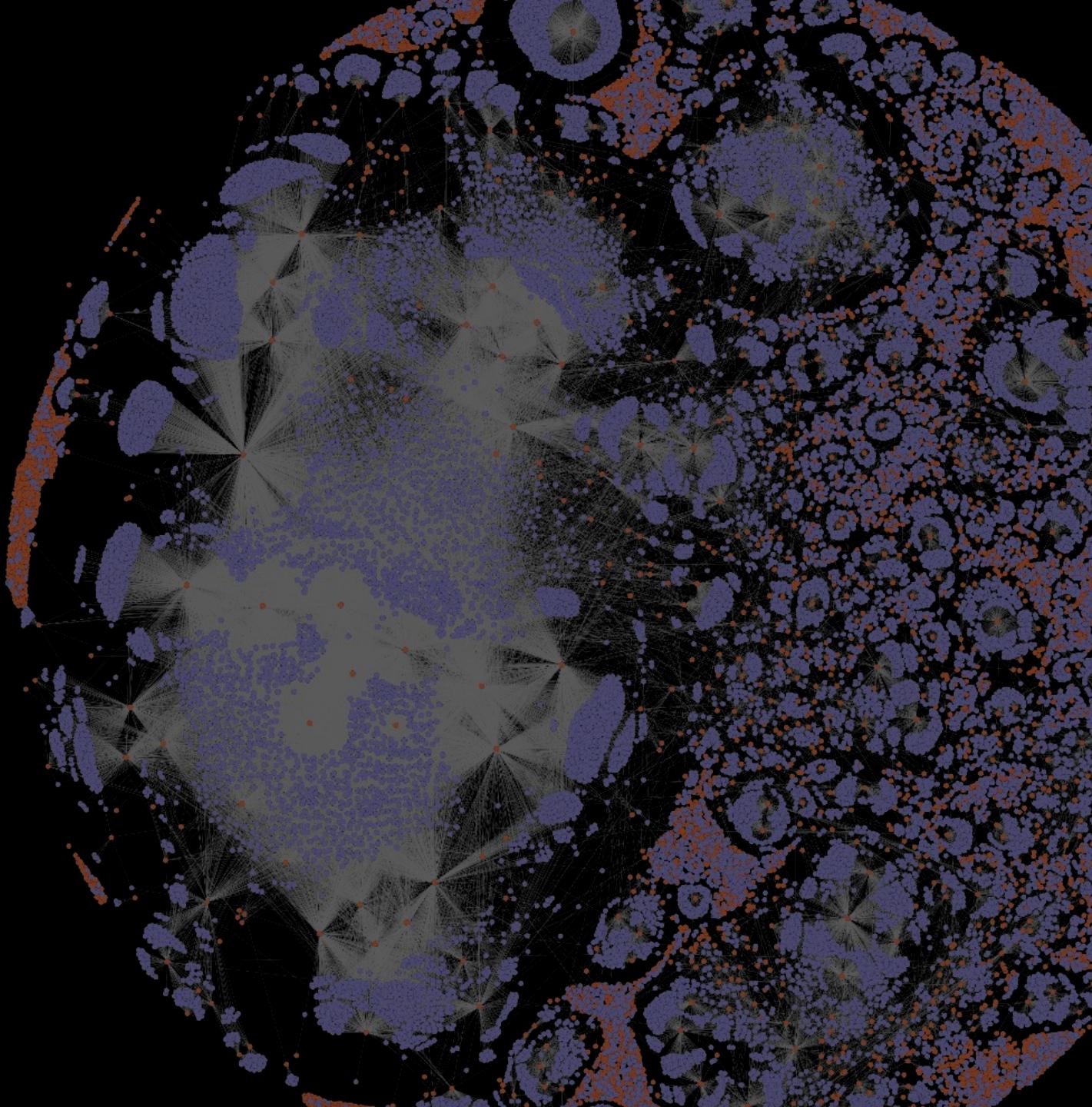


SocialGene: Large Scale Knowledge Graphs for Microbial Based Drug Discovery

Chase Clark
Postdoctoral Research Associate
Computation and Informatics in
Biology and Medicine Training Program

Jason Kwan Lab
School of Pharmacy
University of Wisconsin-Madison

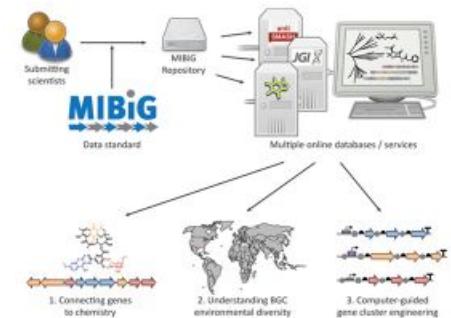


SocialGene

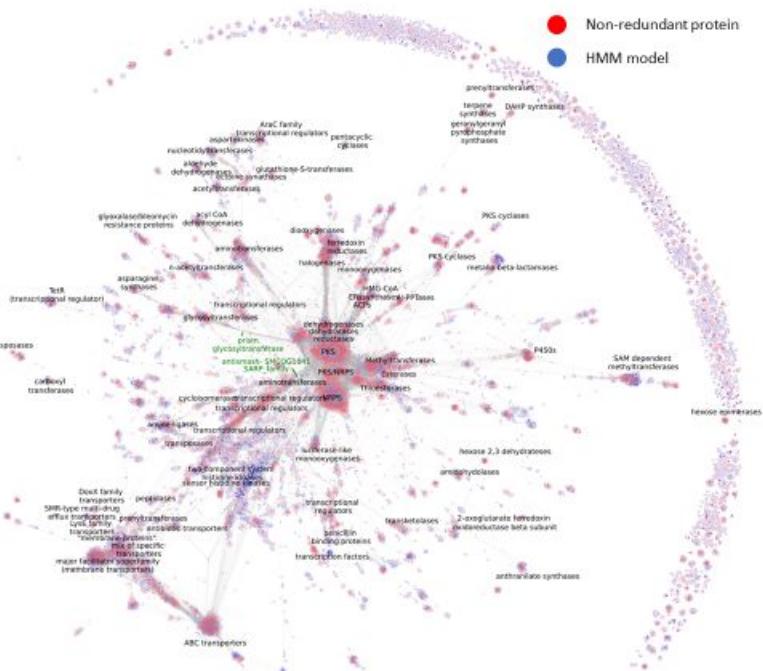
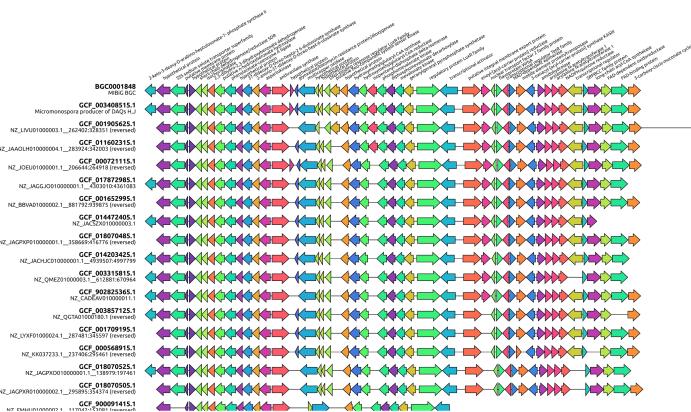
Knowledge graphs for drug discovery
Functional characterization of proteins
and searching for similar BGCs

Input BGC

Search 200k
genomes



SocialGene



SocialGene

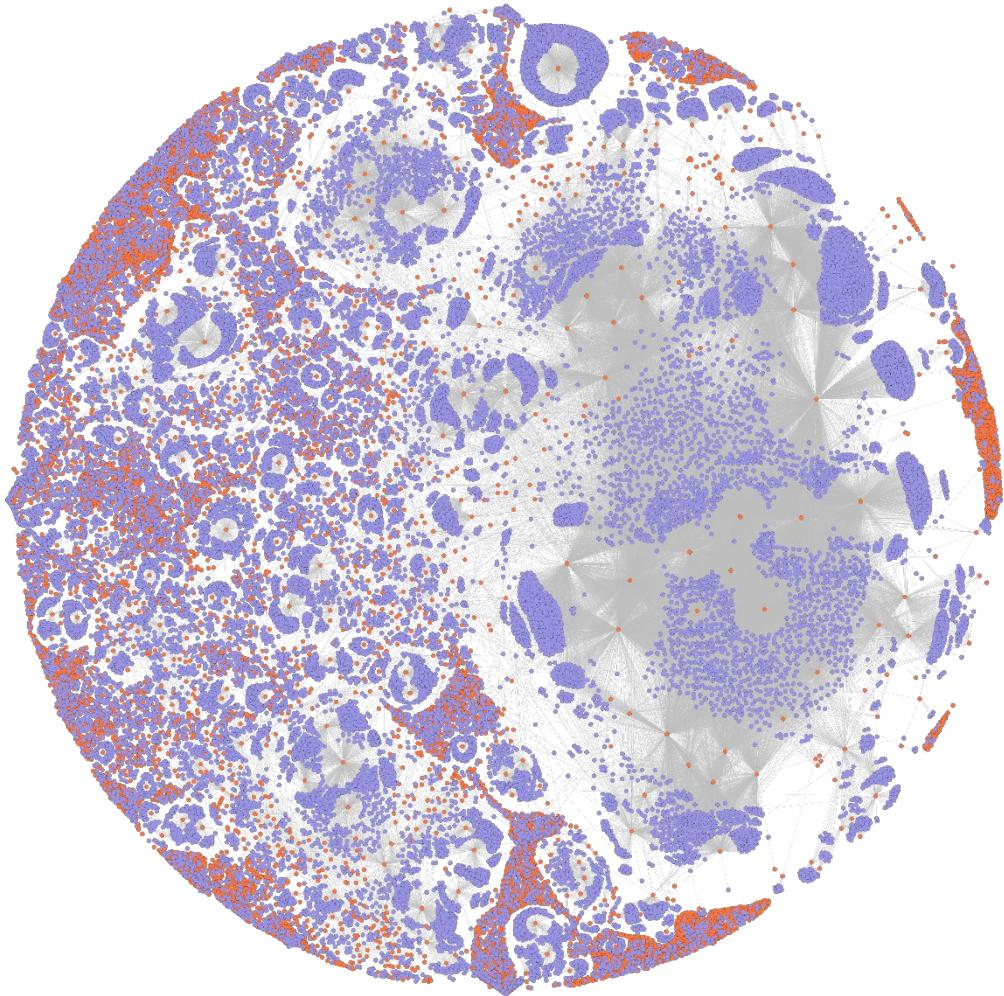
Knowledge graphs for drug discovery
Functional characterization of proteins
and searching for similar BGCs



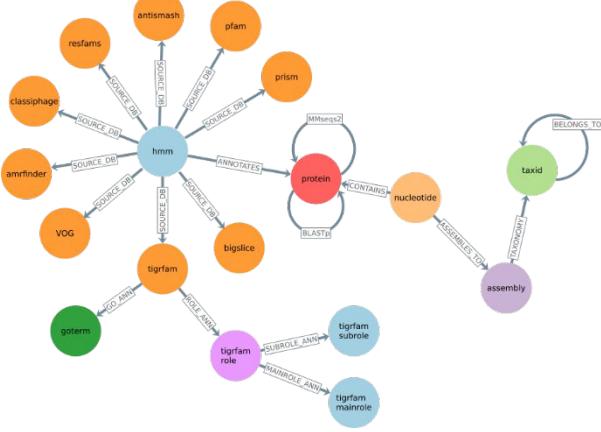
nextflow

django

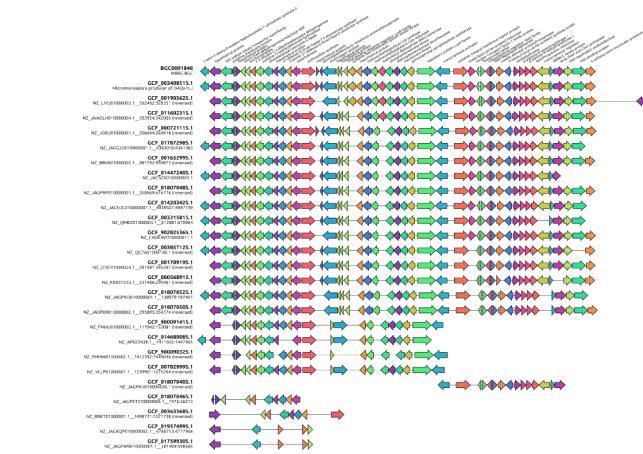
neo4j



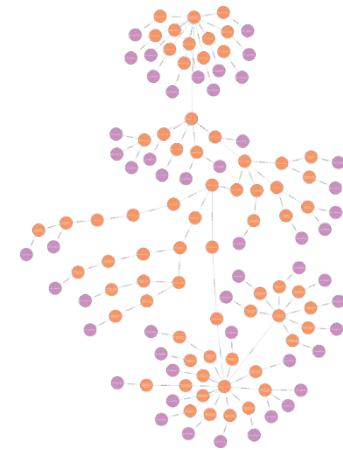
SocialGene



Reproducibly create Graph Databases
of up to 100s of thousands of genomes



Search for remotely-homologous
proteins and biosynthetic gene clusters



Large-scale phylogenetic and distribution patterns
of enzyme occurrence for targeted discovery of enzyme
function

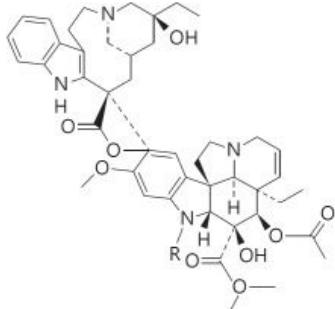
“Drugs”
“Natural Products”
“Specialized Metabolites”

“Drugs”
“Natural Products”
“Specialized Metabolites”

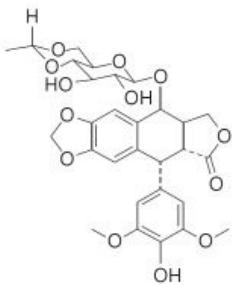
Over 60% of current anticancer drugs are derived in one way or another from natural sources

Cragg, Gordon M, and John M Pezzuto. “Natural Products as a Vital Source for the Discovery of Cancer Chemotherapeutic and Chemopreventive Agents.” *Medical principles and practice : international journal of the Kuwait University, Health Science Centre* vol. 25 Suppl 2,Suppl 2 (2016): 41-59. doi:10.1159/000443404

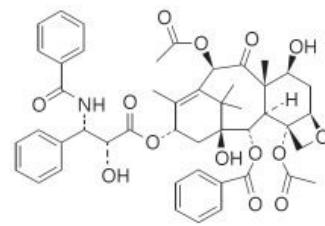
“Natural Products”



1 VBL R = CH₃
2 VCR R = CHO



3 Etoposide

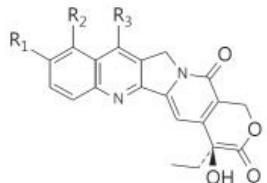


4 Paclitaxel (Taxol™)



5 Docetaxel (Taxotere™) R₁ = R₂ = H

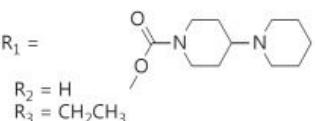
6 Cabazitaxel R₁ = R₂ = CH₃



7 CPT R₁ = R₂ = R₃ = H

8 Topotecan R₁ = OH; R₂ = CH₂NH(CH₃)₂; R₃ = H

9 Irinotecan R₁ =



R₂ = H
R₃ = CH₂CH₃

10 Belotecan R₁ = R₂ = H; R₃ = (CH₂)₂NHCH(CH₃)₂

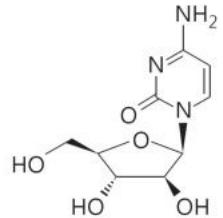
11 Cositecan R₁ = R₂ = H; R₃ = (CH₂)₂Si(CH₃)₃

12 SN-38 R₁ = OH; R₂ = H; R₃ = CH₂CH₃

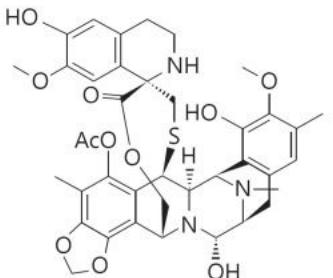
Natural products from plants

- Cragg, Gordon M, and John M Pezzuto. “Natural Products as a Vital Source for the Discovery of Cancer Chemotherapeutic and Chemopreventive Agents.” *Medical principles and practice : international journal of the Kuwait University, Health Science Centre* vol. 25 Suppl 2,Suppl 2 (2016): 41-59. doi:10.1159/000443404
- Cragg, G M et al. “The taxol supply crisis. New NCI policies for handling the large-scale production of novel natural product anticancer and anti-HIV agents.” *Journal of natural products* vol. 56,10 (1993): 1657-68. doi:10.1021/np50100a001
- Stierle, A et al. “The search for a taxol-producing microorganism among the endophytic fungi of the Pacific yew, *Taxus brevifolia*.” *Journal of natural products* vol. 58,9 (1995): 1315-24. doi:10.1021/np50123a002

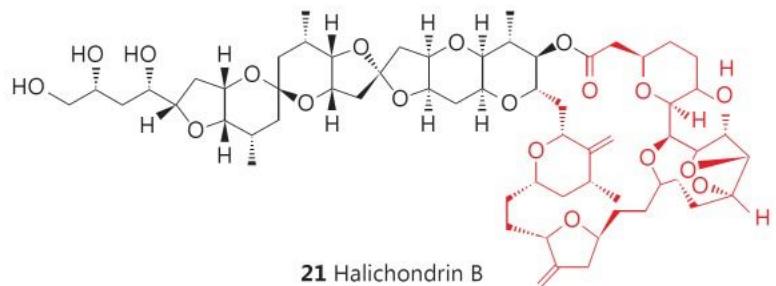
“Natural Products”



19 Cytarabine



20 Trabectedin (ET743; Yondelis®)

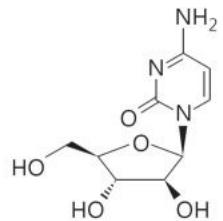


21 Halichondrin B

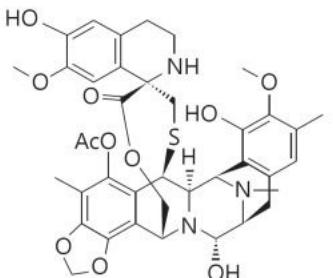
Natural products from
marine organisms

Cragg, Gordon M, and John M Pezzuto. “Natural Products as a Vital Source for the Discovery of Cancer Chemotherapeutic and Chemopreventive Agents.” *Medical principles and practice : international journal of the Kuwait University, Health Science Centre* vol. 25 Suppl 2,Suppl 2 (2016): 41-59. doi:10.1159/000443404

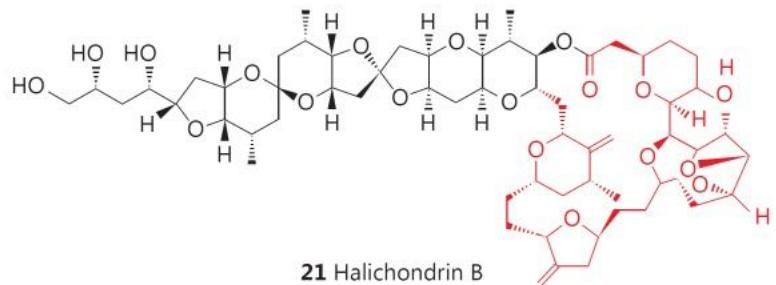
“Natural Products”



19 Cytarabine



20 Trabectedin (ET743; Yondelis®)



21 Halichondrin B

Supply issue can
be a huge problem

Cragg, Gordon M, and John M Pezzuto. “Natural Products as a Vital Source for the Discovery of Cancer Chemotherapeutic and Chemopreventive Agents.” *Medical principles and practice : international journal of the Kuwait University, Health Science Centre* vol. 25 Suppl 2,Suppl 2 (2016): 41-59. doi:10.1159/000443404

“Natural Products”



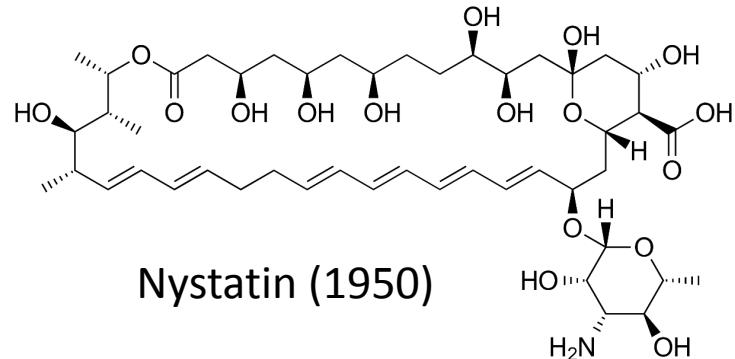
Dr. Elizabeth Lee Hazen (left), microbiologist

Dr. Rachel Fuller Brown (right), chemist



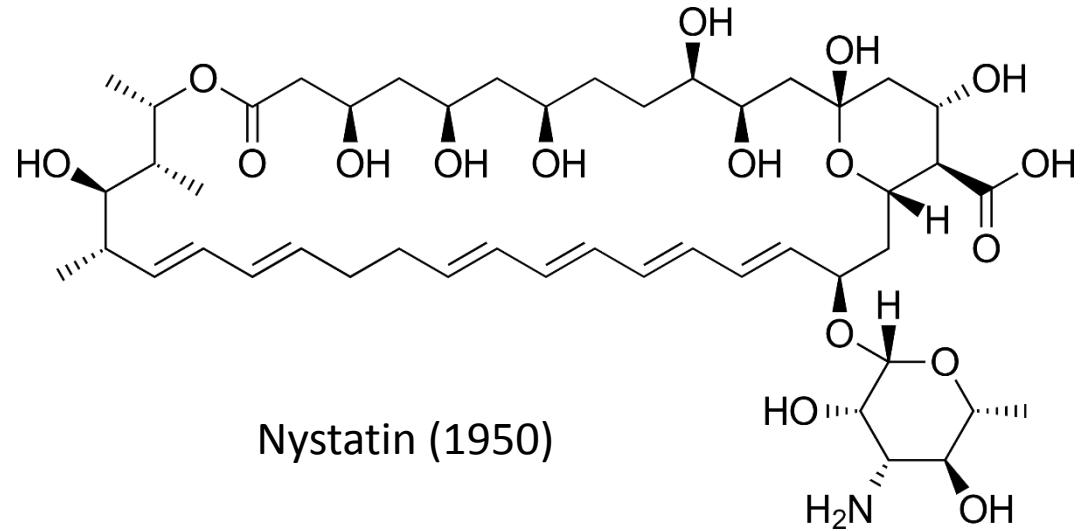
Streptomyces noursei

Natural products from
bacteria and fungi





Dr. Elizabeth Lee Hazen (left), microbiologist
Dr. Rachel Fuller Brown (right), chemist



Streptomyces noursei

MiBiG Repository of Known Biosynthetic Gene Clusters

BGC0000115: nystatin A1 biosynthetic gene cluster from *Streptomyces noursei* ATCC 11455

Location: 1 - 123,580 nt. (total: 123,580 nt).
This entry is originally from NCBI GenBank
AF263912.1.

Download region SVG Download Cluster GenBank file View antiSMASH-generated output

Gene details

Select a gene to view the details available for it

Legend:

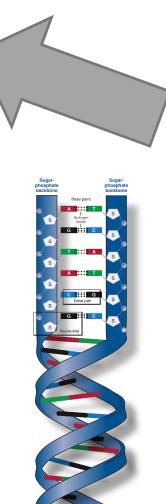
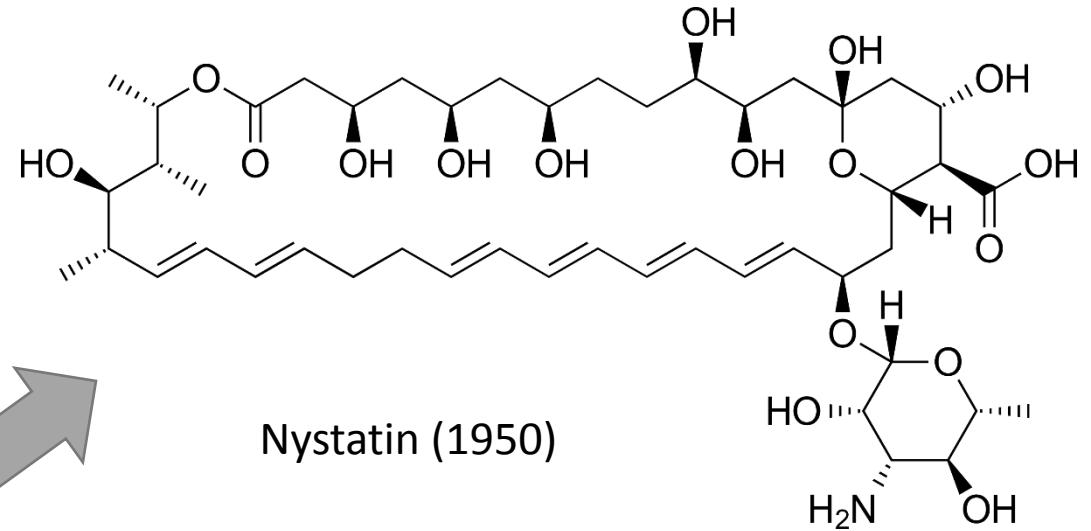
- core biosynthetic genes
- additional biosynthetic genes
- transport-related genes
- regulatory genes
- other genes
- resistance
- TTA codons

reset view zoom to selection

General Compounds Genes Polyketide Saccharide History NRPS/PKS domains KnownClusterBlast

General information about the BGC

MiBiG accession	BGC0000115
Short description	nystatin A1 biosynthetic gene cluster from <i>Streptomyces noursei</i> ATCC 11455
Status	Minimal annotation: no ? Completeness: complete ?
Biosynthetic class(es)	Polyketide (Polyene) Saccharide (hybrid/tailoring)
Loci	NCBP GenBank: AF263912.1
Compounds	nystatin A1
Species	<i>Streptomyces noursei</i> ATCC 11455 [taxonomy]
References	Biosynthesis of the polyene antifungal antibiotic nystatin in <i>Streptomyces noursei</i> ATCC 11455: analysis of the gene cluster and deduction of the biosynthetic pathway. Brautaset T et al., Chem Biol (2000) PMID:10873841



Streptomyces noursei

MISSING TRANSITION

We've barely scratched the surface...

