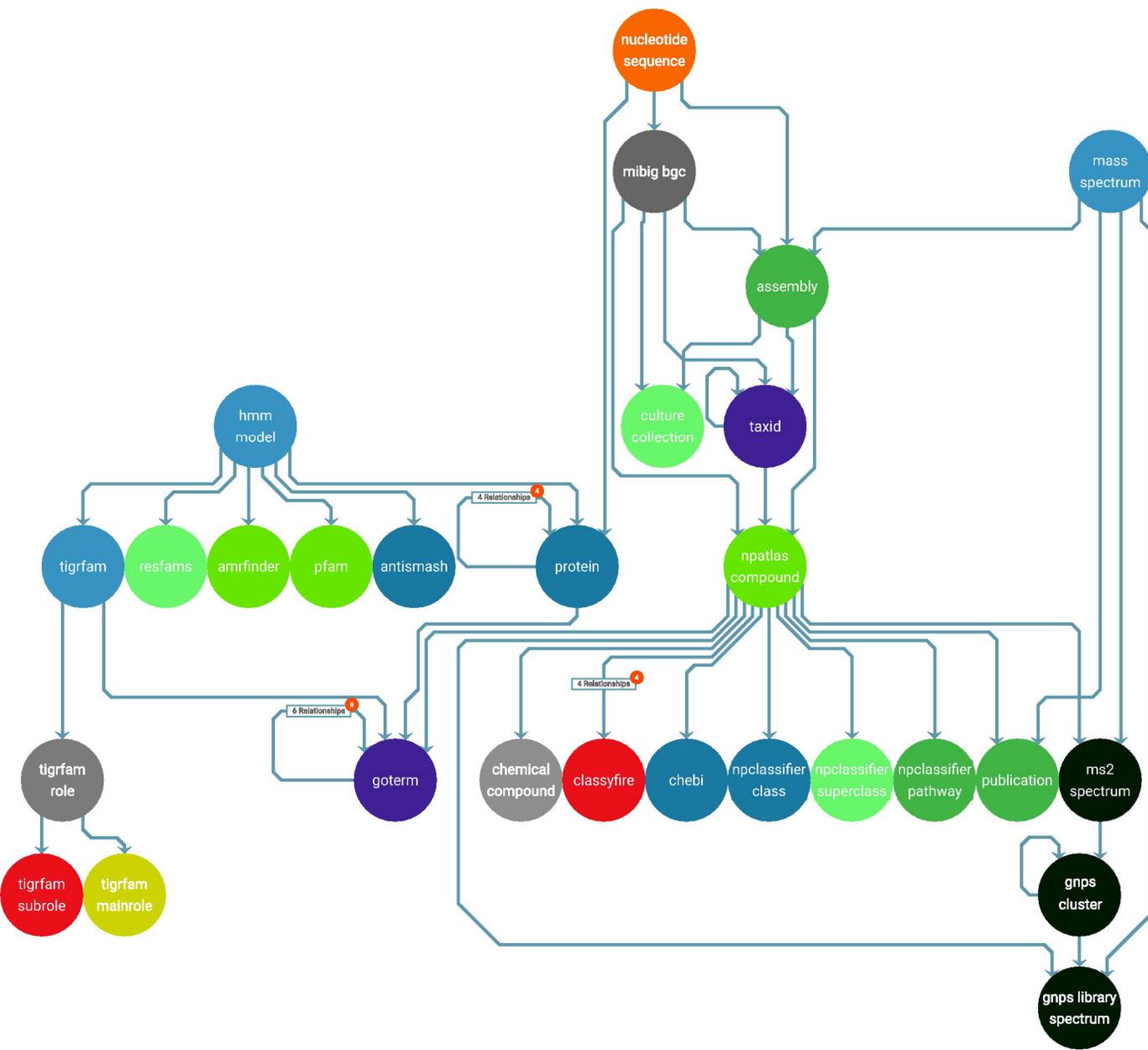
948

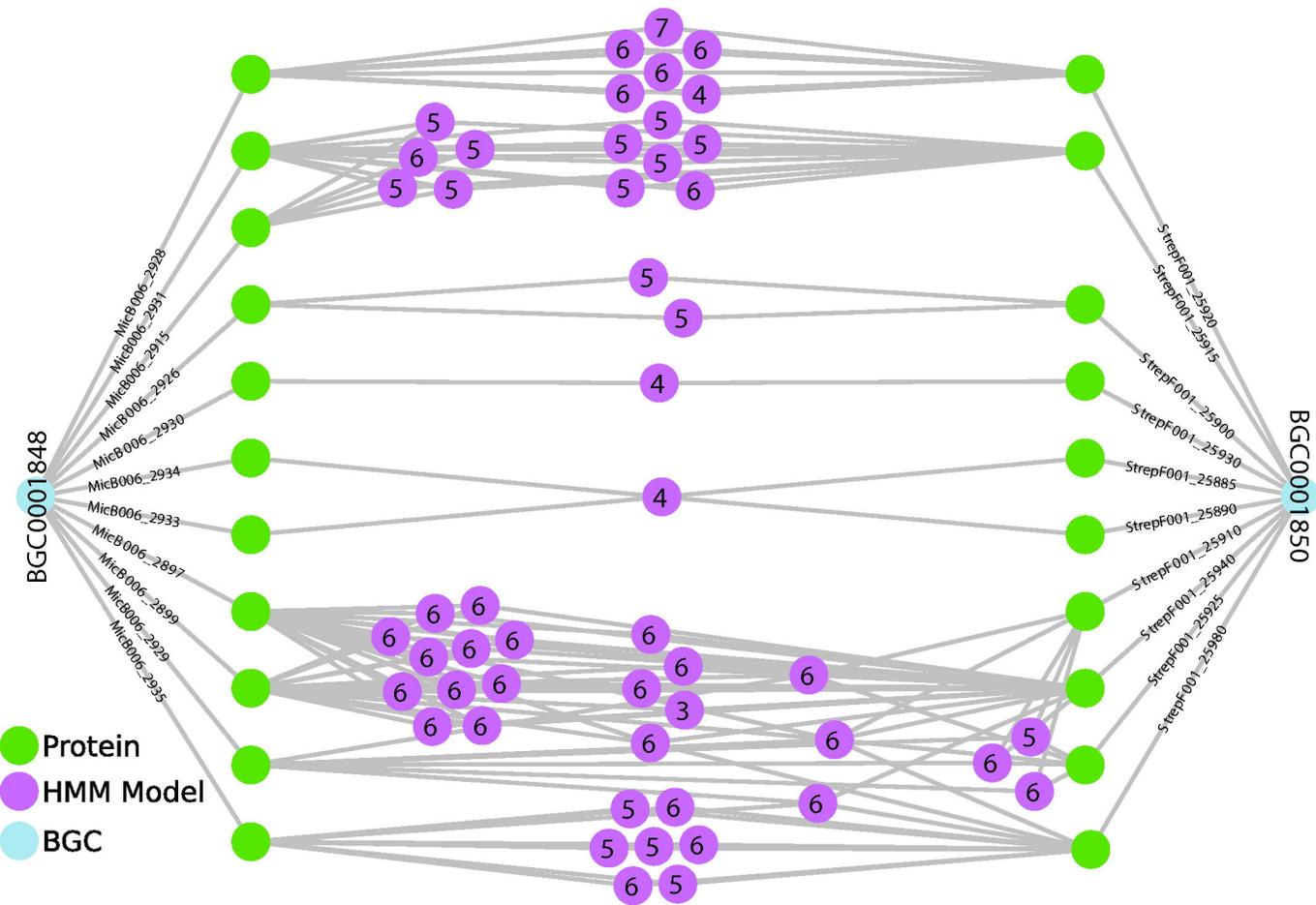
949 **Figure 4.** A SocialGene database containing more than 340,000 RefSeq genome assemblies
950 and MIBiG BGCs was searched for gene regions complete and syntenic to MIBiG BGC000946
951 (encoding for vibrioferrin). The 6,581 points in the graph represent the resulting target BGCs
952 and the y-axis represents the median protein identity of a BGC's reciprocal best hits (RBHs) to
953 BGC000946 proteins. Target BGCs were sorted in the x-axis by median RBHs to BGC000946
954 proteins. All MIBiG BGCs were labeled and highlighted in red and are known vibrioferrin
955 producing BGCs, except BGC0001408 for which an associated chemical structure has not been
956 reported.

