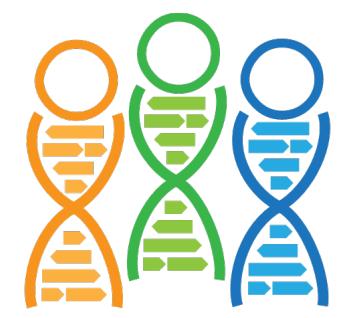


Towards near-instant, repository-scale searching for homologous BGCs with SocialGene

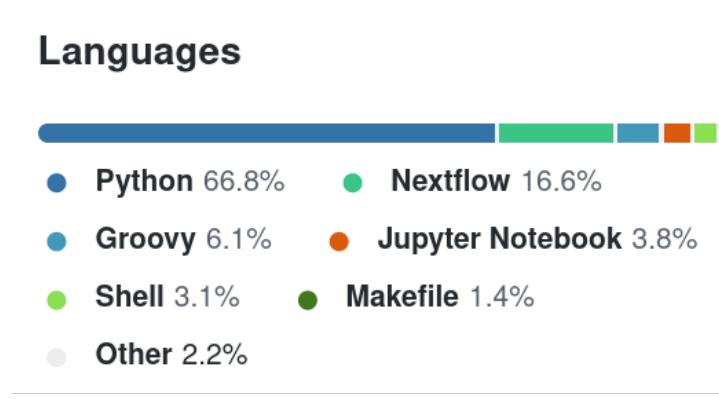


Chase M. Clark, Jason C. Kwan
Division of Pharmaceutical Sciences, School of Pharmacy
University of Wisconsin-Madison,
777 Highland Avenue, Madison, WI 53705, USA

A common workflow in many genomic-based discovery groups is to take a known biosynthetic gene cluster (BGC) and search its individual proteins against NCBI's non-redundant database (nr) using nucleotide (BLASTn) or protein based alignments (BLASTp). A number of tools have been developed to aid this process (e.g. MultiGeneBlast¹, cblaster², clusterfinder³ etc.) However, when trying to include the genomic context of results when searching large databases (RefSeq >230 thousand genomes; GenBank >1.1 million genomes) the results returned from a single protein search often exceed what these programs and/or the NCBI BLAST API can handle (most searches silently limit the number of results).

Socialgene was built to provide a fast and easy way to search large numbers of genomes for similar individual proteins, and BGCs. It is alignment free, relying instead searching through pre-annotated protein domains via a graph database (other tools using domain similarity include micropan, RAMPAGE, ENTS, CO-ED, ClusterFinder, etc.).³⁻⁵ Additionally, there are no inbuilt assumptions of what a BGC should look like, which allows the flexibility of finding BGCs that are broken/split due to biology or incomplete sequencing/assembly. Two common alignment-based tools, BLAST (through Diamond⁶) and MMseqs2⁷, are also incorporated and future searching may take advantage of pre-clustering by MMseqs2.

27,580 HMM models Pfam, TIGRFAM, AMRFinder, PRISM, antiSMASH, BiG-SLiCE, and more?



Mapping MIBiG

As we're ultimately interested in BGCs, it seemed like a good sanity check would be to see how a SocialGene database "looked" when created with only known BGCs.

Therefore, 1,910 BGCs were downloaded from the Minimum Information about a Biosynthetic Gene cluster (MIBiG) database and the graph database created using SocialGene's Nextflow workflow:

```
nextflow run nextflow \
--gbk_input 'mibig_gbk_2.0.tar.gz' \
-profile simple_custom_input,conda
```

Duration: 22m
CPU Hours: 6.2

The resulting database was queried and all proteins and their HMM annotations exported as a graph/network for visualization with Gephi (the graph to the right). Red circles represent the non-redundant proteins present in MIBiG BGCs, while blue circles represent HMM models (protein domains/motifs). Resulting clusters were fairly homogenous in protein function and some have been labeled.

Interactive version:
socialgene.github.io/mibigmap

Searching for Diazaquinomycin BGCs

Diazaquinomycin (DAQ) H and J are selective inhibitors of *Mycobacterium tuberculosis* that stalled in preclinicals due to their horrible aqueous solubility. While just a simple test case, finding BGCs encoding potential chemical analogs could potentially open new routes to compound druggability. DAQs have been found in multiple genera, including *Micromonospora* spp and *Streptomyces* spp.⁸⁻¹⁰ For this test case I uploaded the MIBiG DAQ BGC to Socialgene's web interface, connected to a socialgene database containing all 226 *Micromonospora* genomes currently available in RefSeq. With its success I then repeated the experiment using >1,200 *Streptomyces* genomes (all assemblies except contig-quality).

Search the database,
using a single BGC as input

(Input expects a gbk or gbff file)

Description:

Document: No file selected.

	<i>Micromonospora</i>	<i>Streptomyces</i>
Genomes	226	1,202
Contigs/Scaffolds	34,856	186,574
Proteins	1,373,641	8,561,870
Non-Redundant Proteins	1,076,690	6,236,718
HMM Annotations	10,656,754	65,467,422
MMSEQS2	1,009,056	6,234,349
Bidirectional BLASTP (45.3 GB)	420,788,981	Not performed
Database Creation (Real Duration)	4h 30m	22h
Database Creation (CPU Hours)	133	524

