



genomics
aotearoa

Metagenomics
Summer School 2023

Day 3

Taxonomic classification

Phylogenetic inference

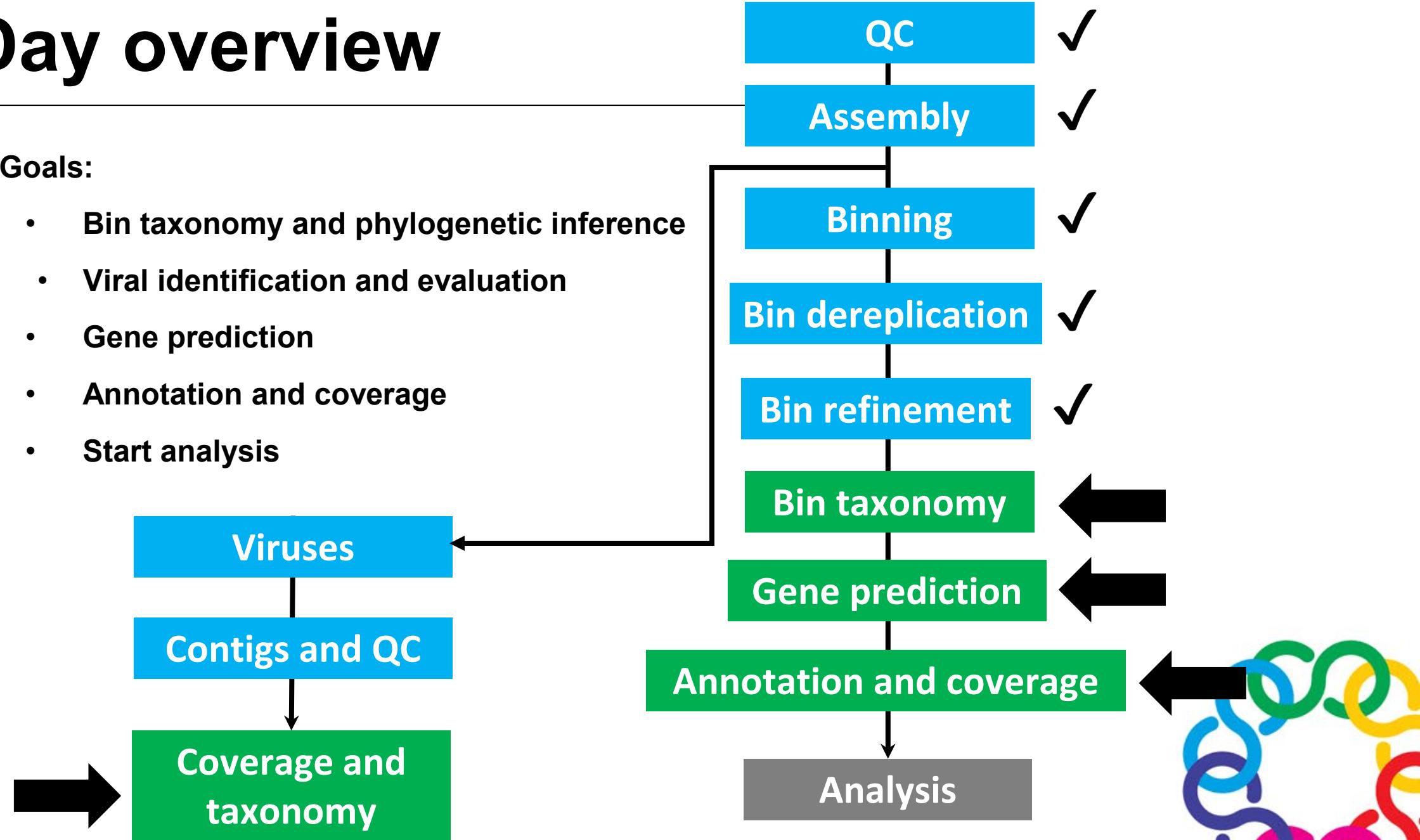
Viruses

Gene prediction & annotation



Day overview

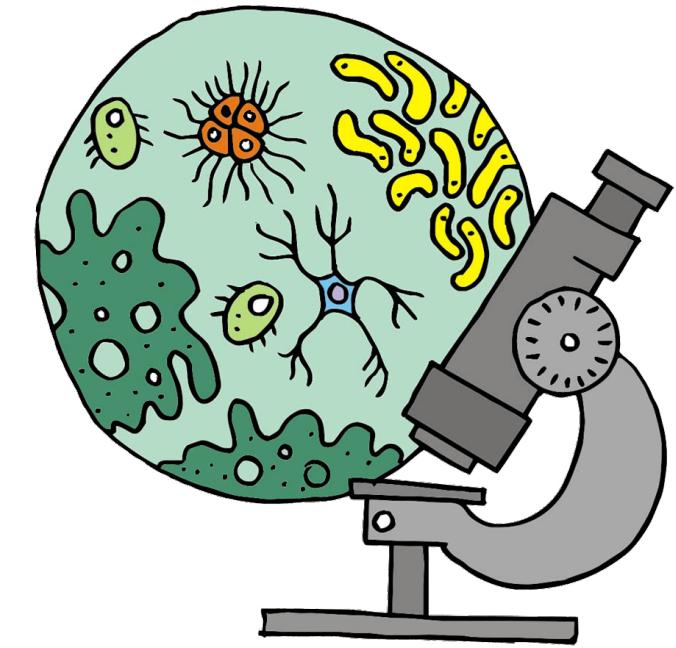
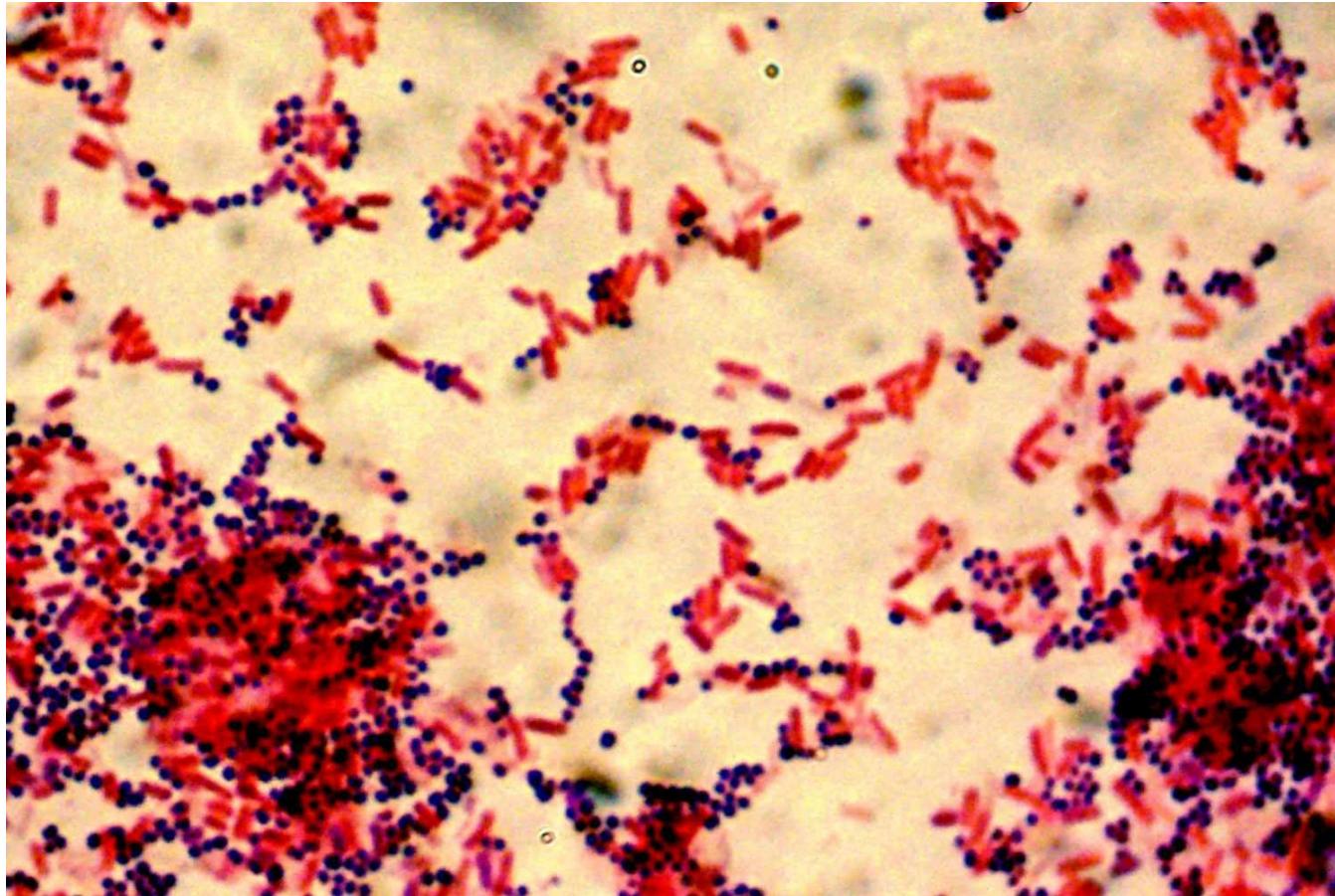
- Goals:
 - Bin taxonomy and phylogenetic inference
 - Viral identification and evaluation
 - Gene prediction
 - Annotation and coverage
 - Start analysis



Bin taxonomic classification

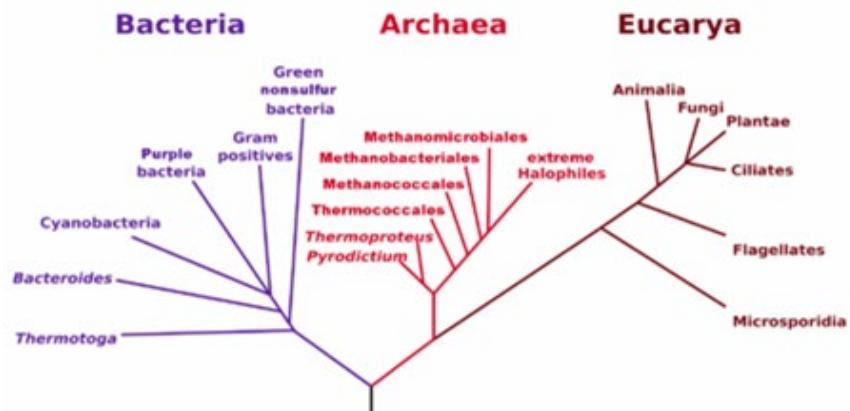


Bin taxonomic classification



Bin taxonomic classification

- Taxonomy is a useful abstraction of the evolutionary process
 - Captures major routes of diversification
 - Not a perfect representation
- Traditionally used a single molecular marker – 16S rRNA gene



Woese et al. (1990) Proc Natl Acad Sci USA 87: 4576



Bin taxonomic classification

- 16S rRNA commonly not recovered by *de novo* assembly
- Can recover 16S and 18S using PhyloFlash and/or EMIRGE
- **Caveat:** can be difficult to assign to genomes in complex communities with many similar taxa

How do we assign taxonomy to bins then?



Bin taxonomic classification

Solution:

- **Use one or more single copy core genes**
- **Universal single-copy genes (USCGs) are marker genes that occur once and only once in almost every genome**

([Wang et al., 2022, Scientific Reports](#))

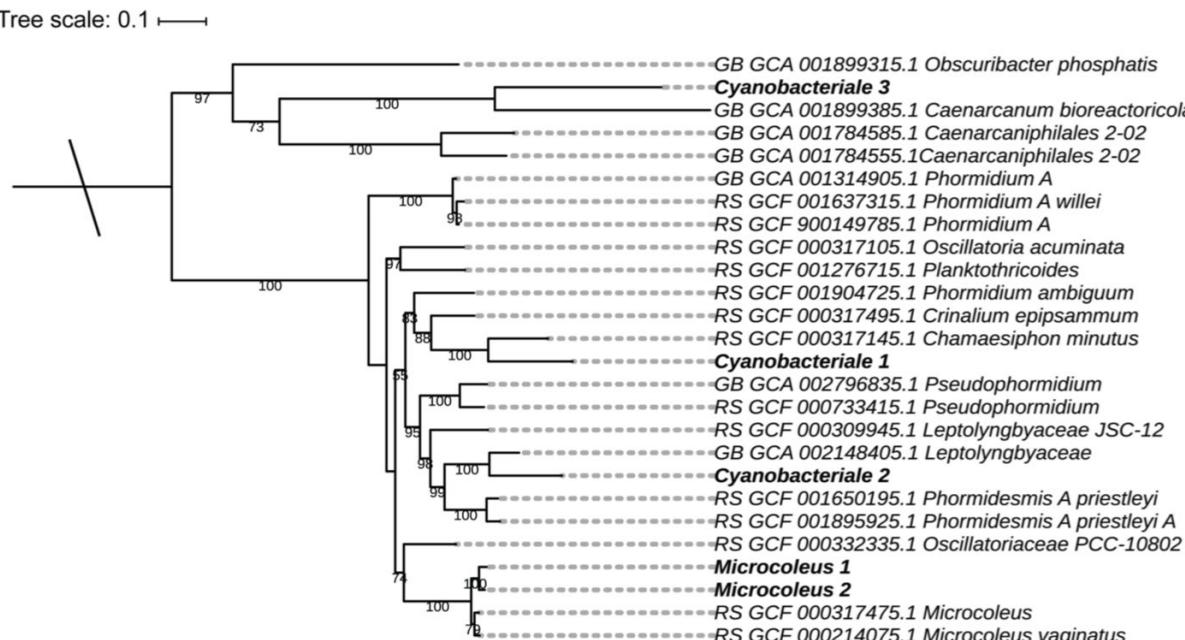


Bin taxonomic classification

Solution:

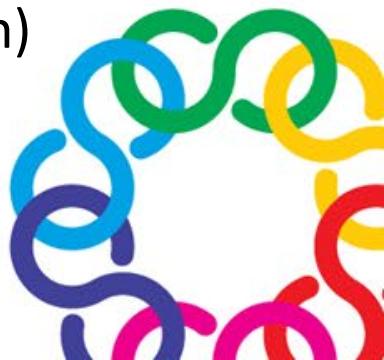
- Use one or more single copy core genes
- Universal single-copy genes (USCGs) are marker genes that occur once and only once in almost every genome
- Concatenate protein sequences of multiple single copy core genes

(Wang et al., 2022, Scientific Reports)



Concatenated protein
sequence tree:

Phylogenetic placement of
cyanobacterial genome bins
(Wai-iti River, Nelson)



Bin taxonomic classification



GTDB builds upon a number of existing public resources in order to provide a taxonomic resource that reflects recently proposed taxa, changes in taxonomic opinion, and the wealth of publicly available genomes.



