

A Multi-Repository Scale Genomic and Chemical Search Engine to Enable the Discovery, Production, and Function of Natural Products

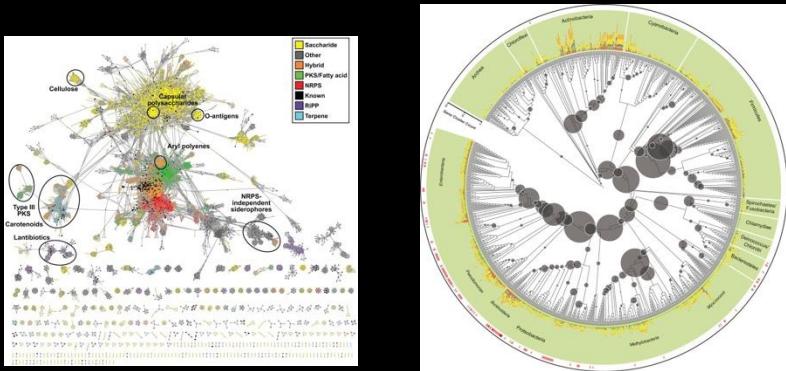
Chase Clark

Postdoctoral Associate

Jason Kwan's Group

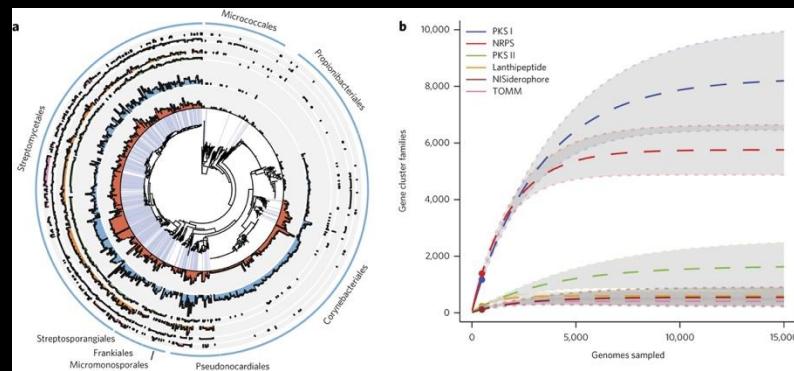
University of Wisconsin-Madison

2014:
1,154 Genomes; 33,351 BGCs; 2,400 GCFs



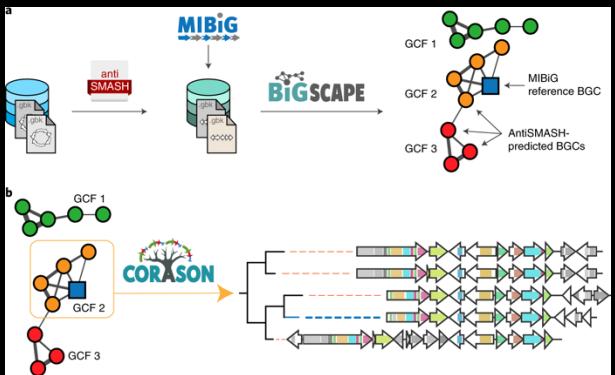
Cimermancic P, Medema MH, Claesen J, Kurita K, Wieland Brown LC, Mavrommatis K, Pati A, Godfrey PA, Koehrsen M, Clardy J, Birren BW, Takano E, Sali A, Linington RG, Fischbach MA. Insights into secondary metabolism from a global analysis of prokaryotic biosynthetic gene clusters. *Cell*. 2014

2014:
830 Genomes; 11,422 BGCs; 4,122 GCFs



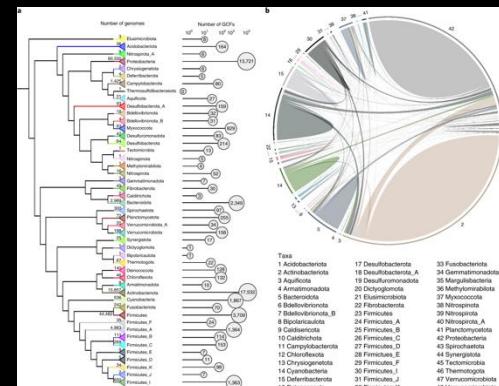
Doroghazi JR, Albright JC, Goering AW, Ju KS, Haines RR, Tchalukov KA, Labeda DP, Kelleher NL, Metcalf WW. A roadmap for natural product discovery based on large-scale genomics and metabolomics. *Nat Chem Biol*. 2014

2020:
3,080 Genomes; 74,652 BGCs

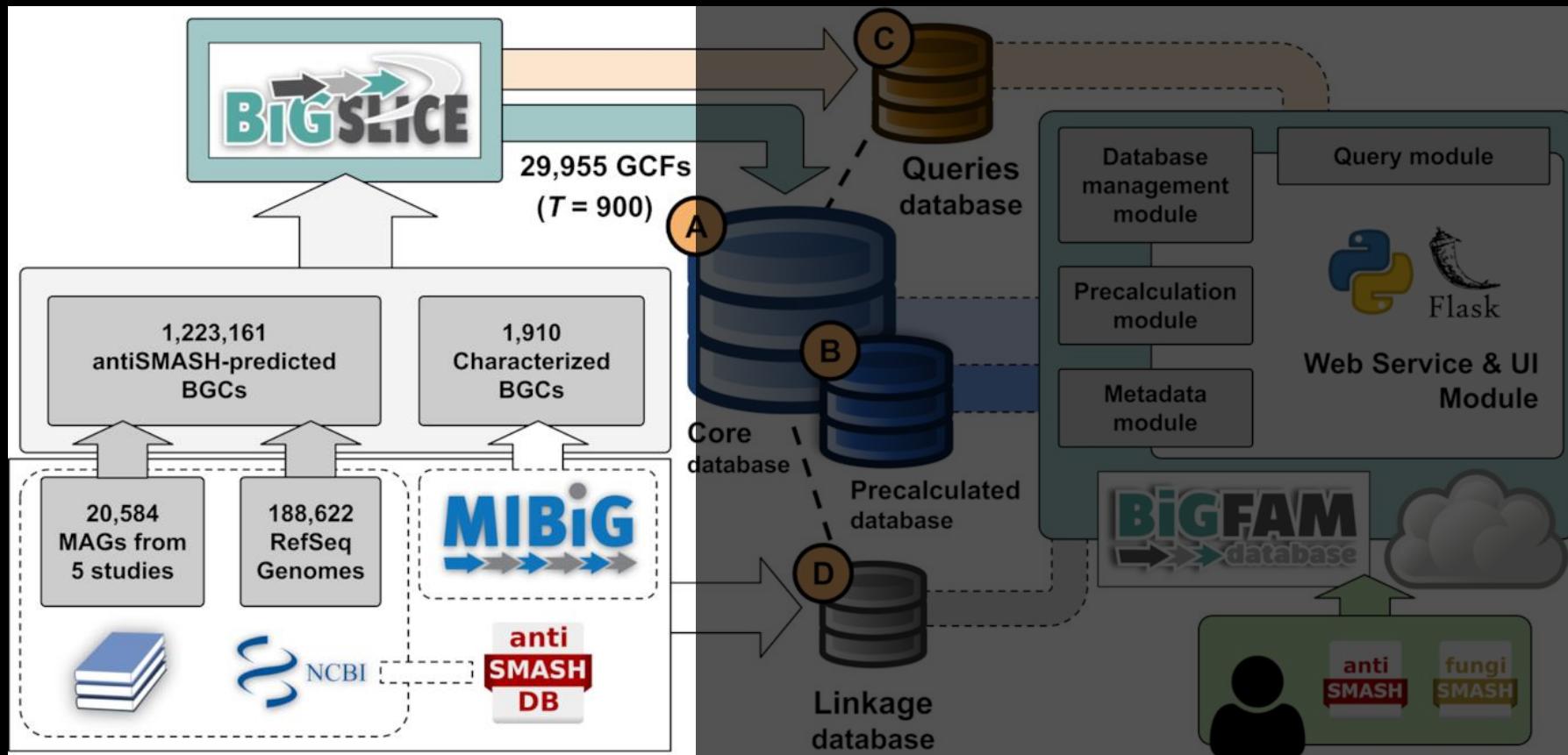


Navarro-Muñoz JC, Selem-Mojica N, Mullowney MW, Kautsar SA, Tryon JH, Parkinson EI, De Los Santos ELC, Yeong M, Cruz-Morales P, Abubucker S, Roeters A, Lokhorst W, Fernandez-Guerra A, Cappelini LTD, Goering AW, Thomson RJ, Metcalf WW, Kelleher NL, Barona-Gomez F, Medema MH. A computational framework to explore large-scale biosynthetic diversity. *Nat Chem Biol*. 2020

2022:
170,000 Genomes, 1,185,995 BGCs; 62,449 GCFs



Gavrilidou A, Kautsar SA, Zaburannyi N, Krug D, Müller R, Medema MH, Ziemert N. Compendium of specialized metabolite biosynthetic diversity encoded in bacterial genomes. *Nat Microbiol*. 2022

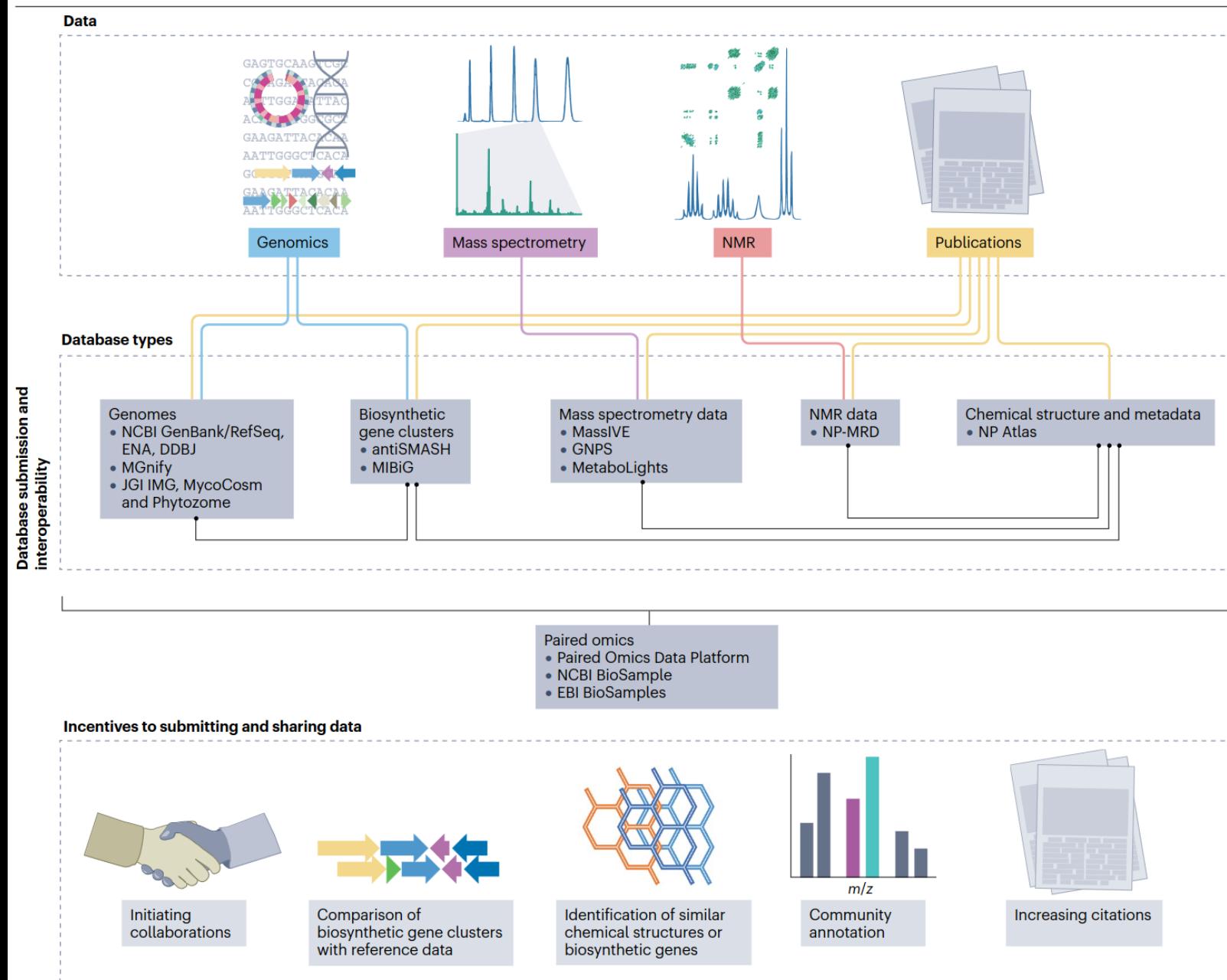


Kautsar SA, Blin K, Shaw S, Weber T, Medema MH. BiG-FAM: the biosynthetic gene cluster families database. Nucleic Acids Res. 2021

Kautsar SA, van der Hooft JJJ, de Ridder D, Medema MH. BiG-SLiCE: A highly scalable tool maps the diversity of 1.2 million biosynthetic gene clusters. Gigascience. 2021

Artificial intelligence for natural product drug discovery

Michael W. Mullowney  ^{1,62}, Katherine R. Duncan  ^{2,62}, Somayah S. Elsayed  ^{3,62}, Neha Garg  ^{1,62}, Justin J. van der Hooft  ^{5,6,62}, Nathaniel I. Martin  ^{2,62}, David Meijer  ^{5,62}, Barbara R. Terlouw  ^{5,62}, Friederike Biermann  ^{5,6,9}, Kai Bin  ¹⁰, Janani Durairaj  ¹¹, Marina Gorostola González  ^{12,19}, Eric J. N. Helfrich  ¹⁹, Florian Huber  ¹⁴, Stefan Leopold-Messer  ¹⁵, Kohulan Rajan  ¹⁶, Tristan de Rond  ¹⁷, Jeffrey A. van Santen  ¹⁸, Maria Sorokina  ^{19,20}, Marcy J. Balunas  ^{21,22}, Mehdi A. Benidir  ²³, Doris A. van Bergeijk  ³, Laura M. Carroll  ²⁴, Chase M. Clark  ²⁵, Djoek-Arné Clevert  ²⁶, Chris A. Dejong  ²⁷, Chao Du  ³, Scarlet Ferrinho  ²⁸, Francesca Grisoni  ^{29,30}, Albert Hofstetter  ³¹, Willem Jespers  ³², Olga V. Kalinina  ^{32,33,34}, Satria A. Kautsar  ³⁵, Hyunwoo Kim  ³⁶, Tiago F. Leao  ³⁷, Joleen Masschalein  ^{38,39}, Evan R. Rees  ²⁹, Raphael Reher  ^{40,41}, Daniel Reker  ^{42,43}, Philippe Schwaller  ⁴⁴, Marwin Segler  ⁴⁵, Michael A. Skinner  ^{22,46}, Allison S. Walker  ^{47,48}, Egon L. Willighagen  ⁴⁹, Barbara Zdravil ⁵⁰, Nadine Ziemert ⁵¹, Rebecca J. M. Goss ⁵², Pierre Guyomard ⁵², Andrea Volkamer ^{34,53}, William H. Gerwick ⁵⁴, Hyun Uk Kim ⁵⁵, Rolf Müller ^{32,56,57,58}, Gilles P. van Wezel ^{3,39}, Gerard J. P. van Westen ^{12,59}, Anna K. H. Hirsch ^{32,56,57,58}, Roger G. Linington ¹⁸, Serina L. Robinson ⁶⁰ & Marnix H. Medema ^{5,61}.



What if we just connected all the things?