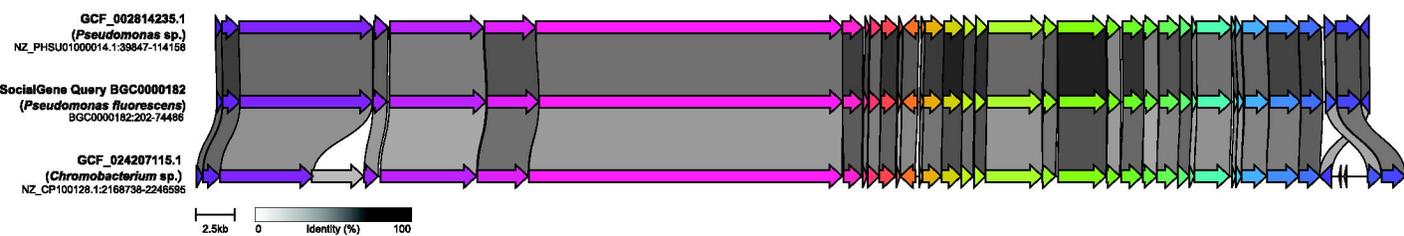
957

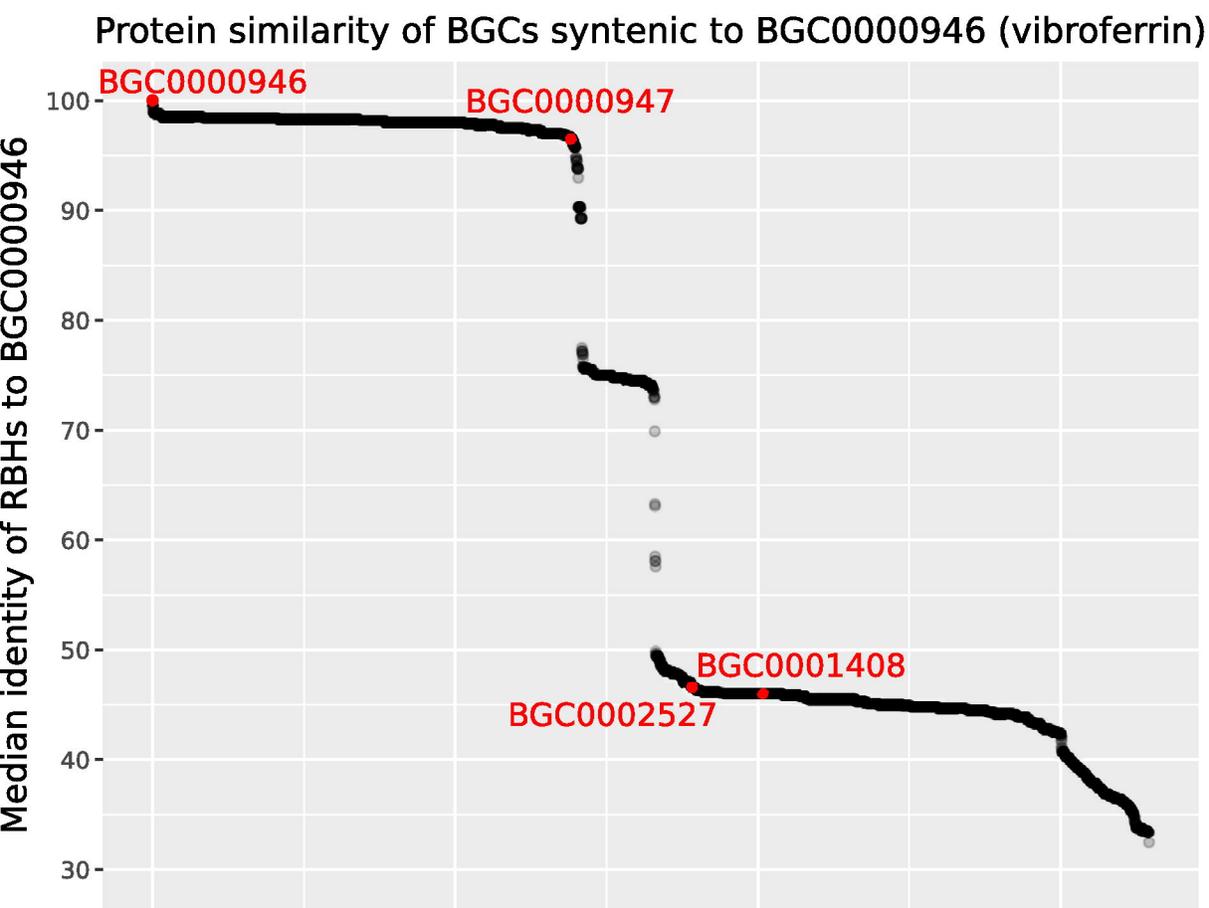
958 **Figure 5.** Using a single Neo4j Cypher statement, new links (edges between nodes) were
959 created between MIBiG BGCs (nodes; some enlarged and labeled by SM product) and
960 genomes assemblies (nodes). New links were created when an assembly contained an
961 antiSMASH predicted BGC with proteins that were at least 70% similar to 70% of a MIBiG
962 BGC's proteins, as determined by traversing MMseqs2 protein cluster relationships. This figure
963 is only a subset of the resulting subgraph (54% of total nodes) and highlights how SocialGene
964 can be used to study complex distributions of BGCs, at scale. See Supplementary Figs. [13-15](#),
965 for the full subgraph.

966









- *Escherichia*
- *Klebsiella*
- *Shigella*
- *Enterobacter*