> [PLOS Biol.](#) 2021 Nov 9;19(11):e3001421. doi: 10.1371/journal.pbio.3001421.

eCollection 2021 Nov.

Exploring bacterial diversity via a curated and searchable snapshot of archived DNA sequences

Grace A Blackwell ^{1 2}, Martin Hunt ^{1 3}, Kerri M Malone ¹, Leandro Lima ¹, Gal Horesh ², Blaise T F Alako ¹, Nicholas R Thomson ^{2 4}, Zamin Iqbal ¹

Affiliations + expand

PMID: 34752446 PMCID: [PMC8577725](#) DOI: [10.1371/journal.pbio.3001421](#)

“639,981 high-quality genomes emphasised the uneven species composition of the ENA/public databases, with just 20 of the total 2,336 species making up 90% of the genomes.”

We've barely scratched the surface...^A

> PLoS Biol. 2021 Nov 9;19(11):e3001421. doi: 10.1371/journal.pbio.3001421.

eCollection 2021 Nov.

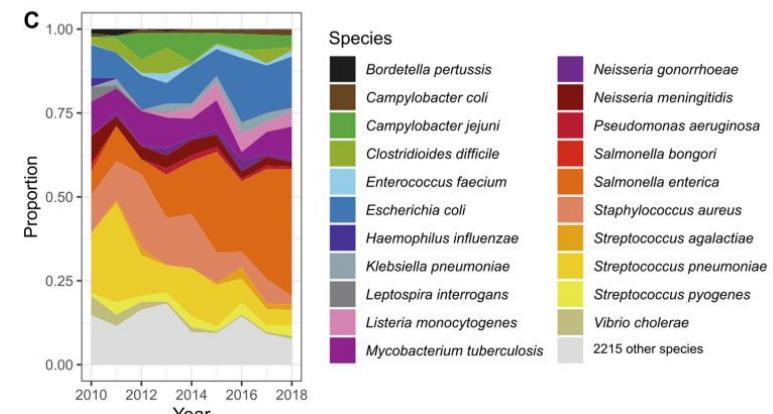
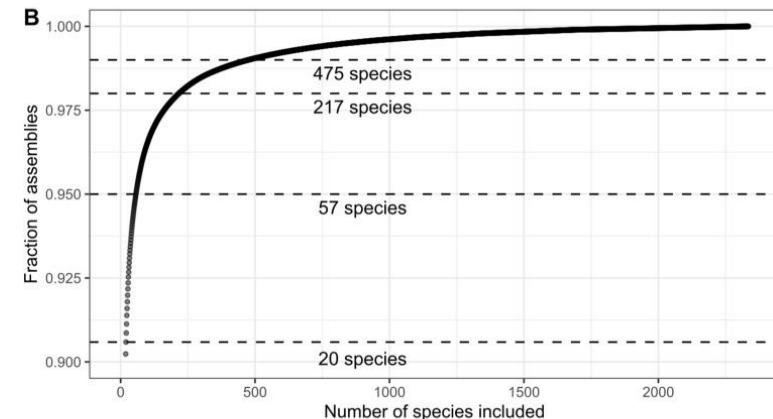
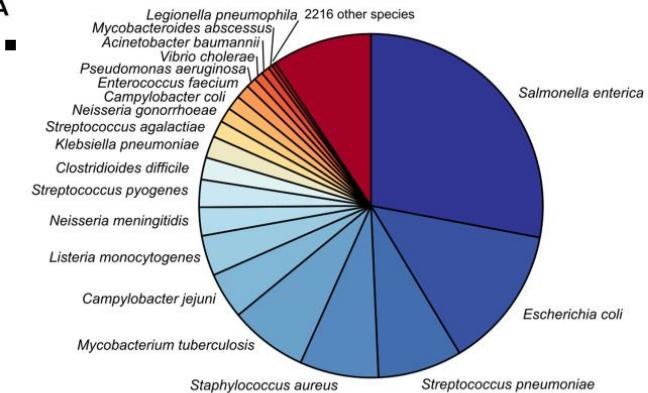
Exploring bacterial diversity via a curated and searchable snapshot of archived DNA sequences

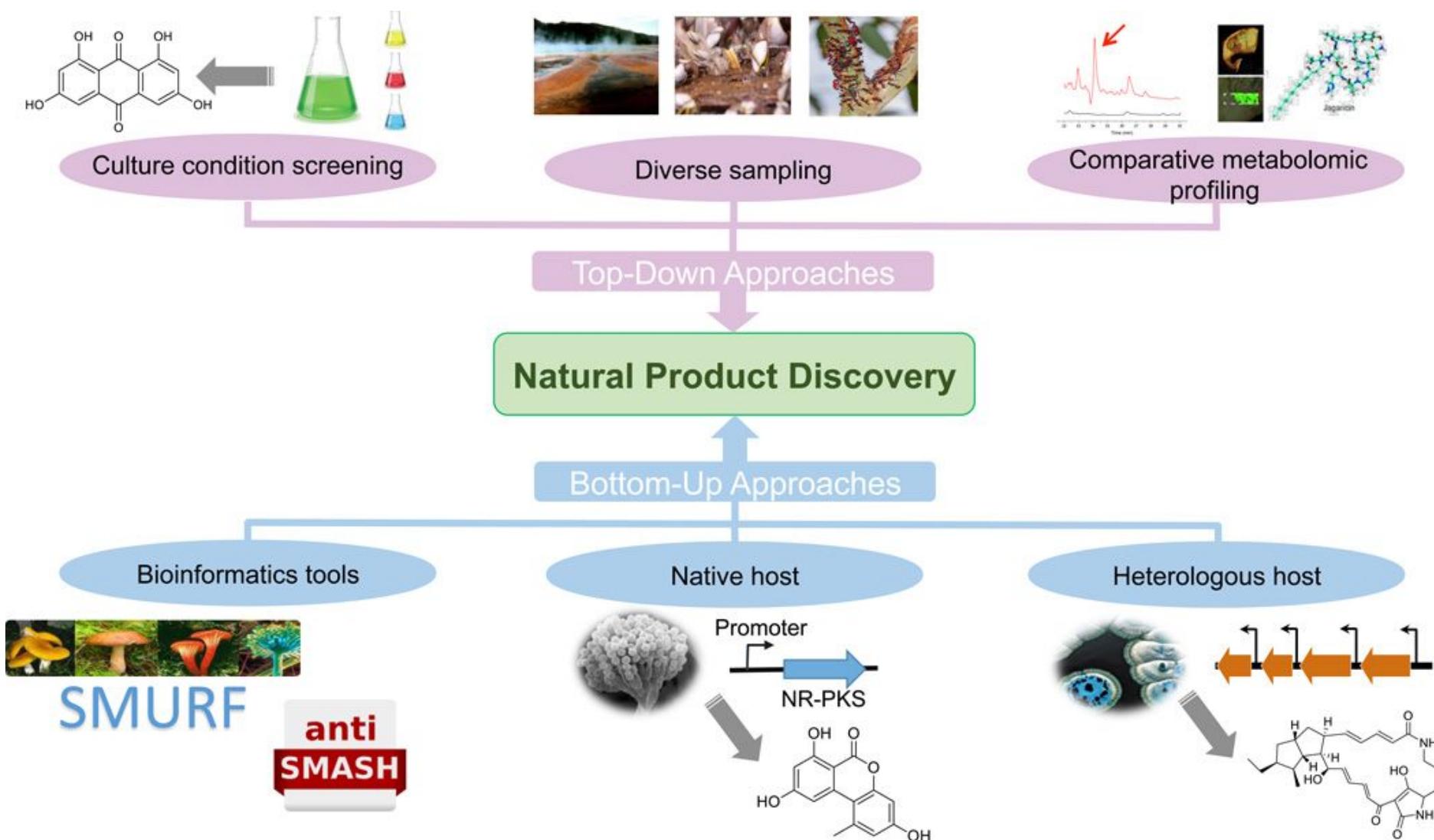
Grace A Blackwell ^{1 2}, Martin Hunt ^{1 3}, Kerri M Malone ¹, Leandro Lima ¹, Gal Horesh ², Blaise T F Alako ¹, Nicholas R Thomson ^{2 4}, Zamin Iqbal ¹

Affiliations + expand

PMID: 34752446 PMCID: [PMC8577725](#) DOI: [10.1371/journal.pbio.3001421](#)

“639,981 high-quality genomes emphasised the uneven species composition of the ENA/public databases, with just 20 of the total 2,336 species making up 90% of the genomes.”





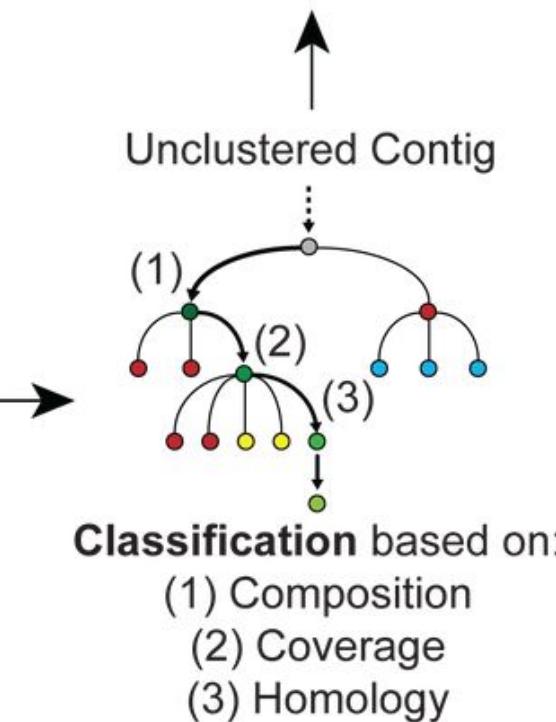
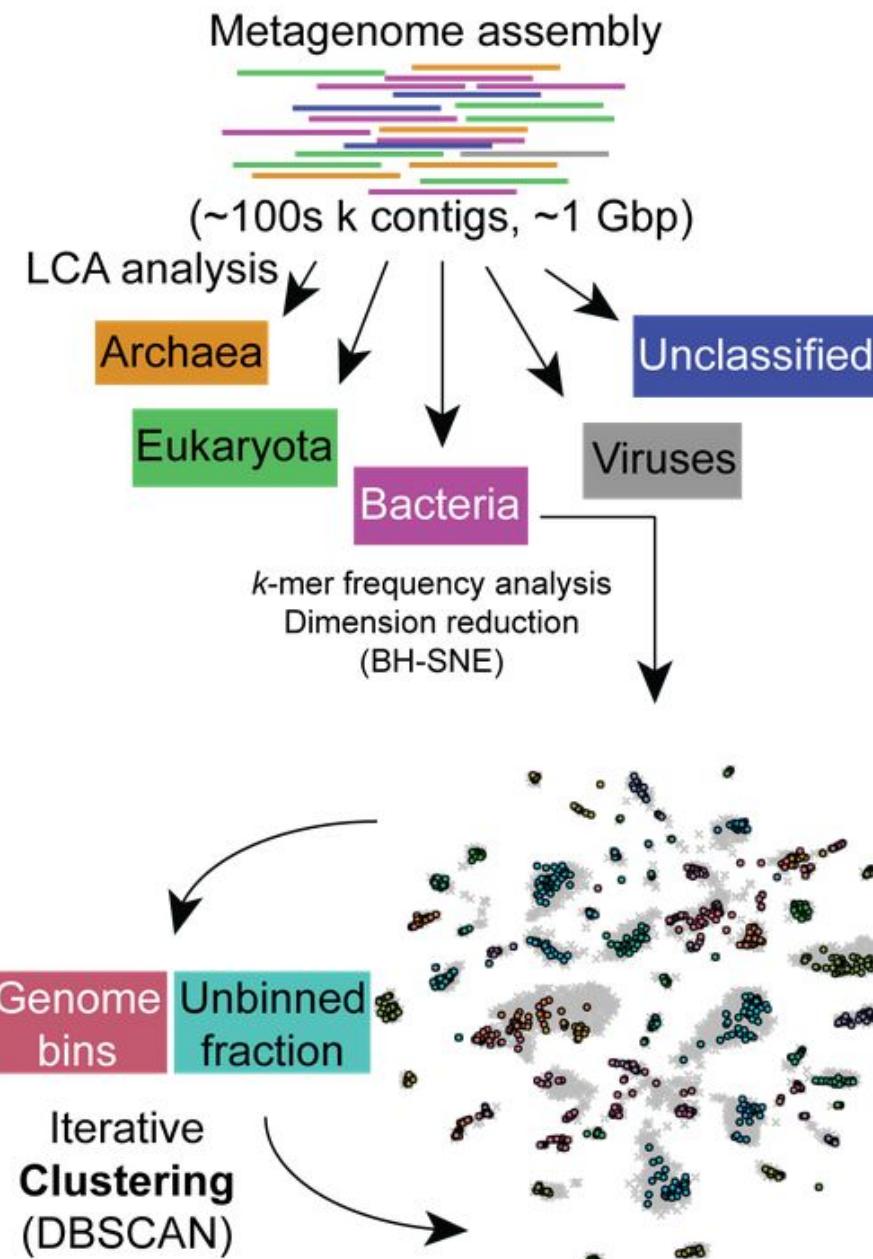
Luo, Yunzi et al. "Recent advances in natural product discovery." *Current opinion in biotechnology* vol. 30 (2014): 230-7. doi:10.1016/j.copbio.2014.09.002

- Bacteria and fungi produce chemical compounds that can be repurposed as drugs
- We've barely genome sequenced bacterial/fungal diversity



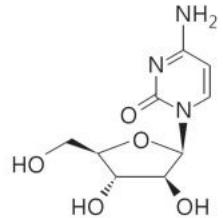
WISCONSIN
UNIVERSITY OF WISCONSIN-MADISON

Postdoctoral Research Bioinformatics, Metagenomics

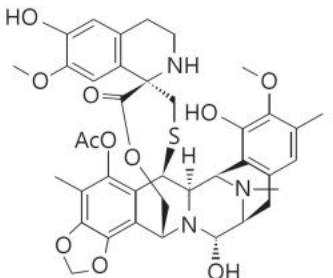


Miller II, et al. Autometa: automated extraction of microbial genomes from individual shotgun metagenomes. Nucleic Acids Res. 2019 Jun 4;47(10):e57.

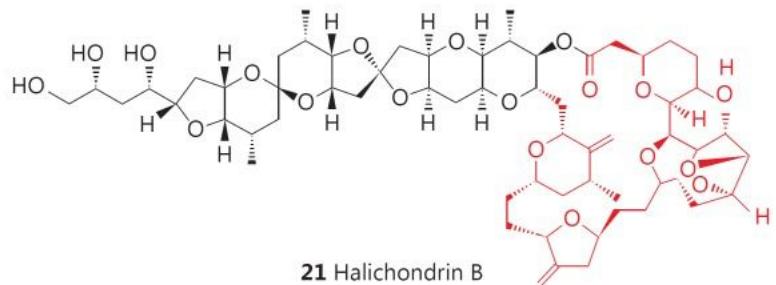
“Natural Products”



19 Cytarabine



20 Trabectedin (ET743; Yondelis®)



21 Halichondrin B

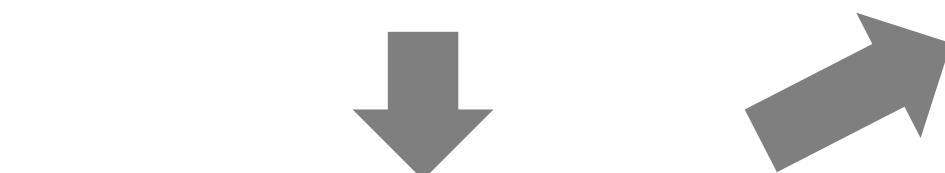
Supply issue can
be a huge problem

Cragg, Gordon M, and John M Pezzuto. “Natural Products as a Vital Source for the Discovery of Cancer Chemotherapeutic and Chemopreventive Agents.” *Medical principles and practice : international journal of the Kuwait University, Health Science Centre* vol. 25 Suppl 2,Suppl 2 (2016): 41-59. doi:10.1159/000443404

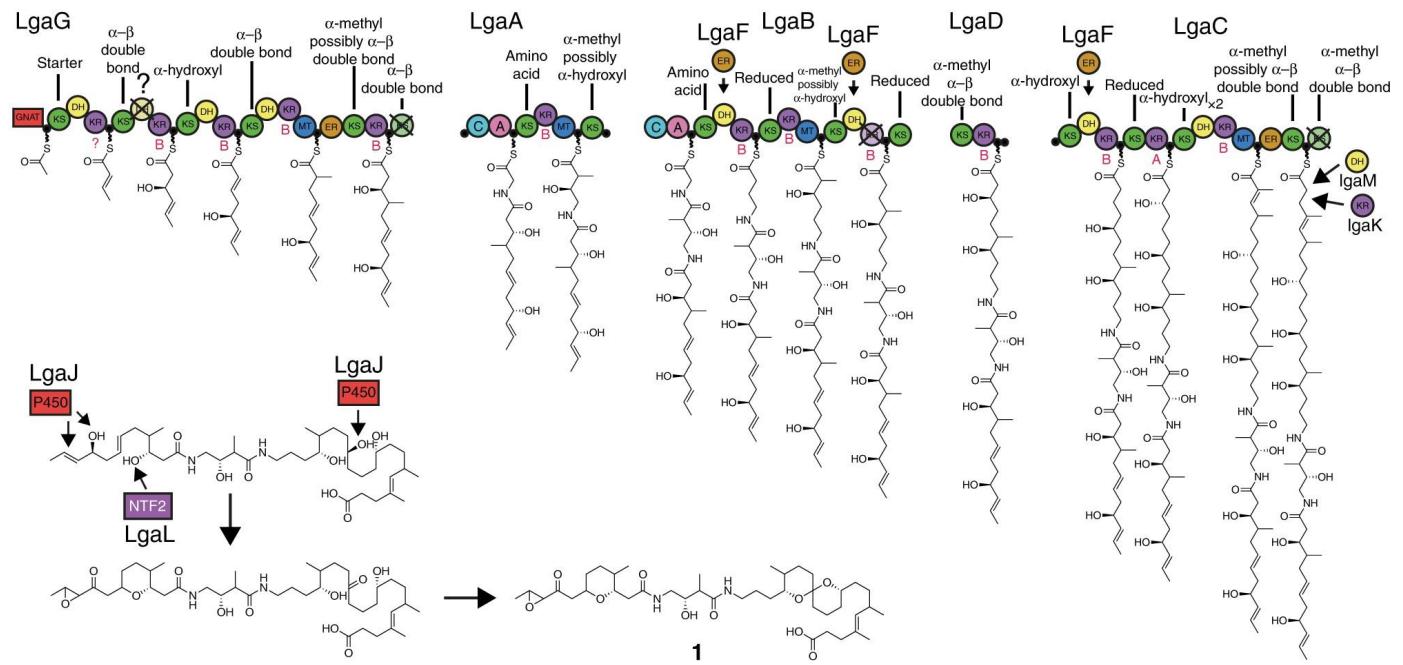


PHOTO BY LAURA FLÓREZ

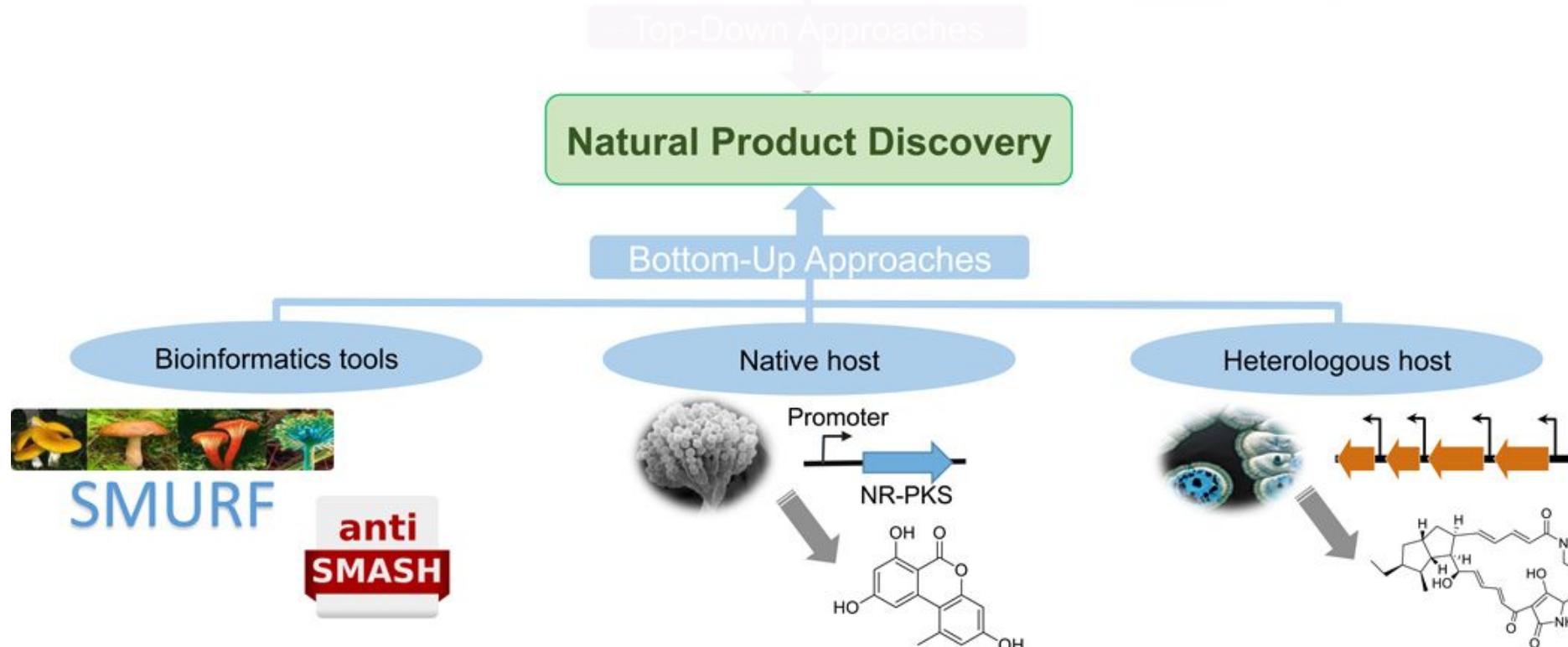
Metagenomics



Metagenomics revealed the biosynthetic gene cluster for lagriamide



Waterworth SC, Flórez LV, Rees ER, Hertweck C, Kaltenpoth M, Kwan JC. Horizontal Gene Transfer to a Defensive Symbiont with a Reduced Genome in a Multipartite Beetle Microbiome. *mBio*. 2020 Feb 25;11(1):e02430-19.