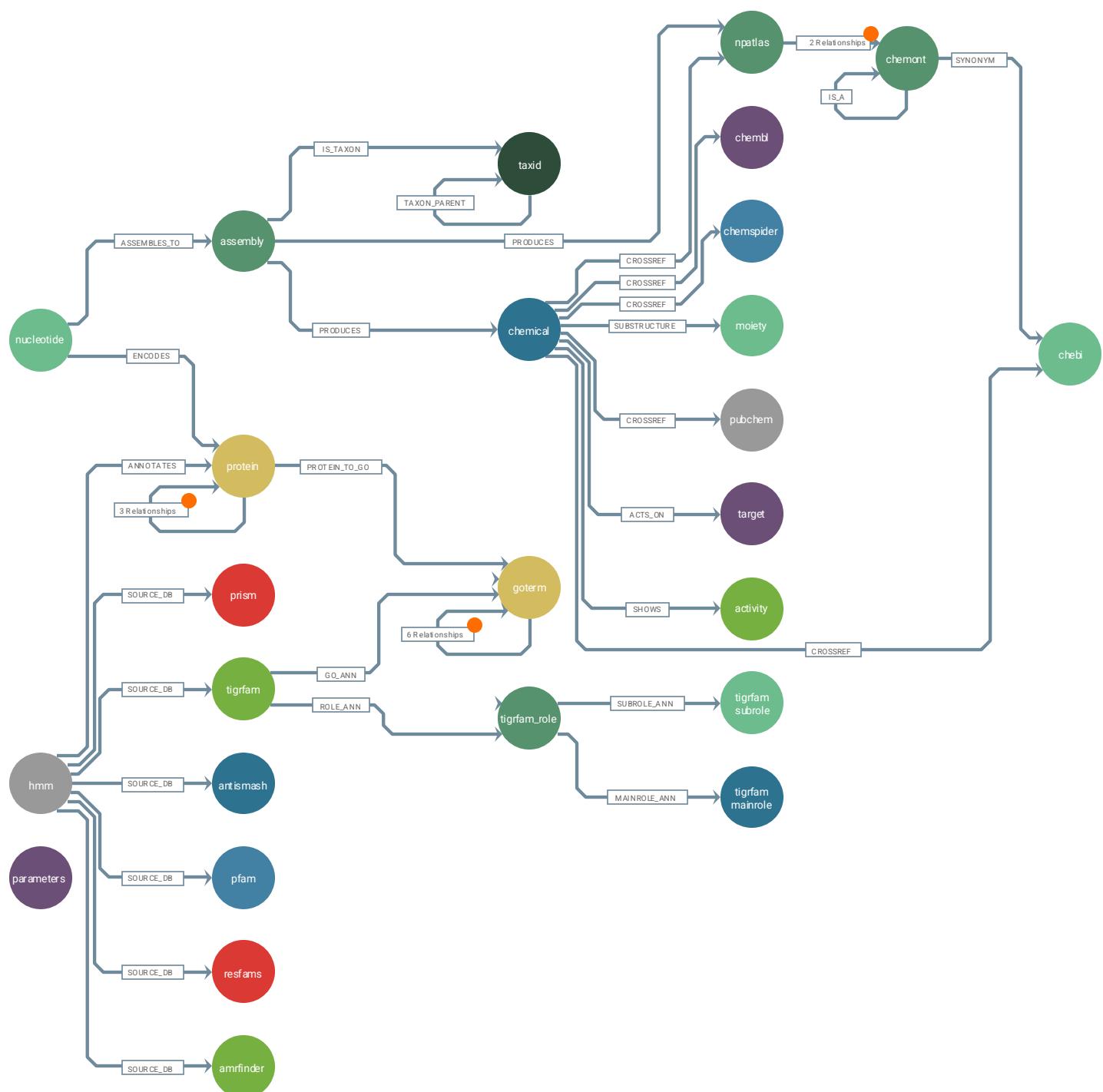


SocialGene

nextflow

python™

neo4j

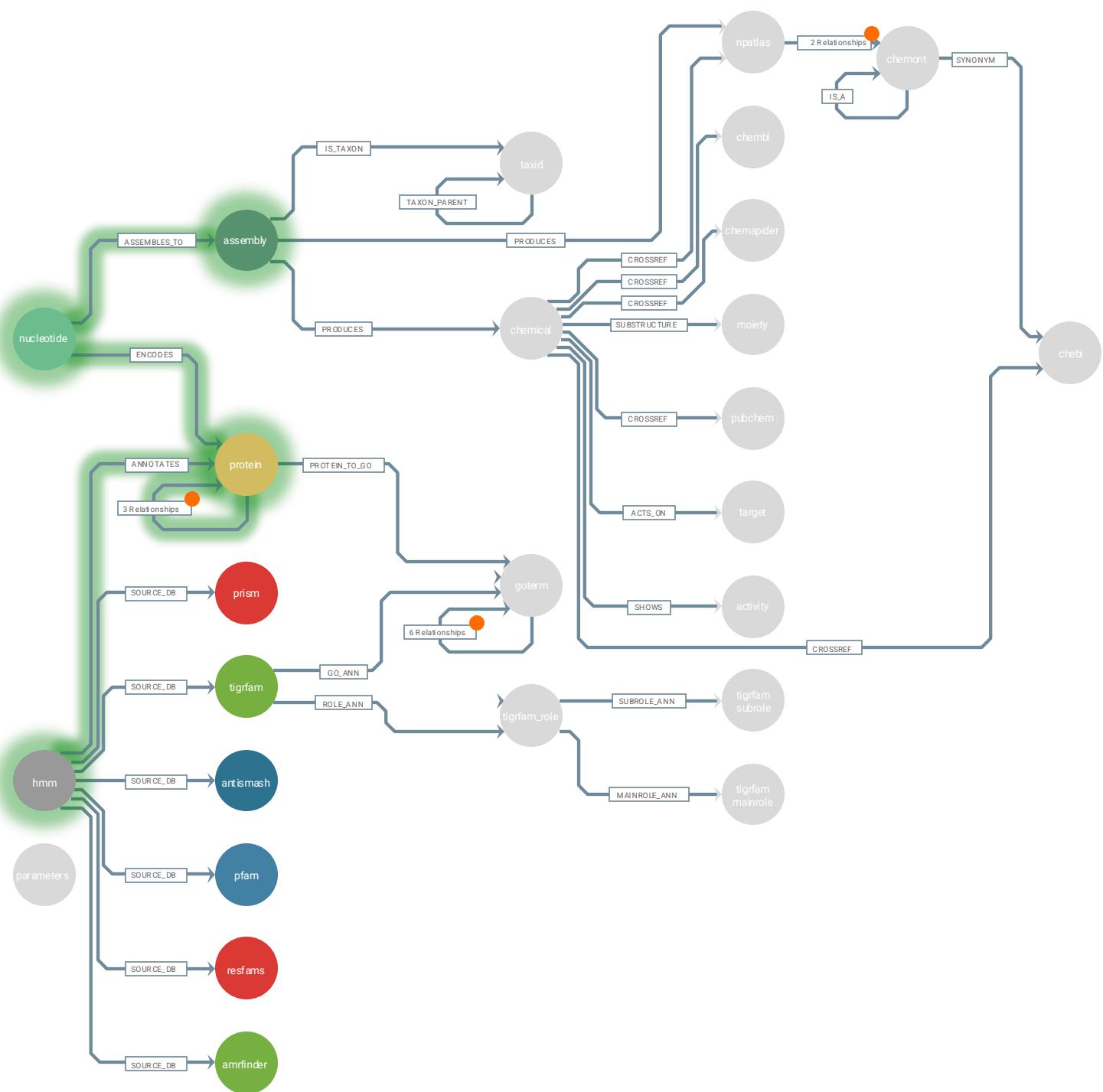


SocialGene

nextflow

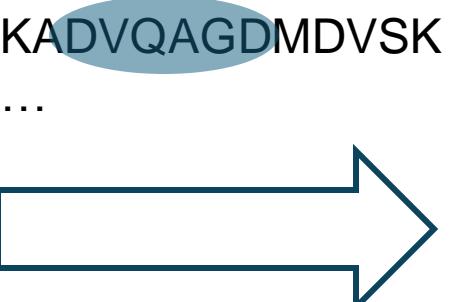
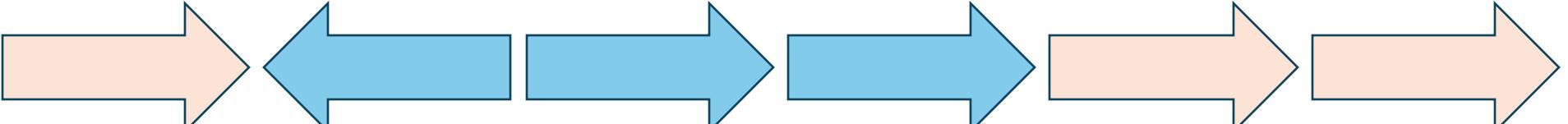
python™

neo4j

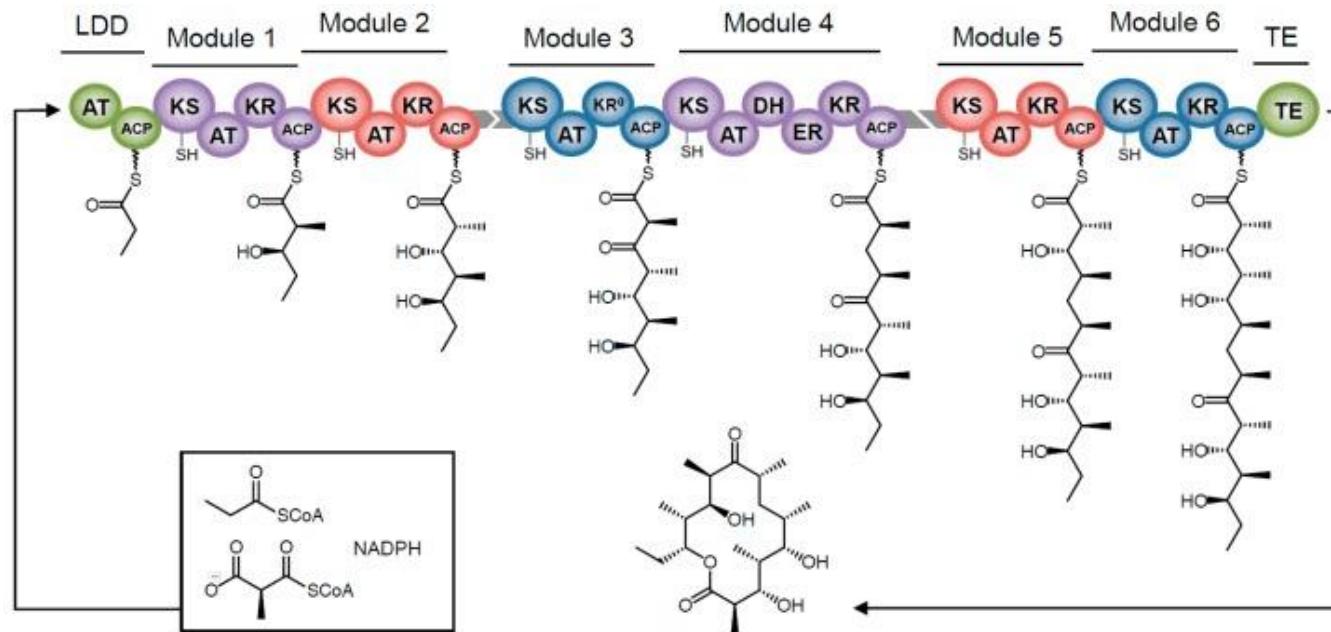


Molecular Biology in < 1 minute

DNA → RNA, RNA → protein

DNA (nucleotides)	ATGTCCAACGCC...	CGCGGCATCCTC...	GGCGCCGTGCTC ...
RNA (nucleotides)	UTG TCC UUC GCC ...	CGC GGC UTC CTC ...	GGC GCC GTG CTC ...
Proteins (amino acids)	MSNARATHLRRGI ...	KADVQAGDMDVSK ...	AKSGPWTFKDDRG T...
Protein have domains	MSNARATHLRRGI ...	KADVQAGDMDVSK ...	AKSGPWTFKDDRG T...
Genes encode proteins			
Gene Cluster			

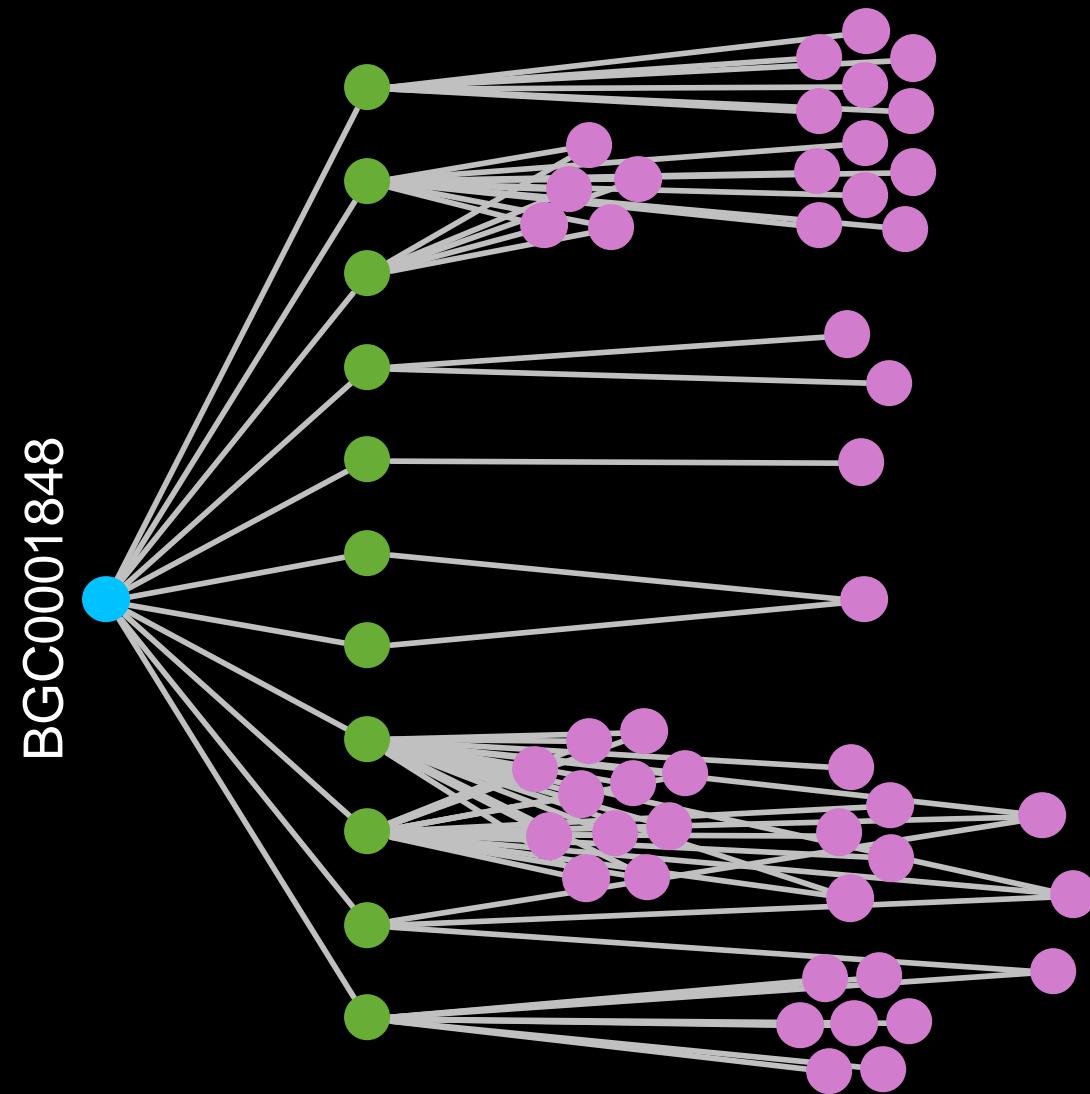
Polyketide synthases are large, multidomain proteins “Beads on a String”



Bayly, Carmen L, and Vikramaditya G Yadav. "Towards Precision Engineering of Canonical Polyketide Synthase Domains: Recent Advances and Future Prospects." *Molecules (Basel, Switzerland)* vol. 22, 2 235. 5 Feb. 2017, doi:10.3390/molecules22020235



Similar domains = similar function



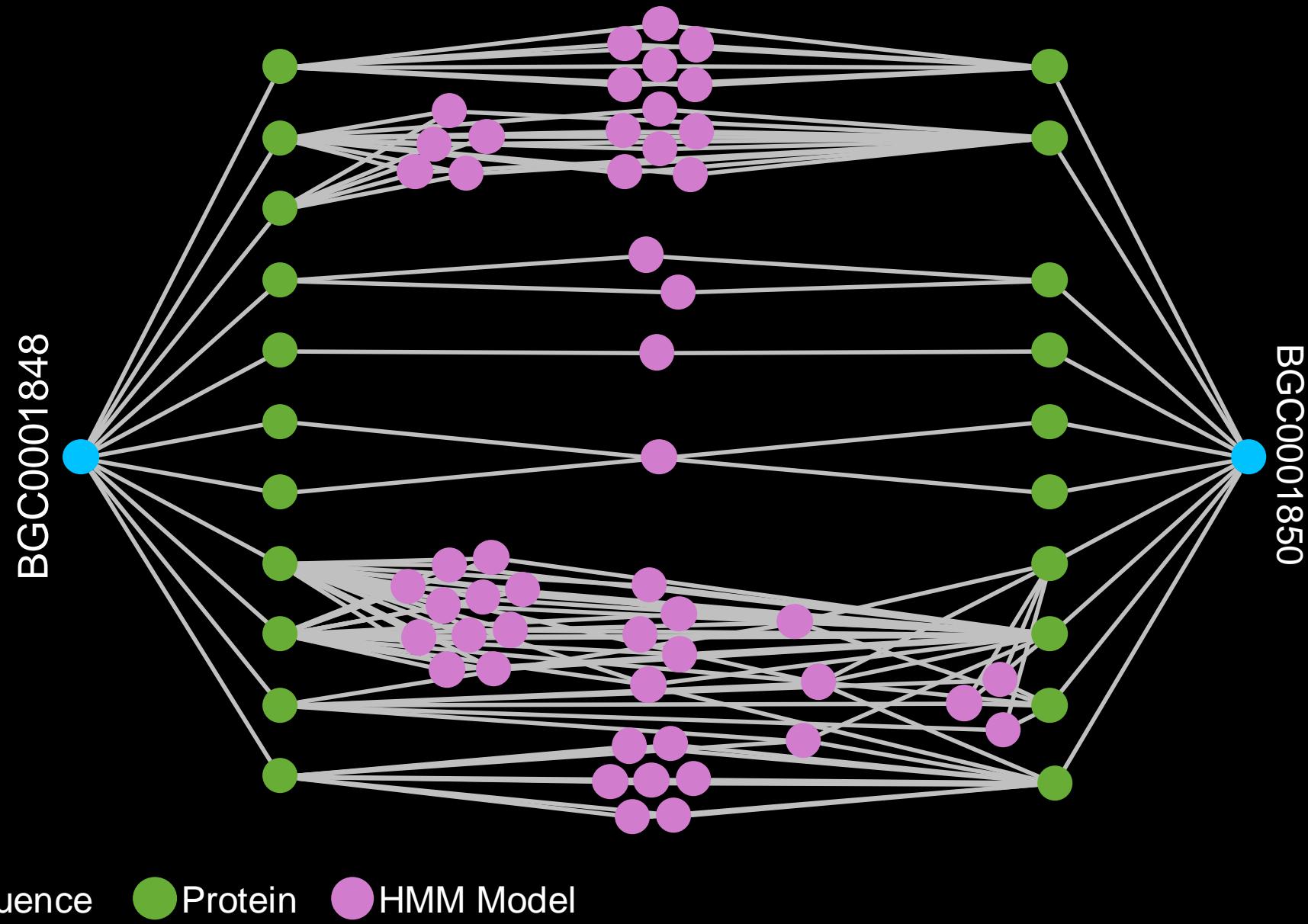
● DNA Sequence

● Protein

● HMM Model

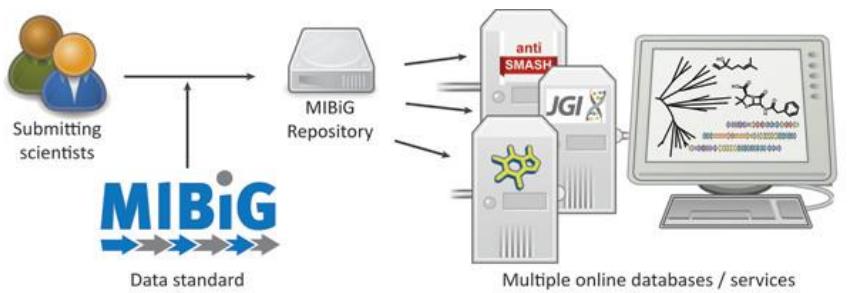


Similar domains = similar function



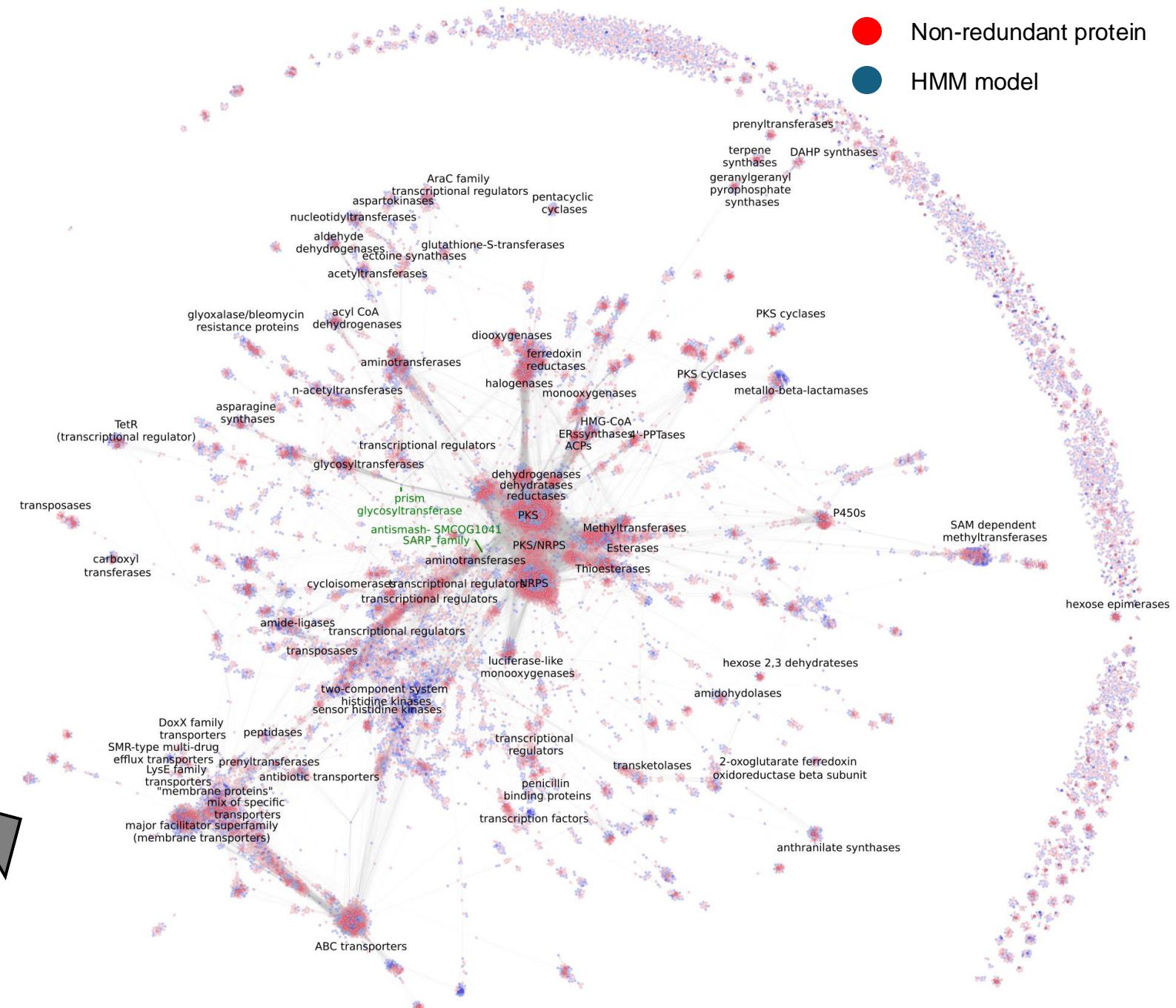
Application

Use domains to cluster proteins by function



1. Connecting genes to chemistry
2. Understanding BGC environmental diversity
3. Computer-guided gene cluster engineering