GTDB - how it works

Genomes are obtained from NCBI and must meet the following criteria to be included in the GTDB reference trees and database:

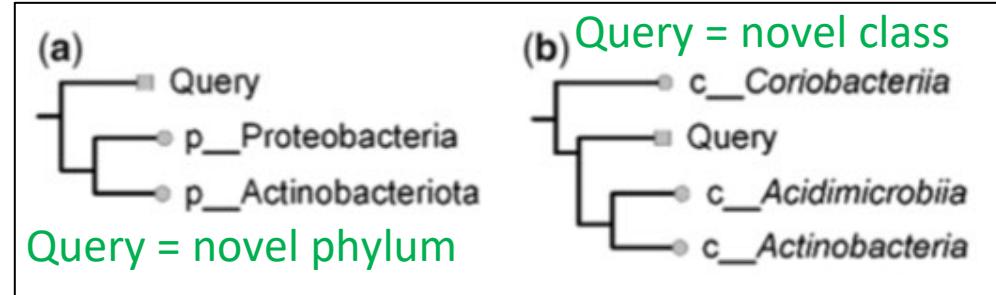
1. CheckM completeness estimate >50%
2. CheckM contamination estimate <10%
3. quality score, defined as completeness - 5*contamination, >50
4. contain >40% of the bac120 or arc53 marker genes
5. contain <1000 contigs
6. have an N50 >5kb
7. contain <100,000 ambiguous bases

Filtered genomes are manually inspected and exceptions are made for genomes of high nomenclatural or taxonomic significance



Bin taxonomic classification

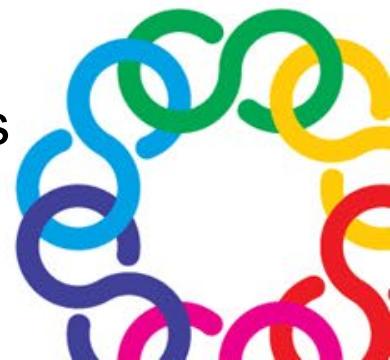
- **Genome Taxonomy Database Toolkit (GTDB-Tk)**
- **Use to classify bins against reference genome trees (GTDB)**
- **Uses set of 120 concatenated protein sequences (of single copy core genes)**



(Chaumeil et al., 2019, Bioinformatics)

Rank assignment based on:

- Tree placement in GTDB ref tree
- Relative Evolutionary Divergence (value between 0=root and 1=tip)
- Species assignment:
 - ANI to reference genomes



Discriminate species

Proxy for DNA-DNA hybridization

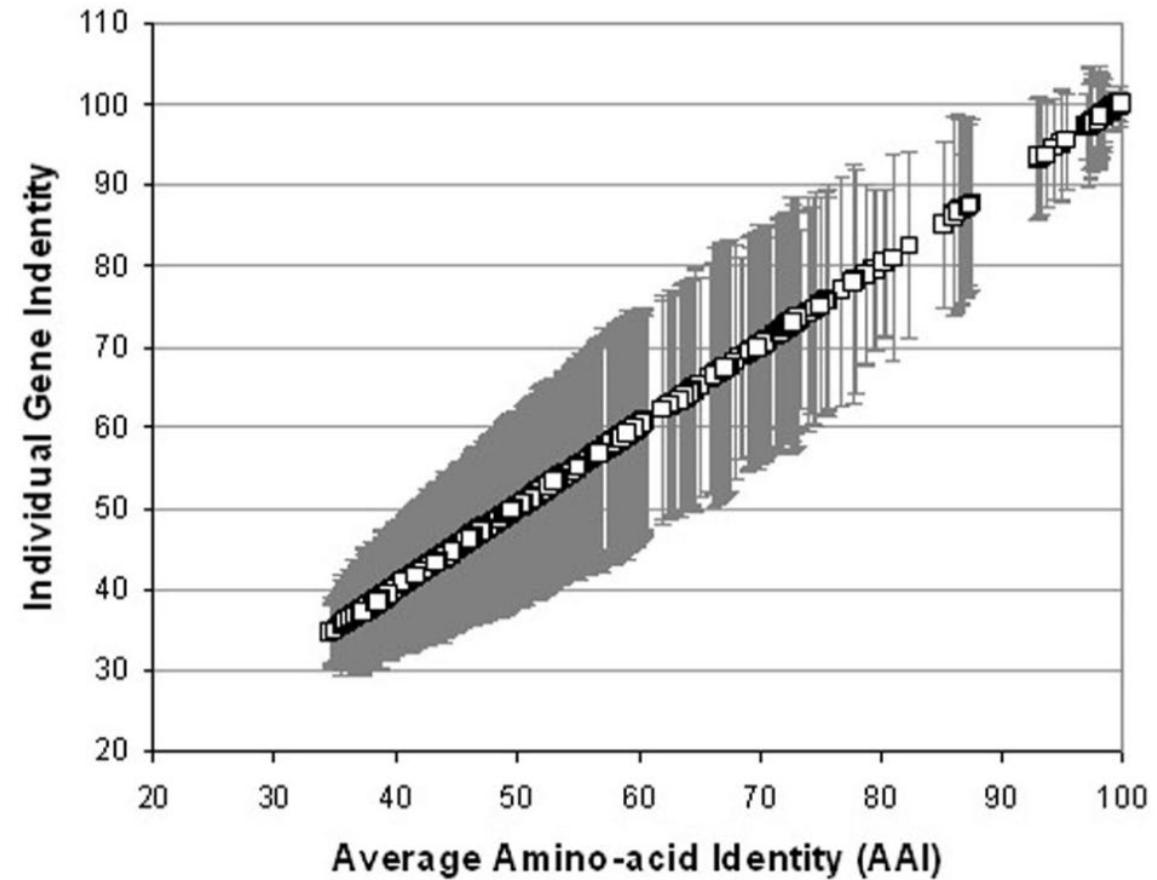
Pairwise genome comparisons:

- Average Nucleotide Identities (ANI)
 - gene comparisons
- Average Amino Acid Identities (AAI)
 - predicted protein comparisons
- Alignable Fraction (AF)
 - proportion of genes that align

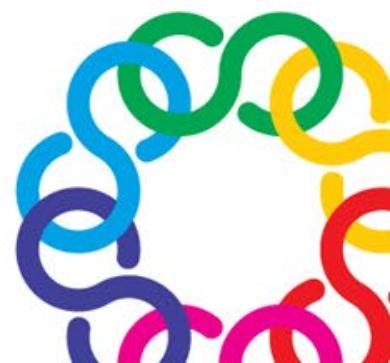
Determine via: Pairwise BLAST-like search



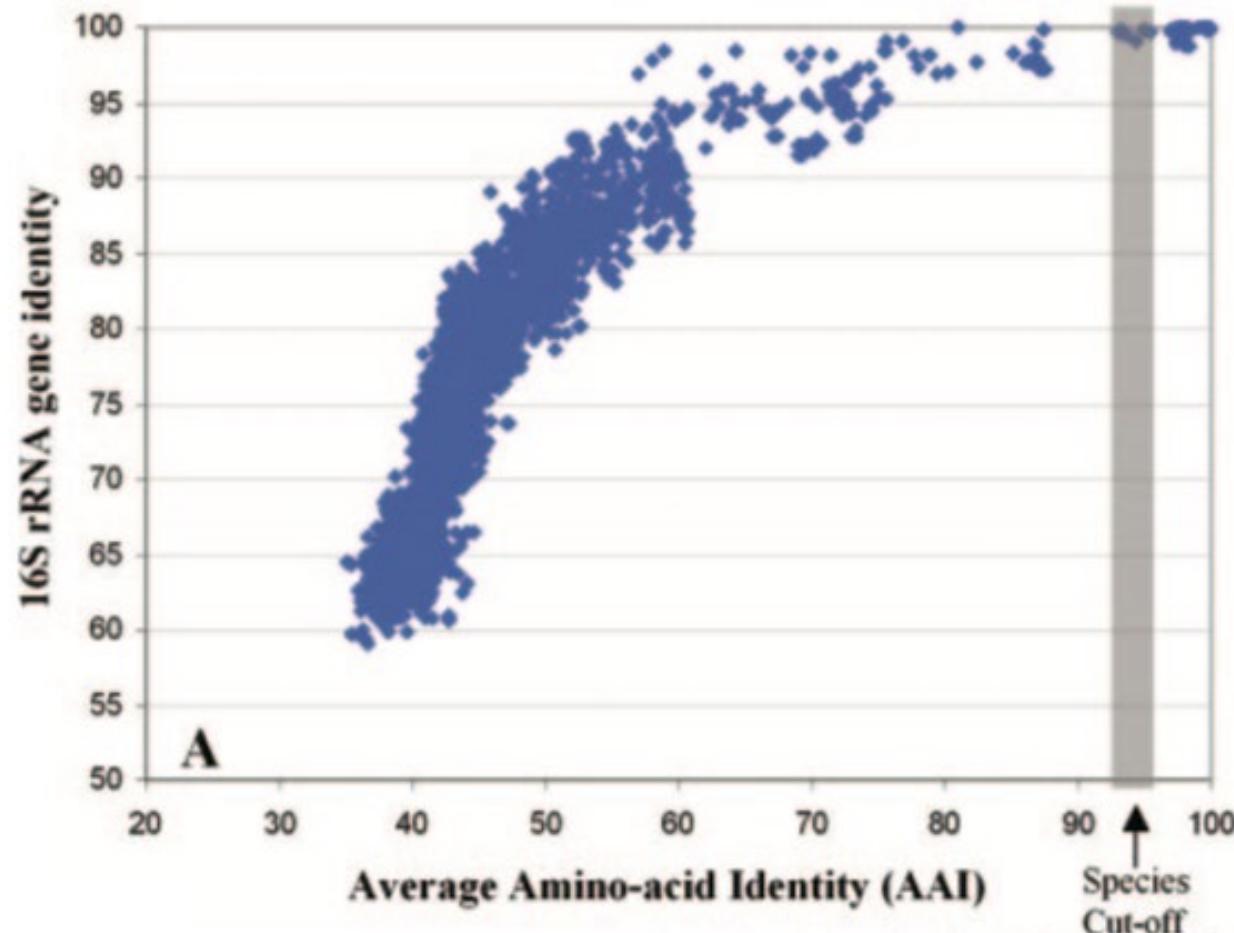
Discriminate species



(Fig. 1, Konstantinidis and Tiedje, 2005, J Bacteriology)



Discriminate species



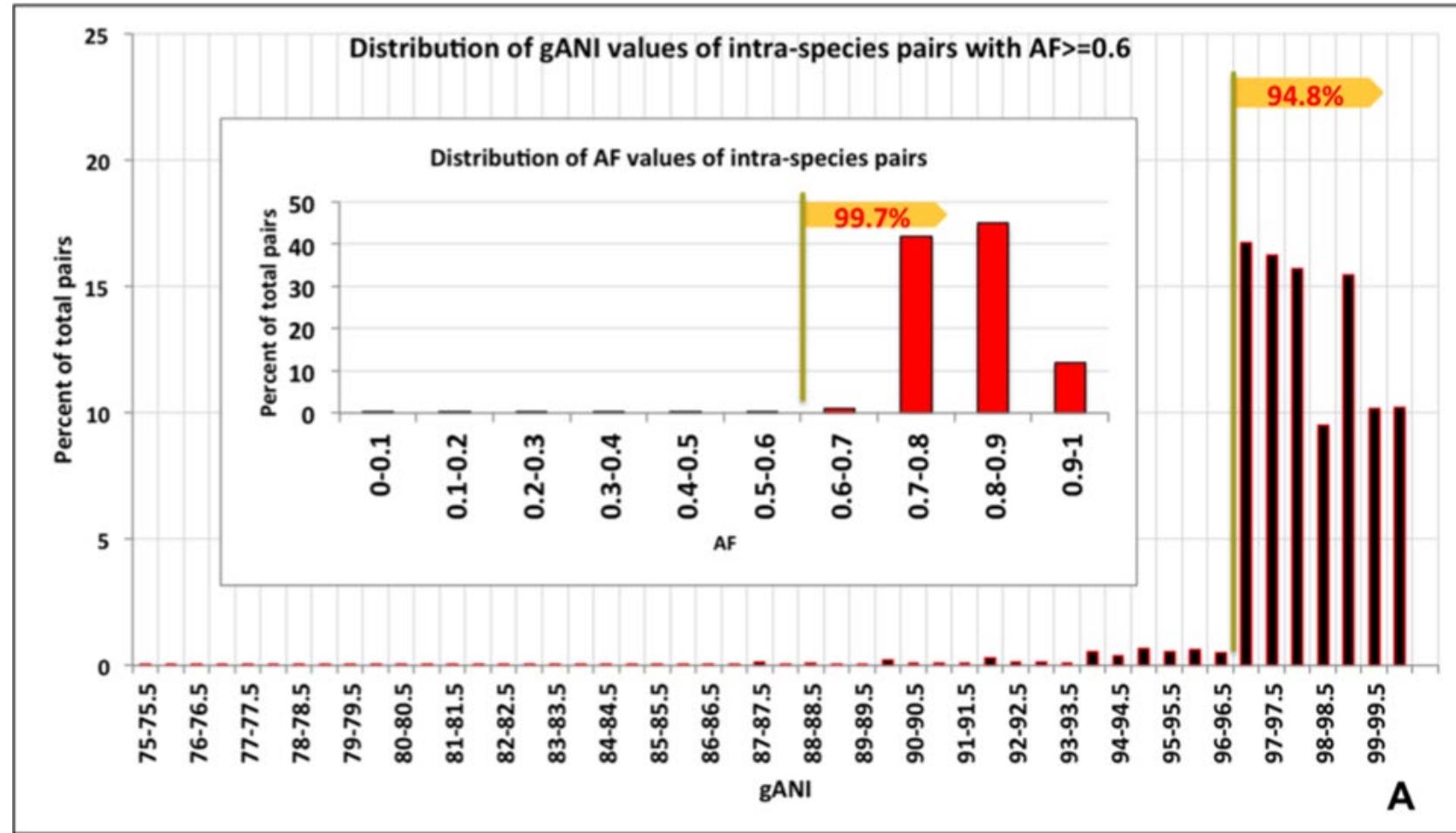
- AAI species cutoff \approx 95-96% ID
- Equivalent to 70% DNA-DNA hybridization threshold for species

(Fig. 3, Konstantinidis and Tiedje, 2005, J Bacteriology)



Discriminate species

10,998 IMG genomes -- 1,130,980 intra-species genome pairs



(Fig. 1A, Varghese et al., 2015, Nucleic Acids Research)



Discriminate species

Table S6. Pairwise average amino acid identities (AAI) shared between genome bins.