Luo, Yunzi et al. "Recent advances in natural product discovery." *Current opinion in biotechnology* vol. 30 (2014): 230-7. doi:10.1016/j.copbio.2014.09.002

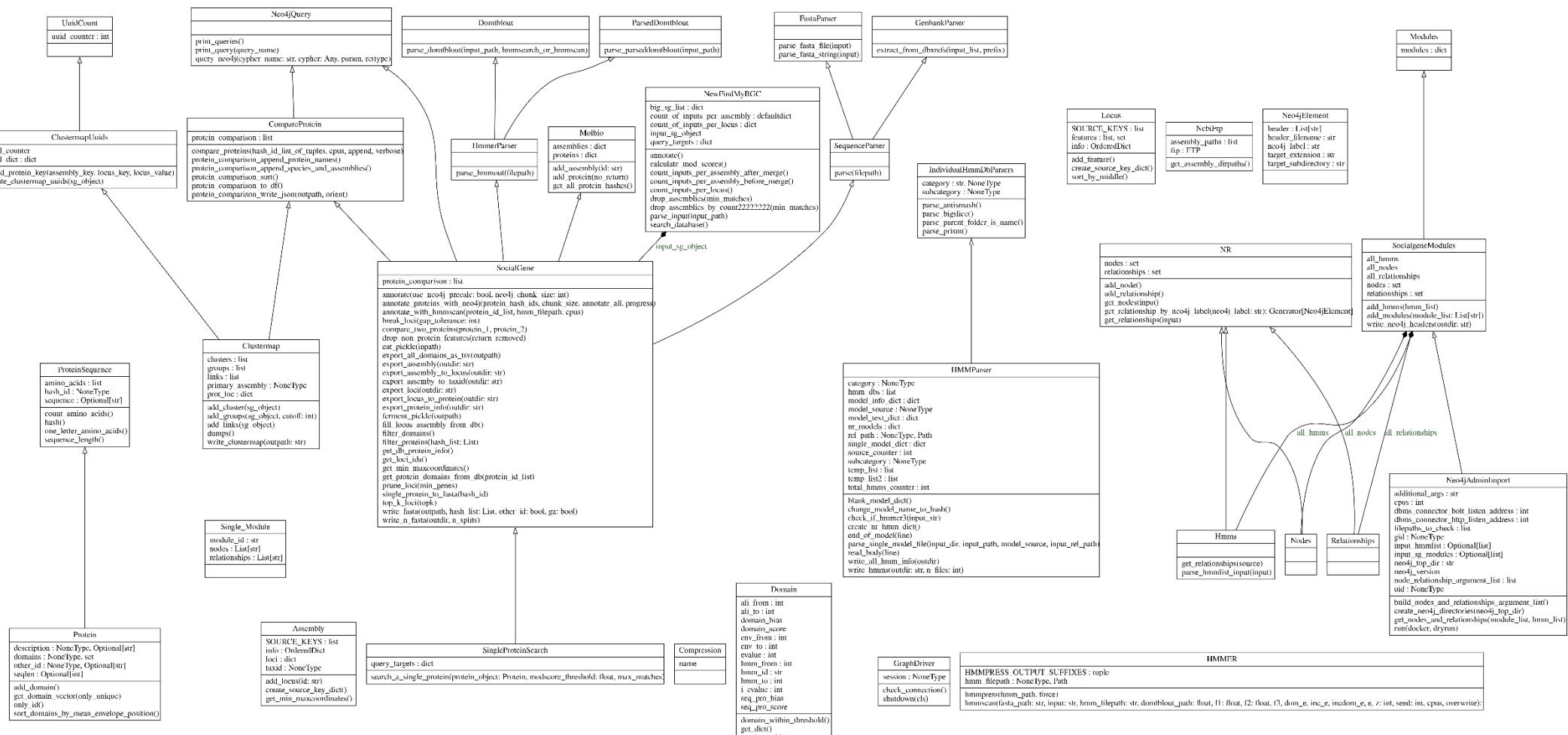
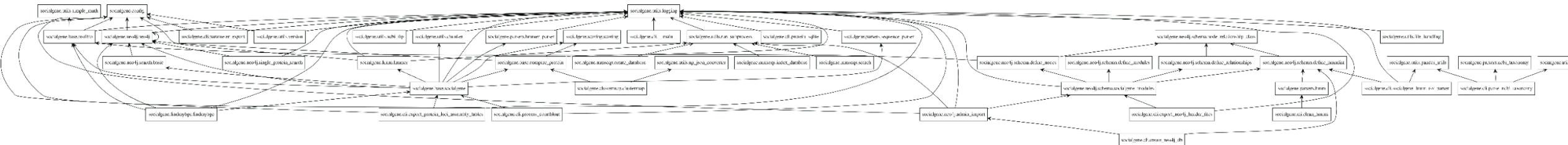
What do we need?

- Collect/process/standardize genomes and the proteins they encode
 - Protein identifiers aren't constant across databases
 - ID conversion files are huge, online API rate-limited, doesn't help with in-house proteins
- All-vs-All sequence comparison too slow for real-time search, results too large to store
- Downloading annotated genomes doesn't because we can't annotate in-house sequences in the same way

A few challenges, of many

Problems	Solutions
Protein identifiers aren't constant across databases	Hash protein sequences and use the hash as the identifier (<i>sgpy</i> and <i>crabhash</i>)
All-vs-All sequence comparison too slow for real-time search, results too large to store	Annotate proteins based on functional domains, use domains as entry-point for search and similarity (<i>HMMER</i> + many, many models, <i>sgnf</i>)
Storing/data-retrieval / machine learning	Graph database (through <i>sgpy</i> , <i>sgnf</i>)
Doing all steps reproducibly and “easily”	(<i>sgpy</i> , <i>sgnf</i>)





<https://github.com/socialgene/sgpy>

Input genome(s)
1-????



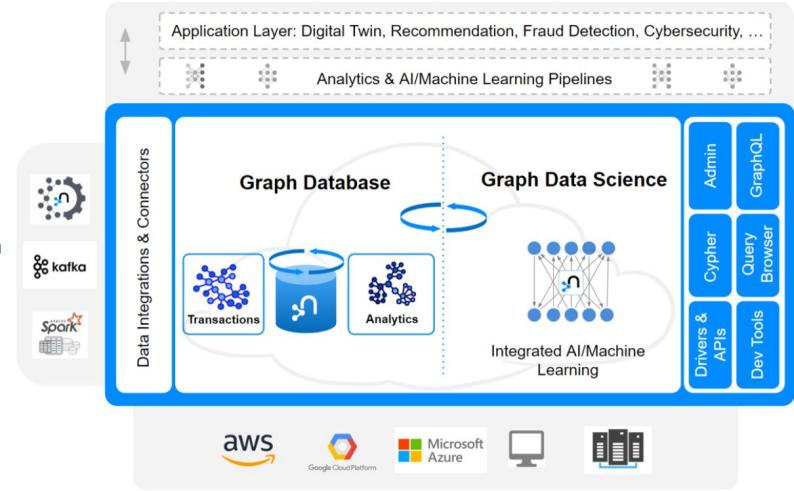
nextflow

Reproducibly run multiple
scripts/software in parallel, across
different machines



 **neo4j**
Graph database

Creating knowledge graphs with Neo4j

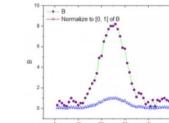


New in Neo4j GDS 1.6

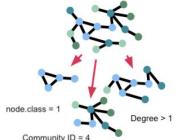
The Fastest Path to Production



Multiple ML models for Community Users (3)



Scaling & Normalization



In-Memory Graph Filtering & Subgraph Projection



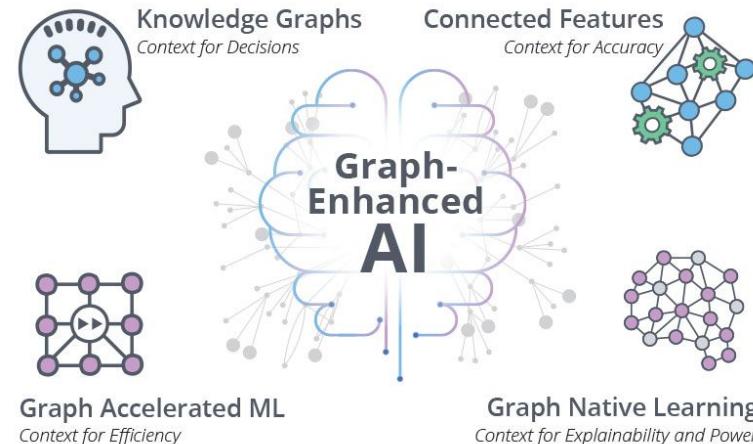
Faster Graph Embeddings, Seeding, ↑ Accuracy



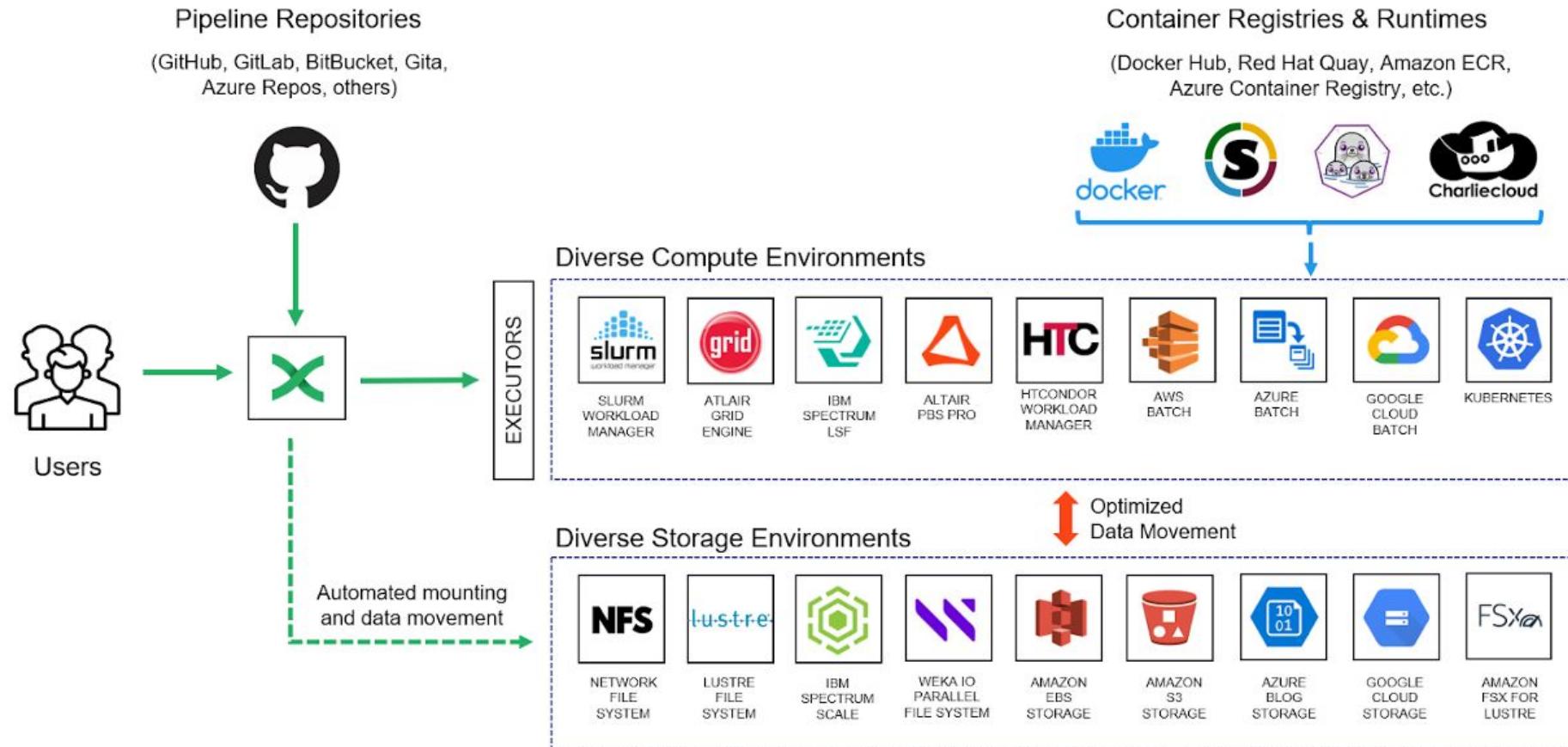
More Supervised ML Features

MLOps

Enterprise Administration

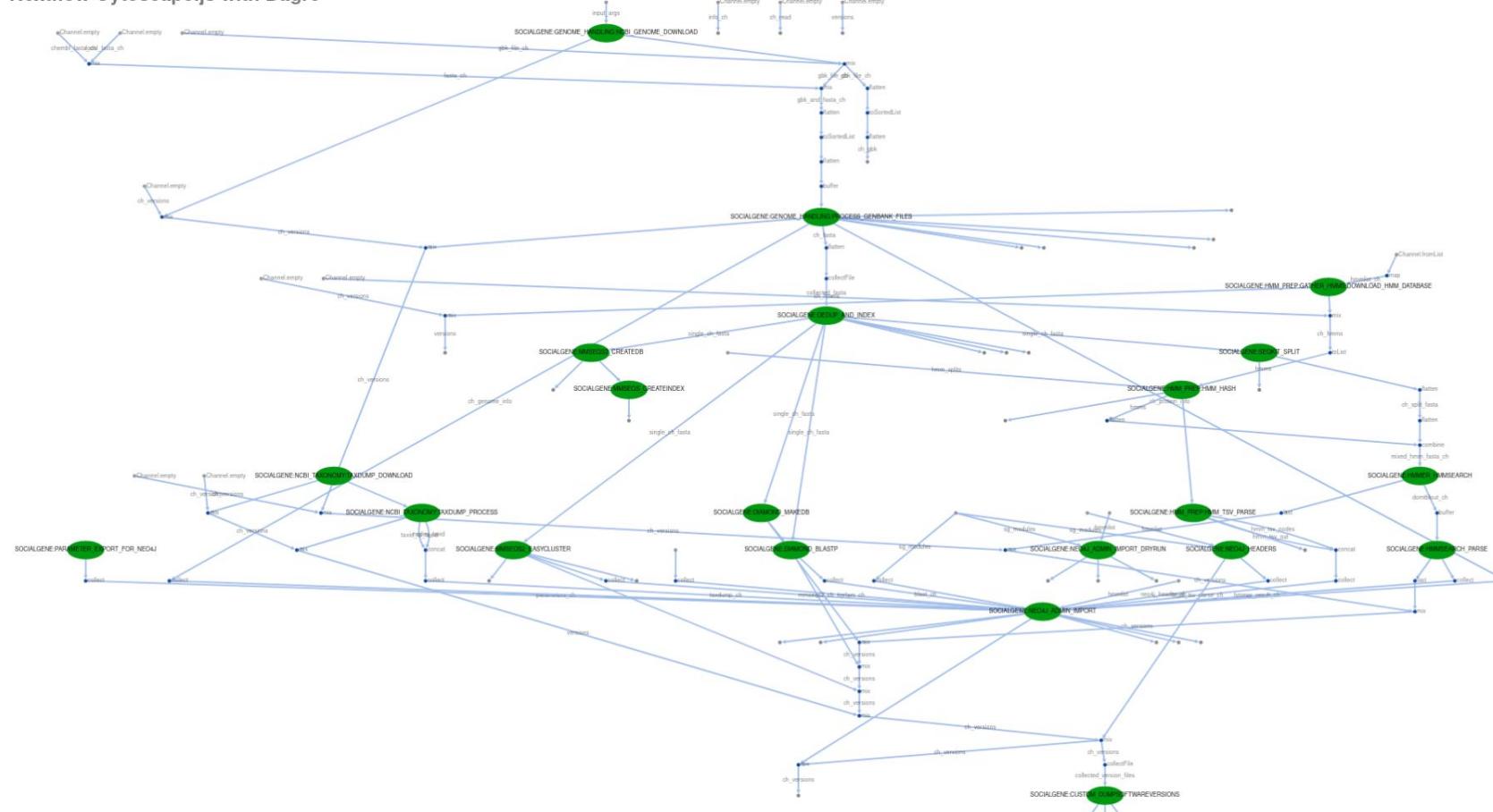


Simplifying complicated bioinformatics workflows with **nextflow**



Why Nextflow? It's complicated

Nextflow Cytoscape.js with Dagre



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

sgnf (Private)
SocialGene's Nextflow Pipeline
Nextflow

<https://github.com/socialgene/sgnf>

Sh#\$
happens...



Chase Clark, PhD
@ChasingMicrobes

...

[99%] 5913 of 5914, failed: 1



3:22 PM · Mar 3, 2023 · 55 Views

View Tweet analytics



Sh#\$
happens...



Chase Clark, PhD @ChasingMicrobes · 1m
[100%] 5914 of 5914, cached: 5913 ✓

...



2



Nextflow workflow report

[trusting_monod] (resumed run)

Workflow execution completed successfully!

Run times

28-Feb-2023 12:09:31 - 28-Feb-2023 12:10:53 (duration: 1m 23s)

37 succeeded

5 cached

Nextflow command

```
nextflow run . -profile ultraquickstart,docker --outdir /home/chase/Documents/socialgene_data/example --outdir_download_cache /home/chase/Downloads/socialgene_data/cache -resume
```

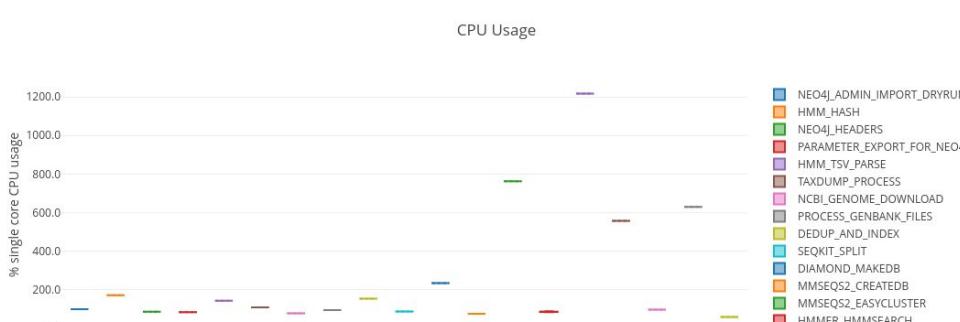
CPU-Hours	0.3 (1.2% cached)
Launch directory	/home/chase/Documents/github/kwan_lab/socialgene/sgnf
Work directory	/home/chase/Documents/github/kwan_lab/socialgene/sgnf/work
Project directory	/home/chase/Documents/github/kwan_lab/socialgene/sgnf
Script name	main.nf
Script ID	a28244b78ec8421acc43ed824425d888
Workflow session	0106669a-7144-4423-82df-a6bd864f98fd
Workflow profile	ultraquickstart,docker
Workflow container	chasemc2/socialgene-nf:0.0.1
Container engine	docker
Nextflow version	version 22.10.6, build 5843 (23-01-2023 23:20 UTC)

Resource Usage

These plots give an overview of the distribution of resource usage for each process.

CPU

Raw Usage % Allocated



Nextflow workflow report

[trusting_monod] (resumed run)

Workflow execution completed successfully!

Run times

28-Feb-2023 12:09:31 - 28-Feb-2023 12:10:53 (duration: 1m 23s)



Nextflow command

```
nextflow run . -profile ultraquickstart_docker --outdir /home/chase/Documents/socialgene_data/example --outdir_download_cache /home/chase /Documents/socialgene_data/cache -resume
```

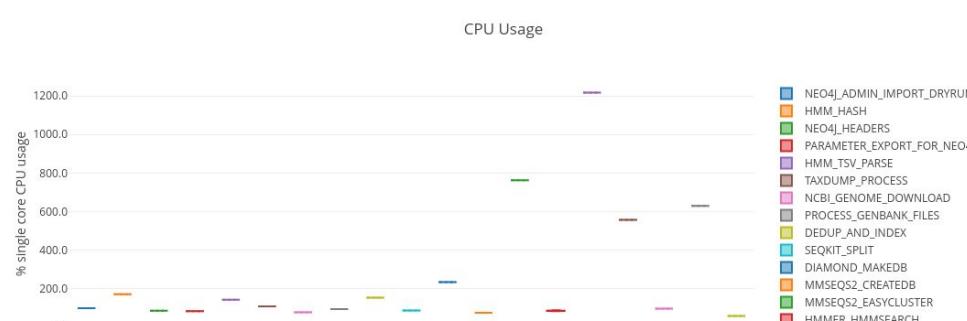
CPU-Hours	0.3 (1.2% cached)
Launch directory	/home/chase/Documents/github/kwan_lab/socialgene/sgnf
Work directory	/home/chase/Documents/github/kwan_lab/socialgene/sgnf/work
Project directory	/home/chase/Documents/github/kwan_lab/socialgene/sgnf
Script name	main.nf
Script ID	a28244b78ec8421acc43ed824425d888
Workflow session	0106669a-7144-4423-82df-a6bd864f98fd
Workflow profile	ultraquickstart_docker
Workflow container	chasemc2/socialgene-nf:0.0.1
Container engine	docker
Nextflow version	version 22.10.6, build 5843 (23-01-2023 23:20 UTC)

Resource Usage

These plots give an overview of the distribution of resource usage for each process.

CPU

Raw Usage % Allocated

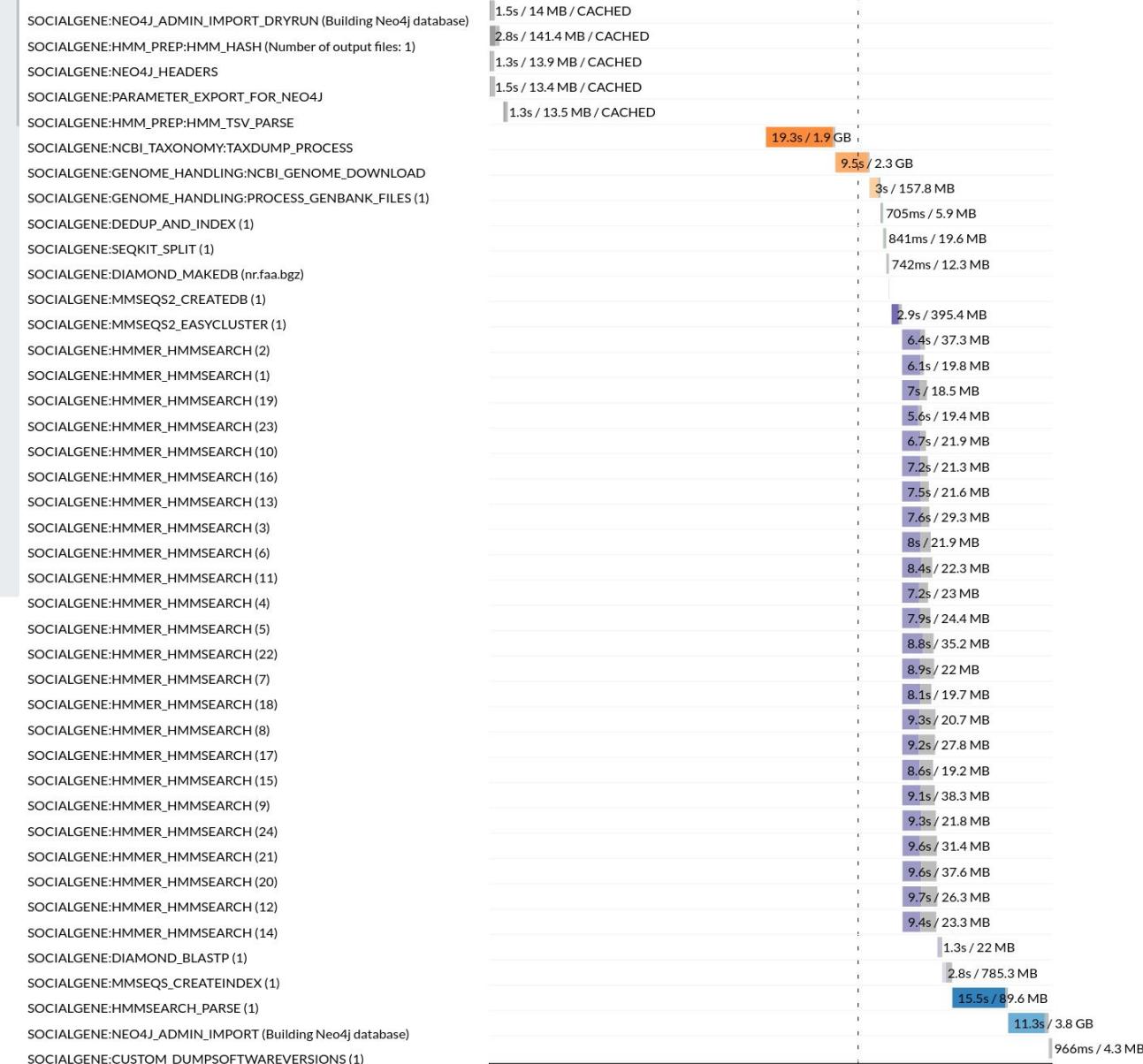


Processes execution timeline

Launch time: 28 Feb 2023 12:08

Elapsed time: 1m 23s

Legend: job wall time / memory usage (RAM)



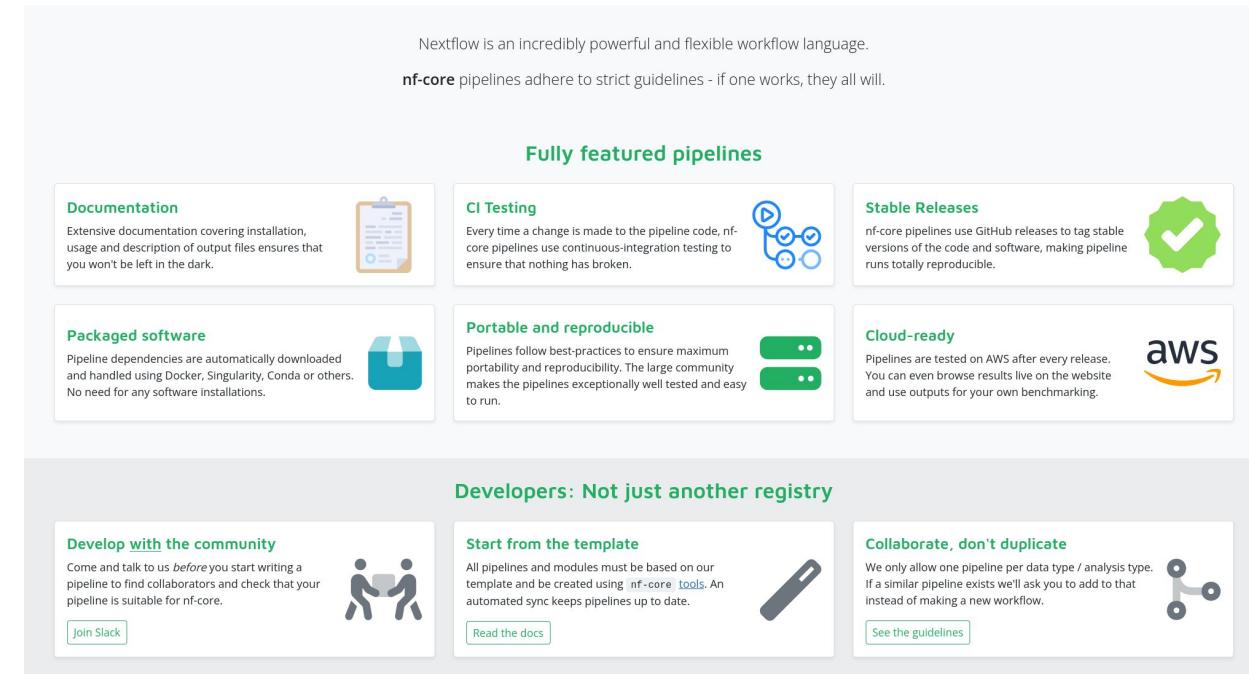
Standardizing the pipeline with **nf-core**



The screenshot shows the nf-core website homepage. At the top, there's a large "nf-core" logo with a green apple icon. Below it is a subtitle: "A community effort to collect a curated set of analysis pipelines built using Nextflow." A prominent "VIEW PIPELINES" button is centered. Below this, there are two search bars and a "VIEW PIPELINES" button. The main content area is divided into three sections: "For facilities", "For users", and "For developers". Each section has an icon and a brief description. At the bottom, there's a footer with the Nature Biotech logo and a citation: "nf-core is published in Nature Biotechnology Nat Biotechnol 38, 276–278 (2020) 8".