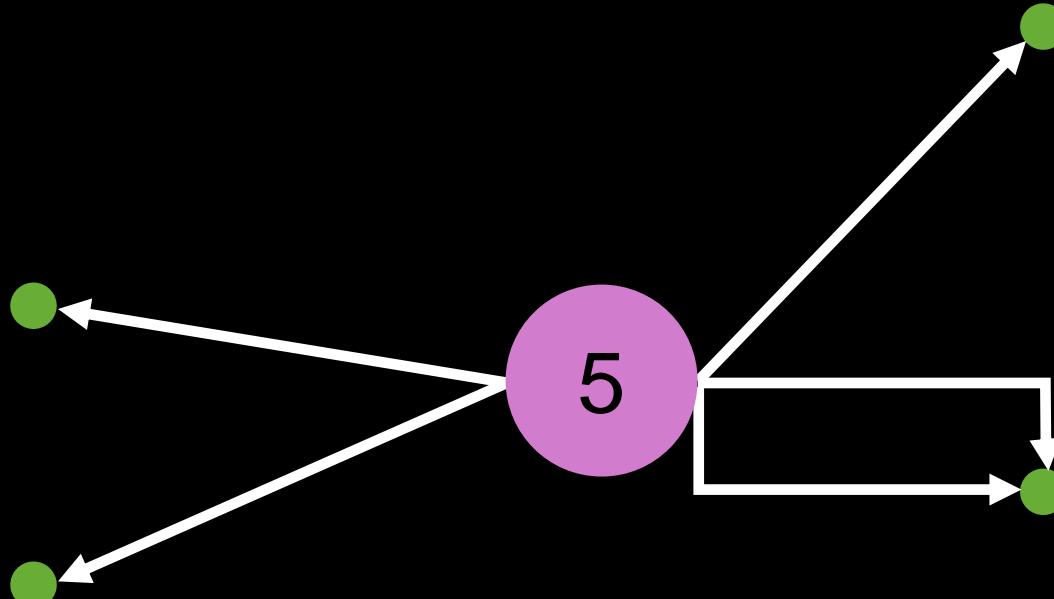
SocialGene ↗



Application

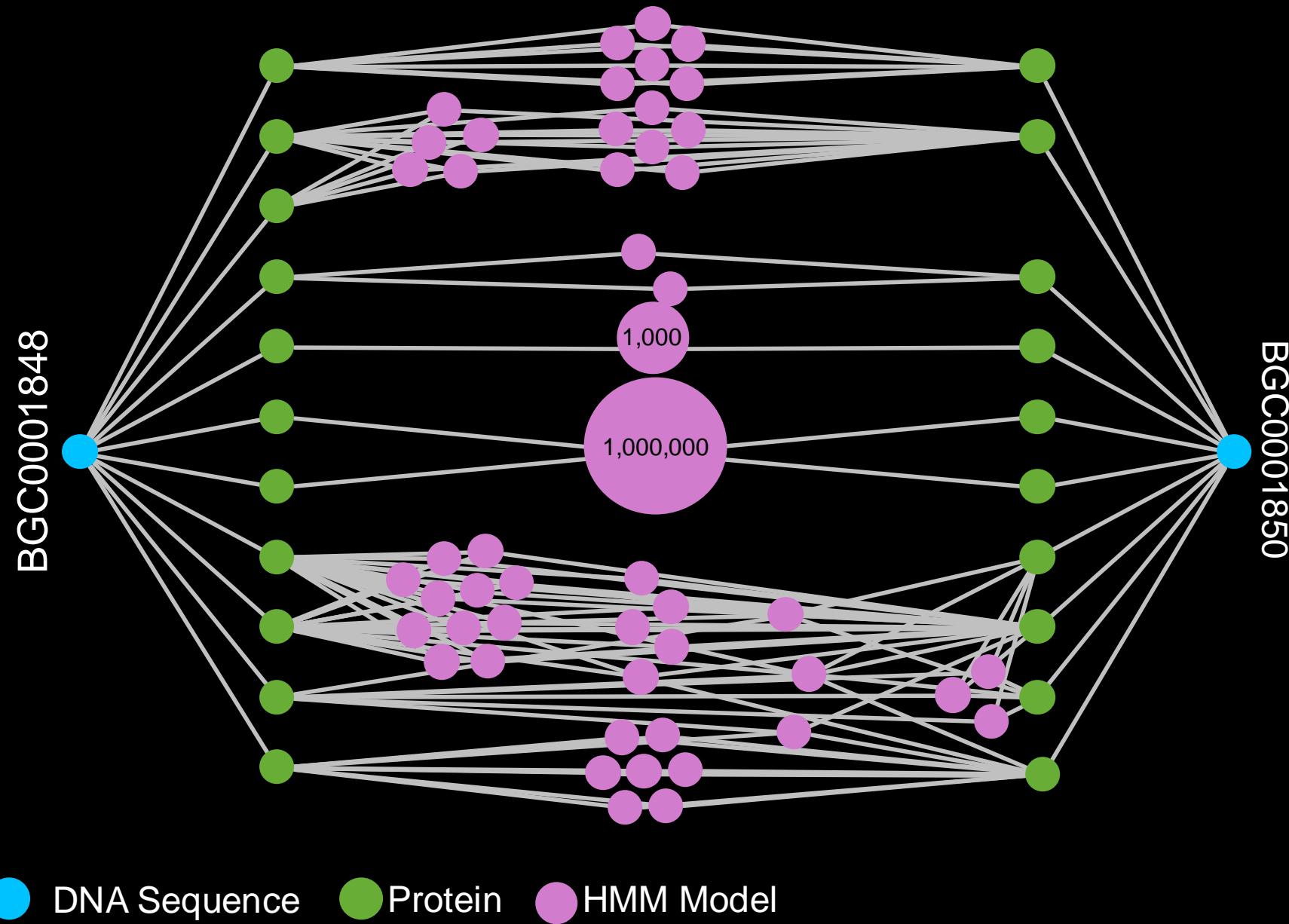
Fast Search for BGCs
with
Remote Homology
at
Really-Big-Scale™

Outdegree = number of annotations that have to be looked at



● Protein ● HMM Model

Outdegree = number of annotations that have to be looked at



Search genomes for remotely- homologous BGCs

- Annotate input BGC domains

Search genomes for remotely- homologous BGCs

- Annotate input BGC domains
- Prioritize domains

Search genomes for remotely- homologous BGCs

- Annotate input BGC domains
- Prioritize domains
- Search database for domains

Search genomes for remotely- homologous BGCs

- Annotate input BGC domains
- Prioritize domains
- Search database for domains
- Find proteins with those domains

Search genomes for remotely- homologous BGCs

- Annotate input BGC domains
- Prioritize domains
- Search database for domains
- Find proteins with those domains
- Find protein locations on genomes

Search genomes for remotely- homologous BGCs

- Annotate input BGC domains
- Prioritize domains
- Search database for domains
- Find proteins with those domains
- Find protein locations on genomes
- Create putative BGCs

Search genomes for remotely- homologous BGCs

- Annotate input BGC domains
- Prioritize domains
- Search database for domains
- Find proteins with those domains
- Find protein locations on genomes
- Create putative BGCs
- Reciprocal Best Hit BLASTp

Search genomes for remotely- homologous BGCs

- Annotate input BGC domains
- Prioritize domains
- Search database for domains
- Find proteins with those domains
- Find protein locations on genomes
- Create putative BGCs
- Reciprocal Best Hit BLASTp
- Evaluate and rank
 - # shared genes
 - Synteny
 - BLASTp scores

Search genomes for remotely- homologous BGCs

- Annotate input BGC domains
- Prioritize domains
- Search database for domains
- Find proteins with those domains
- Find protein locations on genomes
- Create putative BGCs
- Reciprocal Best Hit BLASTp
- Evaluate and rank
 - # shared genes
 - Synteny
 - BLASTp scores
- Create Clinker plot

Search genomes for remotely- homologous BGCs

- Annotate input BGC domains
 - Prioritize domains
 - Search database for domains
 - Find proteins with those domains
 - Find protein locations on genomes
 - Create putative BGCs
 - Reciprocal Best Hit BLASTp
 - Evaluate and rank
 - # shared genes
 - Synteny
 - BLASTp scores
 - Create Clinker plot

protein	Locus/Descriptor	Unique HMM models	Mean	Min	25%	50%	75%	Max	Sum
aUclZlfJ196v15TfvLMaxGpnFc6rMs	CAA60468_1 no-locus-tag pteridine-dependent dioxygenase	1	221	221	221	221	221	221	221
ZKLBeQ7BRLGL8XlsYiYKsLi0gngr2	CAA60470_1 no-locus-tag methyltransferase	2	6735	6670	6702	6735	6768	6801	13471
R10ScC0MvCcyeRnnvB1BL6H4s1MztBmm	CAA60475_1 no-locus-tag None	1	30410	30410	30410	30410	30410	30410	30410
M1bcQ0Vi0hy0DxmzS15hWmIppHTN4	CAA60467_1 no-locus-tag lysine cyclodeaminase	3	13983	10392	11031	11670	15779	19889	41951
X16t-PXLtYKdkwFGK7vzacykhB530CN	CAA60464_1 no-locus-tag ferredoxin	4	23639	13134	22087	26542	28094	28337	94556
tJTHFx2UjgexIdgvwsIZD3B70xUeJt	CAA60454_1 no-locus-tag membrane transport protein	1	132073	132073	132073	132073	132073	132073	132073
tWdDCKsDW8ctciyG0dmP9k_GpdrobM	CAA60453_1 no-locus-tag None	1	132073	132073	132073	132073	132073	132073	132073
pCX16XS2A1sfsb01-7z9KqzhqTIR4D	CAA60467_1 no-locus-tag cystathione synthase	3	77232	60272	70668	81065	85712	90360	231697
92Yftht6aohBzGLV18cgr3MdrkH7wq	CAA60469_1 no-locus-tag cytochrome P450	4	156969	134825	155495	163854	165328	165343	627877
MhLyg_V0j3mc_p23jxag6BzG9f9gy	CAA60465_1 no-locus-tag cytochrome P450	4	156969	134825	155495	163854	165328	165343	627877
F_b3n09j1dDrc78mLlsn54Zh42QvJA	CAA60472_1 no-locus-tag helix-turn-helix DNA binding protein	5	143578	92094	149242	157078	159721	159758	717893
kLkZQzqmmMcS1slnJQdy-E Moy-khTlRE	CAA60456_1 no-locus-tag sensory protein kinase	5	193755	83075	167049	179315	237946	301392	968777
Rz-ZXQZPy0IVh695de5KsL97tDP9v	CAA60471_1 no-locus-tag regulator of cholesterol oxidase	4	314414	80322	265902	330347	378860	516642	1257659
c3RvctUetwT3zMBwVj-xYyoE3qu8p	CAA60455_1 no-locus-tag response regulator	4	326690	185874	207945	302123	420867	516642	130676
U08eDj123grdKzB4ZvBZ18Btadxy6y	CAA60450_1 no-locus-tag None	5	279823	3890	279261	340896	373765	401305	1399117
76p9UrLFBWz5GsfTx10403dkhJA1V8	CAA60473_1 no-locus-tag membrane transport protein	3	716873	608660	672791	736982	771010	805039	215062
ZX6AN34Q-k1oq2AeHPT6WnHnZ5dA	CAA60449_1 no-locus-tag monosaccharide transporter	4	644651	323135	615869	725215	753996	805039	2578604
wTA_WvycakL0KA1mNzKwH0pOnNRUBR	CAA60466_1 no-locus-tag methyltransferase	13	203643	8495	126319	208972	255533	360429	2647359
Fr-Iy2L7A3cJou4R1-wKz6ZDM0Nf1o	CAA60463_1 no-locus-tag methyltransferase	15	178774	16799	16799	144653	249502	360429	2681611
PHRwBKt10u6Fw6wcjuCmc4BwXp	CAA60451_1 no-locus-tag regulator of antibiotic transport complexes	6	524097	42715	508847	618048	663888	713677	3144584
RQfbfx-0B83QzIDPPavLGIQV1cHDI	CAA60452_1 no-locus-tag ABC-transporter	4	955107	17624	925889	1256626	1285843	1289552	3820248
8821v1_uw8VOXGLm_Rx3nt7E106H4r	CAA60458_1 no-locus-tag ketoreductase/dehydrogenase	14	546655	323394	520336	560294	602080	685782	7653170
91UrL7zPvLuJxqdG6vut4UppiW6um_R	CAA60461_1 no-locus-tag pieperolate incorporating enzyme	38	246096	2027	127413	208274	296364	944217	9351666
kTy1YDZGAwMV1rd6tDmUD5MAYRkr	CAA60459_1 no-locus-tag polyketide synthase	59	206147	206225	131559	172651	239874	685782	12162718
FQSeRRg5631GJRz	CAA60462_1 no-locus-tag polyketide synthase	62	201669	16567	122501	173462	233673	685782	1250351
Qcou2BemtHGbZ0w	CAA60460_1 no-locus-tag polyketide synthase	69	234345	960	130116	176911	294729	944217	1619689

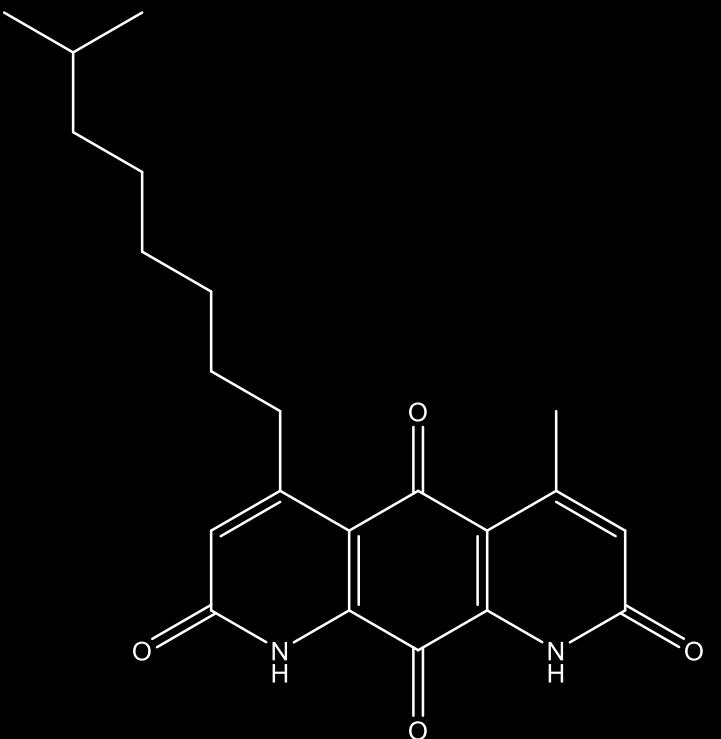
```
INFO Prioritizing input proteins by outdegree
INFO 'max_outdegree' is set to 1,000,000, will remove any domains with a higher outdegree
INFO 'max_outdegree' reduced the total outdegree from 82,446,501 to 78,643,697
INFO 'max_domains_per_protein' is set to 3, will remove domains from proteins from highest to lowest outdegree
INFO 'max_domains_per_protein' reduced the total outdegree from 78,643,697 to 11,544,444
INFO 'query_proteins' is set to 5, will limit search to 5 of 28 input proteins
INFO 'query_proteins' reduced the total outdegree from 11,544,444 to 158,979
INFO Total number of proteins: 28
```

Outdegree of input protein domains											
rotein	Locus/Descriptor		Unique HMM models		Mean	Min	25%	50%	75%	Max	Sum
aUcLZlzfJ196v15TfvLMaxGpnFc6rMs	CAA60468_1	no-locus-tag	pteridine-dependent dioxygenase	1	221	221	221	221	221	221	221
ZKLUBe07BRLGL8XlsY1YKsl1Qngr2	CAA60470_1	no-locus-tag	methyltransferase	2	6735	6702	6735	6768	6801	13471	
PHRwBK10u6uFwFcjuCmnc4BxWm	CAA60452_1	no-locus-tag	ABC-transporter	1	17624	17624	17624	17624	17624	17624	
FQSeRRgS31GJR8_Qco28EmtHbz20w	CAA60460_1	no-locus-tag	polyketide synthase	3	6436	960	1371	1782	9174	16567	19309
R08ScQMVccyErnnVBLQH6H4s1MztBmm	CAA60475_1	no-locus-tag	None	1	30410	30410	30410	30410	30410	30410	
JUr7zPvLuJxqCD6vtu4pPp1W6um_R	CAA60459_1	no-locus-tag	polyketide synthase	59	206147	20625	131559	172651	239874	685782	12167218