- Pairwise AAI comparisons between genomic bins from the Gulf of Mexico seafloor
 - All unique species (i.e. <95-96% AAI)
 - Figures shows clusters of similar genomes
 - Red = more similar
 - Blue = dissimilar

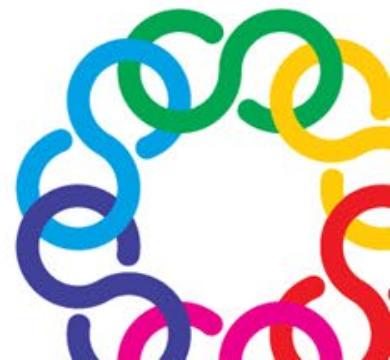
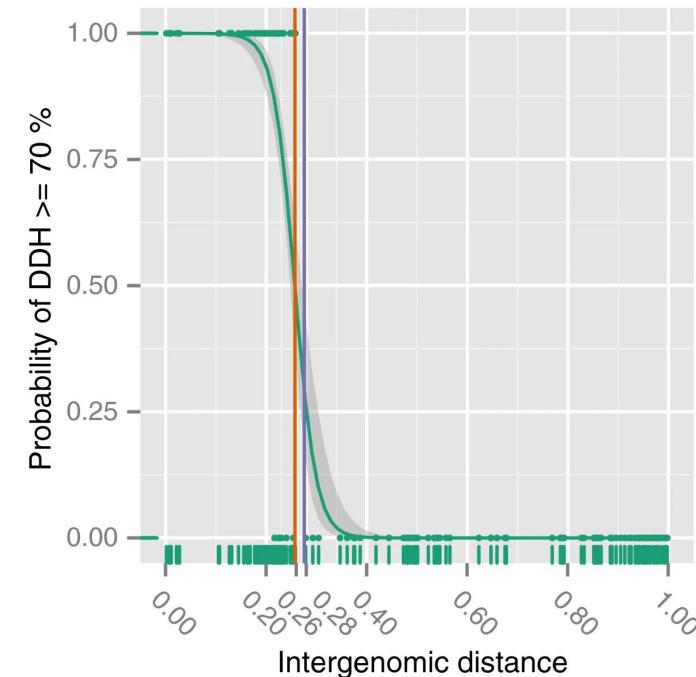


(Handley et al., 2017, ISME J)

Discriminate species

Alternative to gANI or AAI:

- GGDC (Genome Blast Distance Phylogeny) <https://ggdc.dsmz.de/>
 - Makes use of generalised linear models



Task: Genome taxonomic classification

[Go to Github MGSS webpage](#)

Tasks:

- Taxonomic classification of bins using GTDB-Tk



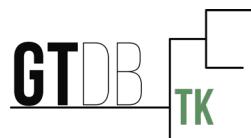
Phylogenetics

Hot take

“There’s no such thing as a genome tree”



Phylogenetics



GTDB-Tk output file: Multiple Sequence Alignment (MSA) of 120 marker genes (protein sequences)

bin_0_Filtered

SVLEAVSLDAWGEAIKRMRESIVKTSAGSEQFRVKTSYKKSLARLVTLNGVEVKIYDKEAIKKASLFNYTIISVNEYAVCTVCKRQTAEDIKIVKKMQUEGEGLESNVEILLEDEIENHSSWADLISDDDISNCNDSSASVSHIIIANARQTTEKTINITKDVVLADAVTIDSLEIKESVGKEIKIVDLNIYVNRQLGIGENIMTTET----EALALDSLGTVNKAAYYKLNLHKYLATGWYNFDILFKIVGEIYQKEDSVAGETGKAFAKQIKLVSAFSIKSFEKVNIDTFLIELIYIRDMSYSDLKFGLYITDQKQARANKSDSKIVDG-KVQLNHDAFLFMNVNLLTVAAAGDIIADEQLIGGGKELGMSI-KDLNKTRGIKATAILALEIEVKKLRAMVDEDIAKENQIEFEDSRCRIMWTRELEELWTSDDWASWIGHTTIG-IVTYSMGNEINRIFTDEK-LLLSPWPPTDEYVAKKRVDELMAKMFAPELLKKHLFIAPPFIFSDVKVKVKE-IAGEASDGSNRYHIVDQLFLIAGRFTAERLIKNAIKSFKATKVNTKI-FAHQKLETD-INASLLENSKVKGQEIEITADMTISRALKIGCIGGIIMEHSIKMNEALDTVDTHGDVIFPMVDELDT-I-EKVLK-FVLAIPDLNVE--LEKVLLASDL-LFAGIDEKIISLMASAEVILQDAMLTKSEIAKALMVLIALYVA-IGKPDNAVNPAFIGNEISQERLEKEDIEADGVETVRVMPHIALLSVRTQHNESEDRGARMVLNLTKLKAESFTTAKAEQILLAGQHRKFSRFNQVGRGHIDIPKEIVTTVKAGFTEVRLVIVKFVFRVVKNTFADSITAIKYNAIPLGAAEIASPVA---RRAAIKLLWVKITVNCSSVIRKIDVIDKTEILRNKNNTIATSTHYTKMEQKAKLKLAGSP--LELRDGVVGTAKGNKGFMPCSQVYRSMRNIGKGRSMIQIKPTKPTKVNIPFVSDEVLLVP-LKIAESNNDGTVVVKLSNDMRTITMDKKMASIAVSTTATSLVFISVQSITLAFTELKGKFQV-IIPAANQIEAEKGVEIYTVKGDNALTKQIKLDKEQVSRSIVETEVIGLYAFEFALNNIGIKVEIRLA-LAVNLFPKASELEIKGYENRKIFAILEGELKLQSHKILEKAGNSNRIENIATTCKDFLIF-EIHYPDIPYEYSDIDK-NVGITTSENLTDEIKYYVLQALIAAMDKEALVFKGSPENIRKETV-FAIAGKDANIAMLTH-SVPVLTKTD-FFTNAAEKTKVKNG-Q

Bin_0

Bin_1

Bin_2

Bin_3

>bin_0.filtered

SVLEAVSLDAWGEAIKRMRESIVKTSAGSEQFRVKTSYKKSLARLVTLNGVEVKIYDKEAIKKASLFNYTIISVNEYAVCTVCKRQTAEDIKIVKKMQUEGEGLESNVEILLEDEIENHSSWADLISDDDISNCNDSSASVSHIIIANARQTTEKTINITKDVVLADAVTIDSLEIKESVGKEIKIVDLNIYVNRQLGIGENIMTTET----EALALDSLGTVNKAAYYKLNLHKYLATGWYNFDILFKIVGEIYQKEDSVAGETGKAFAKQIKLVSAFSIKSFEKVNIDTFLIELIYIRDMSYSDLKFGLYITDQKQARANKSDSKIVDG-KVQLNHDAFLFMNVNLLTVAAAGDIIADEQLIGGGKELGMSI-KDLNKTRGIKATAILALEIEVKKLRAMVDEDIAKENQIEFEDSRCRIMWTRELEELWTSDDWASWIGHTTIG-IVTYSMGNEINRIFTDEK-LLLSPWPPTDEYVAKKRVDELMAKMFAPELLKKHLFIAPPFIFSDVKVKVKE-IAGEASDGSNRYHIVDQLFLIAGRFTAERLIKNAIKSFKATKVNTKI-FAHQKLETD-INASLLENSKVKGQEIEITADMTISRALKIGCIGGIIMEHSIKMNEALDTVDTHGDVIFPMVDELDT-I-EKVLK-FVLAIPDLNVE--LEKVLLASDL-LFAGIDEKIISLMASAEVILQDAMLTKSEIAKALMVLIALYVA-IGKPDNAVNPAFIGNEISQERLEKEDIEADGVETVRVMPHIALLSVRTQHNESEDRGARMVLNLTKLKAESFTTAKAEQILLAGQHRKFSRFNQVGRGHIDIPKEIVTTVKAGFTEVRLVIVKFVFRVVKNTFADSITAIKYNAIPLGAAEIASPVA---RRAAIKLLWVKITVNCSSVIRKIDVIDKTEILRNKNNTIATSTHYTKMEQKAKLKLAGSP--LELRDGVVGTAKGNKGFMPCSQVYRSMRNIGKGRSMIQIKPTKPTKVNIPFVSDEVLLVP-LKIAESNNDGTVVVKLSNDMRTITMDKKMASIAVSTTATSLVFISVQSITLAFTELKGKFQV-IIPAANQIEAEKGVEIYTVKGDNALTKQIKLDKEQVSRSIVETEVIGLYAFEFALNNIGIKVEIRLA-LAVNLFPKASELEIKGYENRKIFAILEGELKLQSHKILEKAGNSNRIENIATTCKDFLIF-EIHYPDIPYEYSDIDK-NVGITTSENLTDEIKYYVLQALIAAMDKEALVFKGSPENIRKETV-FAIAGKDANIAMLTH-SVPVLTKTD-FFTNAAEKTKVKNG-Q

Phylogenetic tree software takes the MSA as input and builds an evolutionary tree



Phylogenetics



Ronald Fisher

In 1928, Fisher was the first to use equations to calculate the distribution of allele frequencies and the estimation of genetic linkage by **maximum likelihood methods** among populations.

R. A. FISHER

89

In the absence of linkage their deviations will be independent, but if linkage is present the mean value of pq may be found to be

$$-3n \frac{1-4x}{3},$$

or, the correlation between p and q is

$$\rho = -\frac{1-4x}{3}.$$

The simultaneous deviation of p and q from zero will therefore be measured by

$$Q^2 = \frac{1}{3n} \left\{ \frac{1}{1-\rho^2} (p^2 - 2\rho pq + q^2) \right\} \\ = \frac{3}{8(1-x)(1+2x)n} \left\{ p^2 + q^2 + \frac{2}{3}(1-4x) pq \right\}.$$

This expression, which of course depends upon x , is a quadratic function of the frequencies; in this it resembles χ^2 , and on comparing term by term the two expressions it appears that

$$\chi^2 = Q^2 + \frac{1}{I} \left\{ \frac{a}{2+x} - \frac{b+c}{1-x} + \frac{d}{x} \right\}^2,$$

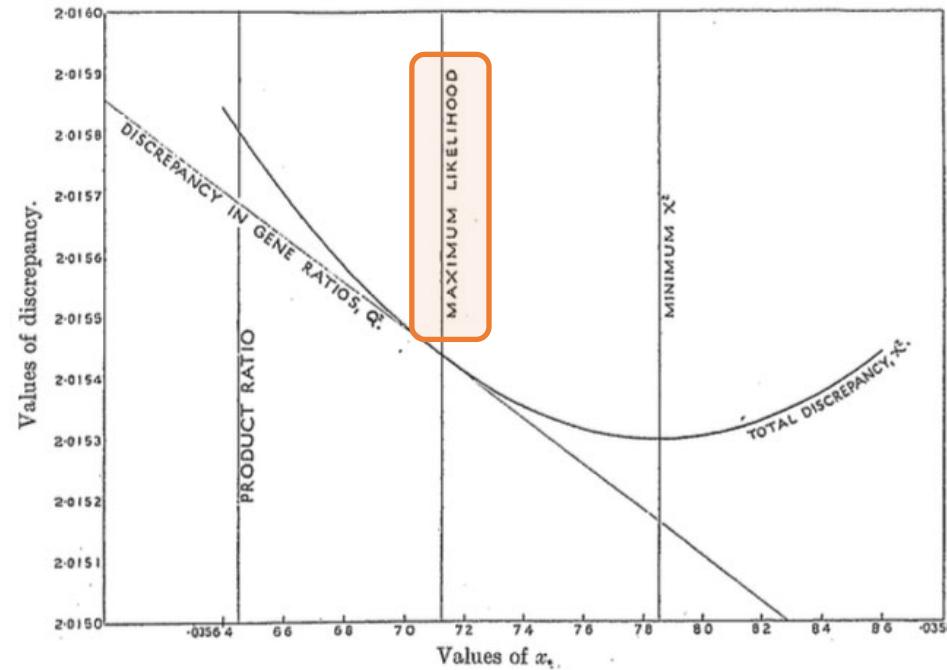


Fig. 2.

Maximum Likelihood Phylogenetic Inference estimates the most likely phylogenetic tree that represents the branching patterns and relationships among a set of biological sequences, usually DNA or protein sequences.

Phylogenetics

Phylogenetic inference using maximum likelihood:



PLOS ONE

OPEN ACCESS PEER-REVIEWED

RESEARCH ARTICLE

FastTree 2 – Approximately Maximum-Likelihood Trees for Large Alignments

Morgan N. Price , Paramvir S. Dehal, Adam P. Arkin

Published: March 10, 2010 • <https://doi.org/10.1371/journal.pone.0009490>



MOLECULAR BIOLOGY AND EVOLUTION

JOURNAL ARTICLE

IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era

Bui Quang Minh , Heiko A Schmidt, Olga Chernomor, Dominik Schrempf, Michael D Woodhams, Arndt von Haeseler, Robert Lanfear [Author Notes](#)

Molecular Biology and Evolution, Volume 37, Issue 5, May 2020, Pages 1530–1534,
<https://doi.org/10.1093/molbev/msaa015>

Published: 03 February 2020



CITATIONS



VIEWS



ALTMETRIC



Phylogenetics

JOURNAL ARTICLE

IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era

Bui Quang Minh , Heiko A Schmidt, Olga Chernomor, Dominik Schrempf, Michael D Woodhams, Arndt von Haeseler, Robert Lanfear [Author Notes](#)

Molecular Biology and Evolution, Volume 37, Issue 5, May 2020, Pages 1530–1534,
<https://doi.org/10.1093/molbev/msaa015>

Published: 03 February 2020

JOURNAL ARTICLE

RAXML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies

Alexandros Stamatakis [Author Notes](#)

Bioinformatics, Volume 30, Issue 9, May 2014, Pages 1312–1313, <https://doi.org/10.1093/bioinformatics/btu033>

Published: 21 January 2014 [Article history ▾](#)

FastTree 2 – Approximately Maximum-Likelihood Trees for Large Alignments

Morgan N. Price , Paramvir S. Dehal, Adam P. Arkin

Published: March 10, 2010 • <https://doi.org/10.1371/journal.pone.0009490>



Phylogenetics



10 bins / MAGs

FastTree
2

CPU: 02
RAM: 64MB
Time: 4 seconds



10 bins / MAGs

CPU: 22
RAM: 80MB
Time: 03:56:34

Pros: Speed, memory efficiency, easy to use

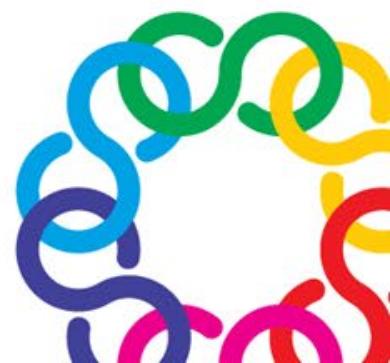
Pros: Provides better accuracy due to its advanced optimization methods (**Maximum Likelihood Optimization**)

Cons: It may not always produce the most accurate phylogenetic trees compared to more computationally intensive methods.