- For facilities**
Highly optimised pipelines with excellent reporting. Validated releases ensure reproducibility.
- For users**
Portable, documented and easy to use workflows. Pipelines that you can trust.
- For developers**
Companion templates and tools help to validate your code and simplify common tasks.

nf-core is published in Nature Biotechnology
Nat Biotechnol 38, 276–278 (2020) 8



This screenshot shows the "Fully featured pipelines" section of the nf-core website. It highlights several key features with icons and descriptions:

- Documentation**: Extensive documentation covering installation, usage and description of output files ensures that you won't be left in the dark.
- CI Testing**: Every time a change is made to the pipeline code, nf-core pipelines use continuous-integration testing to ensure that nothing has broken.
- Packaged software**: Pipeline dependencies are automatically downloaded and handled using Docker, Singularity, Conda or others. No need for any software installations.
- Portable and reproducible**: Pipelines follow best-practices to ensure maximum portability and reproducibility. The large community makes the pipelines exceptionally well tested and easy to run.
- Stable Releases**: nf-core pipelines use GitHub releases to tag stable versions of the code and software, making pipeline runs totally reproducible.
- Fully featured pipelines**
- Cloud-ready**: Pipelines are tested on AWS after every release. You can even browse results live on the website and use outputs for your own benchmarking.
- Developers: Not just another registry**
- Start from the template**: All pipelines and modules must be based on our template and be created using `nf-core tools`. An automated sync keeps pipelines up to date.
- Collaborate, don't duplicate**: We only allow one pipeline per data type / analysis type. If a similar pipeline exists we'll ask you to add to that instead of making a new workflow.

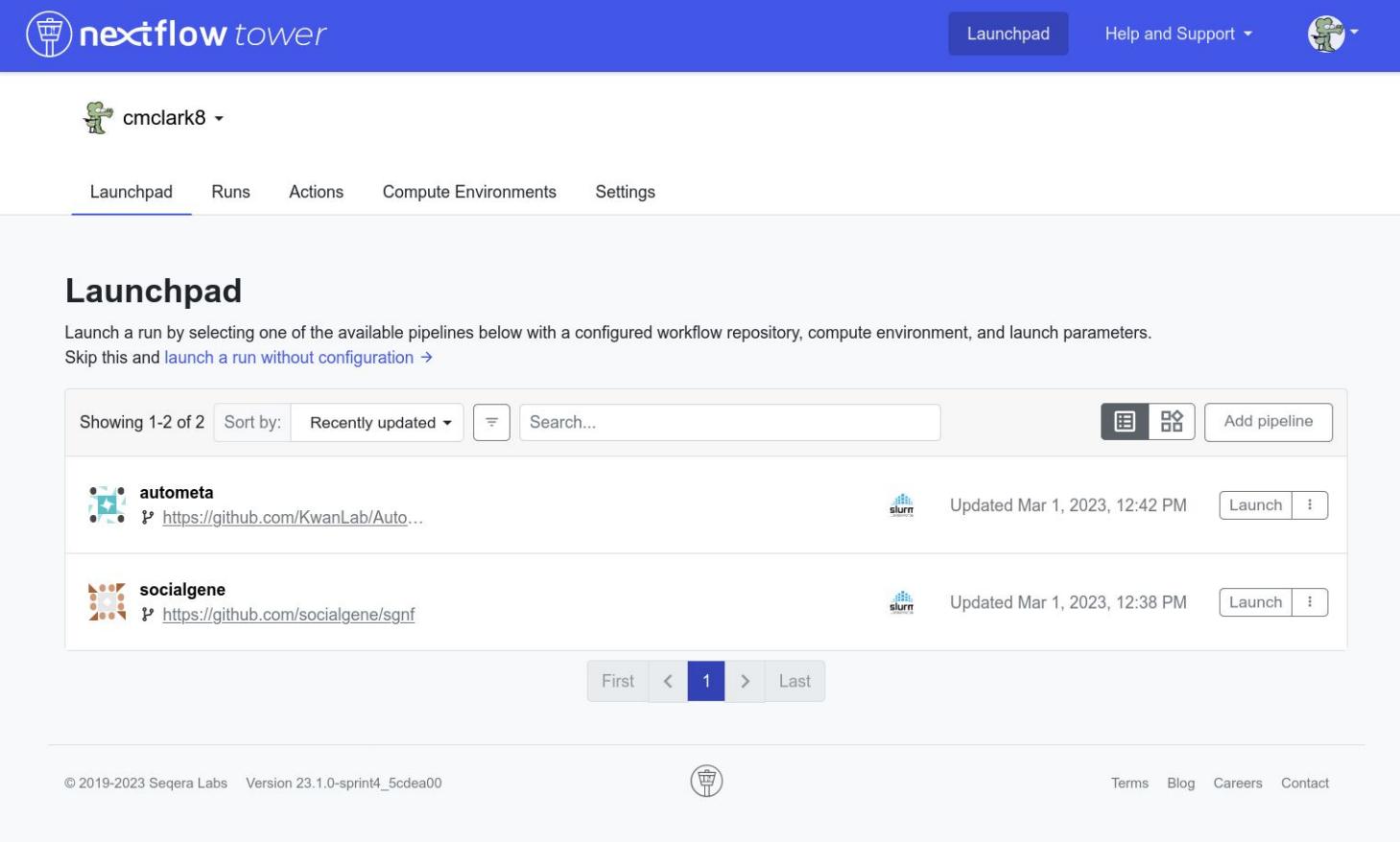
Making the pipeline more accessible with nf-core



```
o (sg) chase@titan:~/Documents/github/kwan_lab/socialgene/sgnf$ nf-core launch .  
  
NF-CORE  
nf-core/tools version 2.7.2 - https://nf-co.re  
  
INFO NOTE: This tool ignores any pipeline parameter defaults overwritten by Nextflow config files or profiles  
INFO [x] Default parameters match schema validation  
INFO [x] Pipeline schema looks valid (found 75 params)  
INFO Would you like to enter pipeline parameters using a web-based interface or a command-line wizard?  
? Choose launch method (Use arrow keys)  
  > Web based  
    Command line
```

A screenshot of the nf-core Launch pipeline web interface. The page is titled 'Launch pipeline' and shows a configuration interface with various input fields and sections. A large gray downward arrow is positioned above the interface, pointing towards it. The interface includes sections for 'Nextflow command-line flags', 'Input Genomes', 'Required Output Directories', 'Input IMeRs', 'SocialGene parameters', 'SocialGene modules', 'Max job resource request options', and 'Unwrapped parameters'. At the bottom, there is a 'Launch' button. The top navigation bar includes links for Home, Pipelines, Modules, Tools, Docs, Events, About, and a 'Join nf-core' button.

Making the pipeline more accessible with **nexiflow tower**

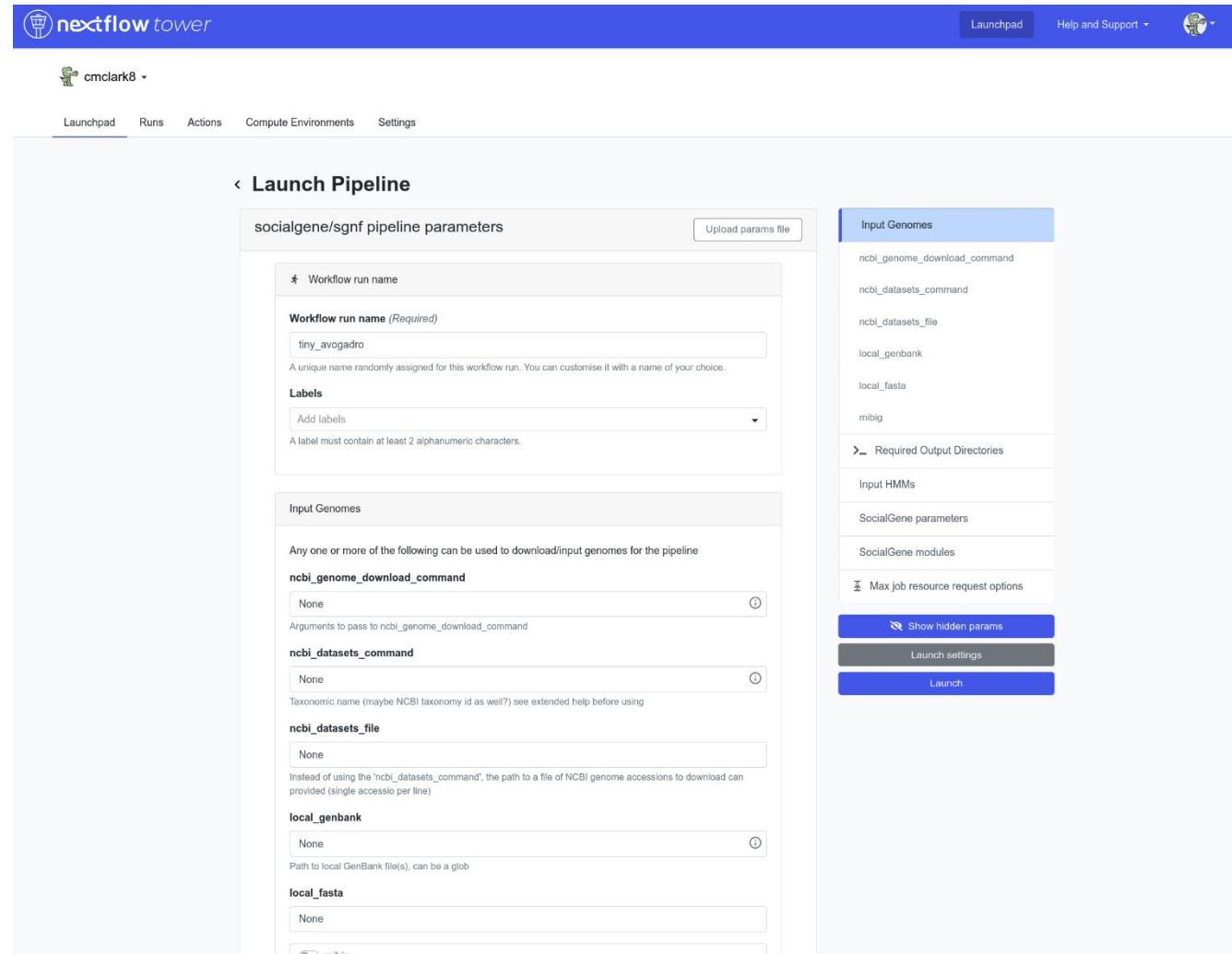


The screenshot shows the nexiflow tower web interface. At the top, there's a blue header bar with the logo on the left, "Launchpad" and "Help and Support" buttons in the center, and a user icon on the right. Below the header, the user "cmclark8" is logged in. A navigation bar with tabs for "Launchpad", "Runs", "Actions", "Compute Environments", and "Settings" is visible. The main content area is titled "Launchpad" and contains instructions: "Launch a run by selecting one of the available pipelines below with a configured workflow repository, compute environment, and launch parameters. Skip this and [launch a run without configuration](#)". There are two pipeline entries:

Pipeline Name	Repository URL	Last Updated	Actions
autometa	https://github.com/KwanLab/Auto...	Updated Mar 1, 2023, 12:42 PM	Launch ⋮
socialgene	https://github.com/socialgene/sgnf	Updated Mar 1, 2023, 12:38 PM	Launch ⋮

At the bottom, there are links for "First", "Last", and page number "1". The footer includes copyright information ("© 2019-2023 Seqera Labs Version 23.1.0-sprint4_5cdea00"), a logo, and links for "Terms", "Blog", "Careers", and "Contact".

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



The screenshot shows the Nextflow Tower interface for launching a pipeline. The top navigation bar includes the Nextflow Tower logo, user profile (cmclark8), Launchpad, Help and Support, and a search bar.

The main area is titled "Launch Pipeline" and shows the "socialgene/sgnf pipeline parameters".

Workflow run name: tiny_avogadro

Labels: Add labels

Input Genomes:

- ncbi_genome_download_command:** None
- ncbi_datasets_command:** None
- ncbi_datasets_file:** None
- local_genbank:** None
- local_fasta:** mihin

Input HMMs: None

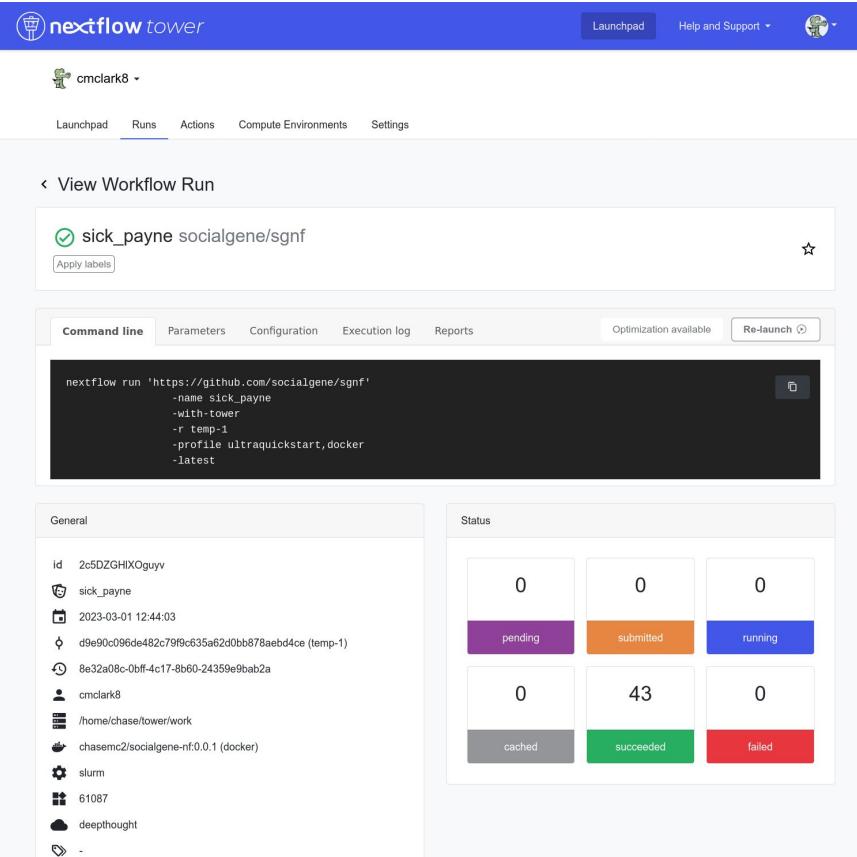
SocialGene parameters: None

SocialGene modules: None

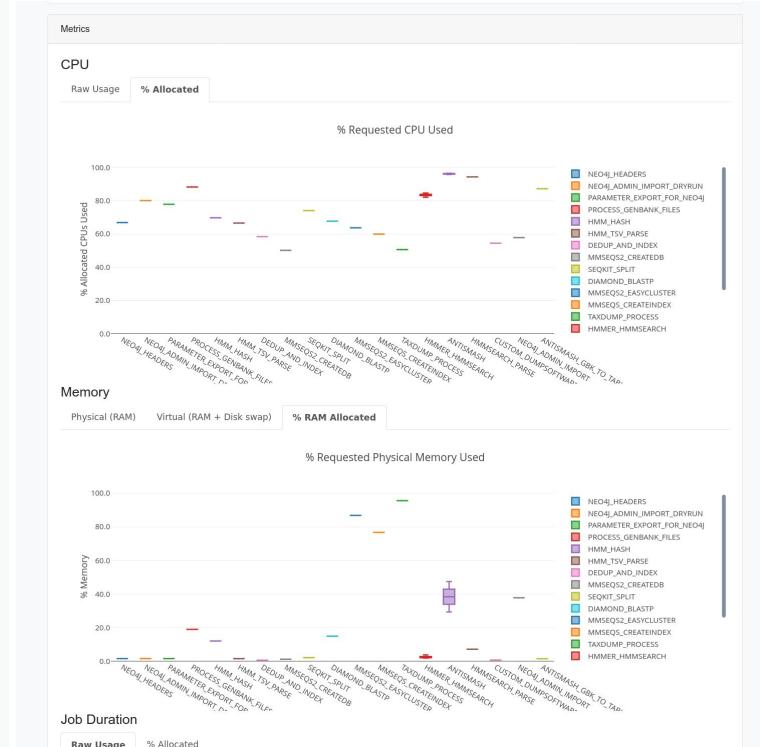
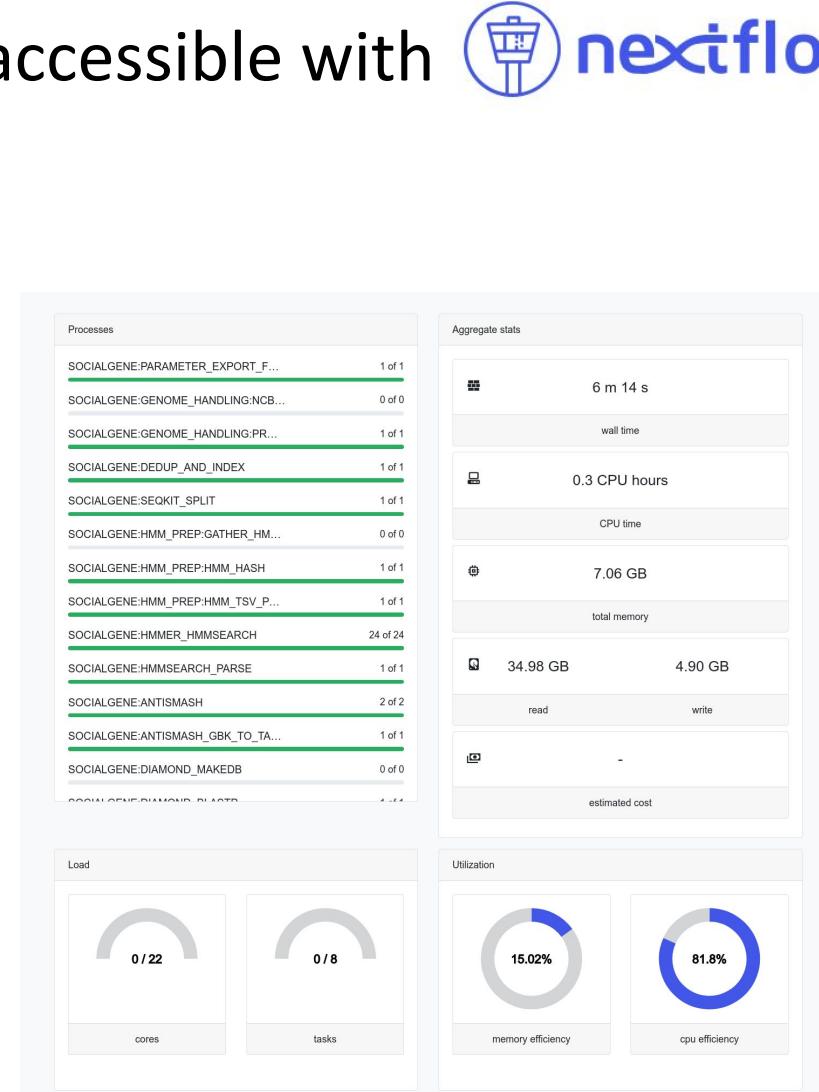
Max job resource request options: None

Buttons at the bottom right include "Show hidden params", "Launch settings", and "Launch".

Making the workflow more accessible with nextflow tower



The screenshot shows the Nextflow Tower interface. At the top, there's a navigation bar with 'Launchpad', 'Help and Support', and a user icon. Below the navigation is a header for the 'View Workflow Run' section, showing the workflow name 'sick_payne' and its source 'socialgene/sgnf'. A 'Command line' tab displays the command used to run the workflow: `nextflow run 'https://github.com/socialgene/sgnf'`. This command includes parameters like '-name sick_payne', '-with-tower', '-r temp-1', '-profile ultraquickstart_docker', and '-latest'. Below the command line are two status panels: 'General' and 'Status'. The 'General' panel lists various system details such as ID, workflow name, date, and user. The 'Status' panel shows counts for pending, submitted, running, cached, succeeded, and failed tasks.



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...

AKSGPWTFKDDRGT...

Protein have domains

MSNARATHL_TRRGI...

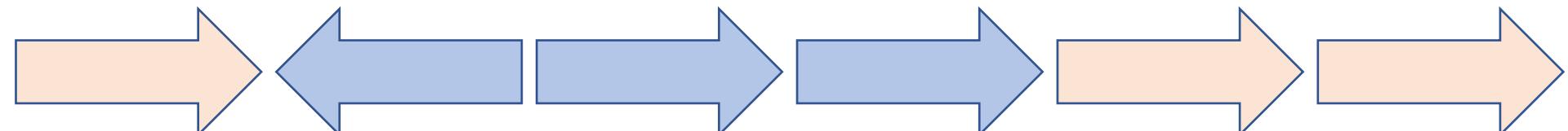
KAD_VVQAGDMDVSK...

AKSGPWT_FKDDRGT...