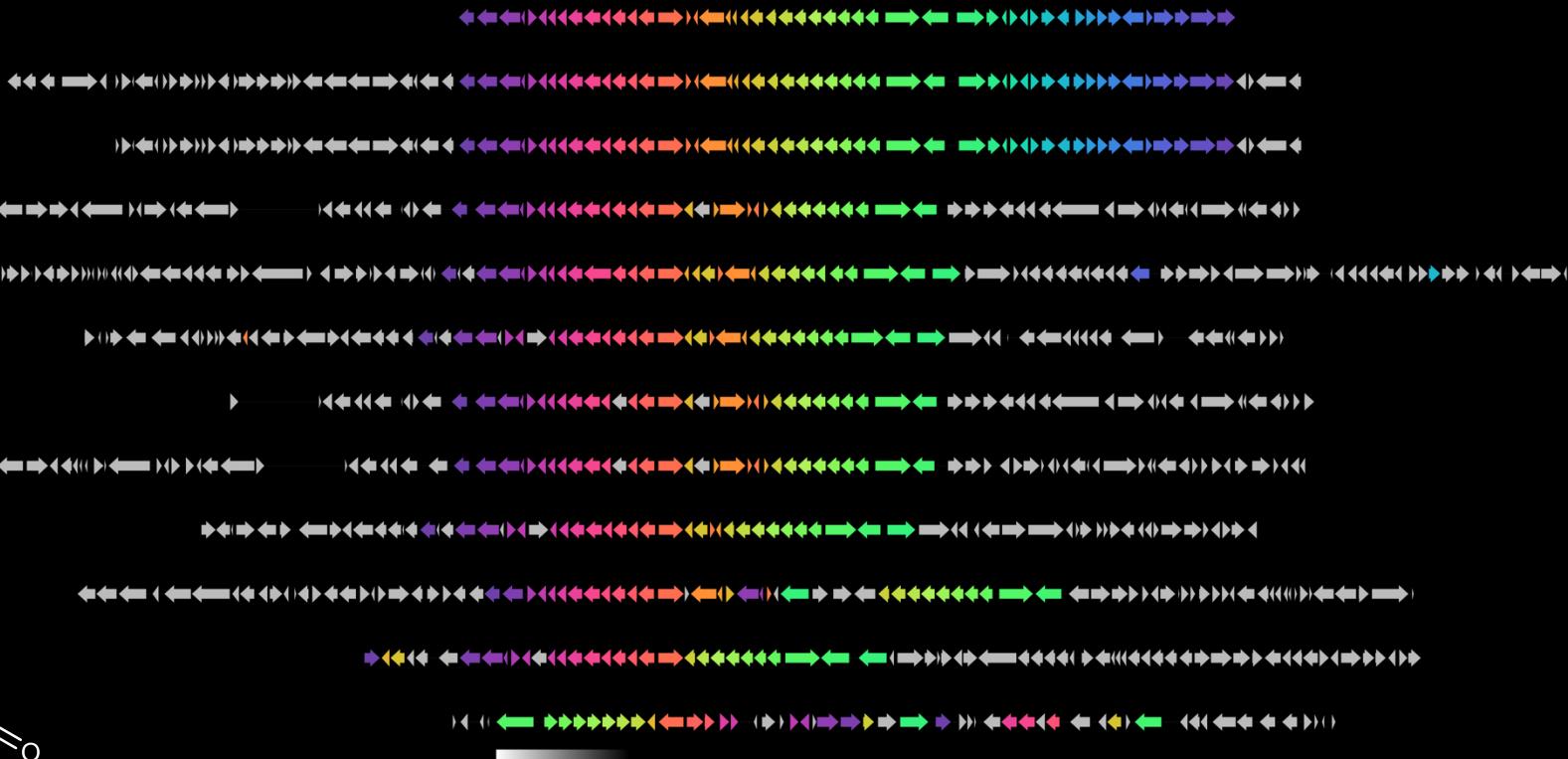
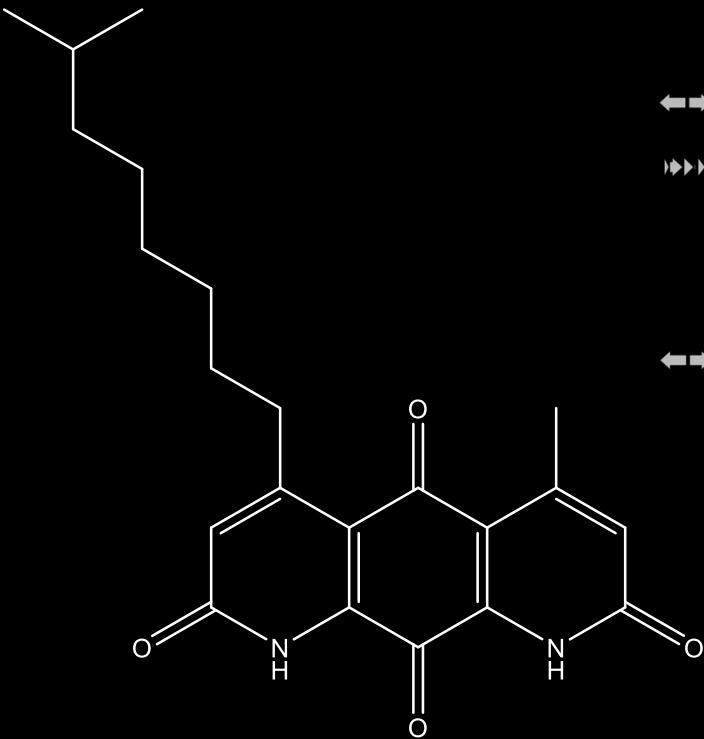
```
INFO    Searching database for proteins with similar domain content and all of the genomes those are found in
4-02-26 12:04:53 INFO    Initial search returned 209,038 proteins, found in 27,879 genomes
4-02-26 12:04:53 INFO    Starting with matches across 28789 genomes
4-02-26 12:04:53 INFO    Filtering on assembly_uid where unique hits >= 0.4
4-02-26 12:04:53 INFO    Counting unique hits per assembly_uid
4-02-26 12:04:53 INFO    14965 assemblies, 80149 nucleotide sequences had assembly_uids with >= 3 unique query hits
4-02-26 12:04:53 INFO    Filtering on nucleotide_uid where unique hits >= 0.4
4-02-26 12:04:53 INFO    Counting unique hits per nucleotide_uid
4-02-26 12:04:54 INFO    4848 assemblies, 4907 nucleotide sequences had nucleotide_uids with >= 3 unique query hits
4-02-26 12:04:54 INFO    Grouping protein hits if less than 20000 bp apart
4-02-26 12:04:55 INFO    Sorting genes by start position
4-02-26 12:04:55 INFO    Pulling data from the database for 28 putative BGCs
4-02-26 12:04:55 INFO    Time to fill: 0 seconds
4-02-26 12:04:56 INFO    Start: Creating links
4-02-26 12:04:56 INFO    Finding reciprocal best hits; protein similarity via Diamond BLASTp
4-02-26 12:05:07 INFO    Finish: Creating links; 279 links produced
4-02-26 12:05:07 INFO    Start: Assigning target BGC proteins to input BGC groups
4-02-26 12:05:07 INFO    Finish: Assigning target BGC proteins to input BGC groups
4-02-26 12:05:07 INFO    Writing clustermap.js output to: /home/chase/Downloads/clinker/plot/data.json
4-02-26 12:05:07 INFO    Creating clustermap.js clusters
4-02-26 12:05:07 INFO    Creating clustermap.js links
4-02-26 12:05:07 INFO    Creating clusterman.js links
```



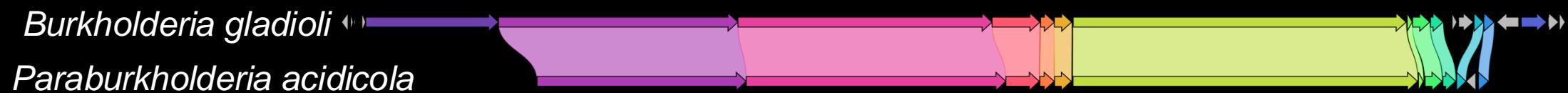
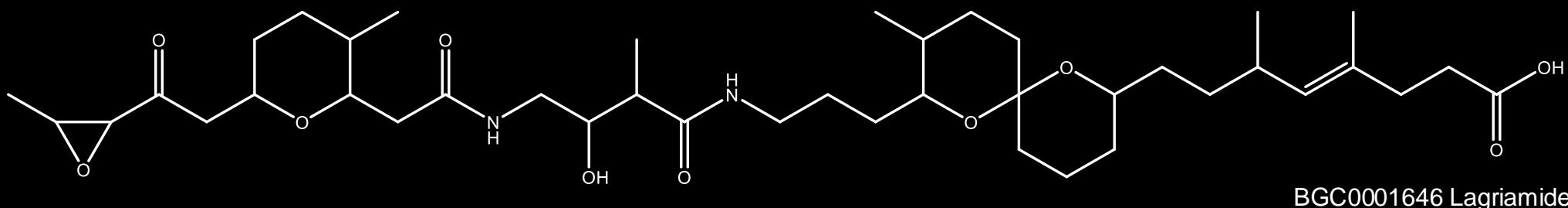
Searching for BGC0001848 in >300k genomes



Searching for BGC0001848
in >300k genomes
filter for organisms in culture collections



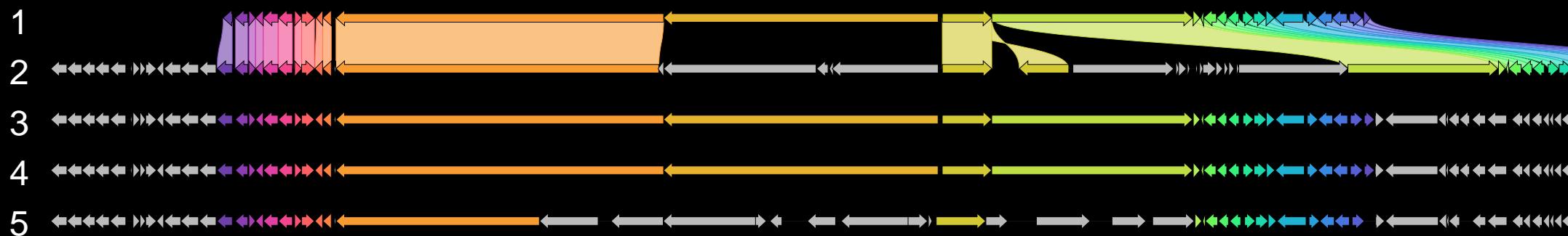
Find isolates containing metagenome-discovered BGCs



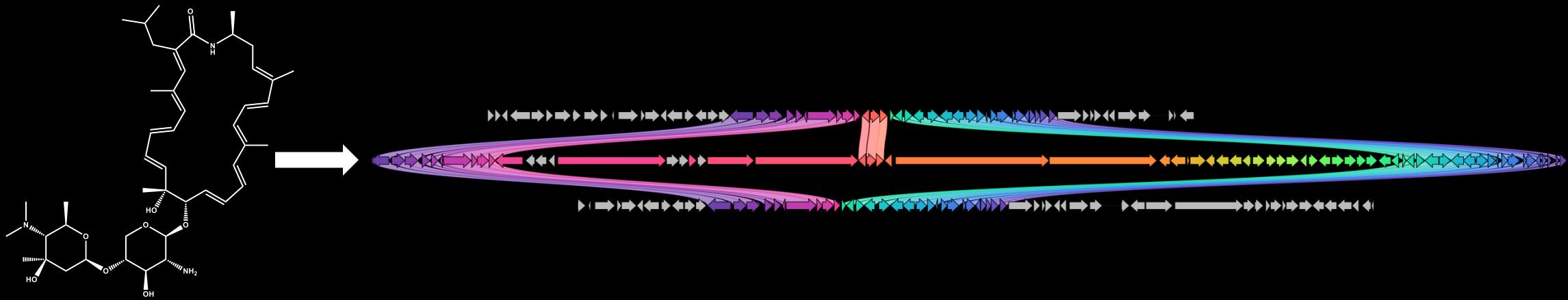
Fergusson, C. H., Saulog, J., Paulo, B. S., Wilson, D. M., Liu, D. Y., Morehouse, N. J., Waterworth, S., Barkei, J., Gray, C. A., Kwan, J. C., Eustaquio, A. S., and Linington, R. (2023) Discovery of the Polyketide Lagriamide B by Integrated Genome Mining, Isotopic Labeling, and Untargeted Metabolomics. ChemRxiv.

Search 30k Actinos for Rapamycin BGCs

Four Genome Assemblies of *Streptomyces rapamycinicus* NRRL 5491



Find... ??



BGC0001522: auroramycin

Application

Can we target specific BGCs
functionality?

Find sequences with a halogenase, near an NRPS, near an antibiotic resistance gene, all on the same strand of DNA

```
:auto MATCH z1=(n:pfam {name:"Trp_halogenase"})-[:SOURCE_DB]-(h1:hmm)
MATCH z2=(h1)-[:ANNOTATES]-(:protein)←[e1:ENCODES]-(n1:mibig_bgc)
CALL {
    WITH n1, e1
    MATCH z3=(an1:antismash)←[:SOURCE_DB]-(:hmm)-[:ANNOTATES]→(p1:protein)←[e2:ENCODES]-(n1)
    MATCH z4=(an2:antismash)←[:SOURCE_DB]-(:hmm)-[:ANNOTATES]→(p1)
    WHERE an1.name = "Condensation"
        AND an2.name in ["AMP-binding", "A-OX"] AND abs(e1.start - e2.start) < 10000 AND e1.strand = e2.strand
    MATCH z5=(:amrfinder)←[:SOURCE_DB]-(:hmm)-[:ANNOTATES]→(p2:protein)←[e3:ENCODES]-(n1)
    WHERE abs(e1.start - e3.start) < 50000 AND e1.strand = e3.strand
    return z3, z4, z5
} in transactions of 1 rows
RETURN z1, z2, z3, z4, z5
```



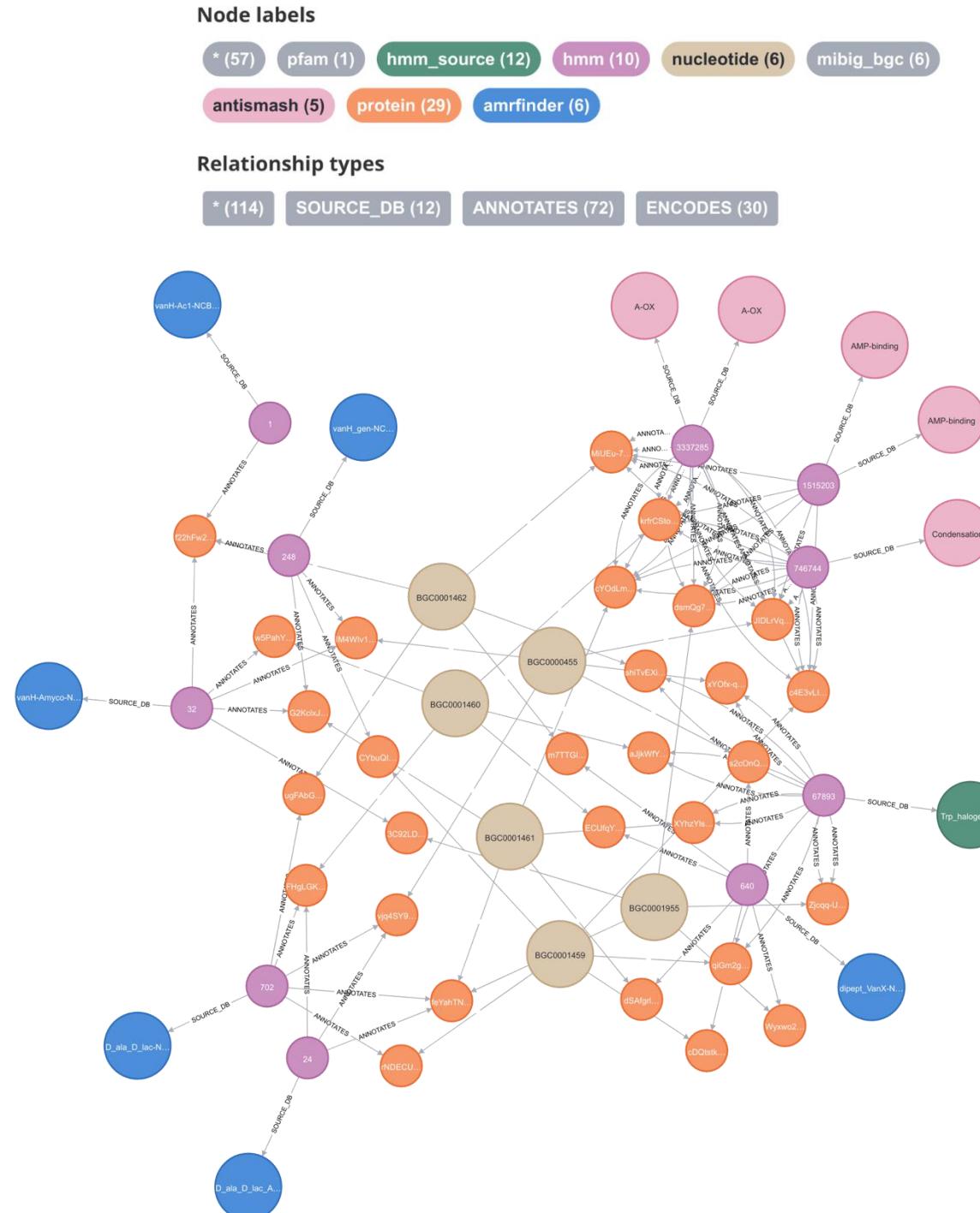
The diagram illustrates the node labels and relationship types used in the query results. Node labels include pfam (57), hmm (1), antismash (5), amrfinder (6), protein (29), nucleotide (6), mibig_bgc (6), hmm_source (12), and hmm (10). Relationship types include SOURCE_DB (12), ANNOTATES (72), and ENCODES (30).

Find sequences with a halogenase, near an NRPS, near an antibiotic resistance gene, all on the same strand of DNA

```

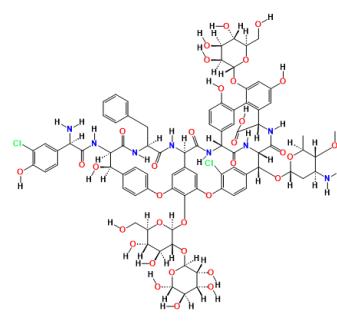
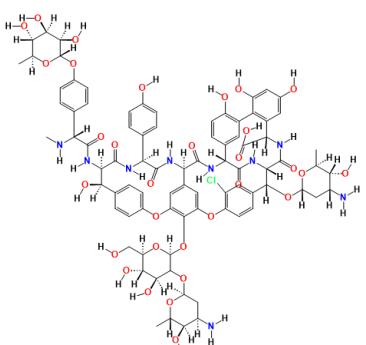
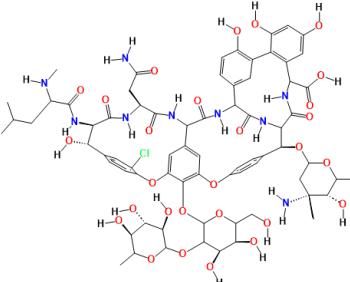
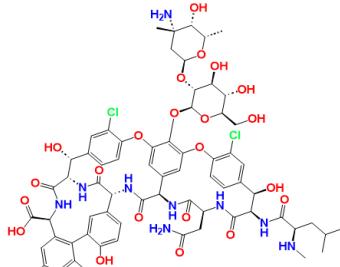
:auto MATCH z1=(n:pfam {name:"Trp_halogenase"})-[:SOURCE_DB]-(h1:hmm)
MATCH z2=(h1)-[:ANNOTATES]-(e1:protein)-[:ENCODES]-(n1:mibig_bgc)
CALL {
    WITH n1, e1
    MATCH z3=(an1:antismash)<[:SOURCE_DB]-[:hmm]-[:ANNOTES]>-(p1:protein)<[:ENCODES]->(n1)
    MATCH z4=(an2:antismash)<[:SOURCE_DB]-[:hmm]-[:ANNOTES]>-(p2:protein)<[:ENCODES]->(n1)
    WHERE an1.name = "Condensation"
        AND an2.name in ["AMP-binding", "A-OX"]
        AND abs(e1.start - e2.start) < 10000 AND e1.strand = e2.strand
    MATCH z5=(:amrfinder)<[:SOURCE_DB]-[:hmm]-[:ANNOTES]>-(p3:protein)<[:ENCODES]->(n1)
    WHERE abs(e1.start - e3.start) < 50000 AND e1.strand = e3.strand
}
} in transactions of 1 rows
RETURN z1, z2, z3, z4, z5

```



Find sequences with a halogenase, near an NRPS, near an antibiotic resistance gene, all on the same strand of DNA