More customizations and flexibilities

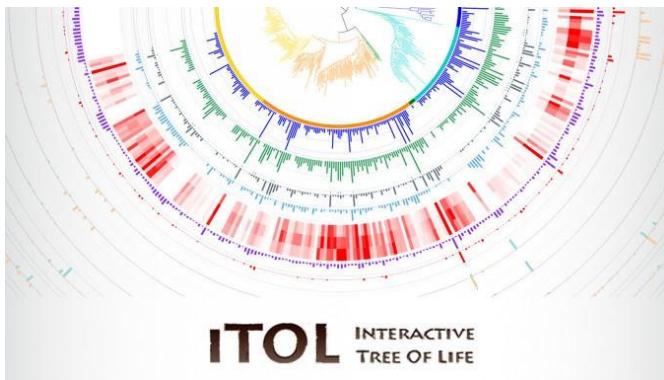
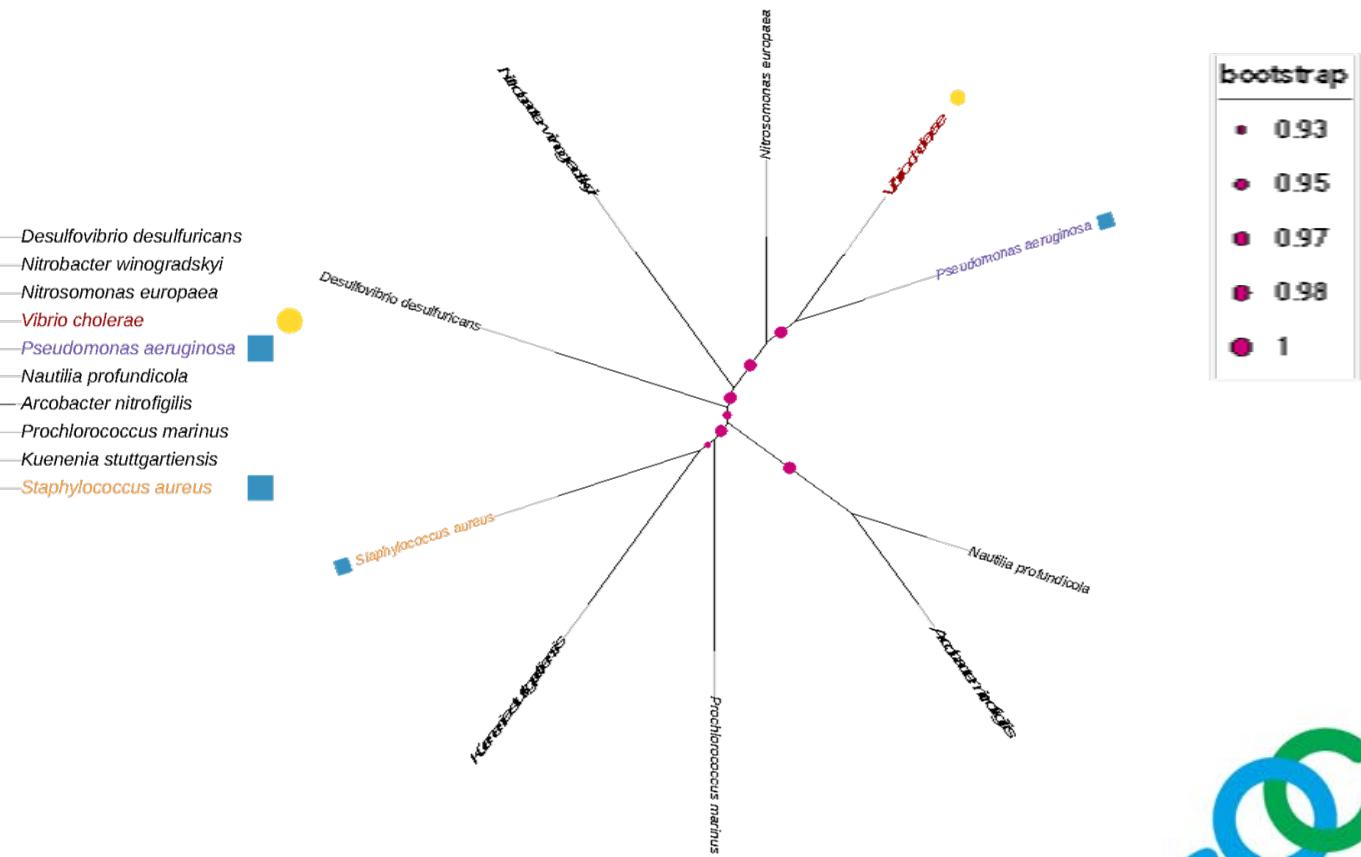
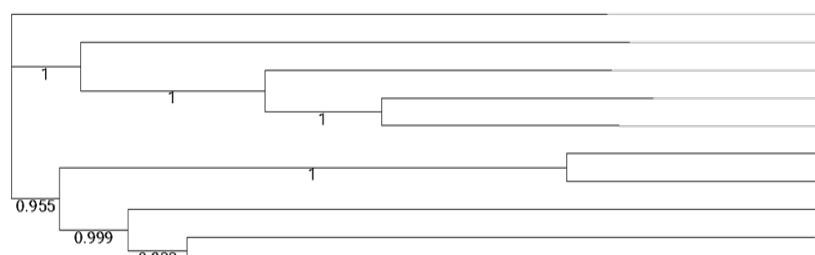
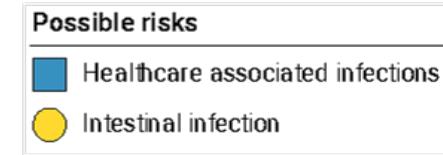
Limited customization and flexibility

Cons: Require more computational resources and time



Phylogenetics

Tree file: (((bin_3.filtered:0.4163372898,(bin_5.filtered:0.2927075220,bin_9.filtered:0.3087018139):0.1195543871)



Task: Build a phylogenetic tree

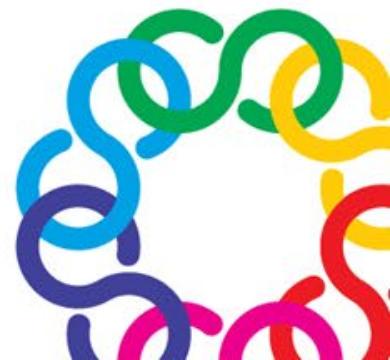
[Go to Github MGSS webpage](#)

Tasks:

- ✓ • Taxonomic classification of bins using GTDB-Tk
 - Obtain MSA from GTDB-Tk output
 - Display phylogenetic tree using iTOL
 - Annotate your tree



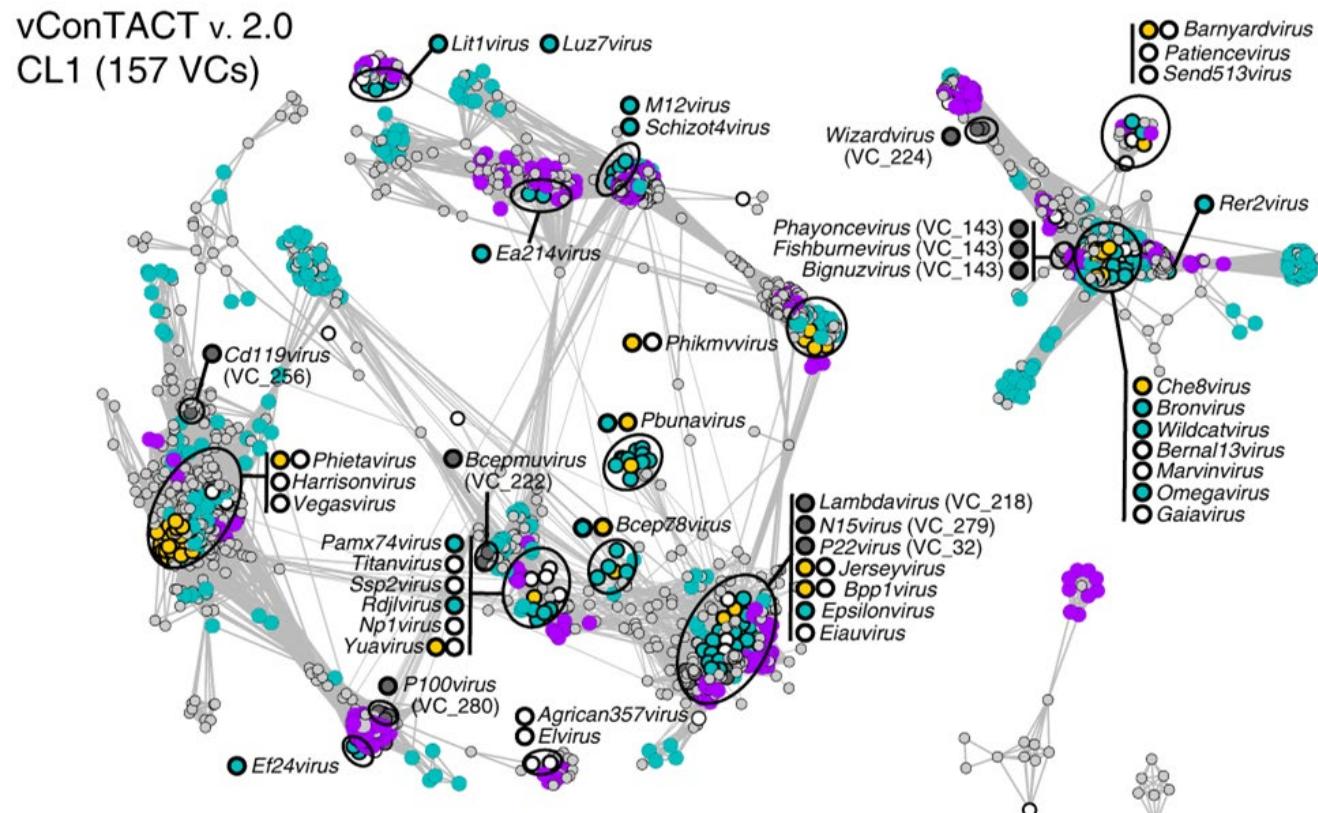
Virus taxonomy



DNA virus taxonomic prediction

vConTACT v.2.0 (<https://bitbucket.org/MAVERICLab/vcontact2/src/master/>)

- Clustering-based: guilt-by-contig-association taxonomic prediction
- Reference database (Viral RefSeq) + identified viral contigs



Viruses - workshop workflow

- Viral identification: **VirSorter2**
- (OPTIONAL: Dereplication across multiple assemblies)
- Viral QC: **CheckV**
- (OPTIONAL: Viral taxonomy and gene-sharing network: **vConTACT2** and **Cytoscape**)
- Viral gene prediction and annotation: **DRAM-v**



Task: QC and taxonomic classification

[Go to Github MGSS webpage](#)

Tasks:

- ✓ • Identifying viral contigs using VirSorter2
- ✓ • QC of viral contigs using CheckV
 - Examine viral output files from VirSorter2 and CheckV
 - Taxonomic classification of viruses using vContact2



Gene prediction



Gene prediction

- Genome annotation is the process of attaching biological information to sequences
- It consists of three main steps:
 - Gene prediction
 - Prediction of protein sequences
 - Functional annotation: Attaching biological information to these elements



Gene prediction

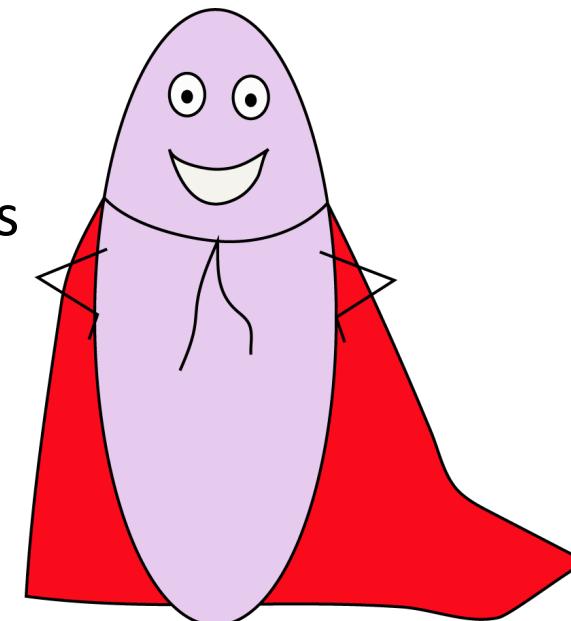
Aim:

- To identify regions of genomic DNA that encode putative genes present in high quality genomes

About 1/1000th of a human genome in size,
but with only 1/10th less coding DNA sequence
⇒ 100 x more power packed!!!