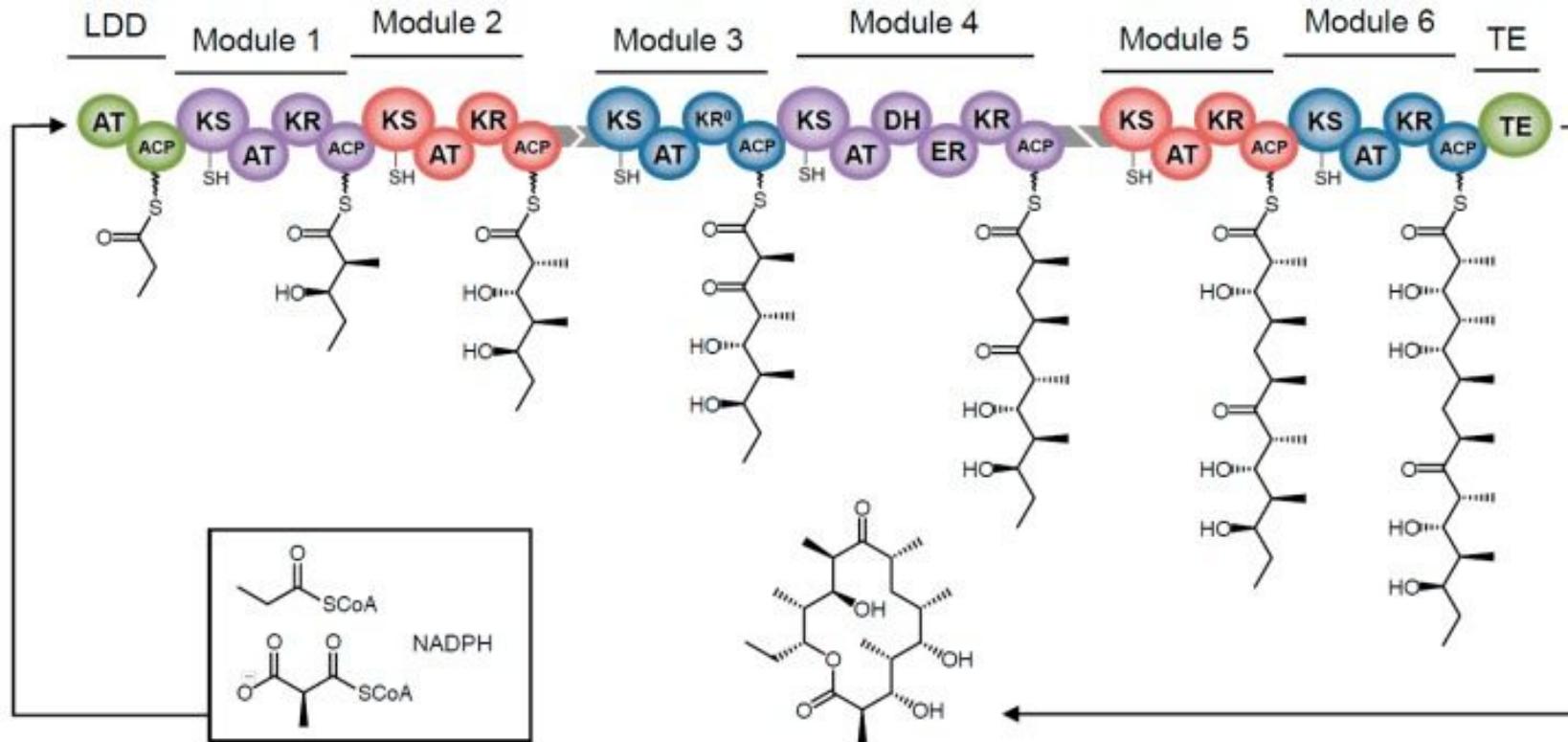
Genes encode proteins



Gene Cluster



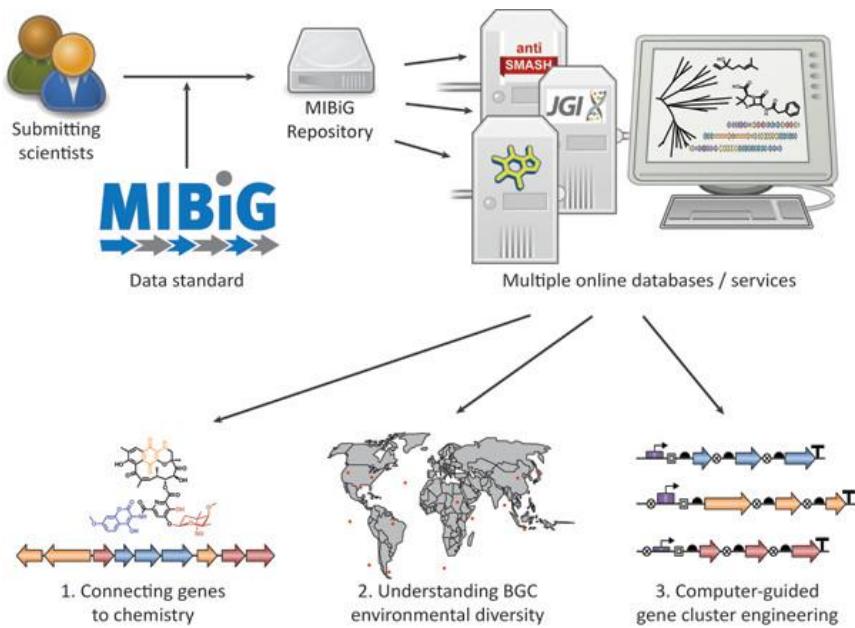
Polyketide synthases are large, multidomain proteins “Beads on a String”



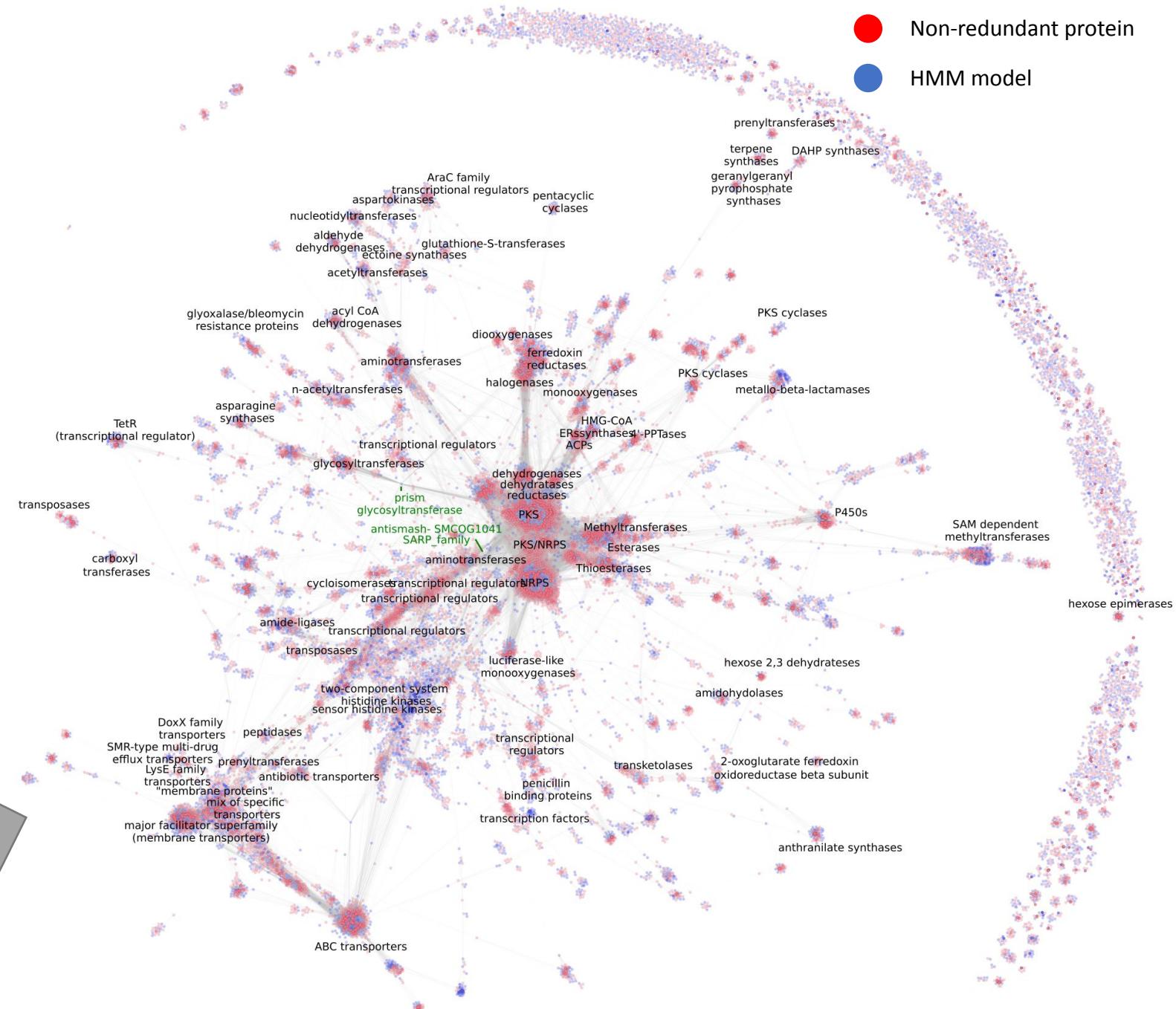
Bayly, Carmen L, and Vikramaditya G Yadav. "Towards Precision Engineering of Canonical Polyketide Synthase Domains: Recent Advances and Future Prospects." *Molecules* (Basel, Switzerland) vol. 22,2 235. 5 Feb. 2017, doi:10.3390/molecules22020235

Vignette 1

Mapping MIBiG

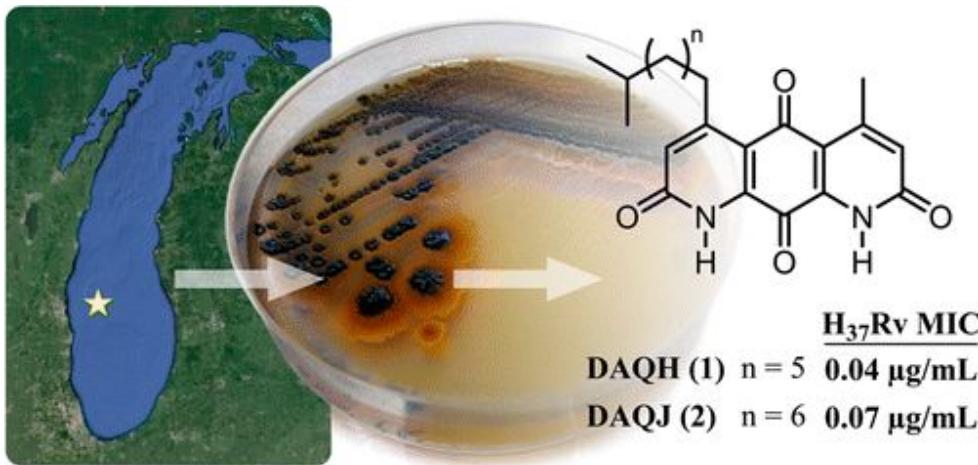


SocialGene



Vignette 2

Looking for diazaquinomycin analogs



Mullowney, Michael W et al. "Diaza-anthracene Antibiotics from a Freshwater-Derived Actinomycete with Selective Antibacterial Activity toward *Mycobacterium tuberculosis*." *ACS infectious diseases* vol. 1,4 (2015): 168-174. doi:10.1021/acsinfecdis.5b00005

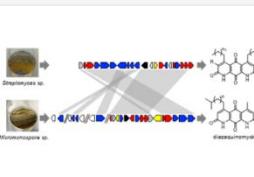
RETURN TO ISSUE | < PREV ARTICLE NEXT >
Diazaquinomycin Biosynthetic Gene Clusters from Marine and Freshwater Actinomycetes
 Jana Braesel, Jung-Ho Lee, Benoit Arnould, Brian T. Murphy, and Alessandra S. Eustáquio*

Cite this: *J. Nat. Prod.* 2019, 82, 4, 937–946
 Publication Date: March 21, 2019
<https://doi.org/10.1016/j.jnp.2019.01.028>
 Copyright © 2019 American Chemical Society and
 American Society of Pharmacogenomics
 RIGHTS & PERMISSIONS

Read Online | PDF (2 MB) | Supporting Info (1) | SUBJECTS: Peptides and proteins, Bacteria, Genetics, ▾

Abstract

Tuberculosis is an infectious disease of global concern. Members of the diazaquinomycin (DAQ) class of natural products have shown potent and selective activity against drug-resistant *Mycobacterium tuberculosis*. However, poor solubility has prevented further development of this compound class. Understanding DAQ biosynthesis may provide a viable route for the generation of derivatives with improved properties. We have sequenced the genomes of two actinomycete bacteria that produce distinct DAQ derivatives. While software tools for automated biosynthetic gene cluster (BGC) prediction failed to detect DAQ BGCs, comparative genomics using MAUVE alignment led to the identification of putative BGCs in the marine *Streptomyces* sp. F001 and in the freshwater *Micromonospora* sp. B006. Deletion of the identified daq BGC in strain B006 using CRISPR-Cas9 genome editing abolished DAQ production, providing experimental evidence for BGC assignment. A complete model for DAQ biosynthesis is proposed based on the genes identified. Insufficient knowledge of natural product biosynthesis is one of the major challenges of productive genome mining approaches. The results reported here fill a gap in knowledge regarding the genetic basis for the biosynthesis of DAQ antibiotics. Moreover, identification of the daq BGC shall enable future generations of improved derivatives using biosynthetic methods.



MIBIG Repository of Known Biosynthetic Gene Clusters

BGC0001848: diazaquinomycin H biosynthetic gene cluster from *Micromonospora* sp. B006

Location: 3,038,865 - 3,097,062 nt. (total: 58,198 nt). This entry is originally from NCBI GenBank CP030865.1.

Download Cluster GenBank file | View antiSMASH-generated output

Select a gene to view the details available for it

Gene details

Legend:

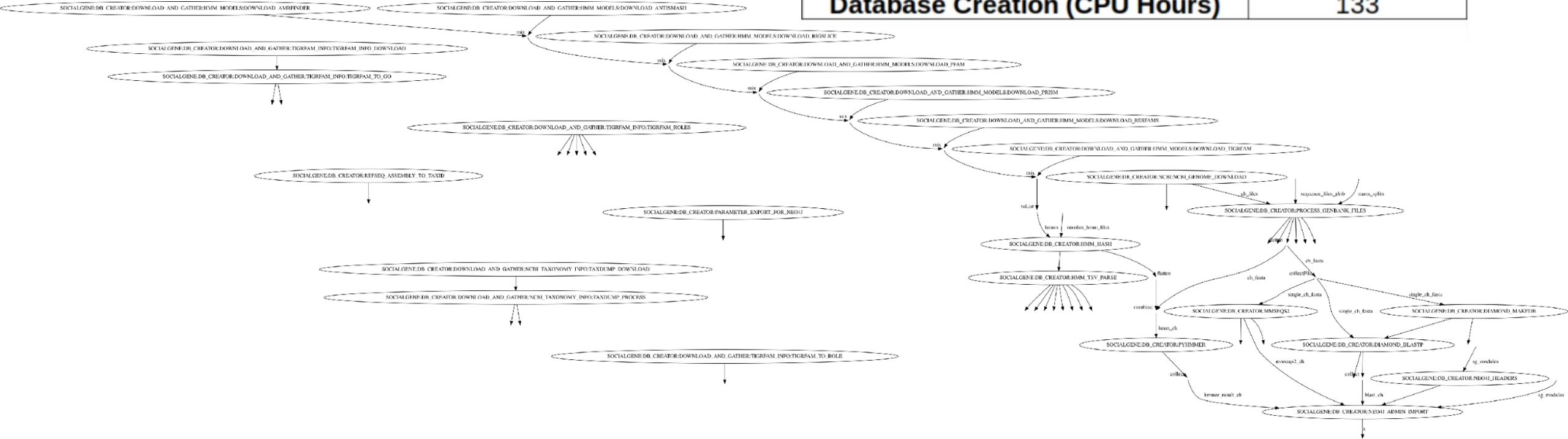
- core biosynthetic genes
- additional biosynthetic genes
- transport-related genes
- regulatory genes
- resistance genes
- other genes

reset view | **zoom to selection**

General | Compounds | Genes | History | KnownClusterBlast | NRPS/PKS domains | List of genes involved in compound(s) production

Identifiers	Position	Product	Functions	Evidence	Extra
MicB006_2892 AX035173.1	3038865 - 3040082 (-)	2-keto-3-deoxy-D-arabino-heptulosonate-7-phosphate synthase II	Precursor biosynthesis	Sequence-based prediction	copy AA seq copy Nt seq
MicB006_2893 AX035174.1	3040216 - 3041814 (-)	hypothetical protein			copy AA seq copy Nt seq
MicB006_2894 AX035175.1	3041874 - 3043553 (-)	SSS sodium solute transporter superfamily			copy AA seq copy Nt seq
MicB006_2895 AX035176.1	3043550 - 3043885 (-)	putative membrane protein			copy AA seq copy Nt seq
MicB006_2896 AX035177.1	3044080 - 3044820 (+)	hypothetical protein			copy AA seq copy Nt seq
MicB006_2897 AX035178.1	3044813 - 3045514 (-)	short-chain dehydrogenase/reductase SDR			copy AA seq copy Nt seq
MicB006_2898 AX035179.1	3045526 - 3046200 (-)	isochorismatase	Precursor biosynthesis	Sequence-based prediction	copy AA seq copy Nt seq
MicB006_2899 AX035180.1	3046233 - 3047018 (-)	2,3-dihydro-2,3-dihydroxybenzoate dehydrogenase	Precursor biosynthesis	Sequence-based prediction	copy AA seq copy Nt seq
MicB006_2900 AX035181.1	3047018 - 3048187 (-)	putative n-hydroxybenzoate hydroxylase			copy AA seq copy Nt seq
MicB006_2901 AX035182.1	3048219 - 3049535 (-)	phenylacetate-coenzyme A ligase			copy AA seq copy Nt seq
MicB006_2902 AX035183.1	3049532 - 3050293 (-)	hypothetical protein			copy AA seq copy Nt seq
MicB006_2903 AX035184.1	3050319 - 3051422 (-)	37-dideoxy-D-threo-hept-2,6-diulosonate synthase			copy AA seq copy Nt seq
MicB006_2904 AX035185.1	3051454 - 3052278 (-)	2-amino-37-dideoxy-D-threo-hept-6-ulosonate synthase			copy AA seq copy Nt seq
MicB006_2905 AX035186.1	3052294 - 3053577 (-)	aspartokinase			copy AA seq copy Nt seq
MicB006_2906 AX035187.1	3053813 - 3055777 (+)	anthranilate synthase	Precursor biosynthesis	Sequence-based prediction	copy AA seq copy Nt seq
MicB006_2907 AX035188.1	3055937 - 3056335 (+)	hypothetical protein			copy AA seq copy Nt seq
MicB006_2908 AX035189.1	3056436 - 3056825 (-)	glyoxalase/bleomycin resistance protein/dioxygenase			copy AA seq copy Nt seq
MicB006_2909 AX035190.1	3056836 - 3058797 (-)	m multicopper oxidase			copy AA seq copy Nt seq
MicB006_2910 AX035191.1	3058818 - 3059273 (-)	hypothetical protein			copy AA seq copy Nt seq
MicB006_2911 AX035192.1	3059338 - 3059733 (-)	hypothetical protein			copy AA seq copy Nt seq
MicB006_2912 AX035193.1	3059926 - 3060588 (-)	DNA-binding response regulator LuxR family			copy AA seq copy Nt seq

SocialGene



	<i>Micromonospora</i>
Genomes	226
Contigs/Scaffolds	34,856
Proteins	1,373,641
Non-Redundant Proteins	1,076,690
HMM Annotations	10,656,754
MMSEQS2	1,009,056
Bidirectional BLASTP (45.3 GB)	420,788,981
Database Creation (Real Duration)	4h 30m
Database Creation (CPU Hours)	133

dj socialgeneweb +

127.0.0.1:8009/findmybgc/ ☆ ☰

socialgeneweb Home About About Single Protein Search BGC Search My Profile Sign Out DJDT

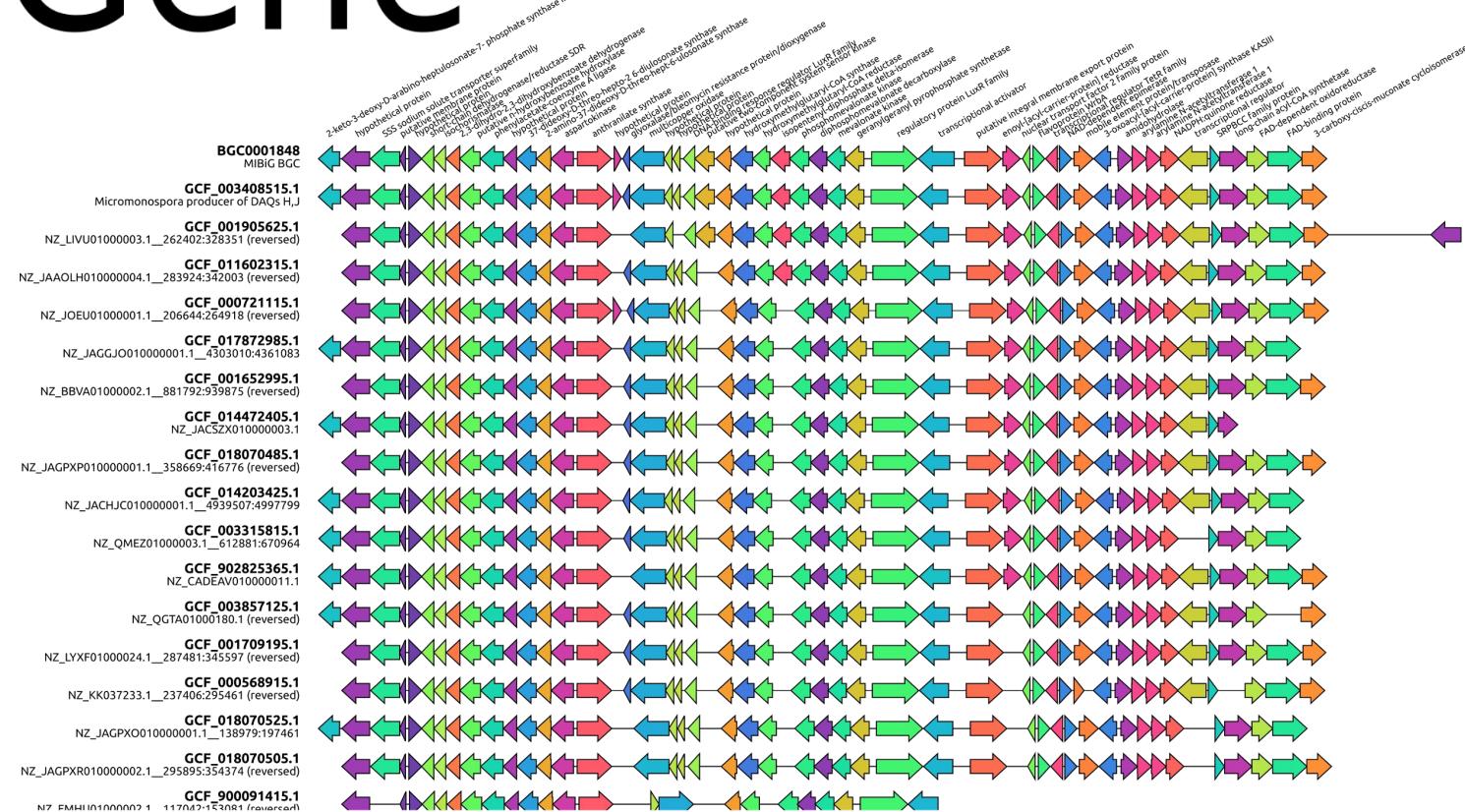
Search the database, using a single BGC as input

(Input expects a gbk or gbff file)

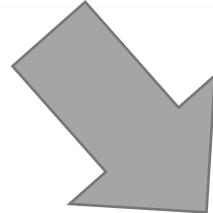
Description:

Document: No file selected.

SocialGene



<i>Micromonospora</i>	
Genomes	226
Contigs/Scaffolds	34,856
Proteins	1,373,641
Non-Redundant Proteins	1,076,690
HMM Annotations	10,656,754
MMSEQS2	1,009,056
Bidirectional BLASTP (45.3 GB)	420,788,981
Database Creation (Real Duration)	4h 30m
Database Creation (CPU Hours)	133



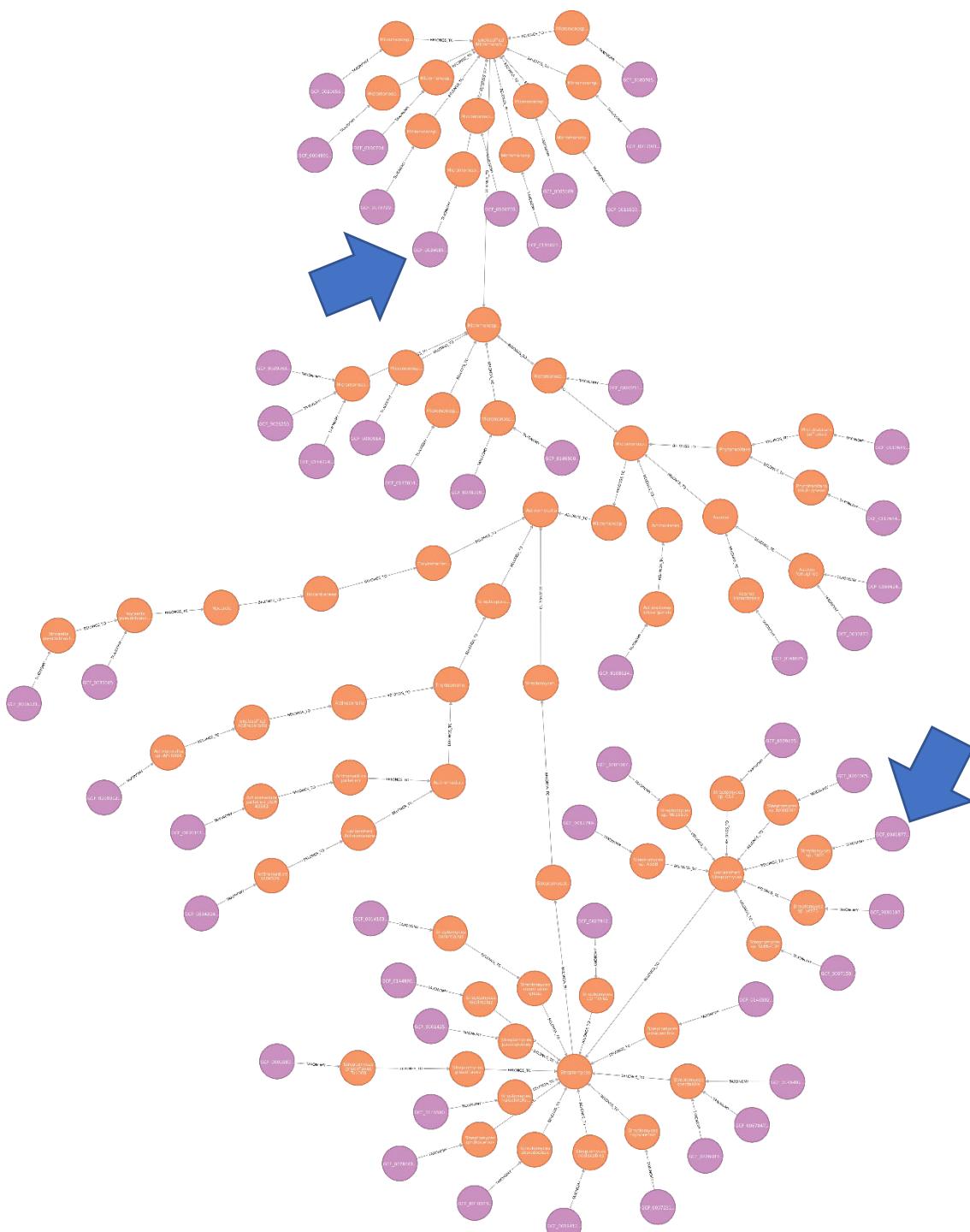
SocialGene



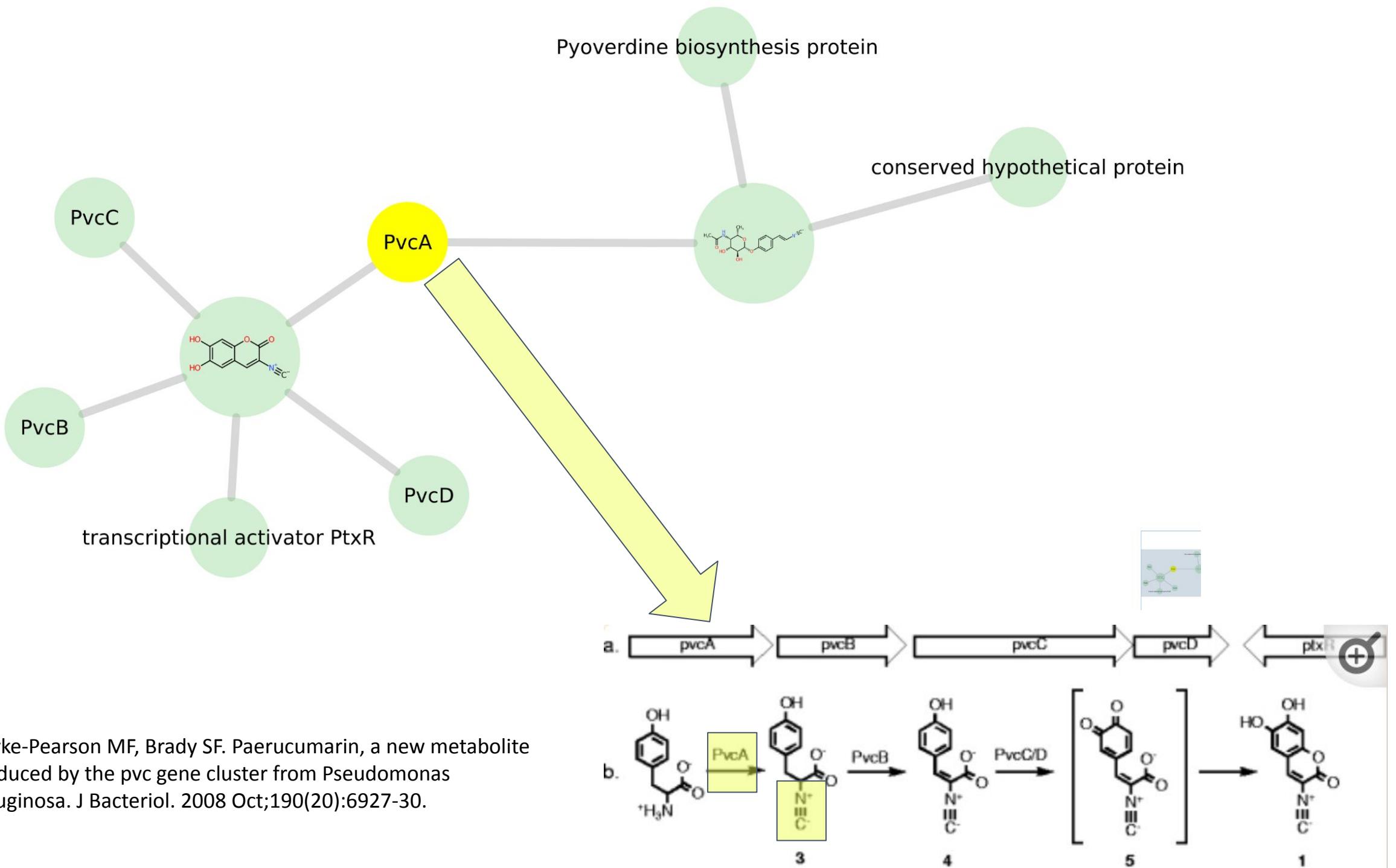
Open Science Grid

HT CENTER FOR
HIGH THROUGHPUT
COMPUTING

	RefSeq
Genomes	266,668
Contigs/Scaffolds	23,941,594
Proteins	188,429,555
HMM models	25,648
HMM annotations	1,403,423,051
MMseqs2 (bug in MMseqs2)	188,327,165
Contigs to Proteins	971,298,319
Species	49,902
Genera	6,460