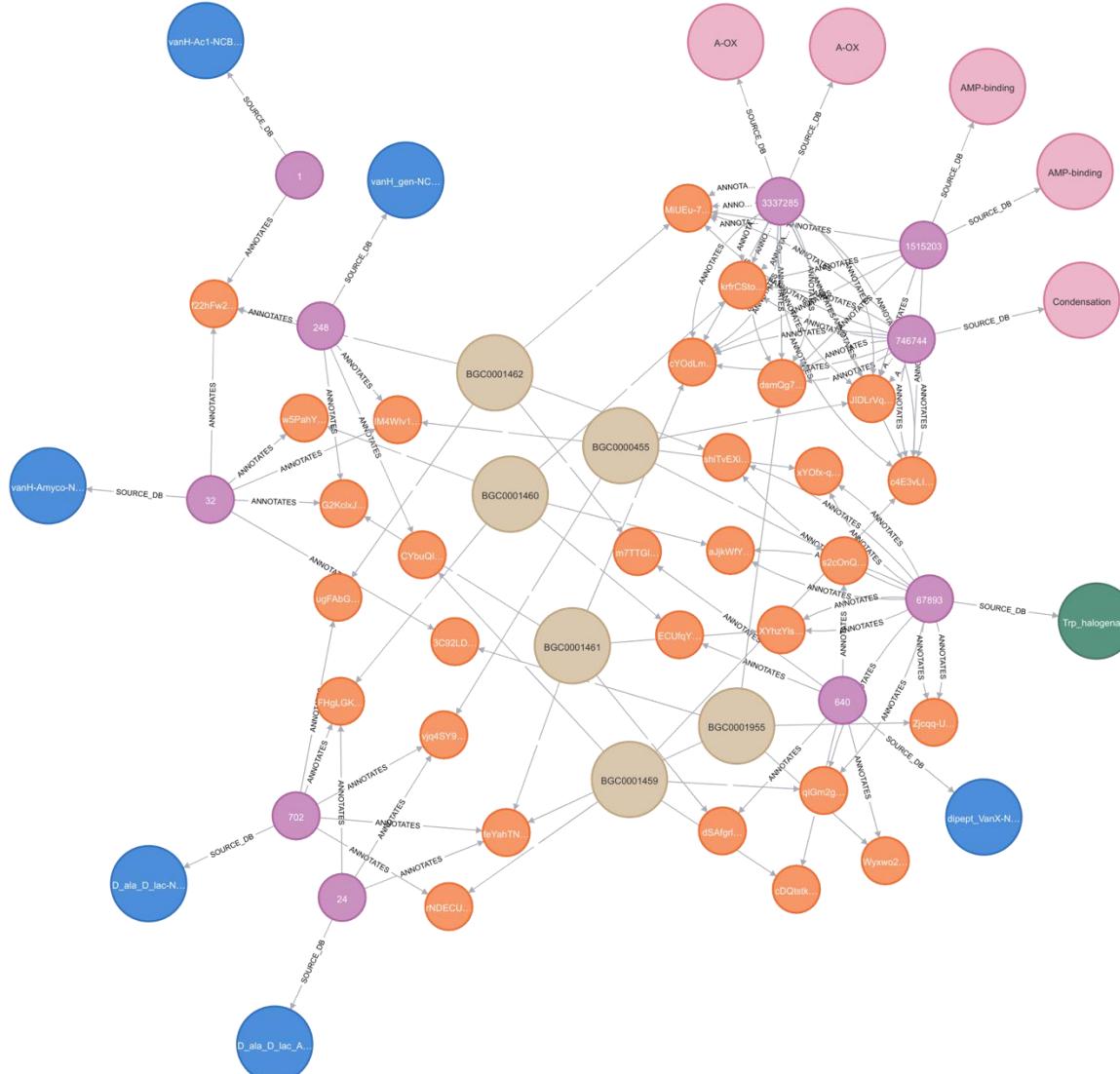
```
:auto MATCH z1=(n:pfam {name:"Trp_halogenase"})-[:SOURCE_DB]-(h1:hmm)
MATCH z2=(h1)-[:ANNOTATES]-(e1:protein)-[:ENCODES]-(n1:mibig_bgc)
CALL {
    WITH n1, e1
    MATCH z3=(an1:antismash) -[:SOURCE_DB]-(:hmm)-[:ANNOTATES]->(p1:protein)-[:ENCODES]->(n1)
    MATCH z4=(an2:antismash) -[:SOURCE_DB]-(:hmm)-[:ANNOTATES]->(p2:protein)-[:ENCODES]->(n1)
    WHERE an1.name = "Condensation"
        AND an2.name in ["AMP-binding", "A-OX"] AND abs(e1.start - e2.start) < 10000 AND e1.strand = e2.strand
    MATCH z5=(amrfinder) -[:SOURCE_DB]-(:hmm)-[:ANNOTATES]->(p3:protein)-[:ENCODES]->(n1)
    WHERE abs(e1.start - e3.start) < 50000 AND e1.strand = e3.strand
} in transactions of 1 rows
RETURN z1, z2, z3, z4, z5
```



Node labels

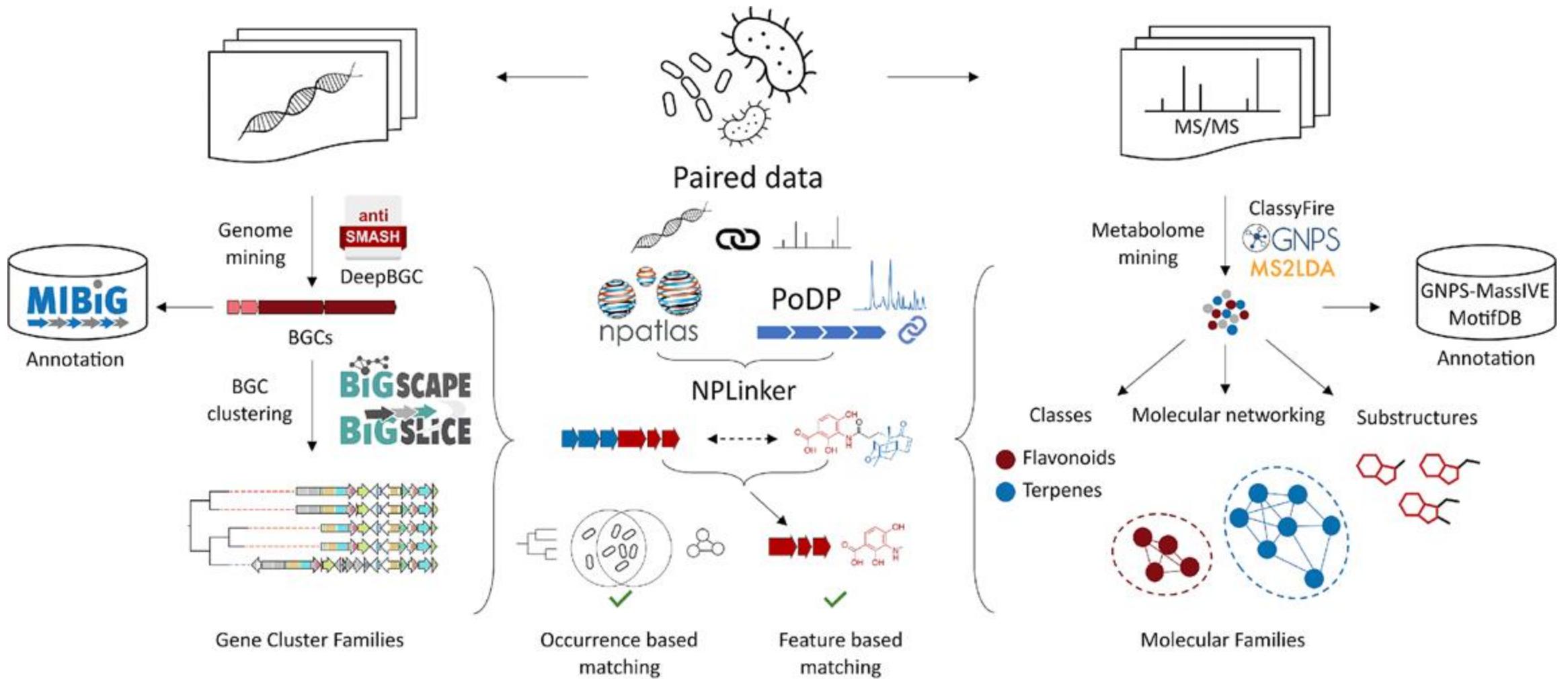


Relationship types



Application

Paired Omics



Louwen JJR, van der Hooft JJJ. Comprehensive Large-Scale Integrative Analysis of Omics Data To Accelerate Specialized Metabolite Discovery. *mSystems*. 2021 Aug 31;6(4):e0072621. doi: 10.1128/mSystems.00726-21. Epub 2021 Aug 24

Paired Omics Data Platform

List

Add

Review

Statistics

Download

Methods

About

Available projects



Example search queries: [Streptomyces](#), [SAMN02603879](#), [scrip*](#), [tubes | flask](#), [Pieter + Carmen](#), [Pieter + -Carmen](#), [Pieter + \(Carmen | Emily\)](#), "acetonitrile with"

Metabolomics project identifier	Principal investigator	Submitter(s)	Nr of (meta)genomes	Nr of proteomes	Nr of growth conditions	Nr of extraction methods	Nr of instrumentation methods	Nr of links between genome and metabolome samples	Nr of links between biosynthetic gene clusters and MS/MS spectra
MSV000084723	Cameron Currie & Tim Bugni	Marc G Chevrette	120	0	1	1	1	122	0

pairedomicsdata.bioinformatics.nl

Metabolomics project identifier	Principal investigator	Submitter(s)	Nr of (meta)genomes	Nr of proteomes	Nr of growth conditions	Nr of extraction methods	Nr of instrumentation methods	Nr of links between genome and metabolome samples	Nr of links between biosynthetic gene clusters and MS/MS spectra
MSV000084723	Cameron Currie & Tim Bugni	Marc G Chevrette	120	0	1	1	1	122	0

```
● (sgpy) chase@titan:~/Downloads$ sg_import_gnps --gnps_dirpath ProteoSAFe-METABOLOMICS-SNETS-V2-927974dc-view_all_clusters_withID_beta --map_path genome_to_file_map.csv
2024-03-08 08:31:47 INFO Connected to Neo4j database at bolt://localhost:7687
2024-03-08 08:31:50 INFO Created 21930 (:ms2_spectrum) nodes, set 131580 properties
INFO Created 63 (:gnps_library_spectrum) nodes, set 1633 properties
INFO Created 0 (:gnps_library_spectrum) nodes, set 0 properties
INFO Created 1214 (:gnps_cluster) nodes, set 25695 properties
2024-03-08 08:31:52 INFO Created 122 (:mass_spectrum_file) nodes, set 366 properties
2024-03-08 08:31:55 INFO 20027 relationships created (:ms2_spectrum)-[:CLUSTERS_TO]->(:gnps_cluster)
2024-03-08 08:31:59 INFO 21930 relationships created (:mass_spectrum_file)-[:HAS]->(:ms2_spectrum)
INFO Created 0 (:gnps_cluster) nodes, set 0 properties
INFO Created 55 (:gnps_library_spectrum) nodes, set 55 properties
INFO 67 relationships created (:gnps_cluster)-[:LIBRARY_HIT]->(:gnps_library_spectrum)
INFO Created 0 (:gnps_cluster) nodes, set 0 properties
INFO Created 0 (:gnps_library_spectrum) nodes, set 0 properties
INFO 67 relationships created (:gnps_cluster)-[:LIBRARY_HIT]->(:gnps_library_spectrum)
2024-03-08 08:32:00 INFO Created 0 (:gnps_cluster) nodes, set 0 properties
INFO Created 0 (:gnps_cluster) nodes, set 0 properties
INFO 2018 relationships created (:gnps_cluster)-[:MOLECULAR_NETWORK]->(:gnps_cluster)
INFO Created 39 (:chemical_compound) nodes, set 1599 properties
INFO 46 relationships created (:gnps_library_spectrum)-[:IS_A]->(:chemical_compound)
INFO Assemblies in GNPS results found in db: 84 of 84
INFO Assemblies in GNPS results not found in db: set()
2024-03-08 08:32:01 INFO 86 relationships created (:mass_spectrum_file)-[:ANALYSIS_OF]->(:assembly)
INFO GNPS molecular network has been integrated into the SocialGene Neo4j database
neo4j_element.py:62
neo4j_element.py:369
parse.py:237
parse.py:240
neo4j_element.py:538
cli.py:146
```

Metabolomics project identifier	Principal investigator	Submitter(s)	Nr of (meta)genomes	Nr of proteomes	Nr of growth conditions	Nr of extraction methods	Nr of instrumentation methods	Nr of links between genome and metabolome samples	Nr of links between biosynthetic gene clusters and MS/MS spectra
MSV000084723	Cameron Currie & Tim Bugni	Marc G Chevrette	120	0	1	1	1	122	0

```
(sgpy) chase@titan:~/Downloads$ sg_import_npatlas --input NPAtlas_download.json  
33372 Processing NPAtlas entries... 100% 0:02:30  
2024-03-08 08:35:13 INFO Creating/Merging npatlas nodes in neo4j  
INFO Connected to Neo4j database at bolt://localhost:7687  
2024-03-08 08:35:16 INFO Created 33372 (:npatlas) nodes, set 467208 properties  
2024-03-08 08:35:18 INFO Creating/Merging nodes linked to npatlas entries in neo4j  
2024-03-08 08:35:19 INFO Created 13058 (:publication) nodes, set 78348 properties  
INFO Created 8 (:taxid) nodes, set 8 properties  
2024-03-08 08:35:20 INFO Created 27196 (:gnps_library_spectrum) nodes, set 27196 properties  
2024-03-08 08:35:21 INFO Created 96 (:assembly:mibig) nodes, set 96 properties  
INFO Created 1632 (:classyfire) nodes, set 1632 properties  
INFO Created 509 (:npclassifier_class) nodes, set 509 properties  
INFO Created 7 (:npclassifier_pathway) nodes, set 7 properties  
INFO Created 76 (:npclassifier_superclass) nodes, set 76 properties  
2024-03-08 08:35:28 INFO Created 33339 (:chemical_compound) nodes, set 1366899 properties  
INFO Created 911 (:chebi) nodes, set 1822 properties  
INFO Linking npatlas entries and related nodes in neo4j  
2024-03-08 08:35:31 INFO 33372 relationships created (:npatlas)-[:HAS]->(:publication)  
2024-03-08 08:35:33 INFO 31481 relationships created (:taxid)-[:PRODUCES]->(:npatlas)  
2024-03-08 08:35:36 INFO 31887 relationships created (:npatlas)-[:HAS]->(:gnps_library_spectrum)  
INFO 2511 relationships created (:assembly:mibig)-[:PRODUCES]->(:npatlas)  
2024-03-08 08:35:38 INFO 32888 relationships created (:npatlas)-[:LOWEST_CLASS]->(:classyfire)  
2024-03-08 08:35:40 INFO 32888 relationships created (:npatlas)-[:DIRECT_PARENT]->(:classyfire)  
2024-03-08 08:35:41 INFO 10694 relationships created (:npatlas)-[:INTERMEDIATE_NODES]->(:classyfire)  
2024-03-08 08:36:15 INFO 444892 relationships created (:npatlas)-[:ALTERNATIVE_PARENTS]->(:classyfire)  
2024-03-08 08:36:18 INFO 31400 relationships created (:npatlas)-[:IS_A]->(:npclassifier_class)  
2024-03-08 08:36:21 INFO 34822 relationships created (:npatlas)-[:IS_A]->(:npclassifier_pathway)  
2024-03-08 08:36:23 INFO 28759 relationships created (:npatlas)-[:IS_A]->(:npclassifier_superclass)  
2024-03-08 08:36:26 INFO 33371 relationships created (:npatlas)-[:IS_A]->(:chemical_compound)  
2024-03-08 08:37:37 INFO 801979 relationships created (:npatlas)-[:IS_A]->(:chebi)
```



```

1 MATCH z1=(:assembly)←[r:ANALYSIS_OF]-(msf:mass_spectrum_file)-[:HAS]→(ms2:ms2_spectrum)-[:CLUSTERS_TO]→(gc:gnps_cluster)-[:LIBRARY_HIT]->(:gnps_library_spectrum)-[:HAS]-(np:npatlas)
2 MATCH (np)←[:PRODUCES]-(:taxid {name: "Streptomyces"})
3 MATCH z2=(gc)-[:MOLECULAR_NETWORK*1 .. 2]-(:gnps_cluster)
4 RETURN z1, z2 limit 20

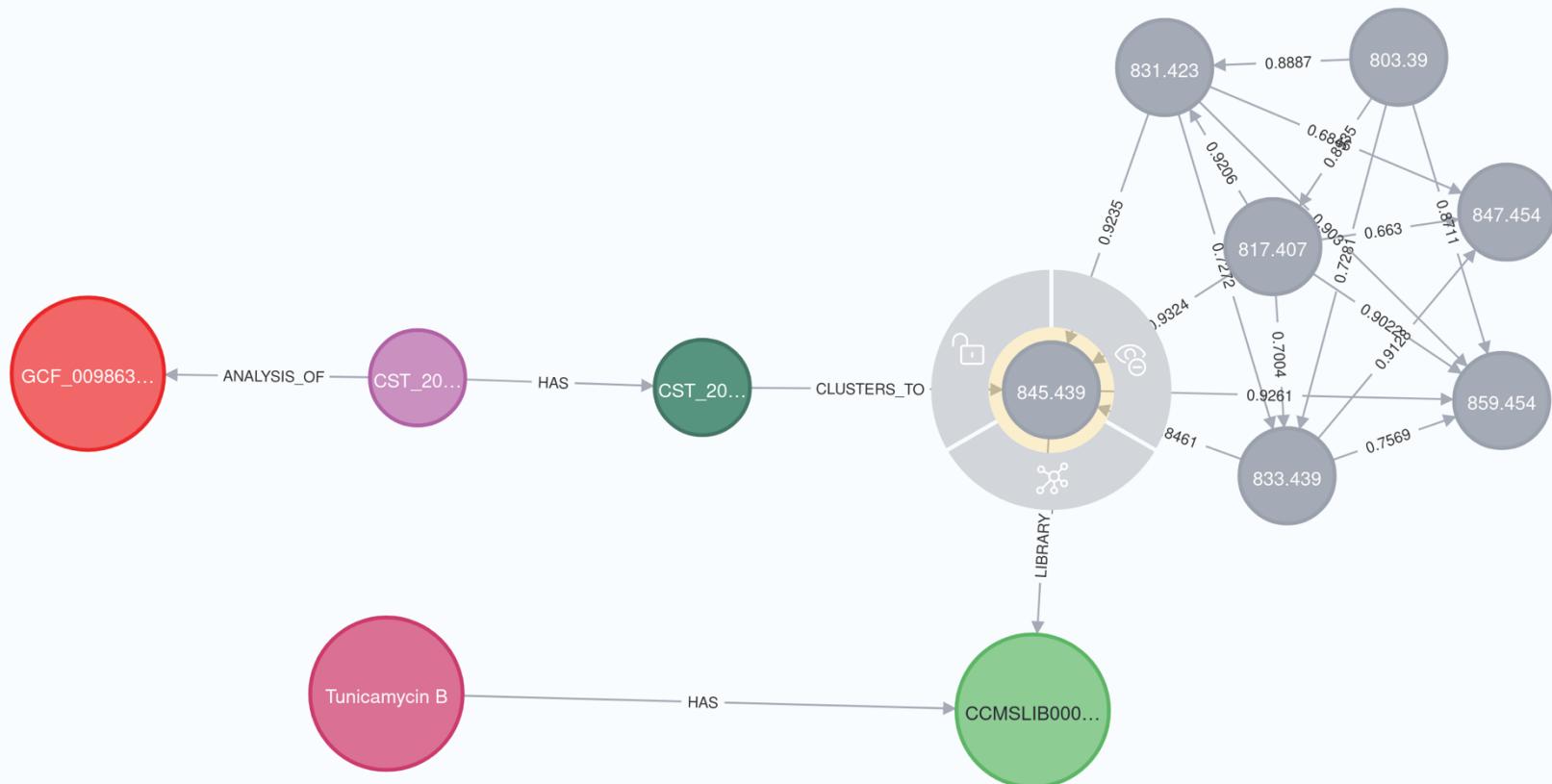
```