Prokaryote genomes:

- High gene density
- Genes = continuous stretches of coding DNA
- Absence of introns in the protein coding regions



Gene prediction

Gene finding algorithms for prokaryotes

- Homology:
 - Search by sequence similarity to homologous sequences
 - Based on the assumption that functional regions are more conserved evolutionarily than non-functional regions
- *Ab initio*:
 - Search by content: find genes by statistical properties that distinguish protein-coding DNA from non-coding DNA
 - Search by signals/sites, e.g. promoters, start and stop codons



Gene prediction

Homology: Sequence similarity searches

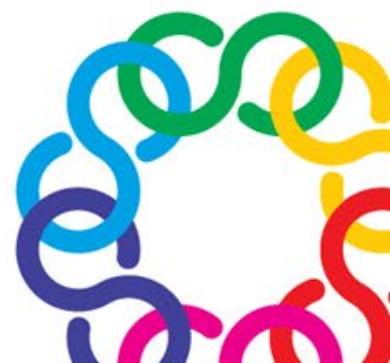
- Finding similarity in gene sequences between expressed sequence tags (ESTs), proteins, or other genomes to the input genome
- Local alignment:
 - BLAST family tools: <https://blast.ncbi.nlm.nih.gov/Blast.cgi>
 - Global alignment
 - GeneWise: <https://www.ebi.ac.uk/jdispatcher/psa/genewise>



Gene prediction

***Ab initio* search by content:** Markov Model Based Algorithms

- Most widespread algorithms for gene finding in prokaryotes are based on Markov Models
- Aim is to capture compositional differences among coding regions, “shadow” coding regions (coding on the opposite DNA strand) and non-coding DNA



Gene prediction

Markov Model Based Algorithms: Glimmer

- <https://ccb.jhu.edu/software/glimmer/index.shtml>
- Interpolated Markov model (IMM) DNA discriminator
- Log-likelihood that a given interval on a DNA sequence was generated by a model of coding versus non-coding DNA



Gene prediction

Markov Model Based Algorithms: GeneMark/GeneMarkHMM/MetaGeneMark

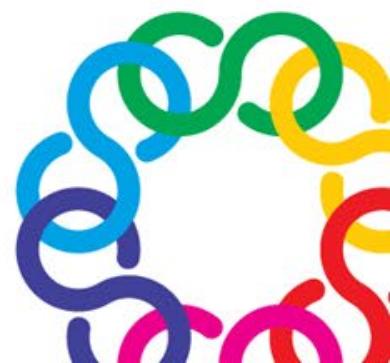
- <https://genemark.bme.gatech.edu/>
- GeneMark is a family of gene prediction tools
- Genomic sequences can be analysed either by the self-training program GeneMarkS-2 (sequences >50 kb) or using Heuristic Models by GeneMark.hmm
- Pre-trained model parameters are available for many species
- Metagenomics sequences can be analysed with MetaGeneMark



Gene prediction

Prodigal (PROkaryotic Dynamic Programming Genefinding ALgorithm)

- <https://github.com/hyattpd/Prodigal>
- Based on Dynamic Programming, not Markov Models
- Gene-finding algorithm for prokaryote genomes developed to predict translation initiation sites more accurately.
- High accuracy in high GC content genomes
- Tends to predict longer genes rather than more genes (minimising number of false positives)



Gene prediction

Prodigal for metagenomics:

- Use anon (meta) mode with metagenomic data (or short sequence data)
 - Copes with diverse genomes
 - Unlike normal mode, it does not attempt to study the input sequence, and predict based on these assumptions
 - Uses pre-calculated training files, and predicts genes based on the best results
- Alternatively, use normal mode on each individual genome bin

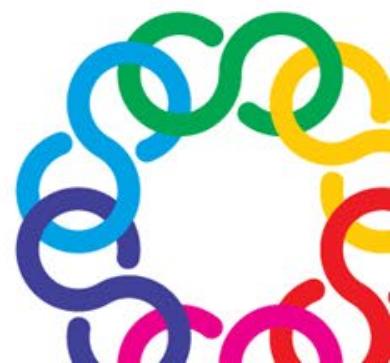


Gene prediction

Prodigal for metagenomics:

Caveat: unusual genetic codes

- First uses genetic code 11 (stop codons TAA, TGA, TAG)
 - Standard for prokaryotes and plant plastid
- If genes are too short, uses alternative code 4 (TGA not a stop codon)
 - Mycoplasma/Spiroplasma (also mold and protozoa)
- Will not try code 25 (must manually select), but will issue warning if genes are short
 - Candidate Phyla Radiation/Patescibacteria and Candidate Division SR1 (uncultured aquatic/host-associated microbiome)



Gene prediction

Prodigal for metagenomics:

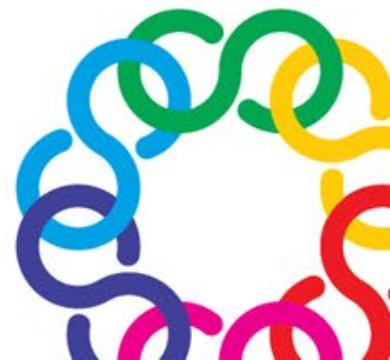
- Important note:
 - Prodigal predicts coding DNA sequence ONLY
 - Provides nucleic acid (.fna) and amino acid (.faa) files
 - **DOES NOT** identify other features (e.g. rRNA, tRNA)
 - Combine with other prediction tools



Gene prediction

Predicting RNA features and non-coding regions:

- MeTaxa2: predicts ribosomal RNA sequences in a genome
- Aragorn: predicts tRNA and tmRNA sequences



Gene prediction

Predicting protein coding sequences in unassembled (short) reads

- FragGeneScan:
 - Tuning parameters for short sequences (and hence incomplete genes)
 - Model sequence error

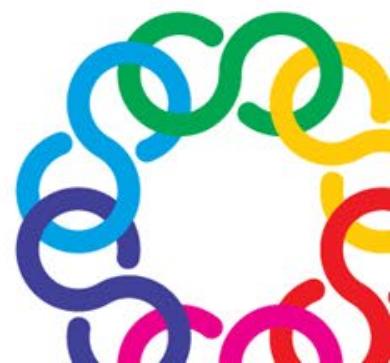


Task: Gene prediction

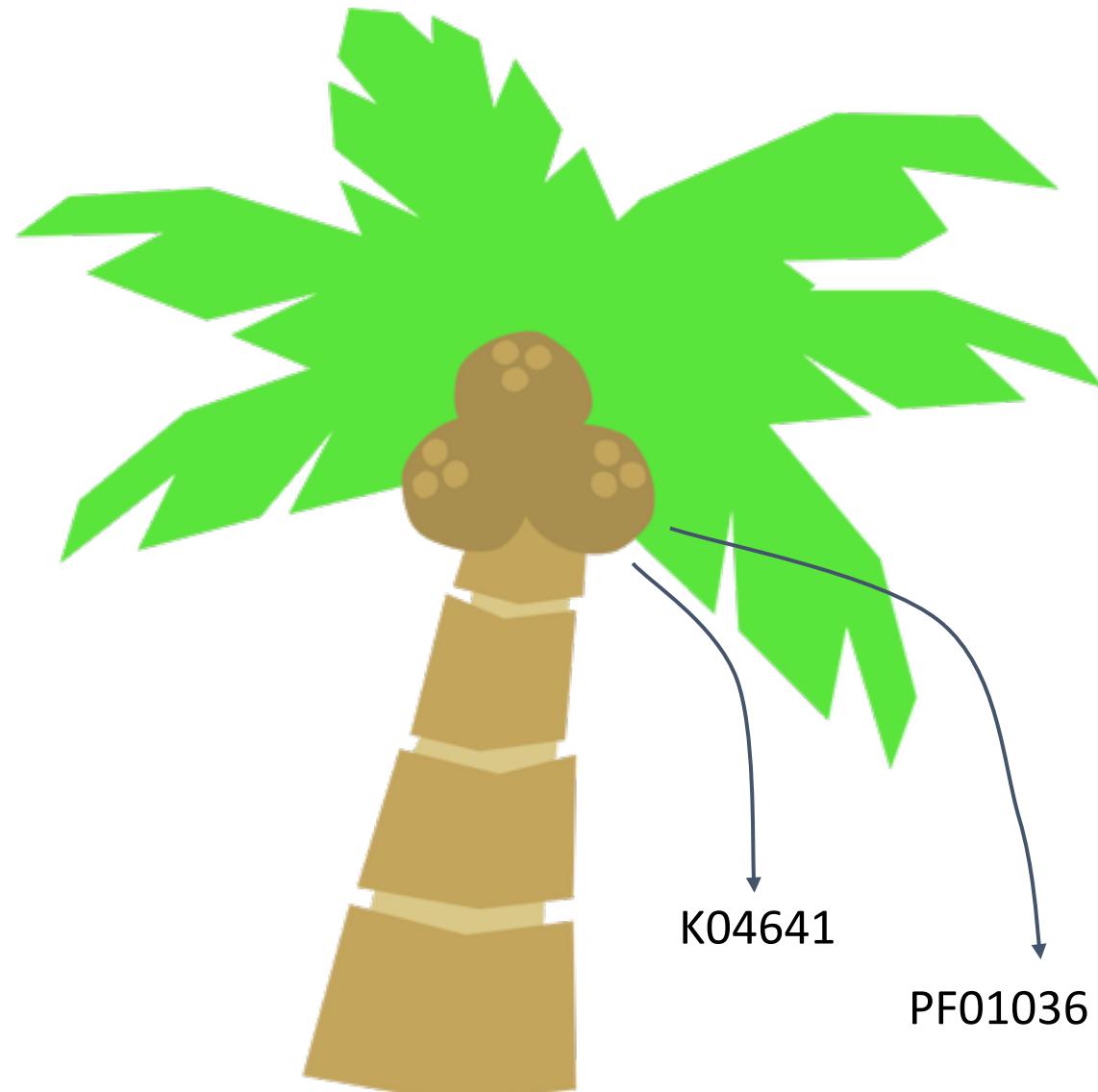
[Go to Github MGSS webpage](#)

Tasks:

- Predict open reading frames and protein sequences



Gene annotation (part 1)



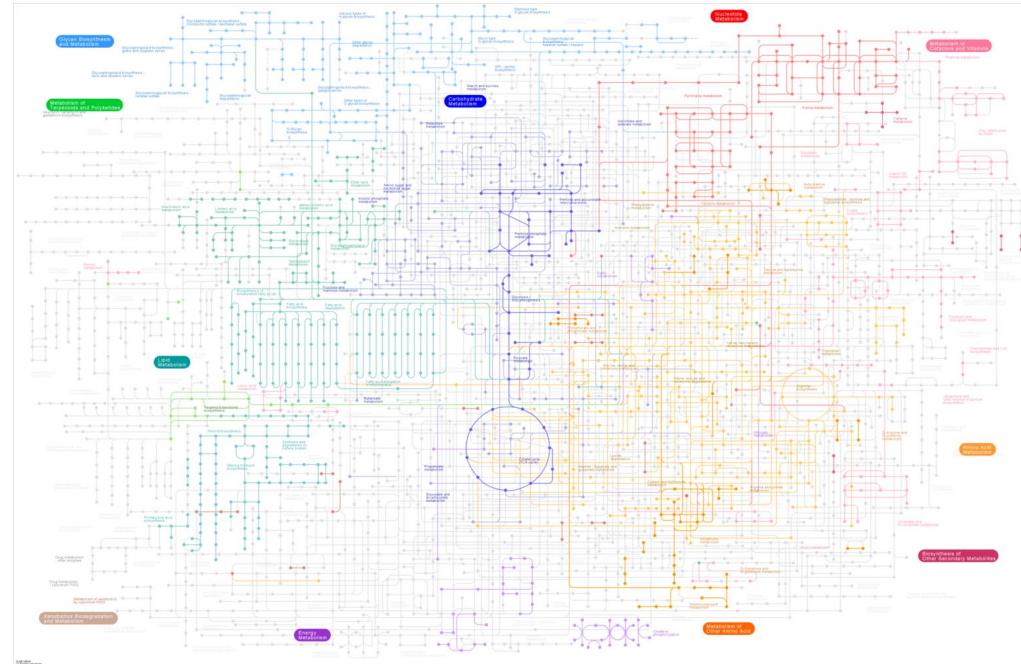
“You exist in the context of
all in which you live and what
came before you”

KH



Gene annotation

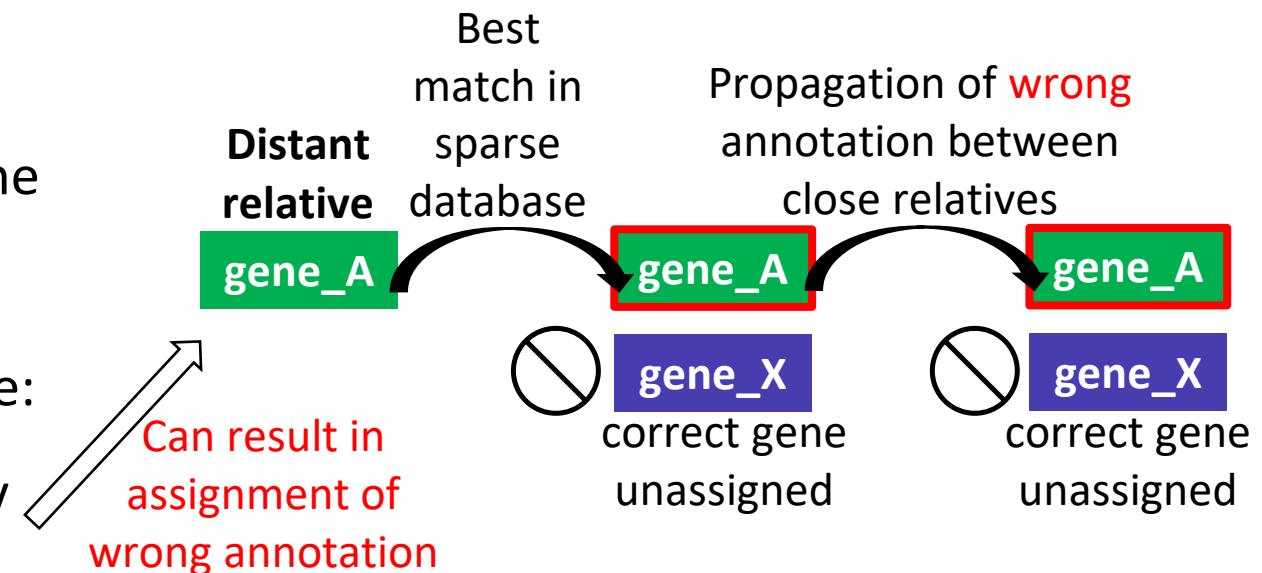
- Genome annotation attempts to predict gene function
- Predicted genes or protein sequences are compared against a curated set of reference sequences for which function is known, or is strongly suspected