Teasers- What Next?



Clarke-Pearson MF, Brady SF. Paerucumarin, a new metabolite produced by the pvc gene cluster from *Pseudomonas aeruginosa*. J Bacteriol. 2008 Oct;190(20):6927-30.

BGC0001848_diazaquinomycin_Micromonospora_B006
BGC0001848 (reversed):58198-1

GCF_016862515_NZ_BONC01000022.region001.Asanoa_iriomotensis
NZ_BONC01000022:1-57188

GCF_016862495_NZ_BONB01000010.region001.Asanoa_ferruginea
NZ_BONB01000010:1-59328

GCF_003387075_NZ_QUMQ01000001.region012.Asanoa_ferruginea
NZ_QUMQ01000001:1-58632

GCF_011764565_NZ_AP022871.1.region028.Phytohabitans_suffuscus
NZ_AP022871:1-85472

GCF_011764425_NZ_BLPF01000001.region004.Phytohabitans_houttuyniae
NZ_BLPF01000001 (reversed):73880-1



BLAST® » blastp suite » results for RID-EU52MYW501R

Home Recent Results Saved Strategies Help

◀ Edit Search Save Search Search Summary ▶

How to read this report? BLAST Help Videos Back to Traditional Results Page

Job Title refWP_203703487.1|

RID EU52MYW501R Search expires on 08-06 21:00 pm Download All ▾

Program Quick BLASTP Citation ▾

Database nr See details ▾

Query ID WP_203703487.1

Description hypothetical protein [Asanoa iriomotensis]

Molecule type amino acid

Query Length 509

Other reports Distance tree of results Multiple alignment MSA viewer ?

Filter Results

Organism only top 20 will appear

Type common name, binomial, taxid or group name

+ Add organism

Percent Identity E value Query Coverage
to to to to

Filter Reset

Compare these results against the new Clustered nr database ? BLAST

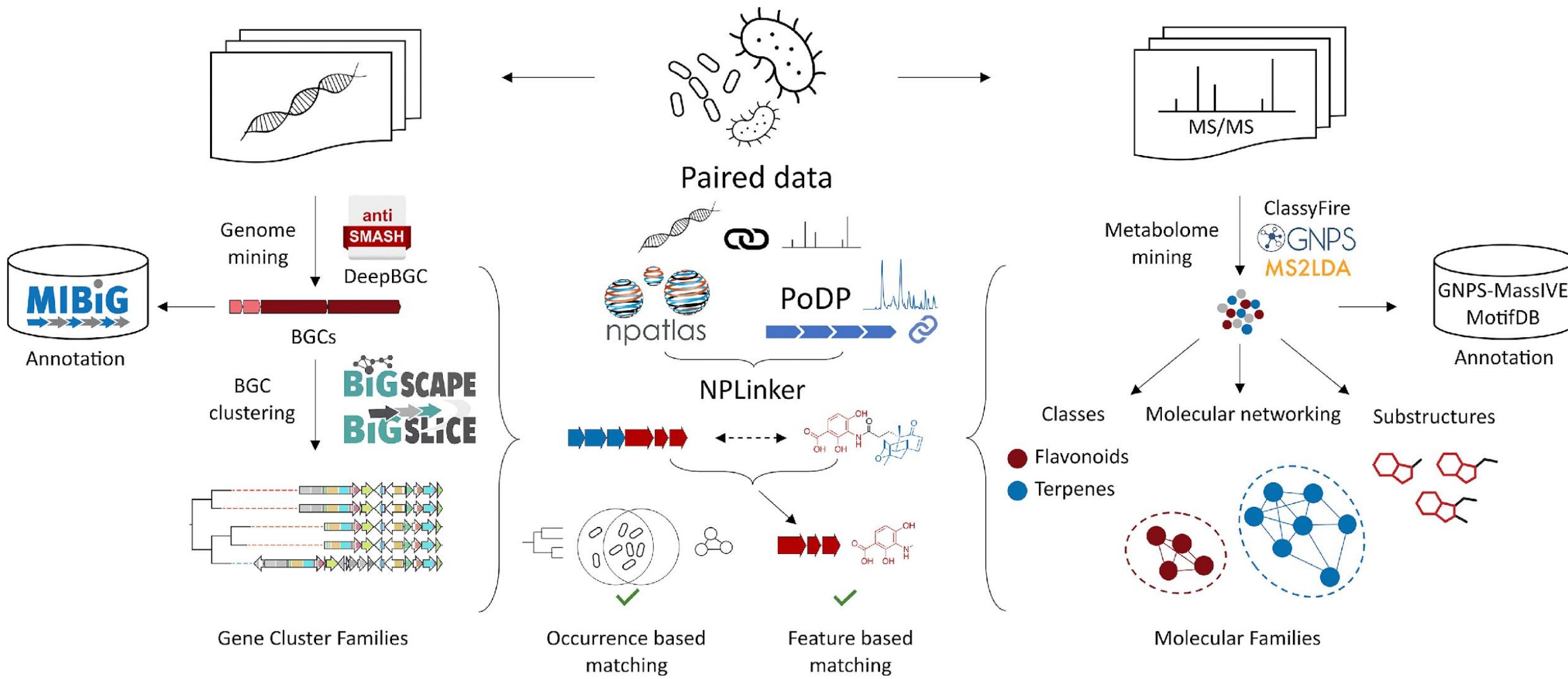
Descriptions Graphic Summary Alignments Taxonomy

Sequences producing significant alignments

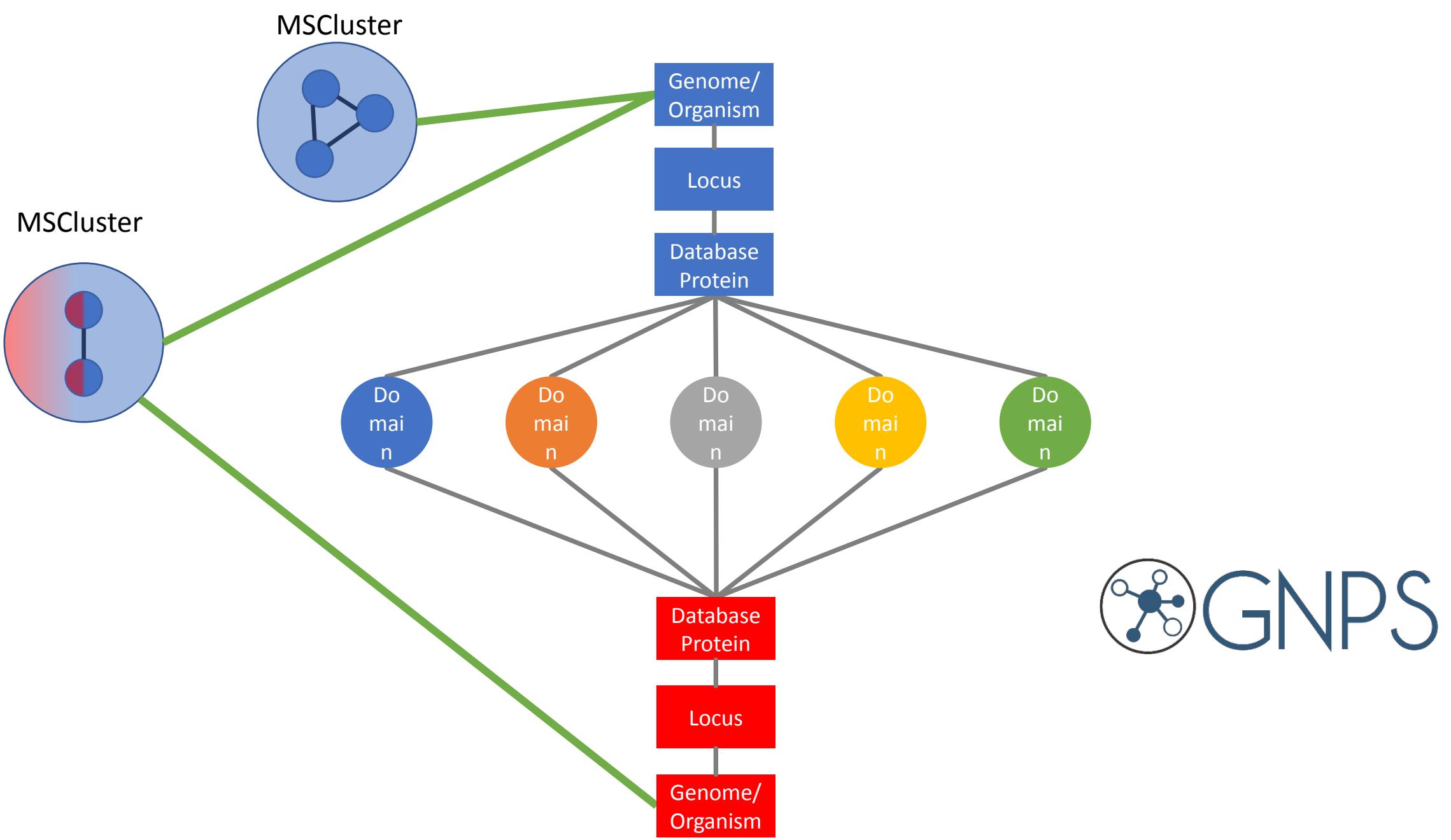
Download Select columns Show 100 ▾ ?

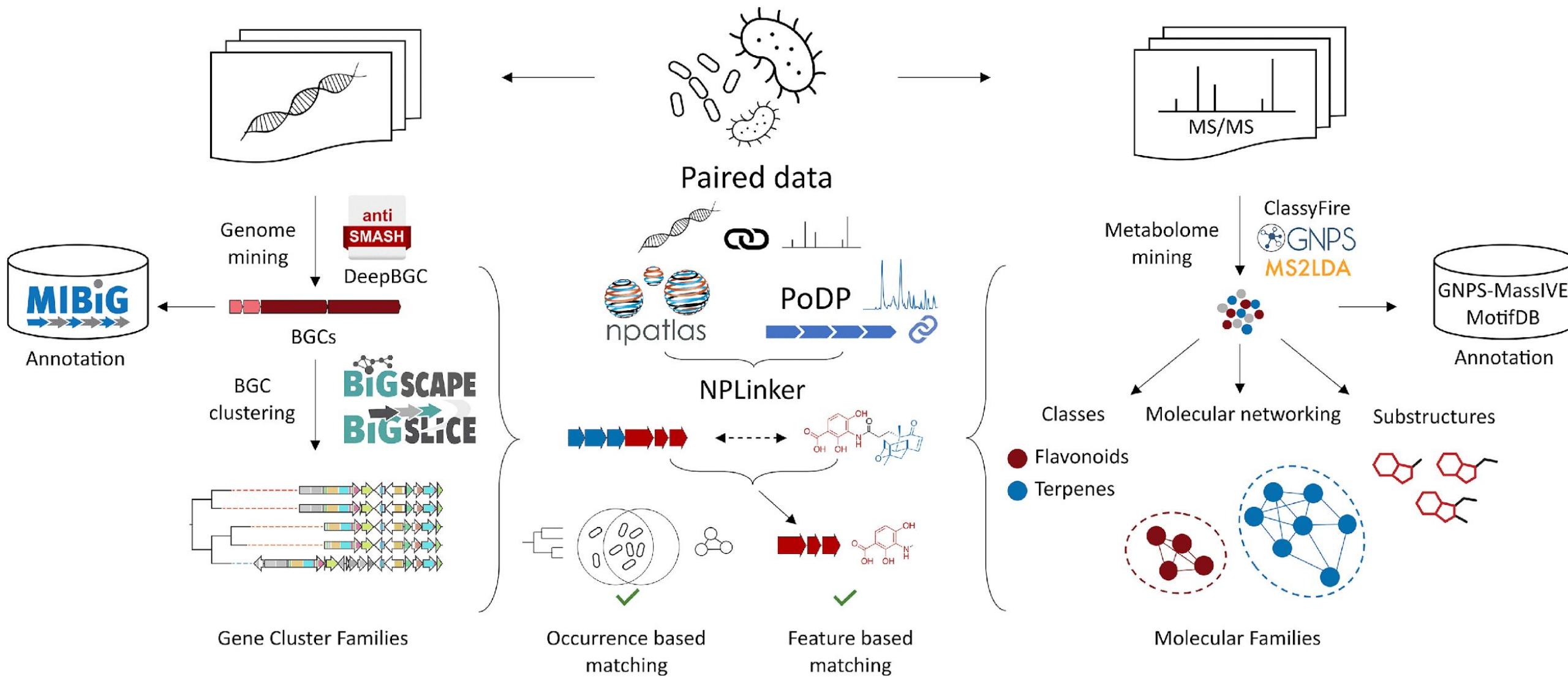
<input checked="" type="checkbox"/> select all 88 sequences selected	GenPept	Graphics	Distance tree of results	Multiple alignment	MSA Viewer				
	Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident.	Acc. Len	Accession
<input checked="" type="checkbox"/>	hypothetical protein [Asanoa iriomotensis]	Asanoa iriomotensis	926	926	100%	0.0	100.0%	509	WP_203703487.1
<input checked="" type="checkbox"/>	vanadium-dependent haloperoxidase [Asanoa ferruginea]	Asanoa ferruginea	831	831	90%	0.0	89.80%	518	WP_203783445.1
<input checked="" type="checkbox"/>	hypothetical protein [Phytohabitans houttuyniae]	Phytohabitans houttu... yiae	654	654	90%	0.0	73.81%	518	WP_218578634.1
<input checked="" type="checkbox"/>	hypothetical protein Phou_003280 [Phytohabitans houttuyniae]	Phytohabitans houttu... yiae	652	652	90%	0.0	73.81%	499	GFJ76148.1
<input checked="" type="checkbox"/>	hypothetical protein C1193_02165 [Micromonospora endophytica (Xie et al. 2001) Li et al. 2019]	Micromonospora end... phytica	643	643	91%	0.0	70.66%	506	PZG00572.1
<input checked="" type="checkbox"/>	hypothetical protein [Streptomyces sp. SID5910]	Streptomyces sp. SID5910	642	642	91%	0.0	70.66%	495	WP_212808405.1
<input checked="" type="checkbox"/>	hypothetical protein [Streptomyces sp. CMB-SIM0423]	Streptomyces sp. CMB-SIM0423	640	640	90%	0.0	69.70%	500	MYR44148.1
<input checked="" type="checkbox"/>	vanadium-dependent haloperoxidase [Streptomyces sp. SBT349]	Streptomyces sp. SBT349	618	618	90%	0.0	67.17%	503	WP_101425741.1
<input checked="" type="checkbox"/>	vanadium-dependent haloperoxidase [Actinomadura alba]	Actinomadura alba	601	601	90%	0.0	63.85%	516	WP_187245424.1

Vanadium-dependent Haloperoxidases



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





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



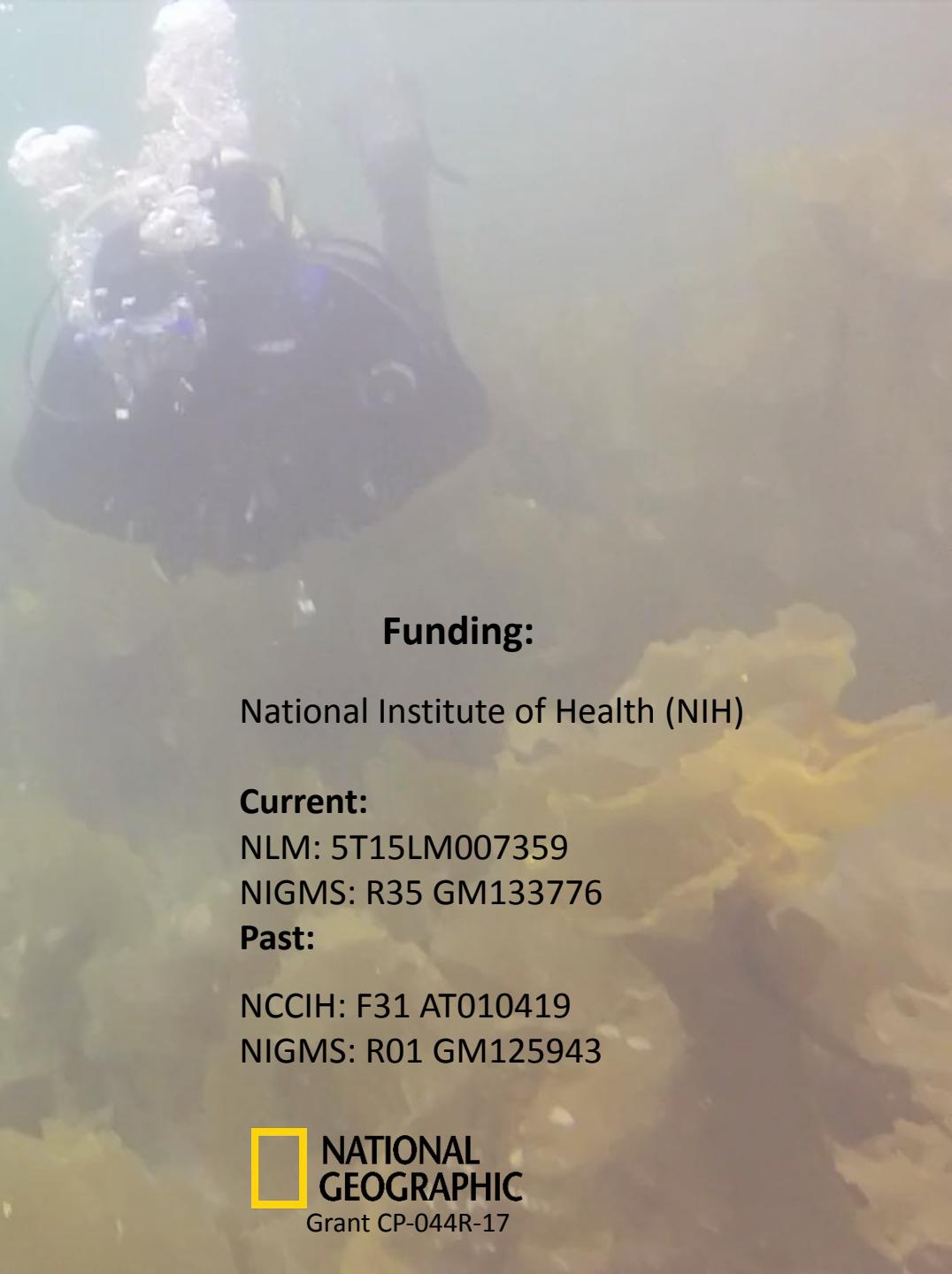
WISCONSIN
UNIVERSITY OF WISCONSIN-MADISON



**PHARMACEUTICAL
SCIENCES
COLLEGE
OF PHARMACY**



**Center for
Biomolecular Sciences**



Funding:

National Institute of Health (NIH)

Current:

NLM: 5T15LM007359

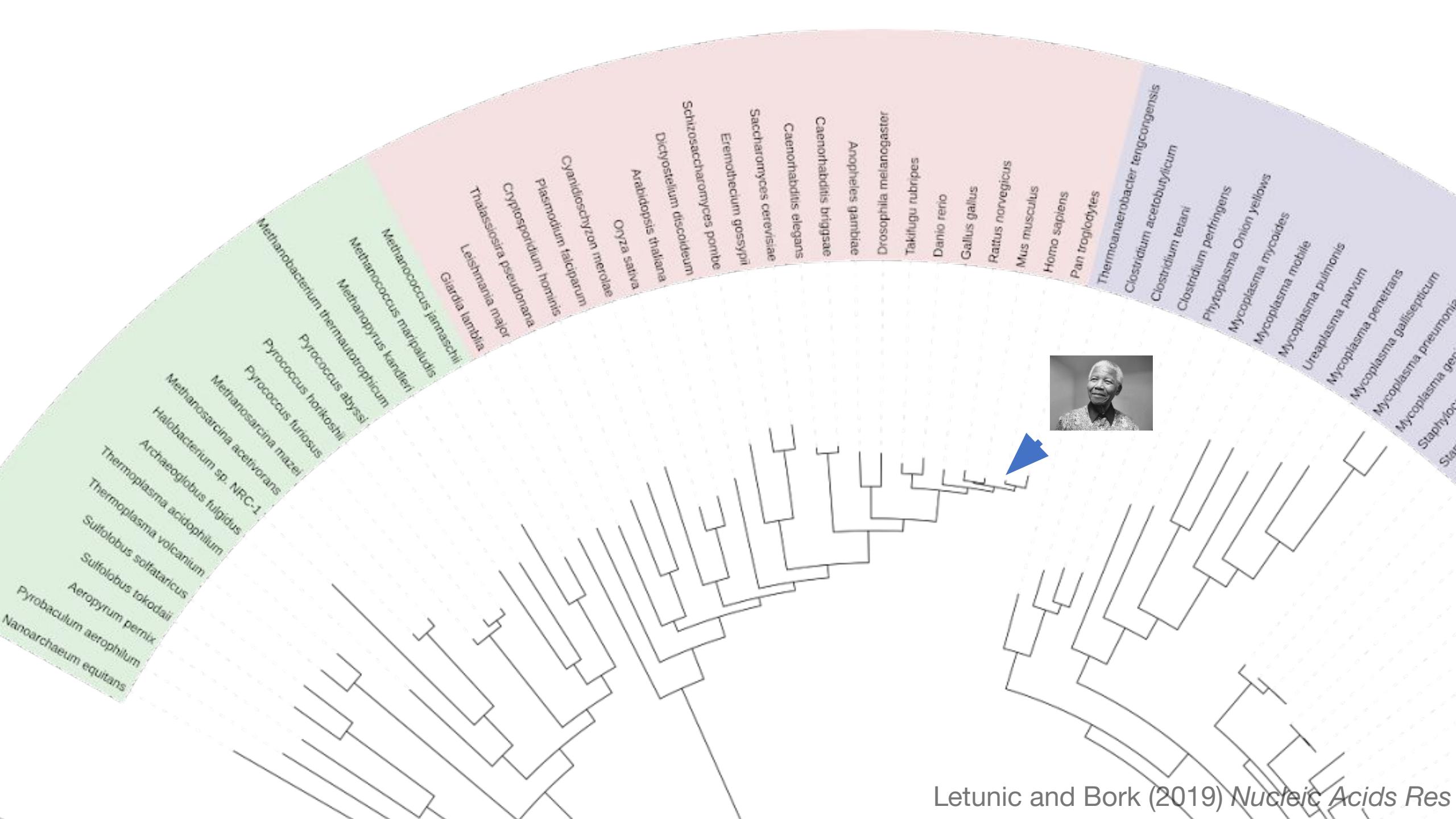
NIGMS: R35 GM133776

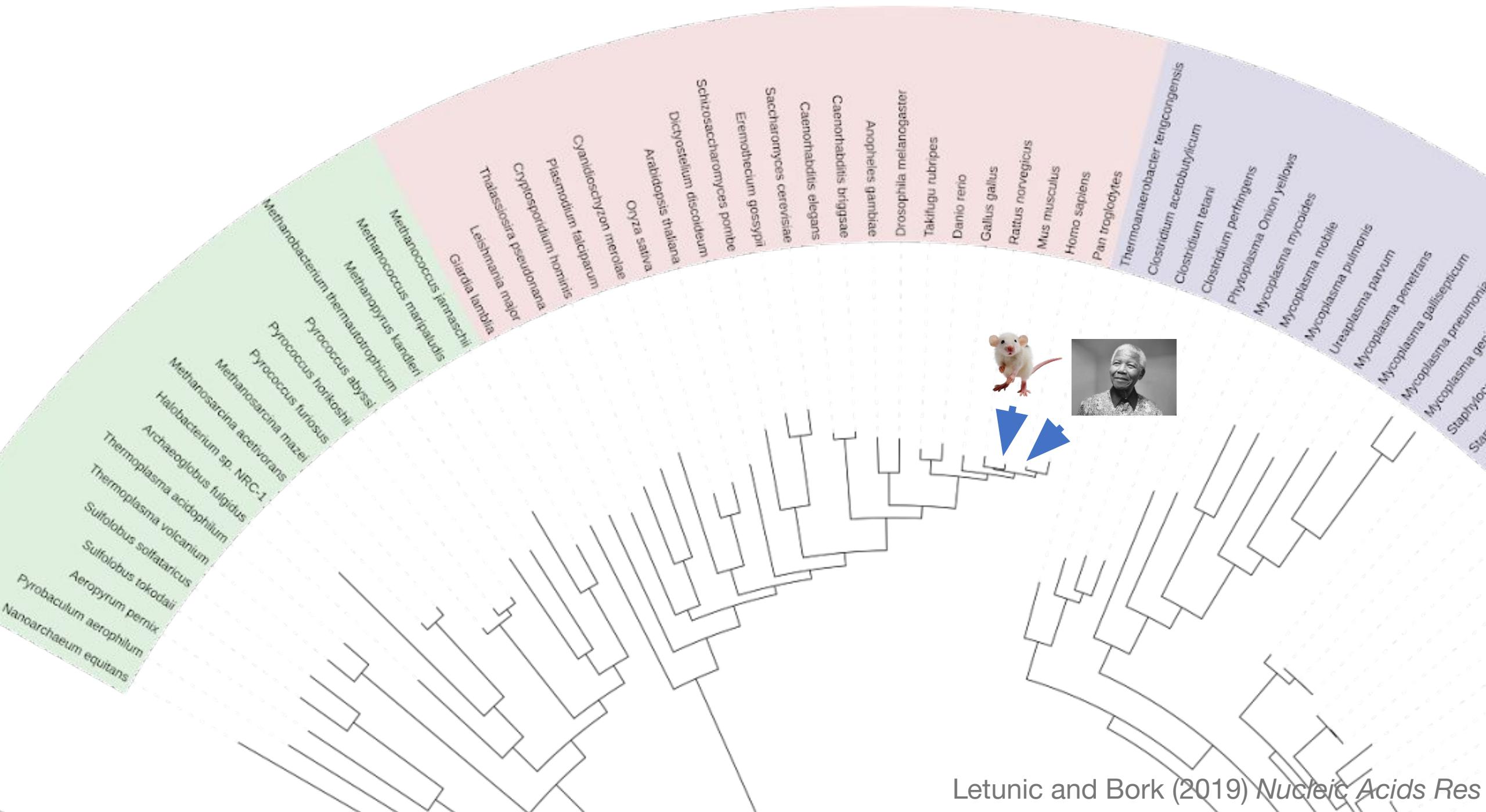
Past:

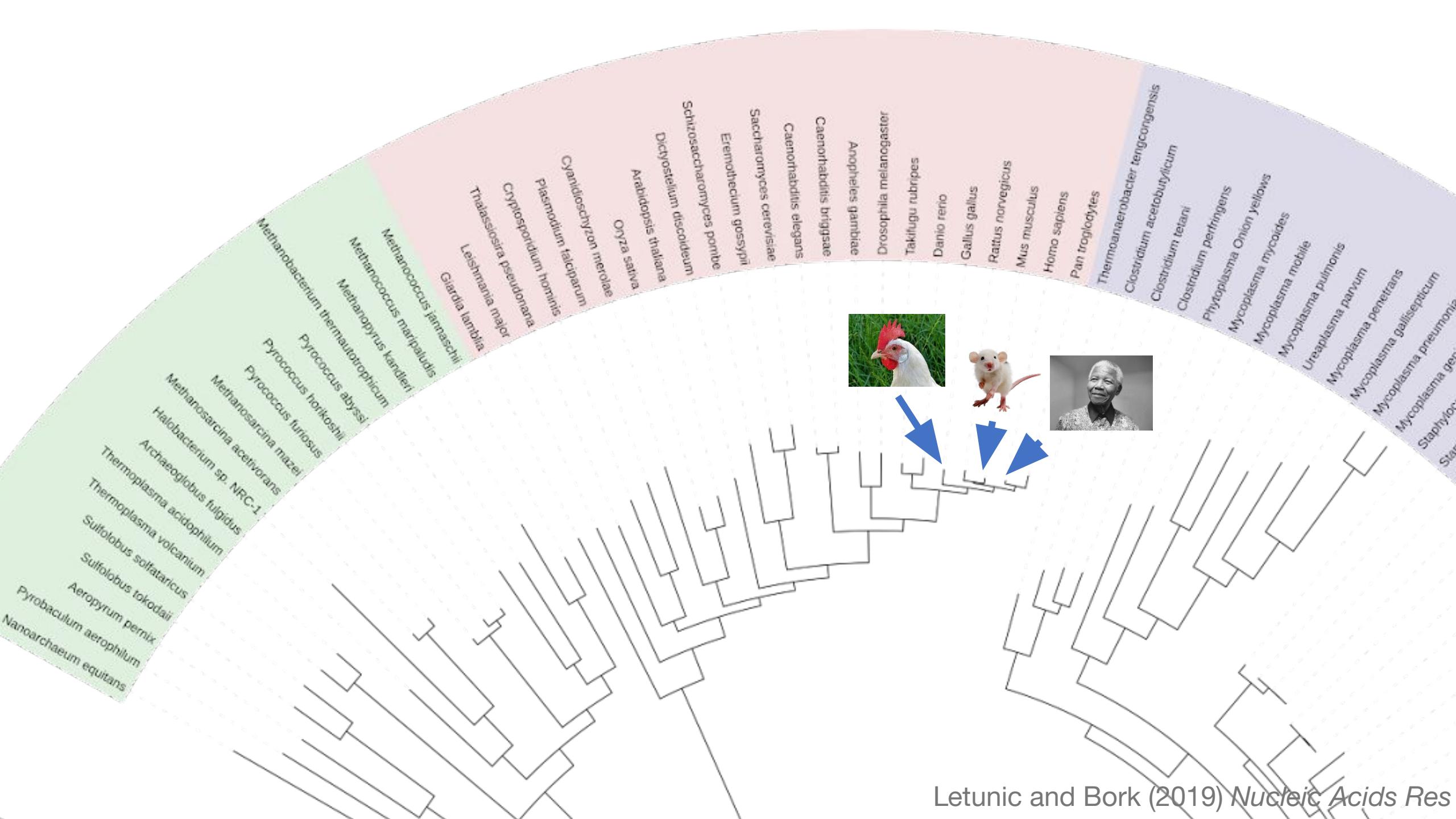
NCCIH: F31 AT010419

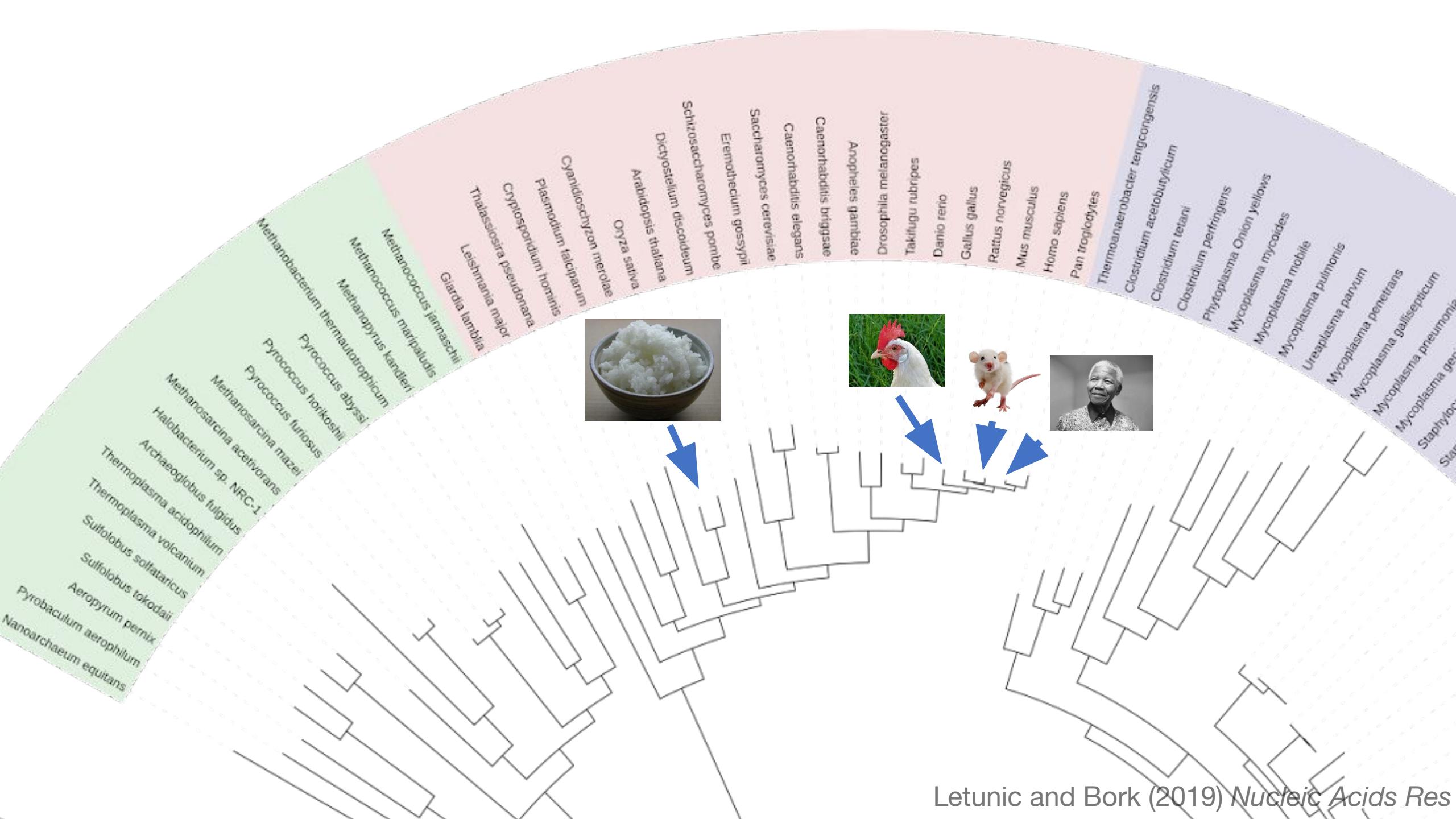
NIGMS: R01 GM125943



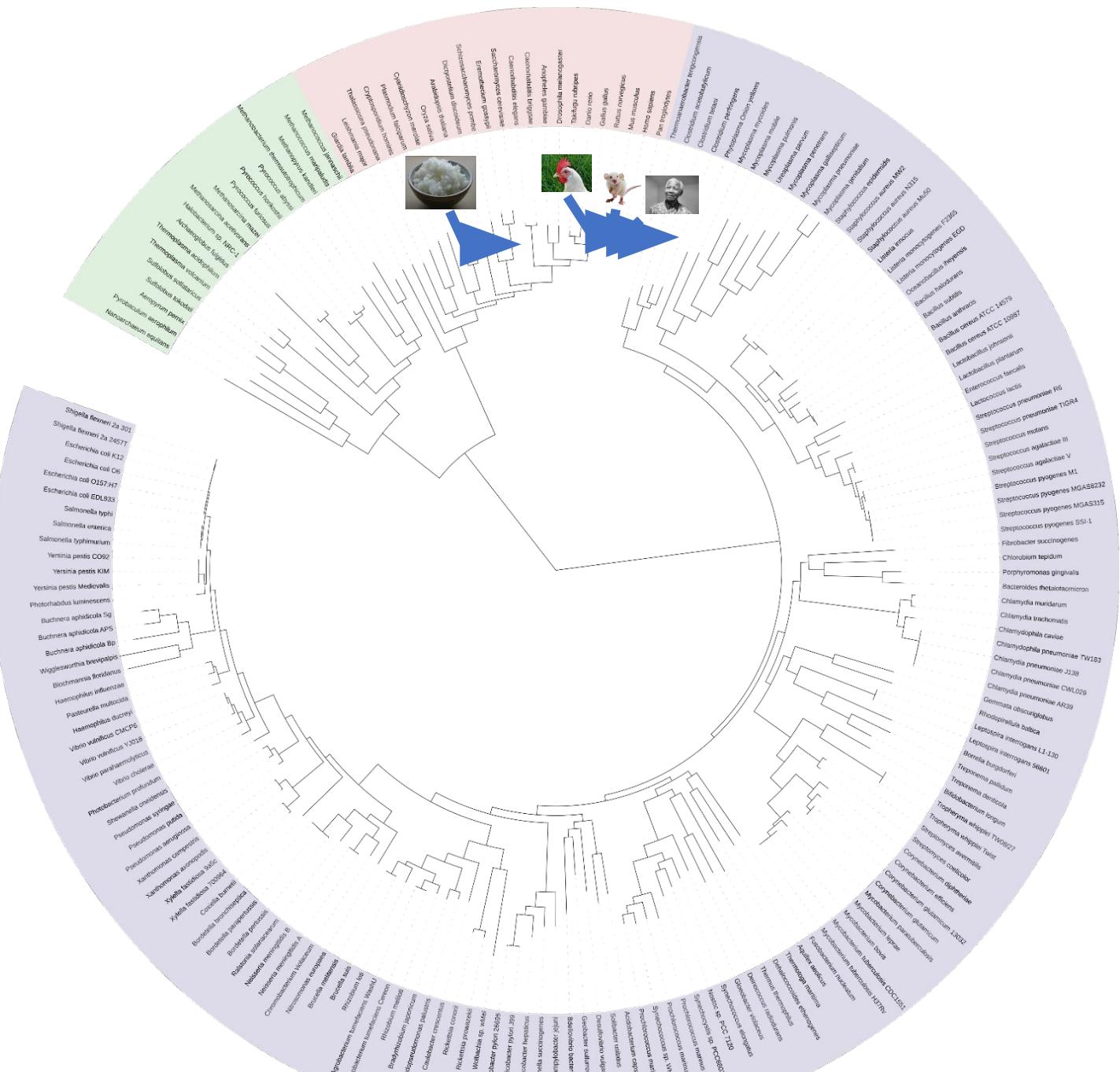
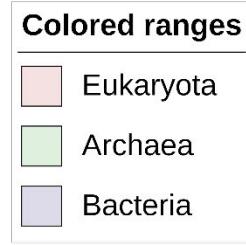








Tree scale: 1

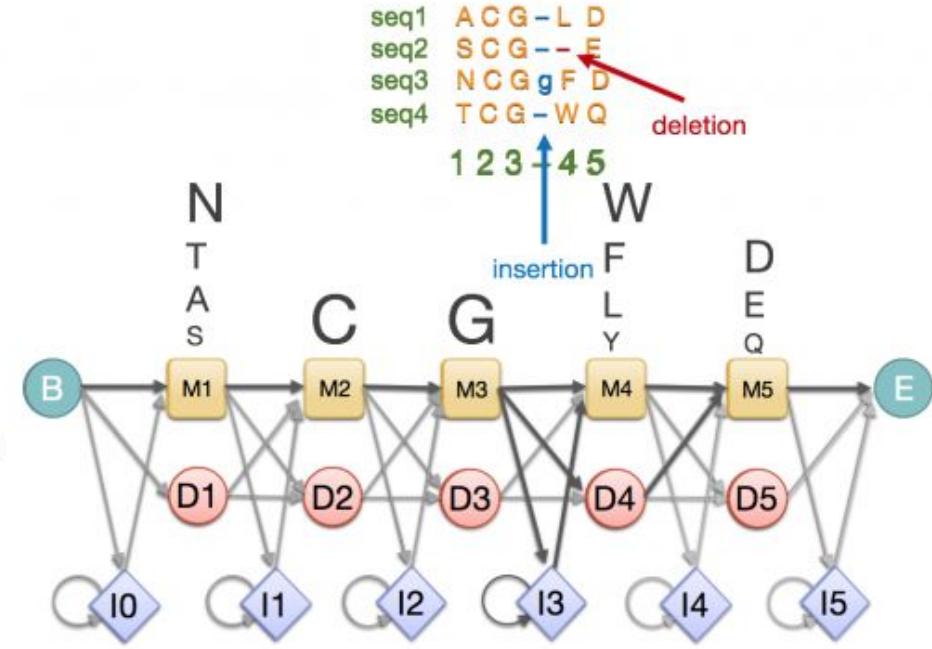


Letunic and Bork (2019) Nucleic Acids Res



pHMM models are a popular tool for classifying sequences

- Start with a multiple sequence alignment
- ↓
- Insertions / deletions can be modelled
- ↓
- Occupancy and amino acid frequency at each position in the alignment are encoded
- ↓
- Profile created





pHMM models are
a popular tool for
classifying
sequences

Input: Query Sequence Set

...SKEAEYLVKQLNTVME...
...SKEAKYLIQQQLDTVMK...
...SKERYAAISMFMK...
...AKEGEYLYSNMLNAV MK...

?

Multiple Alignment

...SKEAEYLVK-QLNTVME...
...SKEAKYLIQ-QLDTVMK...
...SKERYAA---ISMFMK...
...AKEGEYLYSNMLNAV MK...

Input: Target Sequence Set

...CMSDKPDLSLEVTFDKSKLTIQQEKEYNQRS...
...SCALEEHV**SKEAEYLVKMLNAVMKV**TGSFDP...
...DRSQNPPQSKGCCFVTFYTRKAALEAQNALH...
...KMPKDKERSLNAAAQRKLDKQSLKKKGKAE...
...

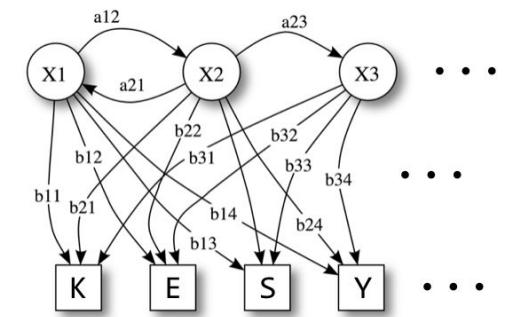
hmmsearch



SKEAEYLVKMLNAVMKV

Output: Resulting Match

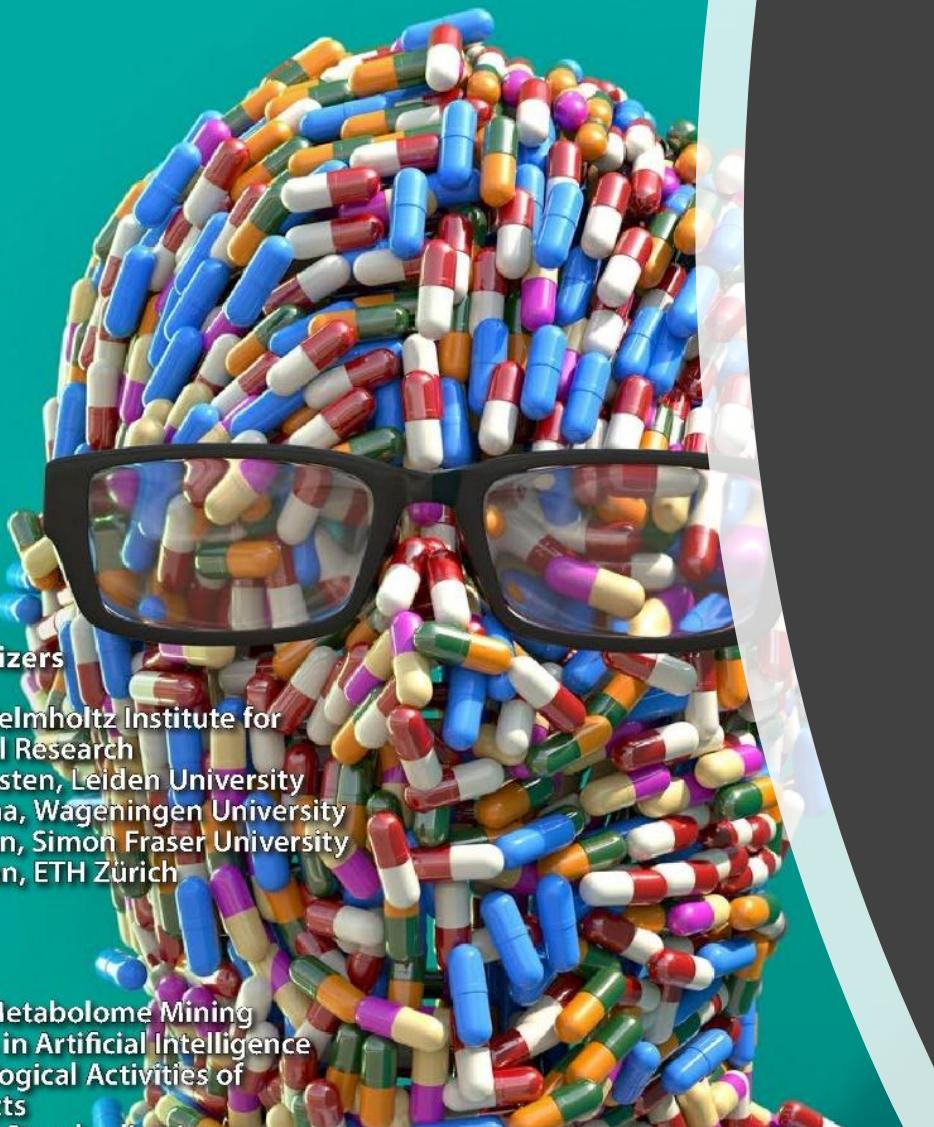
hmmbuild



HMM Profile

Artificial Intelligence for Natural Product Drug Discovery

27 September - 1 October 2021, Leiden, the Netherlands



Scientific Organizers

- Anna Hirsch, Helmholtz Institute for Pharmaceutical Research
- Gerard van Westen, Leiden University
- Marnix Medema, Wageningen University
- Roger Linington, Simon Fraser University
- Serina Robinson, ETH Zürich

Topics

- Genome and Metabolome Mining
- Developments in Artificial Intelligence
- Predicting Biological Activities of Natural Products
- Production and Purification

Artificial Intelligence Approaches to Natural Product Drug Discovery

Nature Reviews Drug Discovery - submitted
(58 authors)

The Threat of Antibiotic Resistance in the United States

Antibiotic resistance—when germs (bacteria, fungi) develop the ability to defeat the antibiotics designed to kill them—is one of the greatest global health challenges of modern time.