Graph

Table

Text

Code

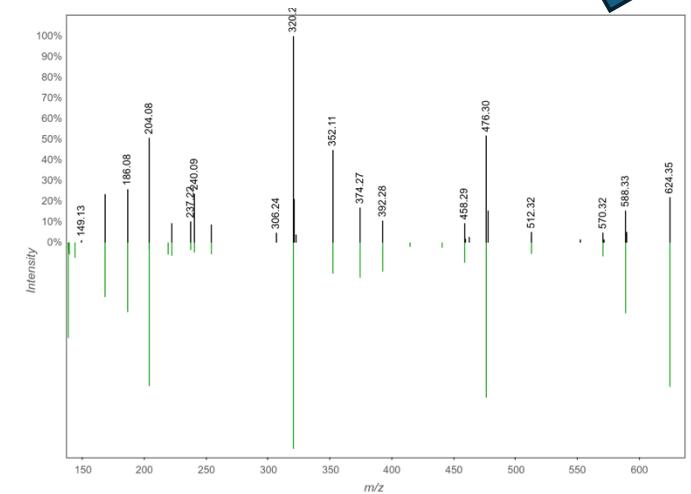
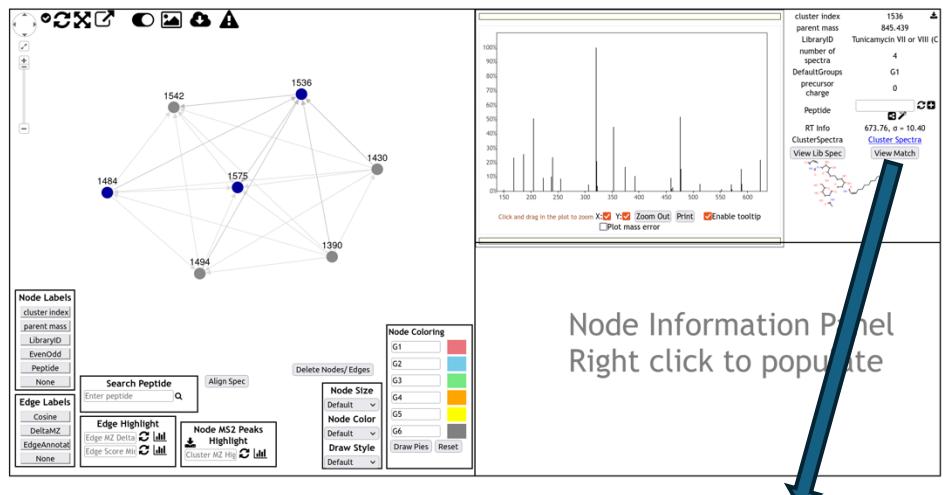
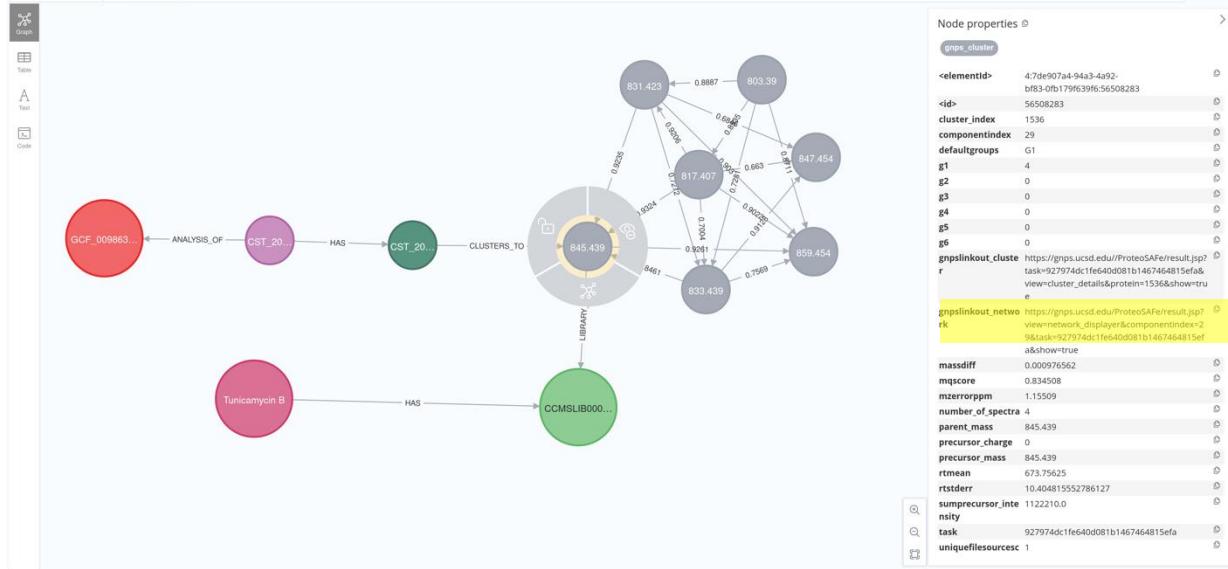


gnps_cluster

<elementId>	4:7de907a4-94a3-4a92-bf83-0fb179f639f6:56508283
<id>	56508283
cluster_index	1536
componentindex	29
defaultgroups	G1
g1	4
g2	0
g3	0
g4	0
g5	0
g6	0
gnpslinkout_cluster	https://gnps.ucsd.edu//ProteoSAFe/result.jsp?task=927974dc1fe640d081b1467464815efa&view=cluster_details&protein=1536&show=true
gnpslinkout_network	https://gnps.ucsd.edu//ProteoSAFe/result.jsp?view=network_displayer&componentindex=29&task=927974dc1fe640d081b1467464815efaa&show=true
massdiff	0.000976562
mqscore	0.834508
mzerrorppm	1.15509
number_of_spectra	4
parent_mass	845.439
precursor_charge	0
precursor_mass	845.439
rtmean	673.75625
rtstderr	10.404815552786127
sumprecursor_intensity	1122210.0
task	927974dc1fe640d081b1467464815efa
uniquefilesources	1



```
1 MATCH z1{(:assembly)-[r:ANALYSIS_OF]-(msf:mass_spectrum_file)-[:HAS]→(ms2:ms2_spectrum)-[:CLUSTERS_TO]→(gc:gnps_cluster)-[:LIBRARY_HIT]->(:gnps_library_spectrum)-[:HAS]→(np:nptlaps)
2 MATCH (np)-[:PRODUCES]-(taxid {name: "Streptomyces"})
3 MATCH z2{gc}-[:MOLECULAR_NETWORK*1..2]-(gnps_cluster)
4 RETURN z1, z2 limit 20
```



Node Information Panel

Right click to populate



```

1 MATCH z1=(:assembly)←[r:ANALYSIS_OF]-(msf:mass_spectrum_file)-[:HAS]→(ms2:ms2_spectrum)-[:CLUSTERS_TO]→(gc:gnps_cluster)-[:LIBRARY_HIT]->(:gnps_library_spectrum)-[:HAS]-(np:npatlas)
2 MATCH (np)←[:PRODUCES]-(:taxid {name: "Streptomyces"})
3 MATCH z2=(gc)-[:MOLECULAR_NETWORK*1 .. 2]-(:gnps_cluster)
4 RETURN z1, z2 limit 20

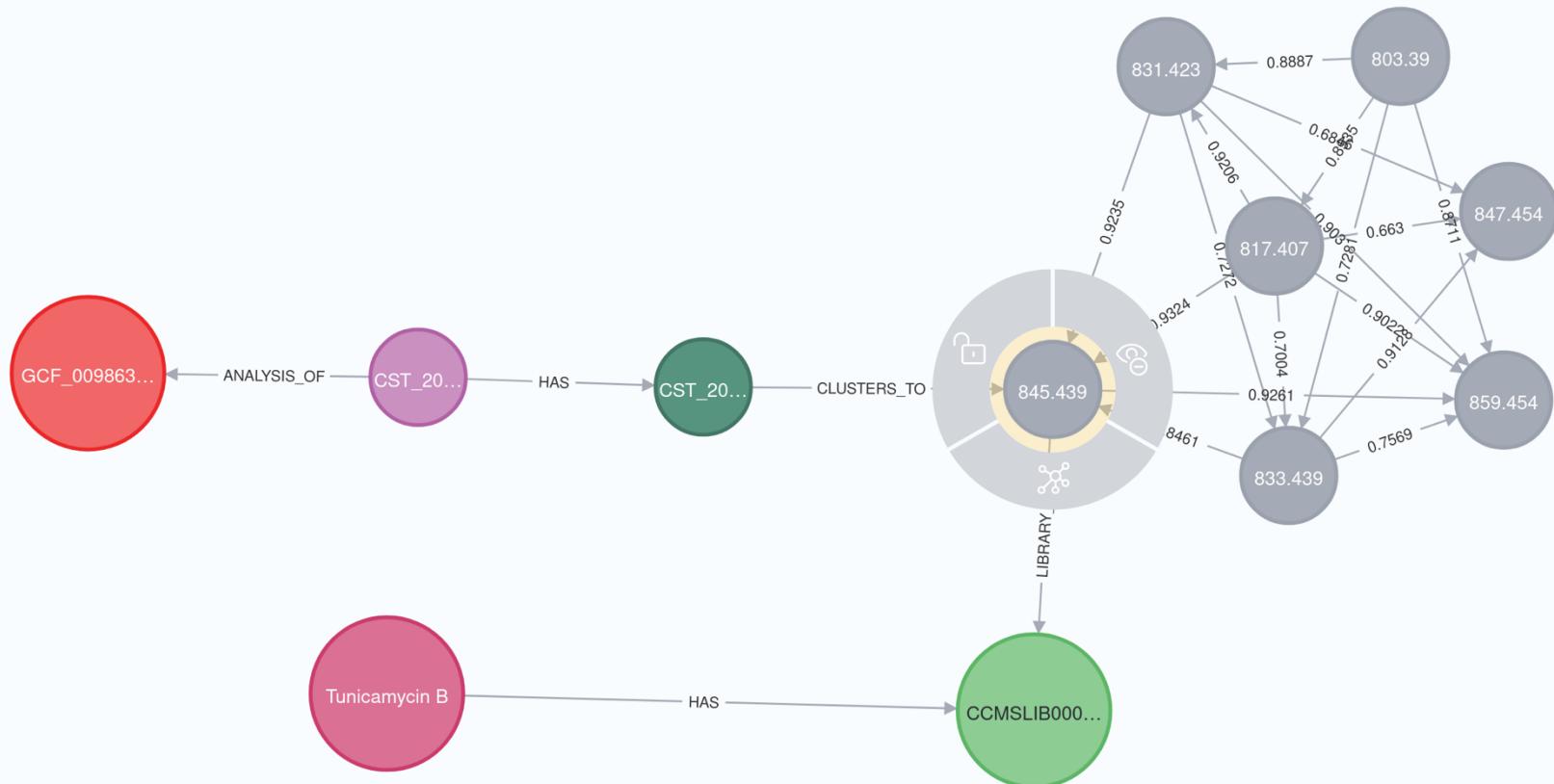
```

Graph

Table

Text

Code



Node properties

gnps_cluster

<elementId>	4:7de907a4-94a3-4a92-bf83-0fb179f639f6:56508283
<id>	56508283
cluster_index	1536
componentindex	29
defaultgroups	G1
g1	4
g2	0
g3	0
g4	0
g5	0
g6	0
gnpslinkout_cluster	https://gnps.ucsd.edu//ProteoSAFe/result.jsp?task=927974dc1fe640d081b1467464815efa&view=cluster_details&protein=1536&show=true
gnpslinkout_network	https://gnps.ucsd.edu//ProteoSAFe/result.jsp?view=network_displayer&componentindex=29&task=927974dc1fe640d081b1467464815efa&show=true
massdiff	0.000976562
mqscore	0.834508
mzerrorppm	1.15509
number_of_spectra	4
parent_mass	845.439
precursor_charge	0
precursor_mass	845.439
rtmean	673.75625
rtstderr	10.404815552786127
sumprecursor_intensity	1122210.0
task	927974dc1fe640d081b1467464815efa
uniquefilesources	1



Application



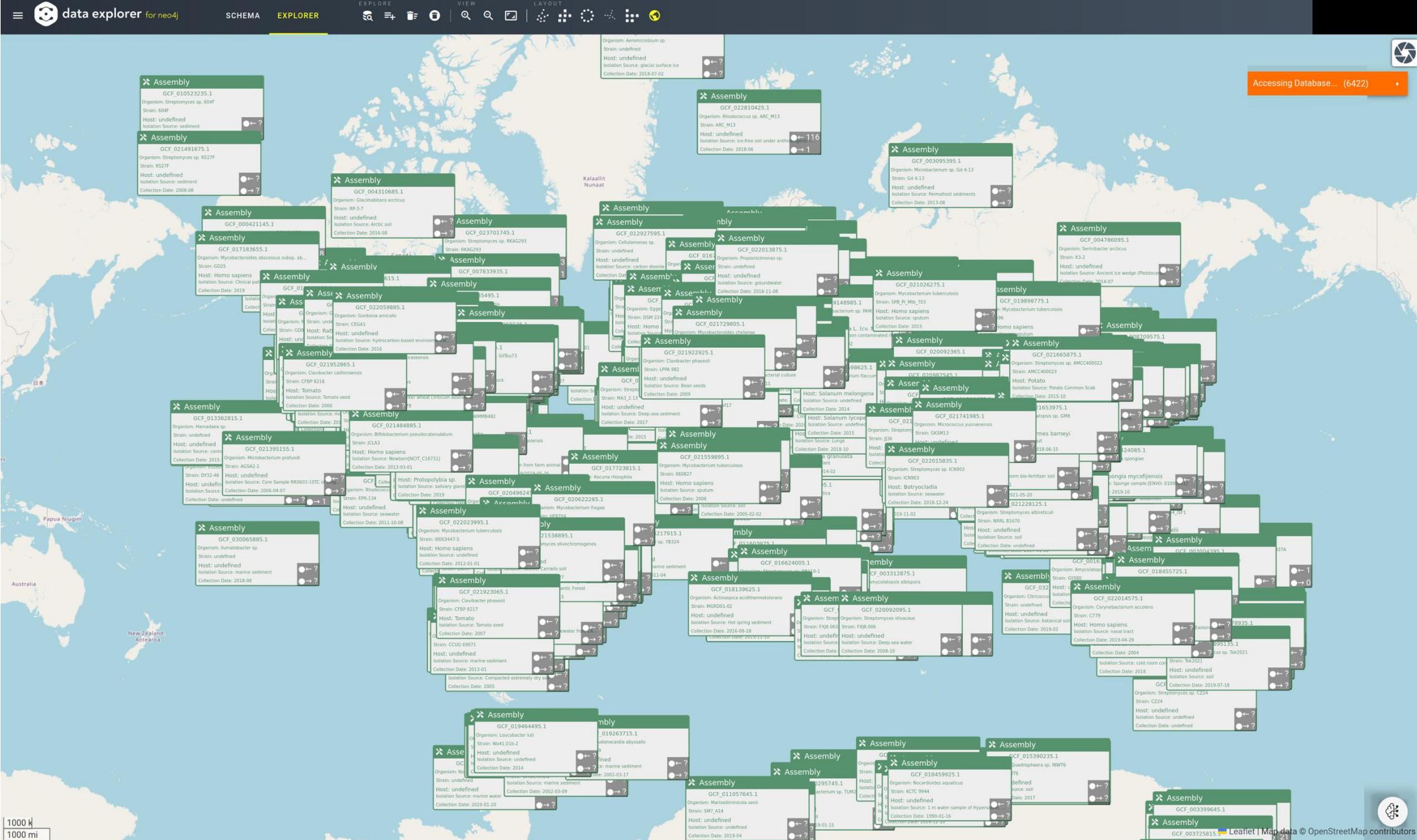
Where in the world are my BGCs?

Pull and store latitude and longitude in-database

```
1 MATCH (n:assembly) where n.lat_lon is not null
2 WITH n, split(n.lat_lon, ' ') AS parts
3 WITH n,toFloat(parts[0]) AS latitude, parts[1] AS lat_direction, toFloat(parts[2]) AS longitude, parts[3] AS lon_direction
4 WITH n,
5   CASE lat_direction WHEN 'N' THEN latitude WHEN 'S' THEN -latitude ELSE null END AS latitude_value,
6   CASE lon_direction WHEN 'E' THEN longitude WHEN 'W' THEN -longitude ELSE null END AS longitude_value
7 SET n.geolocation= point({latitude: latitude_value, longitude: longitude_value})
```

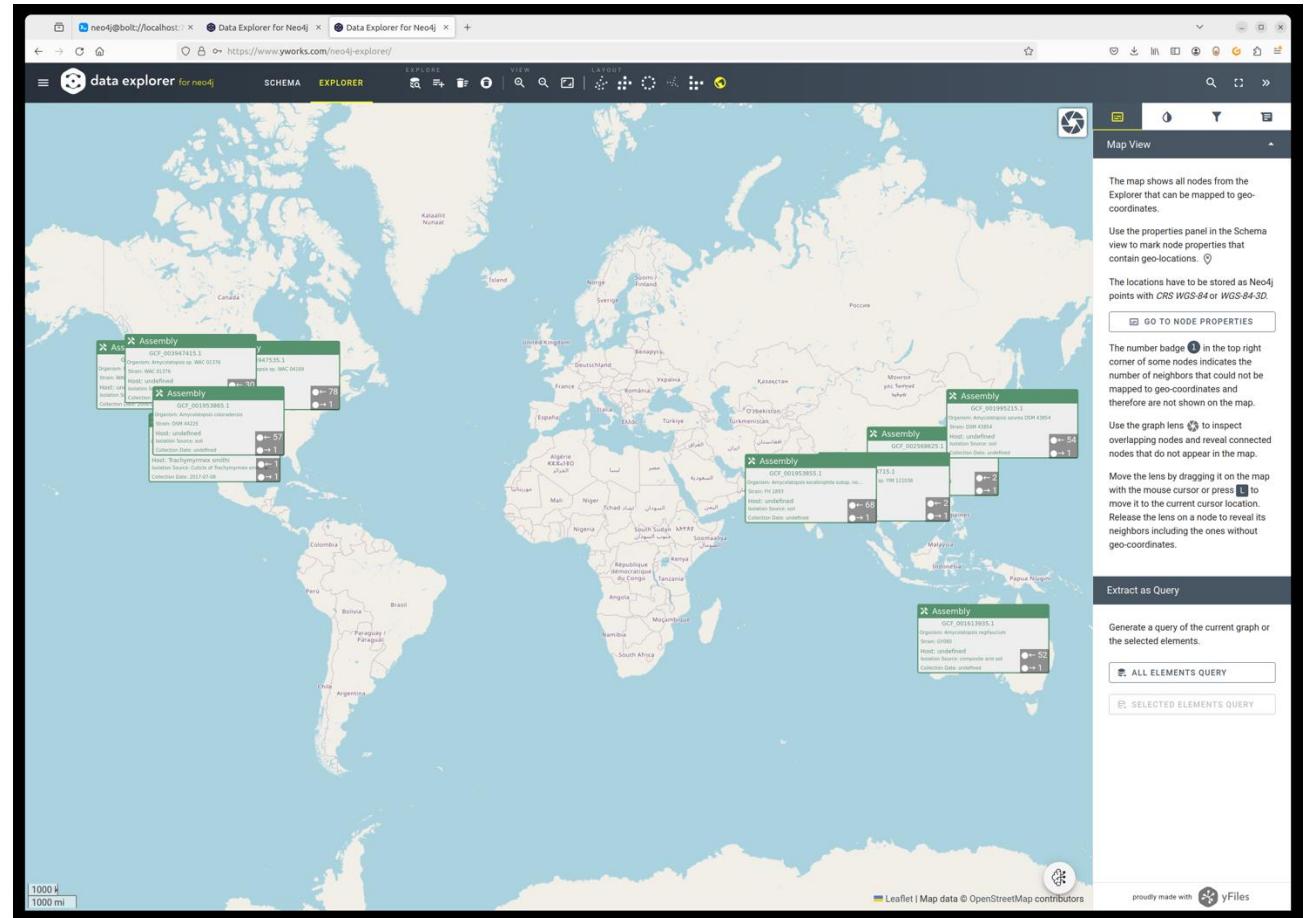


Set 88699 properties, completed after 553 ms.