Gene annotation

Caveat:

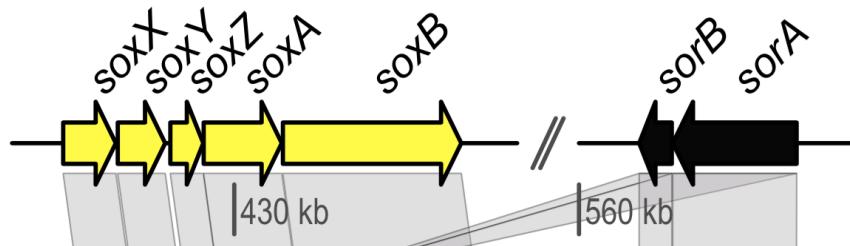
- Annotations are dependent on the reference database
- Environmental genomes can have:
 - Genes with distant homology matches to unrelated taxa
 - Large numbers of “hypothetical” gene annotations (= genes of unknown function)



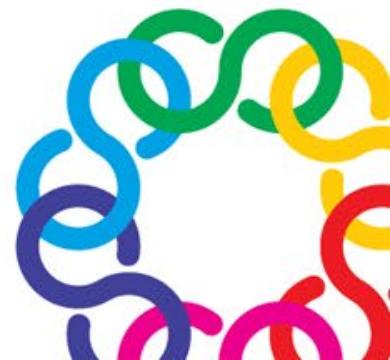
Gene annotation

Caveat:

- Annotations are “**advice**”
- Automated annotations often need to be manually curated
- Interrogate if: expected functional gene is missing from annotations
- Gene synteny is a useful for missing gene discovery, e.g.:
 - check genes co-located in operons for putative functions
 - check for operon truncation (due to contig break)



(Handley et al., 2014, Environmental Microbiology)



Gene annotation

There are two main ways to perform gene annotation with protein sequences:

- BLAST-like gene annotation
- Domain annotation



Gene annotation

BLAST-like gene annotation

- Pairwise local alignment between the gene of interest (query sequence) and the sequences in the database (target sequence)
- Tools:
 - BLAST: web-based and stand alone (usually too slow for metagenomics)
 - USEARCH (64-bit): fast (**subscription needed**)
 - Diamond: fast

Screenshot of a BLAST search results page. The top navigation bar includes 'Descriptions', 'Graphic Summary', **Alignments**, and 'Taxonomy'. The 'Alignments' tab is selected. Below it, the 'Alignment view' dropdown is set to 'Pairwise'. A green banner at the top indicates '100 sequences selected'. The main content area shows a single hit for 'amidase [Gemmatimonadetes bacterium]'. The sequence ID is MBB27982.1, with a length of 447 and one match found. A link to 'See 1 more title(s)' is present. The alignment details show a range from 182 to 446. The 'Related Information' section lists 'Identical Proteins - Identical proteins to MBB27982.1'. The alignment sequence data is presented in a table with columns for Score, Expect, Method, Identities, Positives, and Gaps.

Score	Expect	Method	Identities	Positives	Gaps
268 bits(686)	3e-84	Compositional matrix adjust.	139/265(52%)	175/265(66%)	0/265(0%)

Range 1: 182 to 446 GenPept Graphics ▾ Next Match ▲ Previous Match

Query	1	MGLKPTFGRISLRGILPVSYELDHPGPFTRSVADAALQCLAGKDPPLDPLSADVPVDI +GLKP <small>T</small> GR+S+ G++PVS+ LDHPGP T SV DAA ILQ +AG DP DPLSA
Sbjct	182	VGLKP <small>T</small> LGRVSVHGVVPSFNLDHPGPLTLSVGDARILQVIAGYDPKDPPLSASETTTI
Query	61	RIEPLSRPPRPGIVRTYYPPDNNADETMRAATDDAIERLASEGAEFTDVHMPGSAELHEN PL RPPR+G + Y+ +ADE M +AT AIE L GAE ++ MP SF LHE
Sbjct	242	TPRPLPRPPRIGHLVGYFREQADEMSSATQRAIECLQLAGACIELEMPDSFGCLHEN
Query	121	ALLLAVGAANVLDERYVAHRRDAFPPLCEIMERGRSAGAVDYARARRHQISFKSEVLA +++ A DE++ HR+ +PP L +M+ G + AV YA AR+HQI F+ ++ +
Sbjct	302	RIIMVSEGAAYHDEQFGLHRNEYPPGLRSLMDEGLATSAVTYANARKHQIDFRLQIQSI
Query	181	EGVDLILLTPATPTPAPSGLTSTGNPAFNPSWSYAGLPTIVLPAACSSDGLPAGIQLVAF +D+LLTPAT TPAP L STGNPAFNPSWSY GLPTI LP GLPA IQLV
Sbjct	362	RDLIDLITTPATLTPAKTLESTGNPAFNPSWSYCGLPTISLPVEVGESGLPAAIQLVGE
Query	241	FAEIRLLTVSAWCETRLEWNRTSSI 265 F+E RLL+++ WCE L WN P +
Sbjct	422	FSESRLLSIAQWCEQVLGWNHKPEL 446



Gene annotation

HMM-profiling of domains:

- Considers the query sequences as a collection of independently functioning protein folding domains
- Uses database of Hidden Markov models built from a collection of proteins that share a common domain
- Profiles build from statistical map of the
 - amino acid transitions (from position to position),
 - variations (differences at a position),
 - insertions/deletions between positions
- Tools: HMMer software (<http://hmmer.org/>)



Gene annotation

Common functional databases

- KEGG (Kyoto Encyclopedia of Genes and Genomes) (<https://www.kegg.jp>)
 - Very popular, each entry is well annotated, and often linked into “Modules” or “Pathways”
(Full access now requires a license fee)
- COGs (Clusters of Orthologous Groups of proteins) (<https://www.ncbi.nlm.nih.gov/COG/>)
 - Classify proteins from completely sequenced genomes on the basis of the orthology concept
- PFAM (<https://www.ebi.ac.uk/interpro/>)
 - Focused more on protein domains based on hidden Markov models
- TIGRfam (https://www.ncbi.nlm.nih.gov/genome/annotation_prok/tigrfams/)
 - Database of protein family definitions based on hidden Markov models



Gene annotation

Common functional databases

- The PANTHER (Protein ANalysis THrough Evolutionary Relationships) Classification System (<http://pantherdb.org>)
 - Proteins are classified according to Family and subfamily, molecular function, biological process and pathway
- UniRef (UniProt Reference Clusters) and SwissProt (<https://www.uniprot.org/>)
 - Protein clustering at different levels (e.g. UniRef100, UniRef90, UniRef50)
 - SwissProt is manually curated and experimentally confirmed
- BioCyc (<https://biocyc.org>) and MetaCyc (<https://metacyc.org>) (Subscription required)
 - BioCyc contains organism specific pathway/genome Databases + software tools
 - MetaCyc contains pathways shared between multiple organisms



Gene annotation

Databases for specific functions

- dbCAN3 (<https://bcb.unl.edu/dbCAN2/>) [Carbohydrate-active enzymes, FASTA + HMM]
 - Web server for automated annotation with downloadable databases
- TCDB (<https://tcdb.org/>) [Transporters, FASTA + HMM]
 - Classification system with extensive literature summaries and specialist prediction tools
- MEROPS (<https://www.ebi.ac.uk/merops/>) [Peptidases & inhibitors, FASTA]
 - Also includes peptidase inhibitors
- CARD (<https://card.mcmaster.ca/home>) [Antimicrobial resistance, FASTA]
 - Extensive, curated, and updated AMR gene collection
(Free for academic use, license required for commercial use)

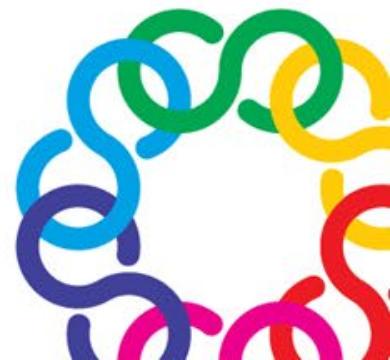


Gene annotation

Interproscan (<https://interproscan-docs.readthedocs.io/en/latest/index.html>)

(Jones et al. 2014, Bioinformatics, 30(9))

- Protein family annotation tool mainly using HMMs
- Annotates across 14 HMM databases
- Uses curated thresholds of member databases for functional assignments
 - E-values reported but often redundant
- Unifies functions and signatures across multiple databases
- Also has additional analyses for protein localisation
(user must possess license and database files)



Gene annotation

Distilling and Refining Annotations of Metabolism

(DRAM; Shaffer et al. 2020. Nucleic Acids Research 48(16))

- Tool for gene prediction and gene annotation of MAGs (DRAM-v for viruses)
 - Functional annotation using KEGG (if provided), UniRef 90, MEROPS, Pfam, dbCAN and VOGDB. tRNAs and rRNAs also detected
 - Genome annotations to metabolic functions in three levels:

1. RAW

Each gene nucleotide and amino acid sequence with annotations

2. DISTILLATE

Taxonomy (GTDB-tk), quality statistics (checkM), and key metabolisms summarized by genome

3. LIQUOR

Genome metabolisms classified by key functional gene, with gene FASTAs output

- **Data compiler:** checkM and GTDB-tk taxonomy summary

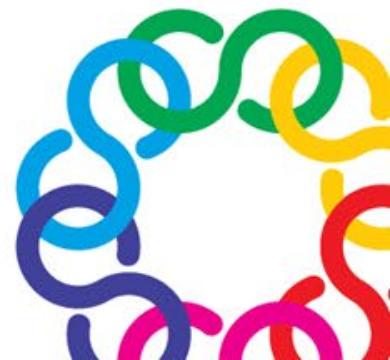


Task: Gene annotation

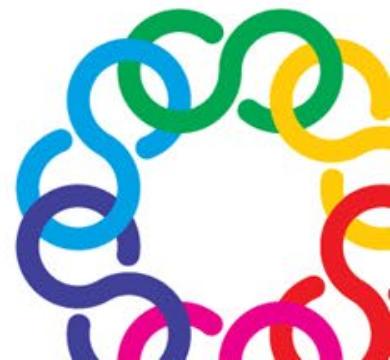
[Go to Github MGSS webpage](#)

Tasks:

- Gene annotation using DIAMOND and HMMER3
- Signal peptide prediction using SignalP6



Online resources and data analysis



Gene annotation

Web-based annotation tools:

- Web BLAST (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>)
- KEGG Automatic annotation and KEGG mapping service
 - BLAST-Koala: BLAST search (<https://www.kegg.jp/blastkoala/>)
 - GHOST-Koala: GHOSTX search (<https://www.kegg.jp/ghostkoala/>)
 - KofamKOALA: HMM search (<https://www.genome.jp/tools/kofamkoala/>)
- IMG/M (The Integrated Microbial Genomes and Microbiomes)
(<https://img.jgi.doe.gov>)



Gene annotation

KEGG: <https://www.genome.jp/kegg/kegg2.html>

KEGG PATHWAY Database
Wiring diagrams of molecular interactions, reactions and relations

Menu PATHWAY BRITE MODULE KO GENES LIGAND NETWORK DISEASE DRUG DBGET
Select prefix Enter keywords Go Help
map Organism [New pathway maps | Update history]

Pathway Maps
KEGG PATHWAY is a collection of manually drawn pathway maps representing our knowledge on the molecular interaction, reaction and relation networks for:

1. Metabolism
Global/overview Carbohydrate Energy Lipid Nucleotide Amino acid Other amino Glycan Cofactor/vitamin Terpenoid/PK Other secondary metabolite Xenobiotics Chemical structure

2. Genetic Information Processing

3. Environmental Information Processing

4. Cellular Processes

5. Organismal Systems

6. Human Diseases

7. Drug Development

KEGG PATHWAY is the reference database for pathway mapping in **KEGG Mapper**.

Pathway Identifiers
Each pathway map is identified by the combination of 2-4 letter prefix code and 5 digit number (see **KEGG Identifier**). The prefix has the following meaning:
map manually drawn reference pathway
ko reference pathway highlighting KOs
ec reference metabolic pathway highlighting EC numbers
rn reference metabolic pathway highlighting reactions
<org> organism-specific pathway generated by converting KOs to gene identifiers
and the numbers starting with the following:
011 global map (lines linked to KOs)
012 overview map (lines linked to KOs)
010 chemical structure map (no KO expansion)
07 drug structure map (no KO expansion)
other regular map (boxes linked to KOs)
are used for different types of maps.