New National Estimate*

Each year, antibiotic-resistant bacteria and fungi cause at least an estimated:



Clostridioides difficile is related to antibiotic use and antibiotic resistance:



2,868,700
infections



223,900
cases



35,900 deaths



12,800 deaths

New Antibiotic Resistance Threats List

Updated urgent, serious, and concerning threats—totaling 18

5 urgent threats

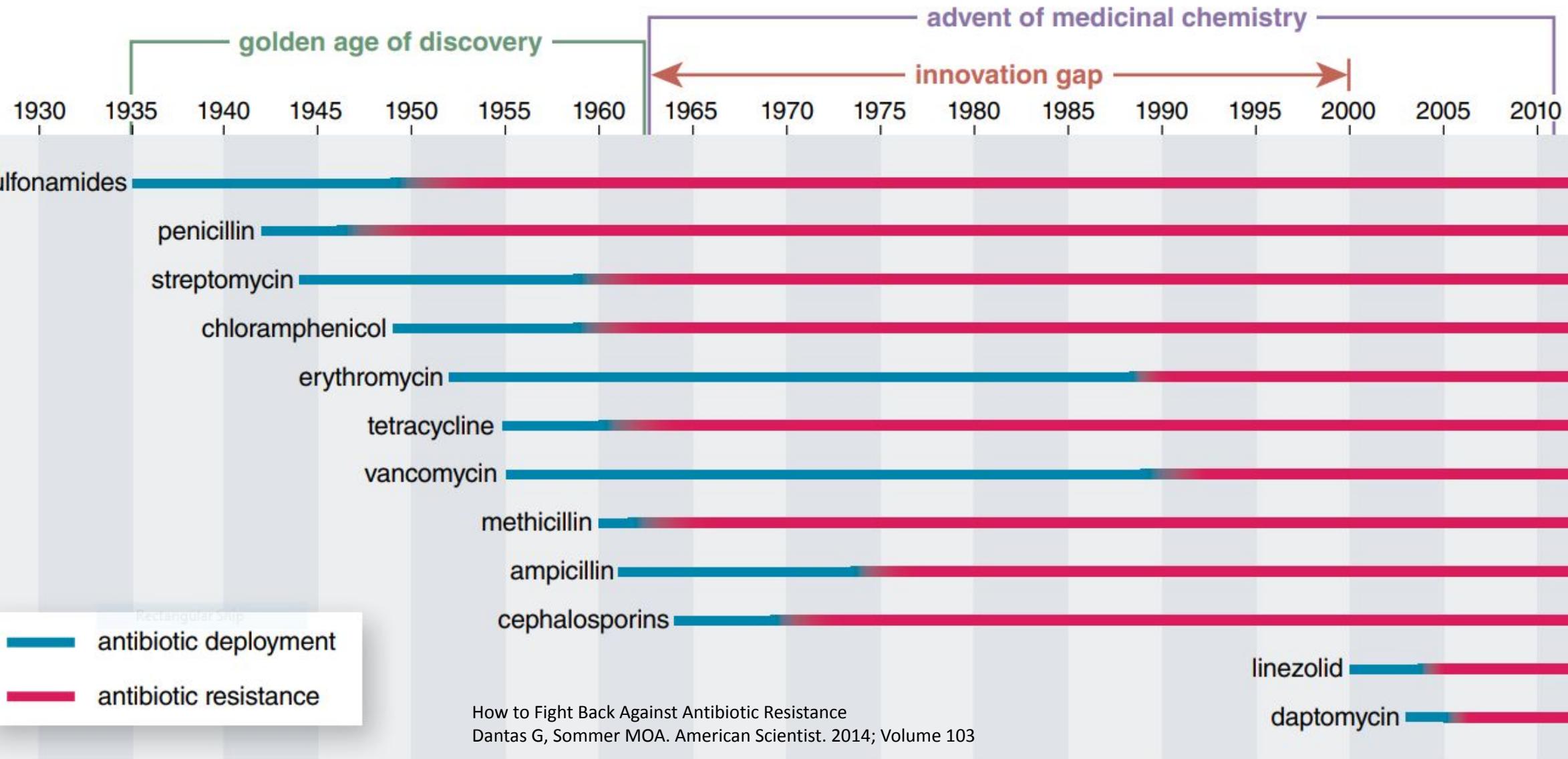
2 new threats

NEW:
Watch List with **3** threats



Antibiotic resistance remains a significant One Health problem, affecting humans, animals, and the environment. Data show infection prevention and control is saving lives—especially in hospitals—but threats may undermine this progress without continued aggressive action now.

Learn more: www.cdc.gov/DrugResistance/Biggest-Threats



Combinatorial Chemistry vs Natural Products



70 HTS campaigns
3 million compounds
19 Hits
5 Leads



65 HTS campaigns
2 million compounds
57 Hits
19 Leads



3 million compounds

— Gram-negative
antibiotics with
cellular activity

Combinatorial Chemistry vs Natural Products



70 HTS campaigns
3 million compounds
19 Hits
5 Leads



65 HTS campaigns
2 million compounds
57 Hits
19 Leads



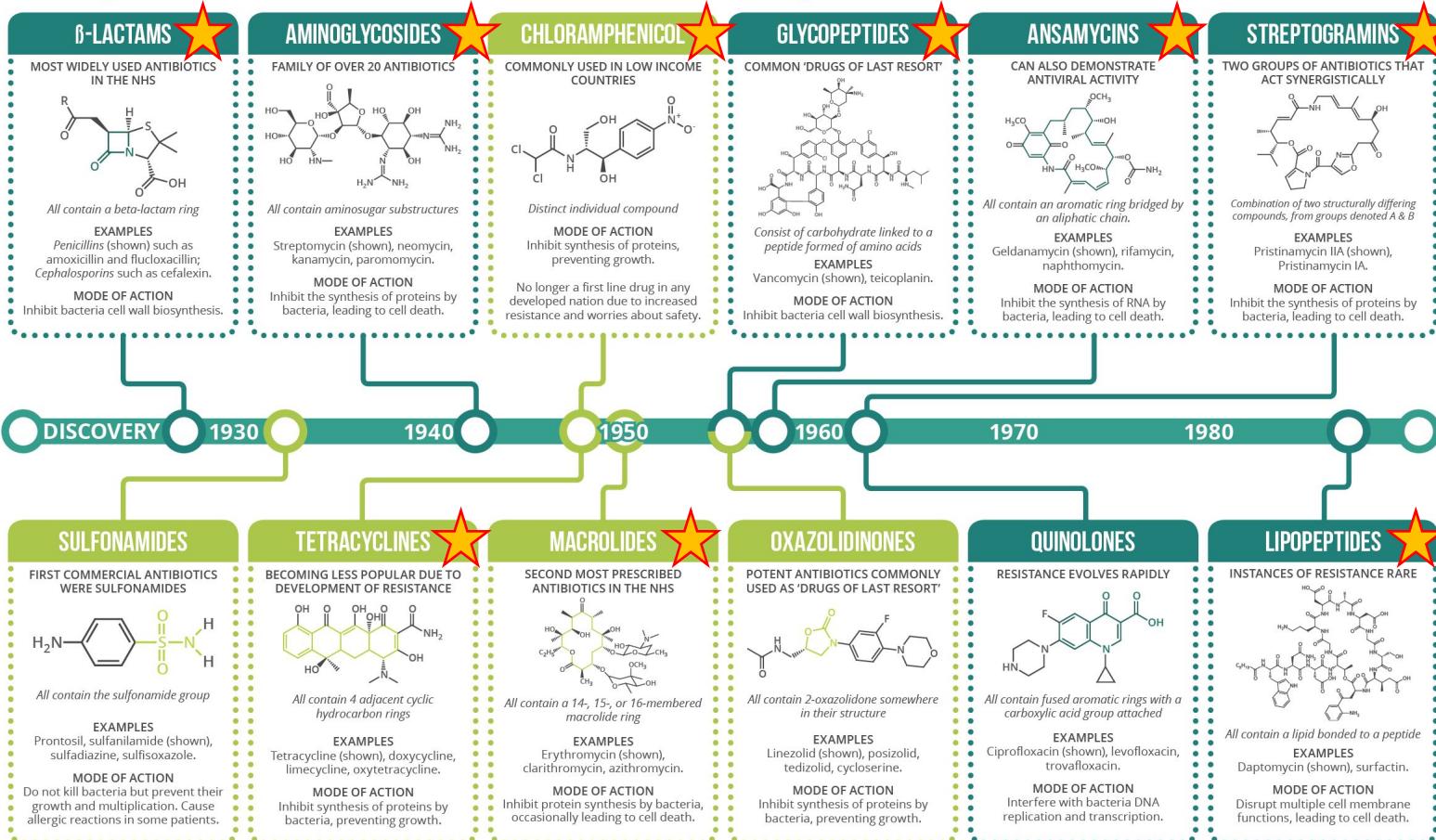
3 million compounds

0 Gram-negative
antibiotics with
cellular activity

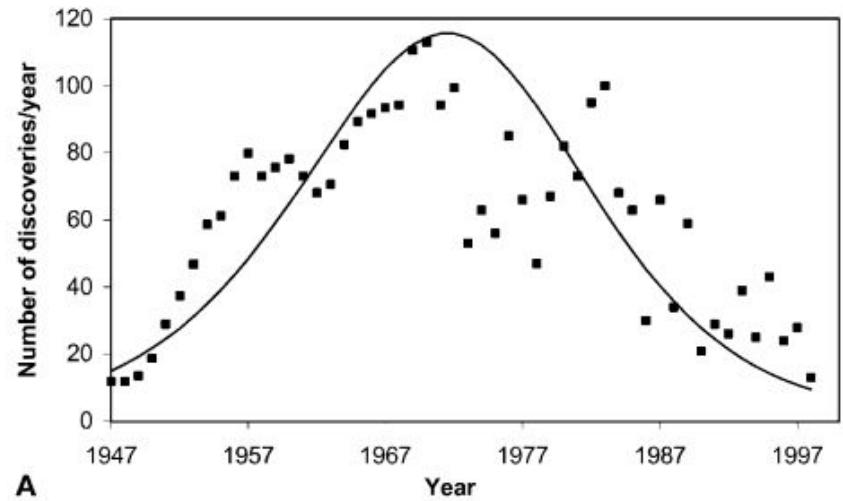
Bacteria (and fungi) as sources of antibiotics

DIFFERENT CLASSES OF ANTIBIOTICS - AN OVERVIEW

Key: ● COMMONLY ACT AS BACTERIOSTATIC AGENTS, RESTRICTING GROWTH & REPRODUCTION ● COMMONLY ACT AS BACTERICIDAL AGENTS, CAUSING BACTERIAL CELL DEATH

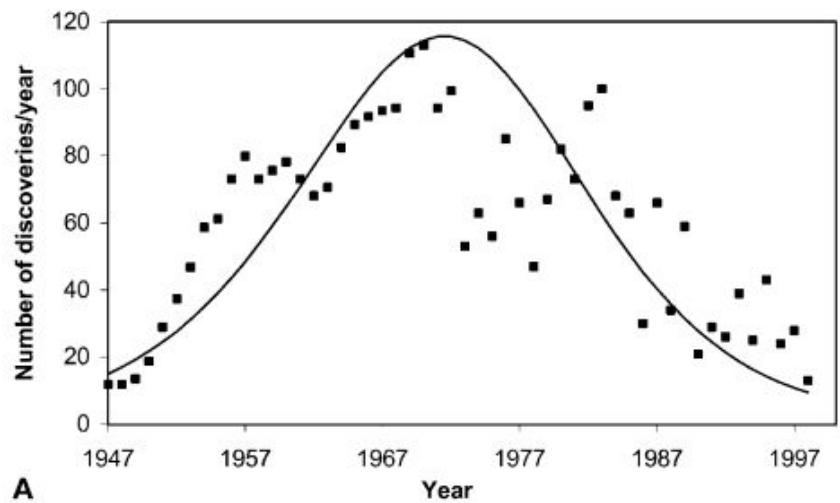


Are there more compounds to find in *Streptomyces*?



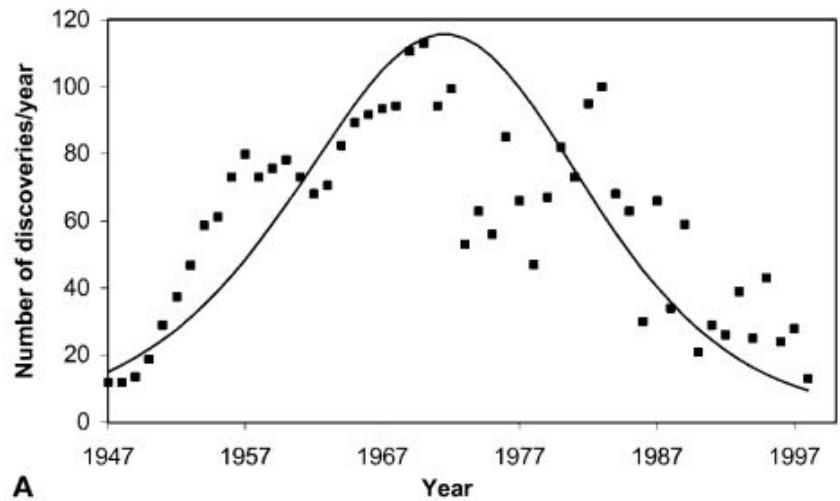
Watve MG, Tickoo R, Jog MM, Bhole BD. How many antibiotics are produced by the genus *Streptomyces*? Archives of Microbiology. 2001. p. 386–390.

Are there more compounds to find in *Streptomyces*?



“...even if we accept the more conservative estimate, only about 3% of the existing compounds have been reported so far.”

Are there more compounds to find in *Streptomyces*?

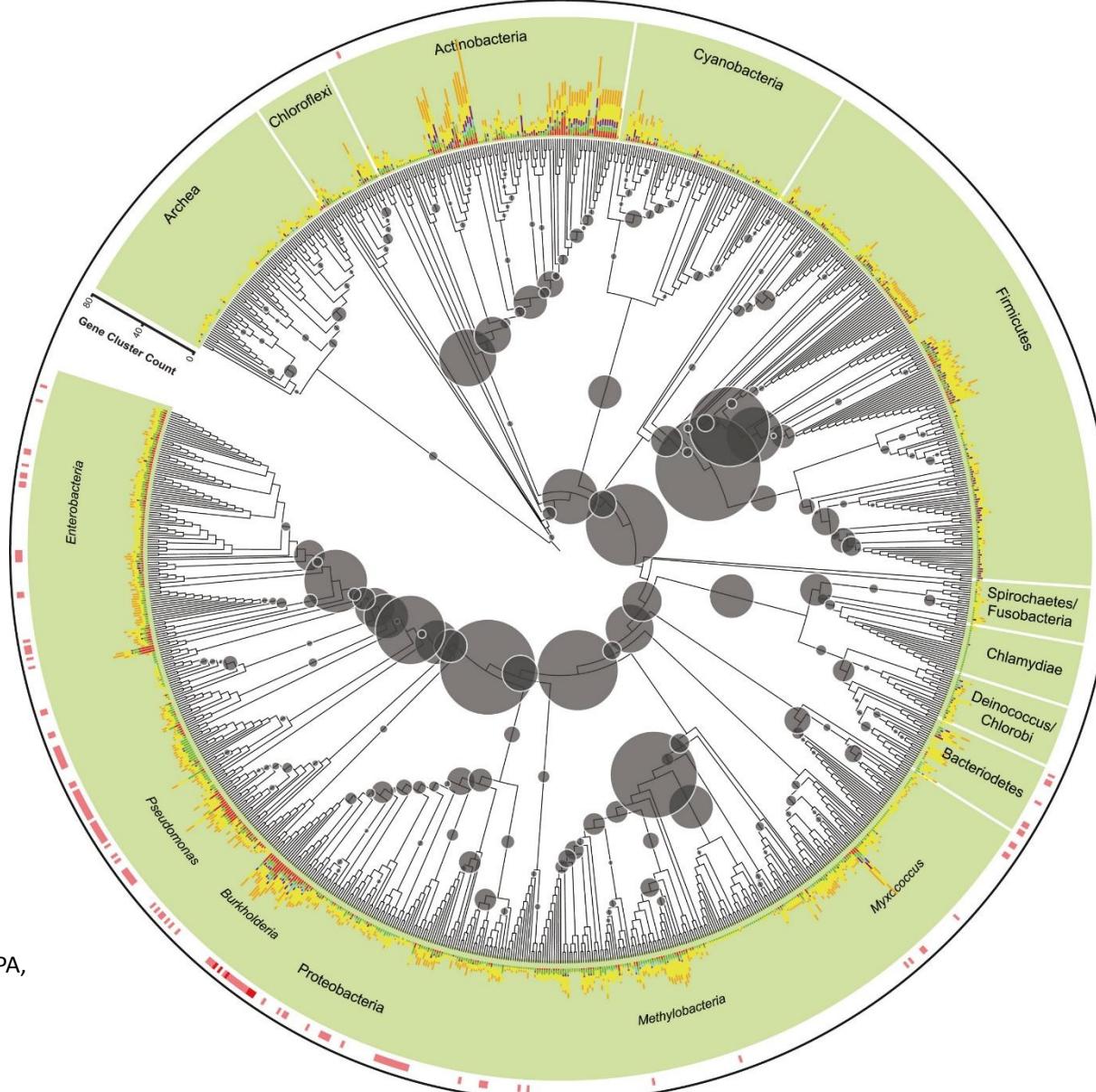


“...even if we accept the more conservative estimate, only about 3% of the existing compounds have been reported so far.”

Problems:

1. Rediscovery of commonly occurring compounds
2. Only tested in limited assays

But... natural products aren't just found in *Streptomyces*

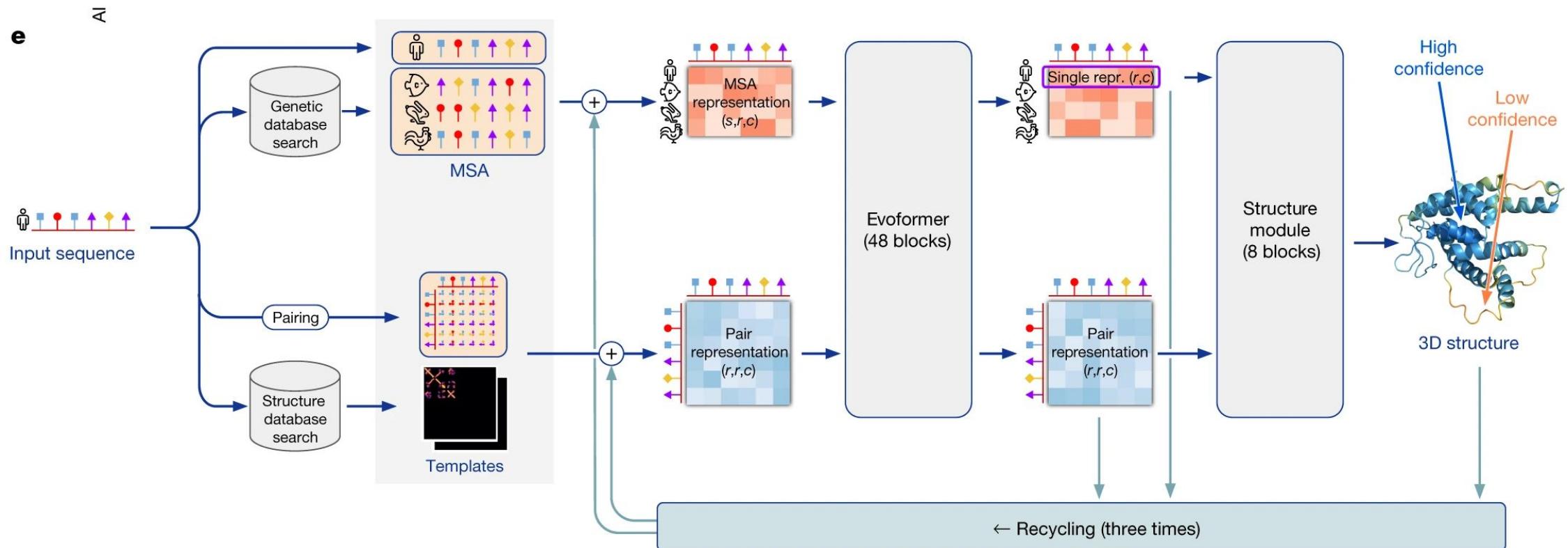


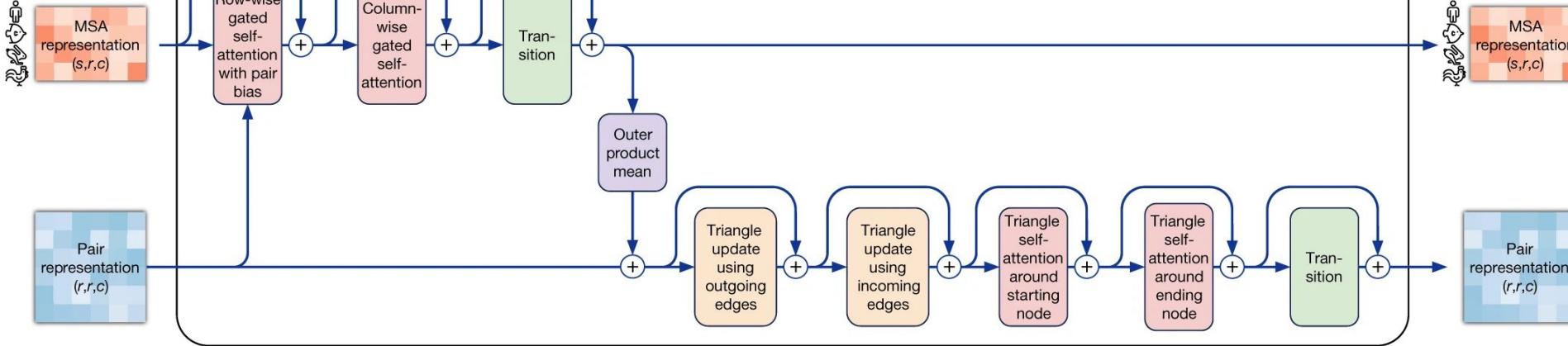
Cimermancic P, Medema MH, Claesen J, Kurita K,
Wieland Brown LC, Mavrommatis K, Pati A, Godfrey PA,
Koehrsen M, Clardy J, Birren BW, Takano E, Sali A,
Lington RG, Fischbach MA. Insights into secondary
metabolism from a global analysis of prokaryotic
biosynthetic gene clusters. *Cell*. Cell Press; 2014 Jul
17;158(2):412–421.

To think about...

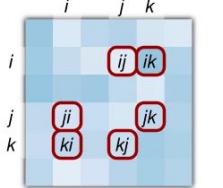
- Where should we look for new bacterial chemistry?
 - New species of bacteria? New phyla?
 - Same species, different location/host organism?
- How to access enough of the chemical compounds?
 - From uncultured organisms?
 - Total synthesis can be slow isn't always an answer

Our methods are scalable to very long protein domains and domain-packing (see Fig. 1c) 2,180-residue protein with no structural homolog

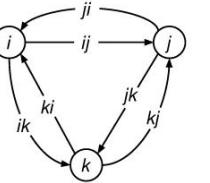




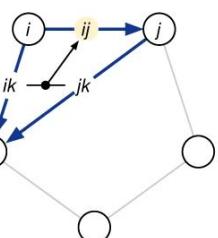
b Pair representation (r, r, c)



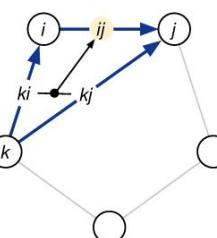
Corresponding edges in a graph



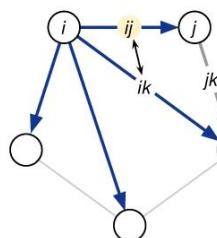
c Triangle multiplicative update using 'outgoing' edges



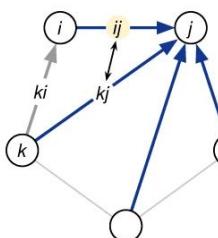
Triangle multiplicative update using 'incoming' edges



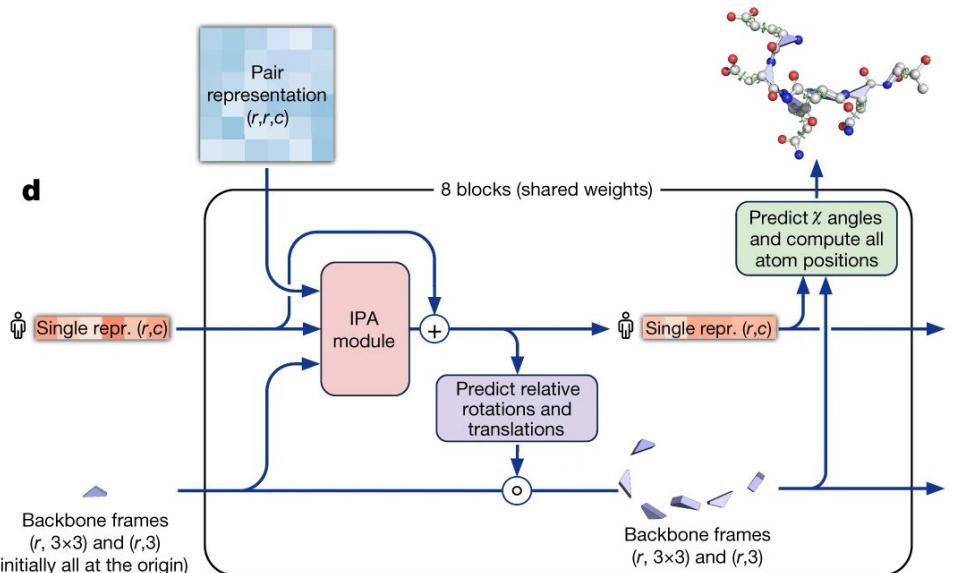
Triangle self-attention around starting node



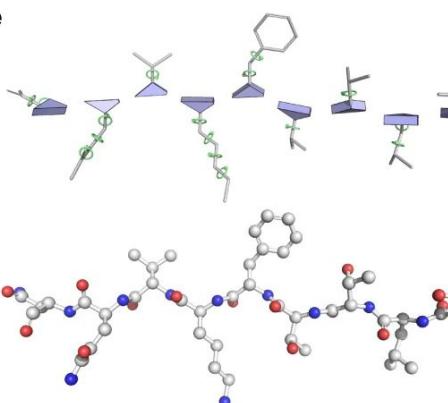
Triangle self-attention around ending node



d



e



f