66 seconds

Geolocate Putative Halogenated NRPS Antibiotics



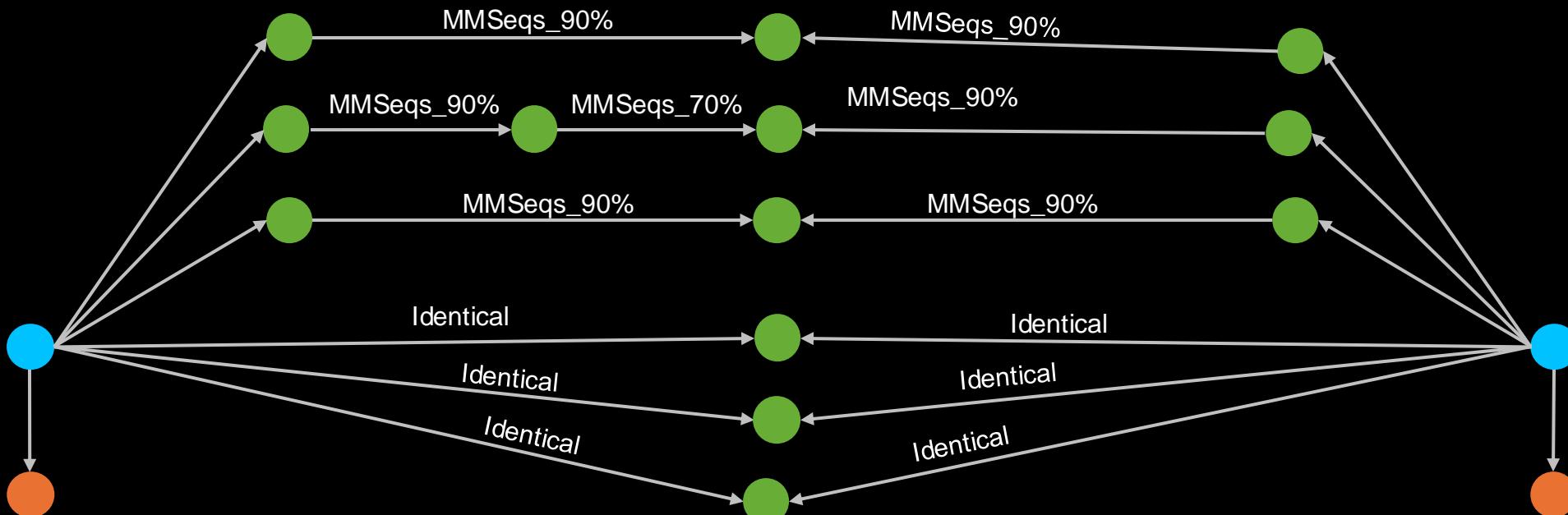
```
1 MATCH z1=(n:pfam {name:"Trp_halogenase"})-[:SOURCE_DB]-(h1:hmm)
2 MATCH z2=(h1)-[:ANNOTATES]-(:protein)<-[:ENCODES]->(n1:nucleotide)
3 CALL {
4   WITH n1, e1
5   MATCH z3=(an1:antismash)<[:SOURCE_DB]-[:hmm]-[:ANNOTATES]>-(p1:protein)<-[:ENCODES]->(n1)
6   MATCH z4=(an2:antismash)<[:SOURCE_DB]-[:hmm]-[:ANNOTATES]>-(p1)
7   WHERE an1.name = "Condensation"
8     AND an2.name in ["AMP-binding", "A-OX"] AND abs(e1.start - e2.start) < 10000 AND e1.strand = e2.strand
9   MATCH z5=(:amrfinder)<[:SOURCE_DB]-[:hmm]-[:ANNOTATES]>-(p2:protein)<-[:ENCODES]->(n1)
10  WHERE abs(e1.start - e3.start) < 50000 AND e1.strand = e3.strand
11  MATCH (n1)-[:ASSEMBLES_TO]-(a1:assembly)
12  WHERE az.lat_lon is not null
13  return az
14 } in transactions of 1 rows
15 RETURN distinct a1
```

Application

Global analyses



Protein sequence similarity with MMSeqs2 clustering



● Genome ● DNA Sequence ● Protein

**Look in 343,381 genomes
For similar BGCs to
Any of 2,502 MIBiG BGCs
Limit to the
2.1 million sequence regions
predicted by antiSMASH 7**

**Look in 343,381 genomes
For similar BGCs to
Any of 2,502 MIBiG BGCs
Limit to the
2.1 million sequence regions
predicted by antiSMASH 7**

- Genome with match



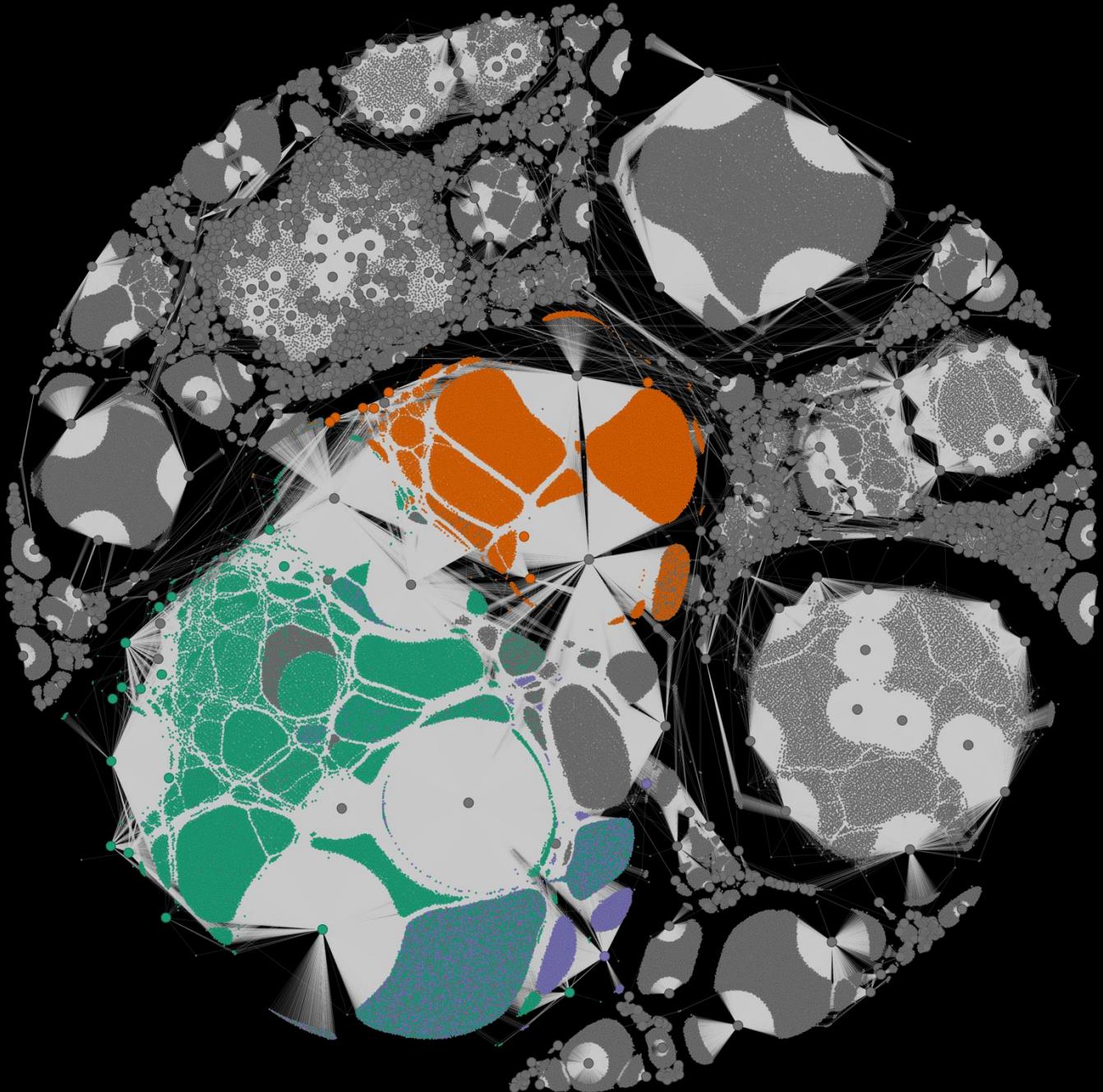
— >70% of proteins, >70% AA identity

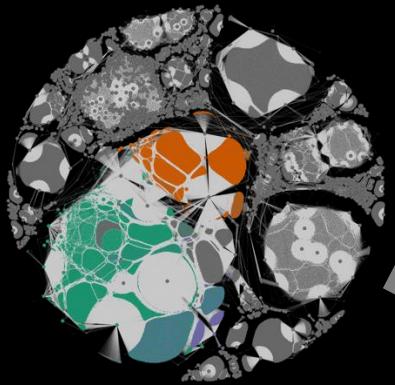
184,571 genomes
1,964 MIBiG BGCs
597,917 relationships
2.5 minutes

Look in 343,381 genomes
For similar BGCs to
Any of 2,502 MIBiG BGCs
Limit to the
2.1 million sequence regions
predicted by antiSMASH 7

- Genome with match
- MIBiG BGC
- *>70% of proteins, >70% AA identity*
- *Klebsiella*
 - *Escherichia*
 - *Salmonella*

184,571 genomes
1,964 MIBiG BGCs
597,917 relationships
2.5 minutes





How does BGC distribution contribute to pathogenicity?