1. Metabolism

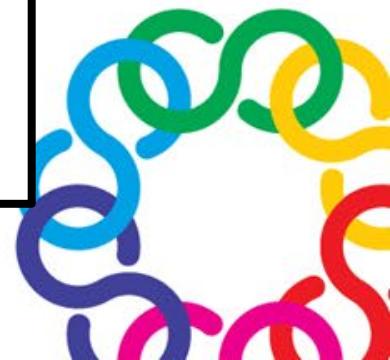
1.0 Global and overview maps
01100 Metabolic pathways
01110 Biosynthesis of secondary metabolites
01120 Microbial metabolism in diverse environments
01130 Biosynthesis of antibiotics
01200 Carbon metabolism
01210 2-Oxocarboxylic acid metabolism
01212 Fatty acid metabolism
01230 Biosynthesis of amino acids
01220 Degradation of aromatic compounds

KEGG Pathway Maps
[Brite menu | Download htext | Download json]
KEGG pathway maps Go

One-click mode

Metabolism

- Global and overview maps
 - 01100 Metabolic pathways
 - 01110 Biosynthesis of secondary metabolites
 - 01120 Microbial metabolism in diverse environments
 - 01130 Biosynthesis of antibiotics
 - 01200 Carbon metabolism
 - 01210 2-Oxocarboxylic acid metabolism
 - 01212 Fatty acid metabolism
 - 01230 Biosynthesis of amino acids
 - 01220 Degradation of aromatic compounds
- Carbohydrate metabolism
- Energy metabolism
 - 00190 Oxidative phosphorylation
 - 00195 Photosynthesis
 - 00196 Photosynthesis - antenna proteins
 - 00710 Carbon fixation in photosynthetic organisms
 - 00720 Carbon fixation pathways in prokaryotes
 - 00680 Methane metabolism
 - 00910 Nitrogen metabolism
 - 00920 Sulfur metabolism
- Lipid metabolism
- Nucleotide metabolism
- Amino acid metabolism
- Metabolism of other amino acids
- Glycan biosynthesis and metabolism
- Metabolism of cofactors and vitamins
- Metabolism of terpenoids and polyketides
- Biosynthesis of other secondary metabolites
- Xenobiotics biodegradation and metabolism
- Chemical structure transformation maps



Gene annotation

KEGG: <https://www.genome.jp/kegg/kegg2.html>

KEGG - Table of Contents

KEGG2 PATHWAY BRITE MODULE KO GENES COMPOUND NETWORK DISEASE DRUG

Search references cited in KEGG PMID, DOI, author, title, journal Go

Number of references (2024/8/27)
total 73,100 pathway 6,683 ko 28,761 glycan 931 network 2,647
brite 445 genome 7,194 reaction 2,082 variant 1,996
module 1,107 agenes 3,065 enzyme 16,187 disease 11,070

Search KEGG for Go

Data-oriented entry points

KEGG databases

Category	Entry point	Database	Content	Classification
Systems information	KEGG PATHWAY	PATHWAY	KEGG pathway maps	Pathway maps
	KEGG BRITE	BRITE	BRITE hierarchies and tables	Brite hierarchies Brite tables
	KEGG MODULE KEGG RModule	MODULE	KEGG modules and reaction modules	Modules Reaction modules
Genomic information	KEGG ORTHOLOGY	KO	Functional orthologs	KO
	KEGG GENES KEGG SeqData	GENES	Genes and proteins	
	KEGG GENOME KEGG Virus	GENOME	Genomes of cellular organisms and viruses	Organisms Viruses
Chemical information	KEGG COMPOUND	COMPOUND	Metabolites and other small molecules	Compounds
	KEGG GLYCAN	GLYCAN	Glycans	
	KEGG REACTION	REACTION RCLASS	Biochemical reactions Reaction class	
	KEGG Enzyme	ENZYME	Enzyme nomenclature with sequence data	EC sequence data
Health information	KEGG NETWORK	NETWORK VARIANT	Disease-related network variations Human gene variants	Network variation maps
	KEGG DISEASE	DISEASE	Human diseases	Human diseases (ICD) Human diseases Infectious diseases
	KEGG DRUG	DRUG DGROUP	Drugs Drug groups	Drugs (ATC) Drugs (target) Antineutraceuticals

KEGG Pathway Maps

Option One-click mode Row border shading Pruning neighbor

Search ID search Join

Metabolism
Global and overview maps
Carbohydrate metabolism
Energy metabolism
00190 Oxidative phosphorylation
00195 Photosynthesis
00196 Photosynthesis - antenna proteins
00710 Carbon fixation by Calvin cycle
00720 Other carbon fixation pathways
00680 Methane metabolism
00910 Nitrogen metabolism
00920 Sulfur metabolism
Lipid metabolism
Nucleotide metabolism
Amino acid metabolism
Metabolism of other amino acids
Glycan biosynthesis and metabolism
Metabolism of cofactors and vitamins
Metabolism of terpenoids and polyketides
Biosynthesis of other secondary metabolites
Xenobiotics biodegradation and metabolism
Chemical structure transformation maps

Genetic Information Processing

Environmental Information Processing

Cellular Processes

Organismal Systems

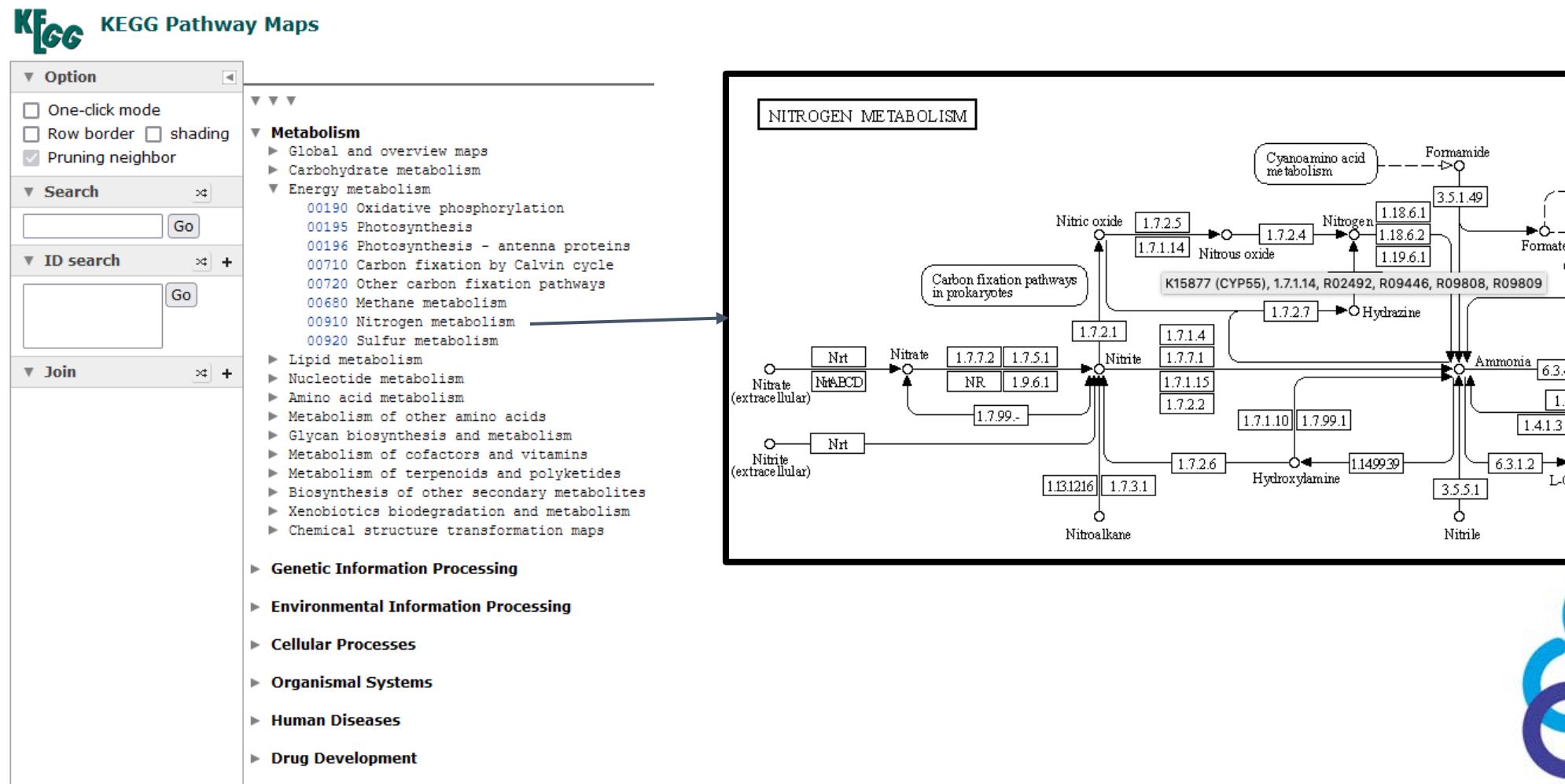
Human Diseases

Drug Development



Gene annotation

KEGG: <https://www.genome.jp/kegg/pathway.html>



Gene annotation

KEGG: <https://www.genome.jp/kegg/pathway.html>

KEGG Pathway Maps

[Brite menu | Download htext | Download json]

KEGG pathway maps

▼ ▼ ▼ One-click mode

▼ Metabolism

- ▼ Global and overview maps
 - 01100 Metabolic pathways
 - 01110 Biosynthesis of secondary metabolites
 - 01120 Microbial metabolism in diverse environments
 - 01130 Biosynthesis of antibiotics
 - 01200 Carbon metabolism
 - 01210 2-Oxocarboxylic acid metabolism
 - 01212 Fatty acid metabolism
 - 01230 Biosynthesis of amino acids
 - 01220 Degradation of aromatic compounds
- Carbohydrate metabolism
- Energy metabolism
 - 00190 Oxidative phosphorylation
 - 00195 Photosynthesis
 - 00196 Photosynthesis - antenna proteins
 - 00710 Carbon fixation in photosynthetic organisms
 - 00720 Carbon fixation pathways in prokaryotes
 - 00680 Methane metabolism
 - 00910 Nitrogen metabolism
 - 00920 Sulfur metabolism
- Lipid metabolism
- Nucleotide metabolism
- Amino acid metabolism
- Metabolism of other amino acids
- Glycan biosynthesis and metabolism
- Metabolism of cofactors and vitamins
- Metabolism of terpenoids and polyketides
- Biosynthesis of other secondary metabolites
- Xenobiotics biodegradation and metabolism
- Chemical structure transformation maps

→

ORTHOLOGY: K15877	
Entry	K15877 KO
Name	CYP55
Definition	fungal nitric oxide reductase [EC:1.7.1.14]
Pathway	ko00910 Nitrogen metabolism ko01100 Metabolic pathways ko01120 Microbial metabolism in diverse environments
Brite	KEGG Orthology (KO) [BR:ko00001] 09100 Metabolism 09102 Energy metabolism 00910 Nitrogen metabolism K15877 CYP55; fungal nitric oxide reductase Enzymes [BR:ko01000] 1. Oxidoreductases 1.7 Acting on other nitrogenous compounds as donors 1.7.1 With NAD+ or NADP+ as acceptor 1.7.1.14 nitric oxide reductase [NAD(P)+, nitrous oxide-fo K15877 CYP55; fungal nitric oxide reductase BRITE hierarchy
Other DBs	RN: R02492 R09446 R09808 R09809 GO: 0016966



Gene annotation

Identification of Biosynthetic Gene Clusters with antiSMASH

antiSMASH bacterial version Submit Bacterial Sequence Submit Fungal Sequence Submit Plant Sequence Download About Help Contact

Server status: working
Running jobs: 0
Queued jobs: 0
Jobs processed: 585148

Nucleotide input Results for existing job

Search a genome sequence for secondary metabolite biosynthetic gene clusters Load sample input Open example output

Incomplete RefSeq annotations
Dear antiSMASH users, it has come to our attention that a recent RefSeq reannotation again broke NRPS/PKS ORFs. If your results look weird, try uploading the corresponding GenBank record or a FASTA file.

Notification settings
Email address (optional): your@email.com

Data input
Upload file Get from NCBI NCBI acc # NCBI accession number of desired sequence

Detection strictness: relaxed
strict relaxed loose
• Detects well-defined clusters containing all required parts.
• Detects partial clusters missing one or more functional parts.

Extra features All off All on
 KnownClusterBlast ClusterBlast SubClusterBlast
 ActiveSiteFinder Cluster Pfam analysis Pfam-based GO term annotation

Submit

Please be considerate in your use of antiSMASH. Help us keep antiSMASH available for everybody by limiting yourself to 5 concurrent jobs. Need to run more? See the [antiSMASH install guide](#) for instructions for getting your own antiSMASH installation.

antibiotics & Secondary Metabolite Analysis SHell Version 4.2.0

Select Gene Cluster: Overview 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40

Identified secondary metabolite clusters

Cluster	Type	From	To	Most similar known
The following clusters are from record c00552_NODE_55.. (original name was: NODE_552_length_5782_cov_0.632764):				
Cluster 1	Other	1	5782	-
The following clusters are from record c00573_NODE_57.. (original name was: NODE_573_length_5701_cov_0.539986):				
Cluster 2	T3pks	1	5701	-
The following clusters are from record c00895_NODE_89.. (original name was: NODE_895_length_4918_cov_0.147126):				
Cluster 3	T1pks	1	4918	-
The following clusters are from record c02406_NODE_24.. (original name was: NODE_2406_length_3433_cov_0.091782):				
Cluster 4	Otherks	1	3433	-
The following clusters are from record c02525_NODE_25.. (original name was: NODE_2525_length_3365_cov_0.089712):				
Cluster 5	Nrps	1	3365	-
The following clusters are from record c02784_NODE_27.. (original name was: NODE_2784_length_3235_cov_0.116390):				
Cluster 6	Terpene	1	3235	-



<https://antismash.secondarymetabolites.org/>

NCBI Conserved Domain Search

Search for Conserved Domains within a protein or coding nucleotide sequence

Enter **protein** or **nucleotide** query as accession, gi, or sequence in [FASTA format](#). For multiple protein queries, use [Batch CD-Search](#).

MAINKHHITPMLDQLESGFWPSFISGIKRLRDEHPEERINKMTNDLLGQLEHSYETRKGYWKGGTVSVFQYGGGIIPRFSEVGHAFPEKFHTLRLVQPAGNHYSTMLRQMADEWKGSLVTFHQQTGNIMF1GTDTEQTHFFDEINDYGWLGGAGPCVRTAMSVCVGAAARCEMSCTNEQKAIRLLVNNFTDDVHRLPALPYKFKFKVSGCGNDQCNAVERADFAVIGTWRDDMNVNQDEFKAYGRKGRQHVIDNIITRCPTNALSINDDDSLEVNNKDCVRCMHCLNVPKALHPGDRGVTLIIGGKRTLKGDLMGTVVVVFKKLDTEEDWEELAEEIIDFWAENALEHERCGERMIERIGLNFLEGVGVEVDPMVNPNPRESSYIRMDGWDDEAEVKWFDRQAEAS|

OPTIONS

Search against database: CDD v3.17 - 52910 PSSMs

Expect Value threshold: 0.01000

Apply low-complexity filter:

Composition based statistics adjustment:

Force live search:

Rescue borderline hits: Suppress weak overlapping hits:

Maximum number of hits: 500

Result mode: Concise Standard Full

Submit Reset

Retrieve previous CD-search result

Request ID: Retrieve

References:

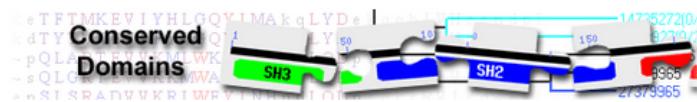
- Marchler-Bauer A et al. (2017), "CDD/SPARCLE: functional classification of proteins via subfamily domain architectures.", *Nucleic Acids Res.* **45**(D)200-3.
- Marchler-Bauer A et al. (2015), "CDD: NCBI's conserved domain database.", *Nucleic Acids Res.* **43**(D)222-6.
- Marchler-Bauer A et al. (2011), "CDD: a Conserved Domain Database for the functional annotation of proteins.", *Nucleic Acids Res.* **39**(D)225-9.
- Marchler-Bauer A, Bryant SH (2004), "CD-Search: protein domain annotations on the fly.", *Nucleic Acids Res.* **32**(W)327-331.