- Genome with match



Klebsiella

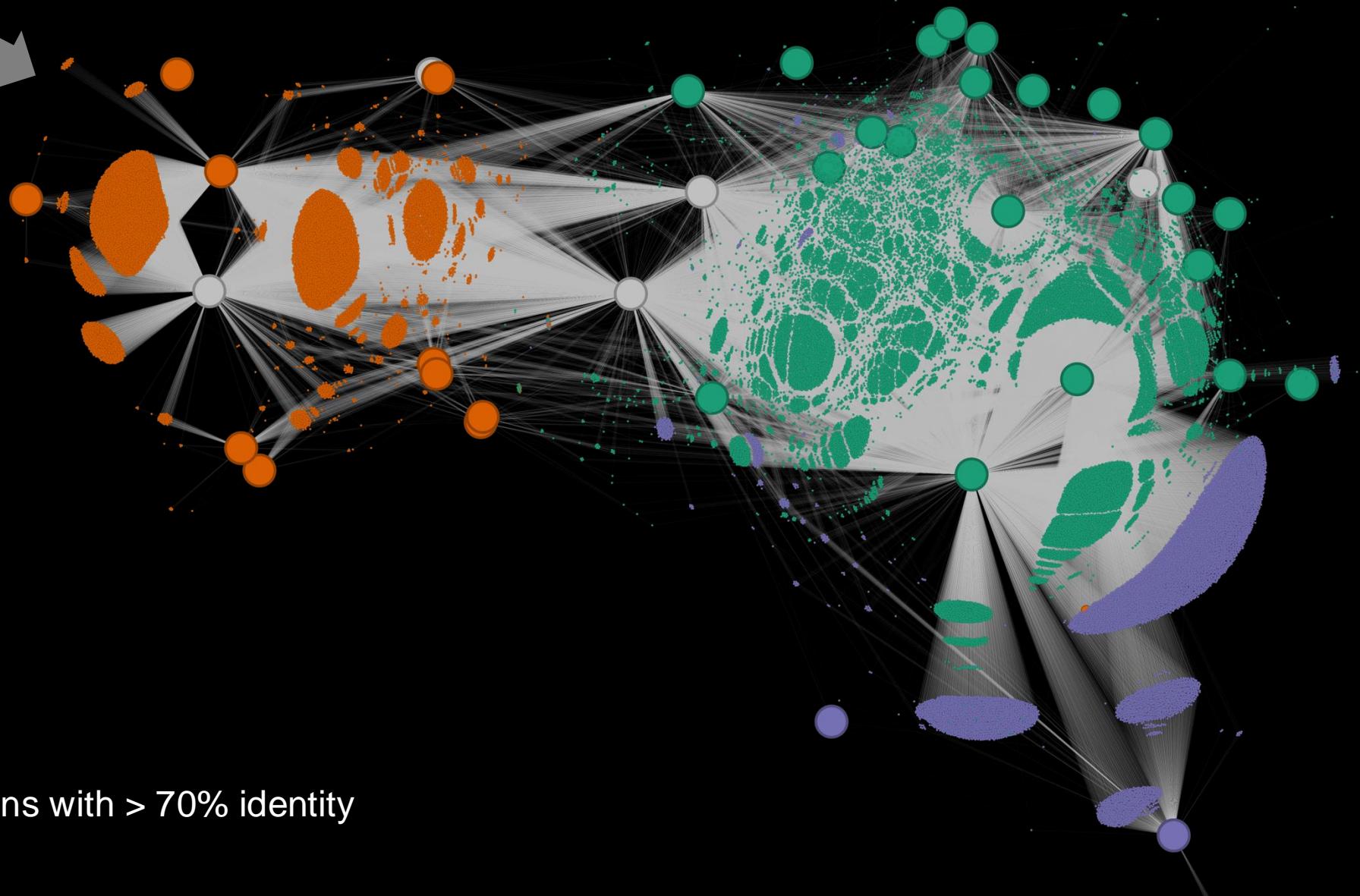
Escherichia

Salmonella

69,880 genomes

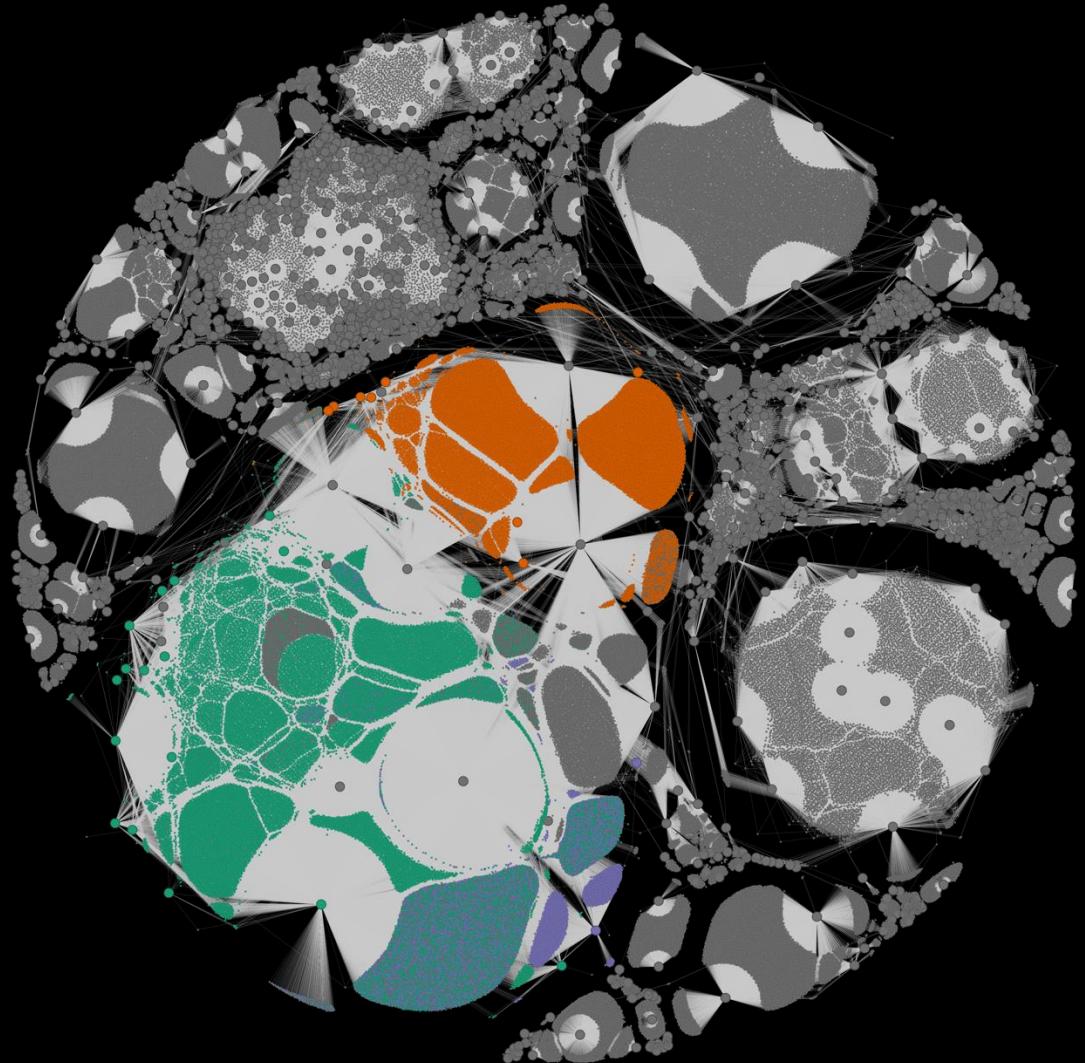
42 MIBiG BGCs

203,092 relationships



Weighted: % shared proteins with $> 70\%$ identity

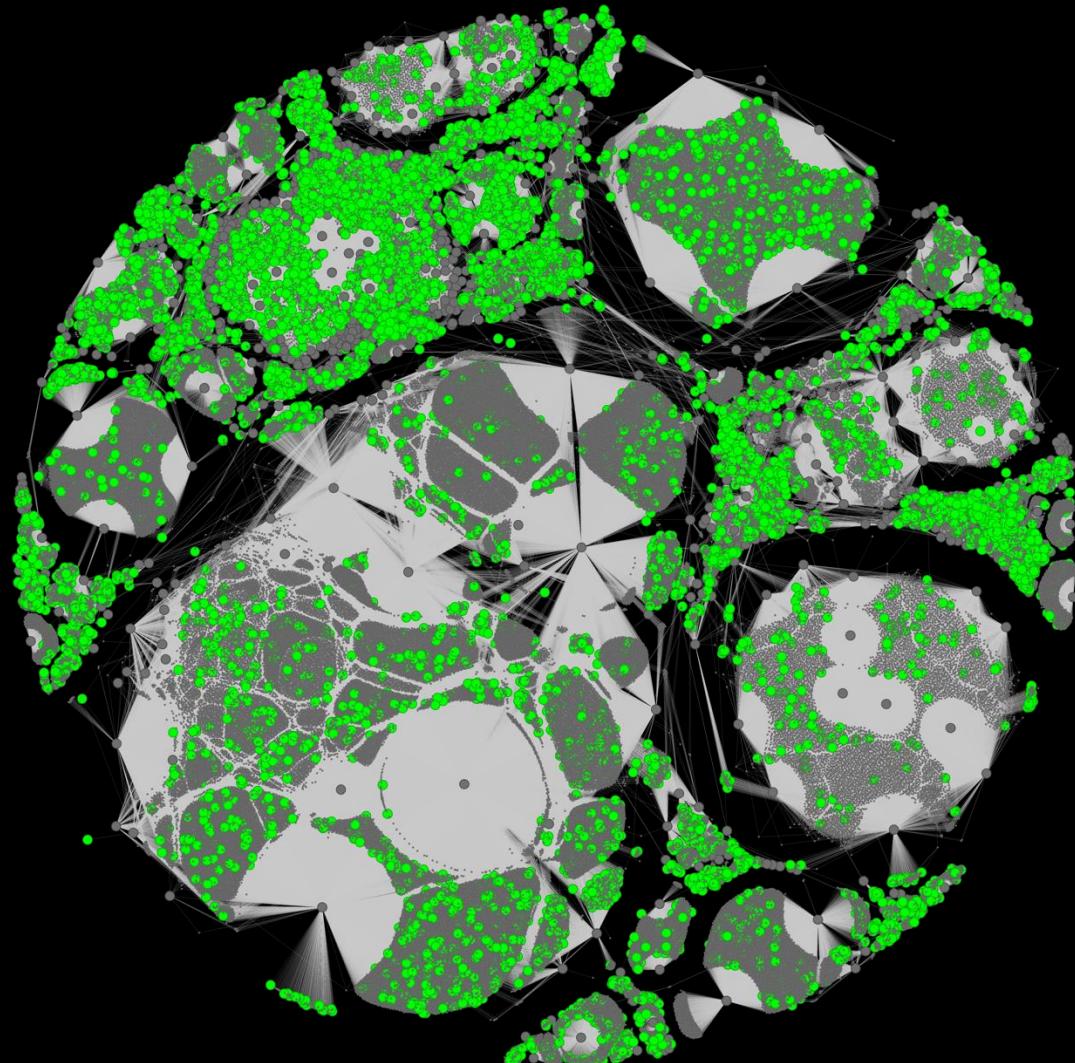
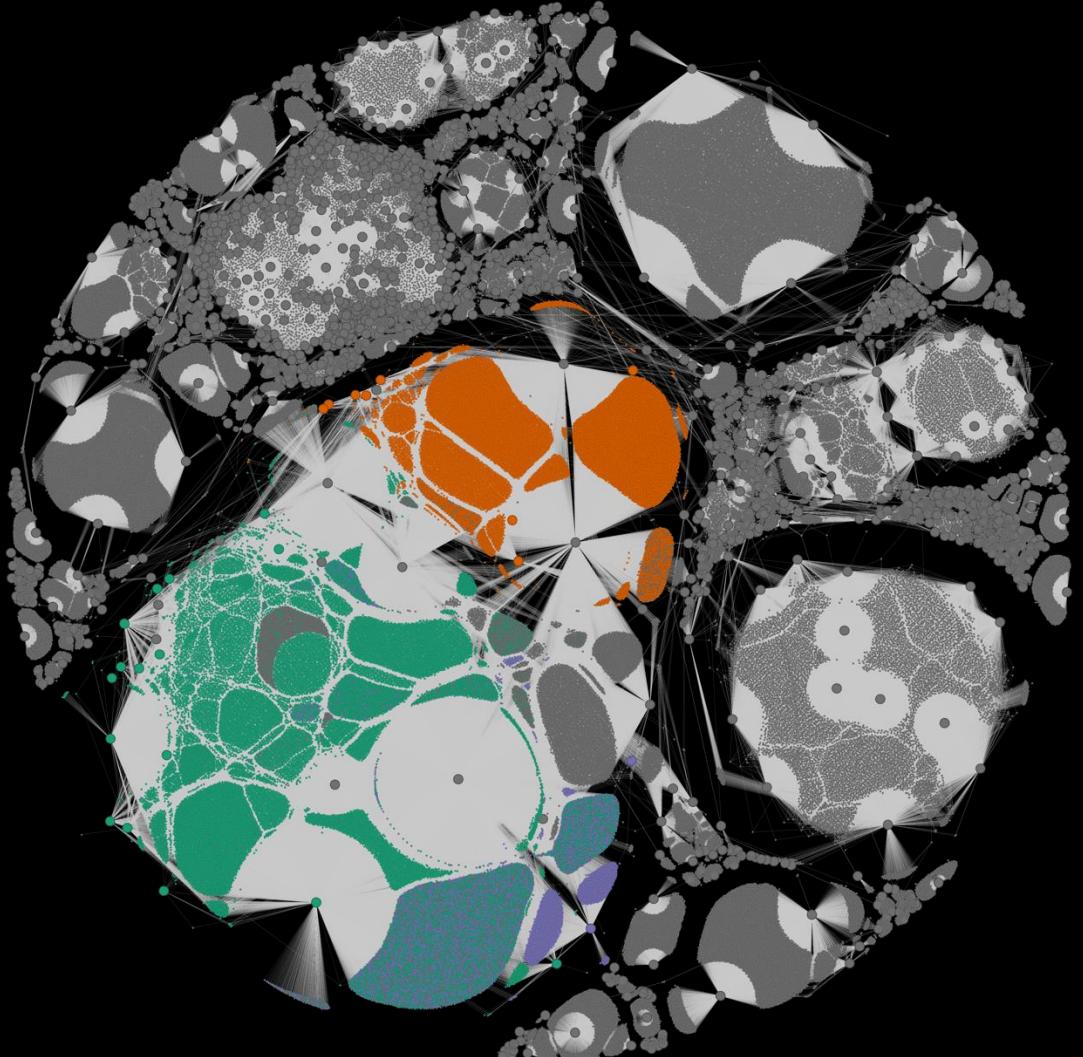
How to test our hypotheses?



8,984 of the 184,571 genomes are in a culture collection



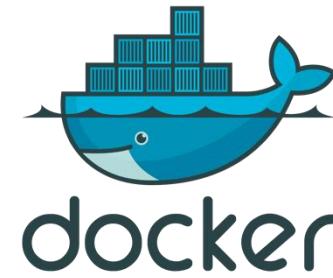
In a culture collection



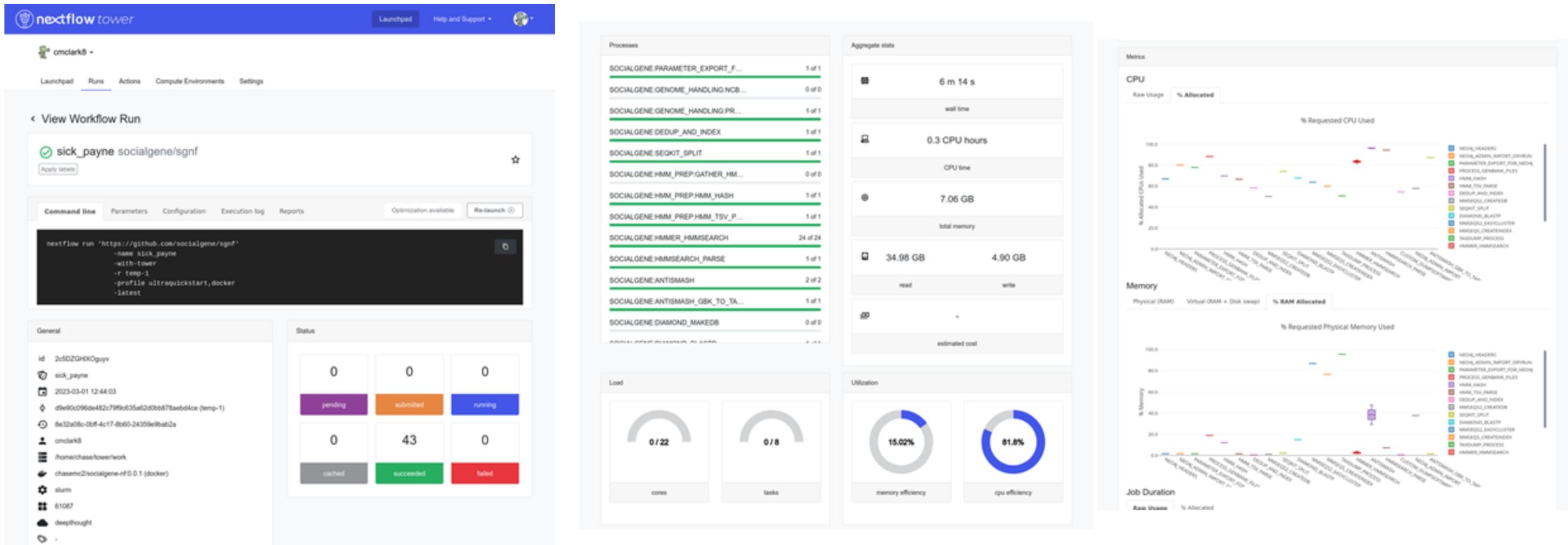
How?

```
(base) chase@titan:~/Documents/github/kwan_lab/socialgene/sgnf$
```

nextflow



Making the workflow more accessible with **nextflow tower**



The image shows the Nextflow Tower web interface. On the left, a sidebar lists user profiles (cmclark8, sick Payne, etc.) and navigation links (Launchpad, Runs, Actions, Compute Environments, Settings). The main area displays a workflow run titled "sick Payne socialgene/sgnf". The command line used is:

```
nextflow run 'https://github.com/socialgene/sgnf'  
-name sick Payne  
-with塔  
-r temp  
-profile ultrarapidstart,docker  
-latest
```

The "Runs" tab is selected, showing the workflow's status with 0 pending, 0 submitted, and 0 running tasks. Below this, it shows 0 cached, 43 succeeded, and 0 failed tasks.

The central part of the interface shows the "Processes" section with a list of tasks and their counts (e.g., 1 of 1 for various HMM-related steps). To the right, the "Aggregate stats" section provides summary information:

- wall time: 6 m 14 s
- CPU time: 0.3 CPU hours
- total memory: 7.06 GB
- read: 34.98 GB
- write: 4.90 GB
- estimated cost: ~

The "Metrics" section contains three line charts for CPU, Memory, and Job Duration, each comparing "Raw Usage" and "% Allocated" across various Nextflow processes. A legend on the right identifies the processes by color.



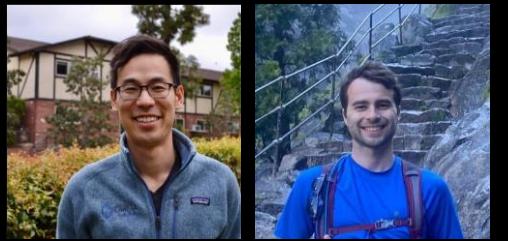
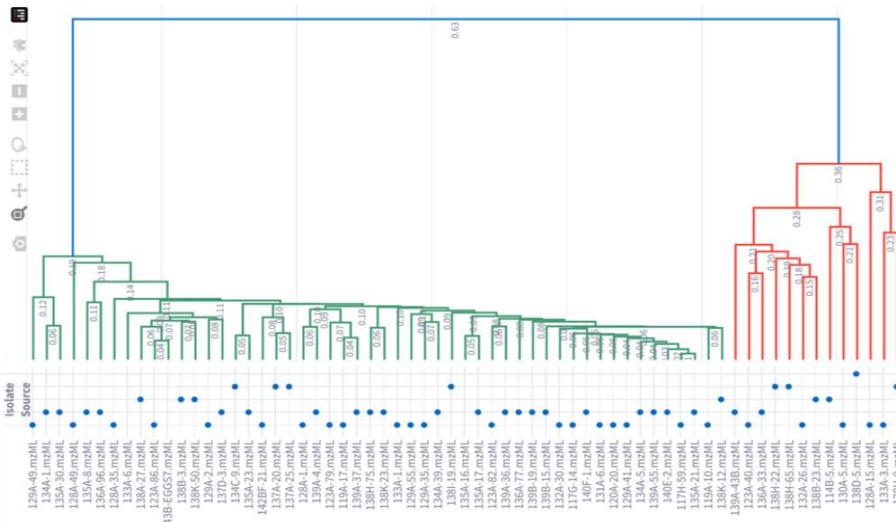
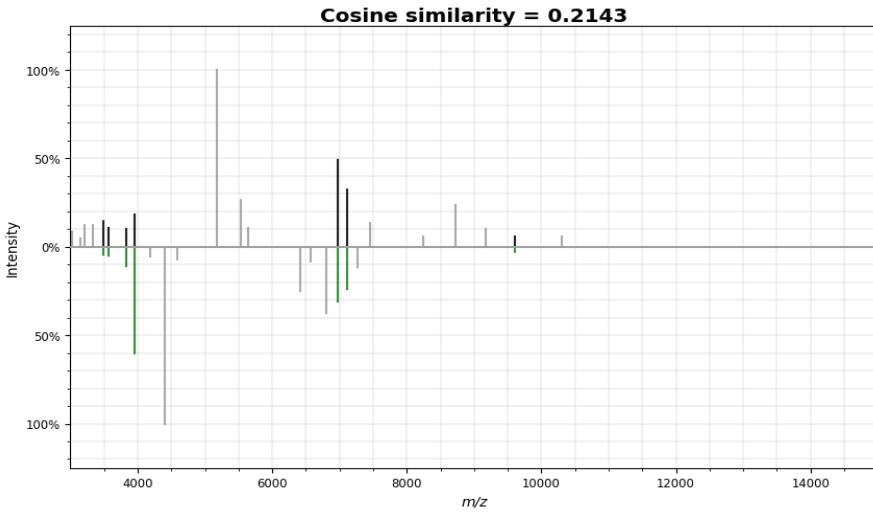
B. Murphy



N. Krull

IDBac

- Available on the web soon!
- *Taxa ID.* Query MS spectrum of ‘unknown’ against a public repository (>4500 seeds)
- Integrate metadata into dendrogram of bacteria



M. Wang



M. Strobel



L. Sanchez



R. Shepherd

code4np.github.io



Chase Clark



Joseph Egan

Computational Pharmacognosy

Posts About Contributors

Computational Pharmacognosy

Feb 4, 2024 Chase M Clark

Part 4: Introduction to Analysis and Plotting Mass Spectrometry Data in R

BEGINNER MASS SPECTROMETRY R

An introduction to working with mass spectrometry data in R

Categories All (4) beginner (4) mass spectrometry (4) r (4)

Feb 3, 2024 Chase M Clark

Part 3: Reading mzXML/mzML into R

BEGINNER MASS SPECTROMETRY R

An introduction to working with mass spectrometry data in R

Feb 2, 2024 Chase M Clark

Part 2: Anatomy of an mzXML/mzML file

BEGINNER MASS SPECTROMETRY R

What mass spectrometry data looks like, using R

Feb 1, 2024 Chase M Clark

Part 1: Mass Spectrometry Data

BEGINNER MASS SPECTROMETRY R

First in a series of introductory posts about working with mass spectrometry data in R

Computational Pharmacognosy

Posts About Contributors

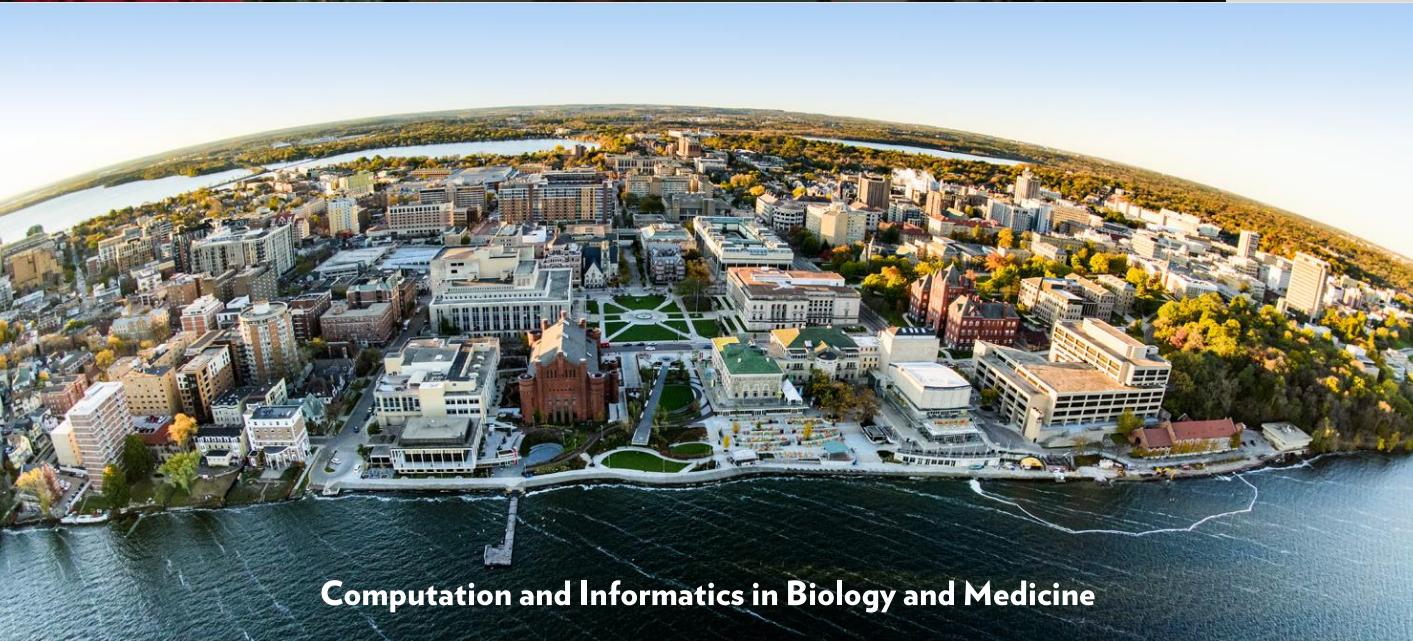
Plot mirror spectra

Rectangles were widened here so they would be easier to see on this blog.

```
transform_df <- function(df){  
  df$x1 = df$mass - 0.75  
  df$x2 = df$mass + 0.75  
  df$y1 = 0  
  df$y2 = df$intensity  
  df2=df  
  df2$mass <- NULL  
  df2$intensity <- NULL  
  return(df2)  
}  
  
centroid_plot <- function(df1, df2, top_color="red", bottom_color="blue"){  
  df2$intensity <- -df2$intensity  
  df1 <- transform_df(df1)  
  df2 <- transform_df(df2)  
  
  p <- ggplot() +  
    scale_x_continuous(name="m/z") +  
    scale_y_continuous(name="Intensity") +  
    geom_rect(data=df1, mapping=aes(xmin=x1, xmax=x2, ymin=y1, ymax=y2), color="grey30", fill=top_color, alpha=0.5)  
  geom_rect(data=df2, mapping=aes(xmin=x1, xmax=x2, ymin=y1, ymax=y2), fill=bottom_color, color="grey30", alpha=0.5)  
  
  return(p)  
}  
  
df <- as.data.table(mzR::peaks(centroid_msfile_handle, 2243))  
colnames(df) <- c("mass", "intensity")  
  
ggplotly(centroid_plot(df, df))
```

Compare with a spectrum in the GNPS library

Download [GNPS library spectrum CCMSLIB00000072054](#) (Acyl desferrioxamine C14).



socialgene.github.io



Funding
NLM: 5T15LM007359
NIGMS: R35GM133776