[Help](#) | [Disclaimer](#) | [Write to the Help Desk](#)
NCBI | NLM | NIH

Search nucleotide/protein sequence(s) for conserved domains

Individual search: <https://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi>
Batch: <https://www.ncbi.nlm.nih.gov/Structure/bwrpsb/bwrpsb.cgi>





Conserved domains on [lcl|1]

NZ_JRAA01000001.1:c722683-721433 Solemya velum gill symbiont strain WH SV_sym_Scaffold_1, whole genome shotgun sequence



[Search for similar domain architectures](#) [?](#) [Refine search](#) [?](#)

List of domain hits

Name	Accession	Description	Interval	E-value
dsrA	TIGR02064	sulfite reductase, dissimilatory-type alpha subunit; Dissimilatory sulfite reductase catalyzes ...	31-1242	0e+00

sulfite reductase, dissimilatory-type alpha subunit; Dissimilatory sulfite reductase catalyzes the six-electron reduction of sulfite to sulfide, as the terminal reaction in dissimilatory sulfate reduction. It remains unclear however, whether thiosulfinate and thiosulfate serve as intermediate compounds to sulfide, or as end products of sulfite reduction. Sulfite reductase is a multisubunit enzyme composed of dimers of either alpha/beta or alpha/beta/gamma subunits, each containing a siroheme and iron sulfur cluster prosthetic center. Found in sulfate-reducing bacteria, these genes are commonly located in an unidirectional gene cluster. This model describes the alpha subunit of sulfite reductase. [Central intermediary metabolism, Sulfur metabolism]

Pssm-ID: 273948 [Multi-domain] Cd Length: 402 Bit Score: 667.31 E-value: 0e+00

	10	20	30	40	50	60	70	80		
1********		
Cdd:TIGR02064	11 LDQLESGPWPSFI	SGIKRLRDEPPEIRINKMTNDLLGOL	HSYETRKGYWKGGTVSVFQYGGGI	PPRFSEVGHAFPFESKE	90					
	1 LDQLEKGPWPSFV	SEIKKTAAYRADYQV	PVPDPELDLGVLELSYDERKTHW	KGGIVSVFQYGGGVIGRY	SQDGKEKPPGVAE	80				
	90	100	110	120	130	140	150	160		
1********		
Cdd:TIGR02064	91 FHTLRVQQPAGNH	YSTDMRLRQMADSWEKYKGS	LVTFHGQTCGNIMPICTDTEQ	TQHFDEINDYCWDLGAGPCV	RVTAMSC	170				
	81 FHTVRVAQPSGKF	YSTDYLRLQLCDVWEKYKGS	LTNFHGQTCGDIVFLGTQTP	QLQEIEFEELTNLGTDLGG	GSGSNLRTPE	160				
	170	180	190	200	210	220	230	240		
1********		
Cdd:TIGR02064	171 VGAARCEMSCTNEQ	KAHRLVNNFTDDVHRPA	LPYKFKPVSGCGND	QCNAVERADFAVIGTW	RDMVNQDEFKAYVGR	250				
	161 VGPARCEFACYDTLK	ACYELTM	YEYQDELHRPA	PKFKFSGCPND	CVAAIASDFAVIGTW	KDDIKV	DQEAVKAYIAG	240		
	250	260	270	280	290	300	310	320		
1********		
Cdd:TIGR02064	251 KGRQHVIDNIITRCPT	NALSLNDDDSLEVNN	KDCVRCMHCLNV	PKALHPGD	DRGVFTILIGGK	RTLKIGDLMGTVVVPPK	330			
	241 WGKFIDIEEVNNRCPT	KAISWDGSKELS	IDNRECVRCMHCI	INKMPKALHPGD	ERGVFTILIGGK	KAPILDGAQMCWVVVPPV	320			
	330	340	350	360	370	380	390	400		
1********		
Cdd:TIGR02064	331 kldT	EEEDWEELAEIIIDFWA	ENALEHERC	GEMIERIGLVNFLEG	VGVVEVDPNMVNNP	RESSYIRMDGWDE	AEAVKWF	410		
	321 --EAEYPDE	IKELEK	KI	IIDWWDE	KGNRERI	GETIKRGLQK	PLEVIGIEPD	DPQMVKEPRTNPYIFF	KVDEVPGGWDA	398
	411 RQAE	414								
1*									
Cdd:TIGR02064	399 DIAE	402								

Conserved Domain Search

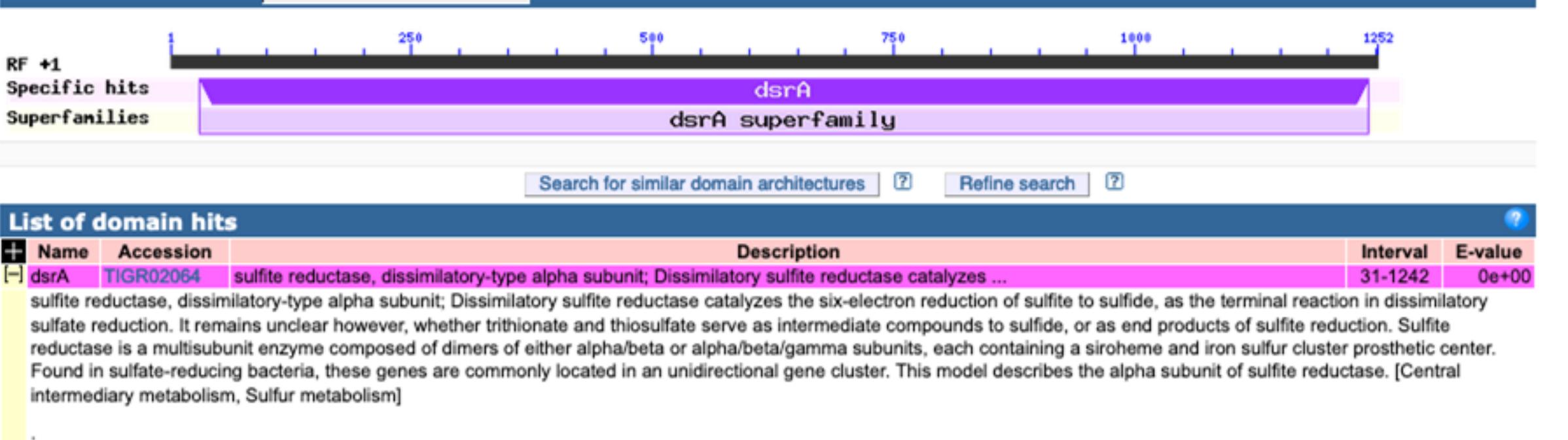
results:

dsrA gene of *Solemya velum* gill symbiont strain WH



References:

- Marchler-Bauer A et al. (2017), "CDD/SPARCLE: functional classification of proteins via subfamily domain architectures.", *Nucleic Acids Res.* 45(D)200-3.



250 260 270 280 290 300 310 320
*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
 1 251 KGRGVIDNRTI RCPTMKALSLNDDDSLVEVN KDKCVRCMHCN LNPVKKALHPGDRGVTLIGKKRKLIGK DLMCTGVVVFK 330
 Cdd:TTGR02064 241 WCKFDIENEVVN RCPMKW SDWGSKETLSDTRCVRCMHCNTK NPKKALHPGDRGVFTLIGKKAPILDAQMGVVVVFV 320

330 340 350 360 370 380 390 400
*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
 1 331 kIDTDEWEEDVVELAEIIDWVAEALAEHERGEMIERIGLVLNFPELEGVVEIDPNVNPNPRESSEYIRMDGWDEEAVKWF 410
 Cdd:TIGR02064 321 --EAEPDYEIKELKEVIIWDNEEGKNRERIGETKLGILOKPLEVIGEIPDPDPVKCPEPRTPNYIYFFKWDVEPGCGDA 398

1 411 RQAE 414
Cdd-TIGR02064 399 DIAE 402



References:

Marchler-Bauer A et al. (2012). "CCD/SPARCLE: functional classification of proteins via subfamily-domain architectures." Nucleic Acids Res 40(D):D200-3.

Gene annotation

Metacyc: experimentally curated metabolic pathways

 **METACYC**
A member of the BioCyc database collection

Two-day Introduction to Pathway Tools Tutorial
Jan 16-17, 2020
Early Registration by Dec 12

Sites ▾ | Search ▾ | Genome ▾ | Metabolism ▾ | Analysis ▾ | SmartTables ▾ | Help ▾

Search Results for **dsra** using database **MetaCyc** what is this?

[Genes \(3\)](#) | [Proteins \(3\)](#) | [EC Numbers \(2\)](#)

Genes Gene/Gene Product pages contain: chromosomal location of gene; depiction of its operon; link to genome browser; detailed summaries; complexes); cofactors, activators, and inhibitors (for enzymes), depiction of regulon (for transcriptional regulators), protein features.

- [dsrA - *Allochromatium vinosum*](#)
- [dsrA - *Desulfovibrio gigas*](#)
- [dsrA - *Archaeoglobus fulgidus*](#)

 [Login](#) to turn into a SmartTable.

Proteins Gene/Gene Product pages contain: chromosomal location of gene; depiction of its operon; link to genome browser; detailed summaries; regulon (for transcriptional regulators), protein features.

- siroheme sulfite reductase, α subunit (*DsrA*) - *Allochromatium vinosum*
- sulfite reductase α subunit (*DsrA*) - *Desulfovibrio gigas*
- sulfite reductase, dissimilatory α subunit (*DsrA*) - *Archaeoglobus fulgidus*

 [Login](#) to turn into a SmartTable.

EC Numbers EC Number pages contain: links to reaction and enzymes associated with the EC number in this database, names, description.

<https://metacyc.org/>



Gene annotation

Metacyc: experimentally curated metabolic pathways

Two-day Introduction to Pathway Tools Tutorial Jan 16-17, 2020 Early Registration by Dec 12

Sites | Search | Genome | Metabolism | Analysis | SmartTables | Help

Search Results for **dsra** using database **MetaCyc** what is this?

Genes (3) | Proteins (3) | EC Numbers (2)

Genes Gene/Gene Product pages contain: chromosomal location of gene; depiction of its operon; complexes; cofactors, activators, and inhibitors (for enzymes); depiction of regulon (for transcriptional regulators).

dsra - Allochromatium vinosum
dsra - Desulfovibrio gigas
dsra - Archaeoglobus fulgidus

Proteins Gene/Gene Product pages contain: chromosomal location of gene; depiction of its operon; complexes; cofactors, activators, and inhibitors (for enzymes); depiction of regulon (for transcriptional regulators), protein features.

- siroheme sulfite reductase, α subunit (**DsrA**) - *Allochromatium vinosum*
- sulfite reductase α subunit (**DsrA**) - *Desulfovibrio gigas*
- sulfite reductase, dissimilatory α subunit (**DsrA**) - *Archaeoglobus fulgidus*

EC Numbers EC Number pages contain: links to reaction and enzymes associated with the EC number.

Sites | Search | Genome | Metabolism | Analysis | SmartTables | Help

gene **dsrA** protein **sulfite reductase α complex**
Desulfovibrio gigas

Accession ID G-385 (MetaCyc)

Reactions thiosulfate + 2 an oxidized unknown electron carrier + 3 H₂O → 2 sulfite + 2 a reduced unknown electron carrier + 3 H⁺ (catalyzed by complex a [DsrC]-trisulfide + an oxidized unknown electron carrier + 3 H₂O → sulfite + a [DsrC protein]-dithiol + a reduced unknown electron carrier + 3 H⁺)
trithionate + an oxidized unknown electron carrier + 3 H₂O → 3 sulfite + a reduced unknown electron carrier + 4 H⁺ (catalyzed by complex)

Pathways dissimilatory sulfate reduction I (to hydrogen sulfide)
dissimilatory sulfate reduction II (to thiosulfate)

Summary GO Terms (1) Reactions (3) References Show All

Subunit Composition [DsrA]₂
Component of dissimilatory sulfite reductase (extended summary available): [(DsrA)₂][(DsrB)₂]

Gene-Reaction Schematic

+ 1.8.99.5 : trithionate + an oxidized unknown electron c...
- 1.8.99.5 : a [DsrC]-trisulfide + an oxidized unknown el...
1.8.99.5 : thiosulfate + 2 an oxidized unknown electron...
1 2 Dg-dsrB
1 2 Dg-dsrA

Gene: dsrA G-385
Product: sulfite reductase α subunit, subunit of sulfite reductase α complex, dissimilatory sulfite reductase
Species: *Desulfovibrio gigas*

GO Terms:
Cellular Component: GO:0005829 - cytosol

Enzymatic activity: sulfite reductase (thiosulfate-forming) (dissimilatory sulfite reductase)

2 sulfite + 2 a reduced unknown electron carrier + 3 H⁺ → thiosulfate + 2 an oxidized unknown electron carrier + 3 H₂O

<https://metacyc.org/>

Gene annotation

Accurate classifier for hydrogenase sequences - HydDB

HydDB  Classify  Browse  Information Pages ▾

Classify

HydDB provides access to an accurate classifier for hydrogenase sequences and a curated database of hydrogenases by known type. The service is provided by the School of Biological Sciences, Monash University and the Bioinformatics Research Centre, Aarhus University.

Please cite! If you use HydDB for research, please cite the following paper: "Søndergaard D, Pedersen CNS, Greening C. **HydDB: A web tool for hydrogenase classification and analysis**. Sci Rep. 2016;6:34212. doi: 10.1038/srep34212.". The preprint is available on [bioRxiv](#). If you have any comments, corrections or questions contact [Chris Greening](#) or [Dan Søndergaard](#).

Classify

HydDB is unable to accurately check whether uploaded sequences correspond to hydrogenases or not. Instead, it is well-suited for functionally-predictive classification of known hydrogenases into different subgroups. Please ensure that all sequences that you upload correspond to catalytic subunits of hydrogenases (e.g. using conserved domain database and phylogenetic trees). Sequences that do not encode catalytic subunits of hydrogenases will still be classified, but the result may be wrong.

Sequences

Sequences File

Choose File no file selected

Instructions

To use the classifier to predict the type of one or more hydrogenases from sequence, either:

- paste your FASTA-formatted protein sequences into the text area, or
- upload a FASTA-formatted file with your protein sequences.

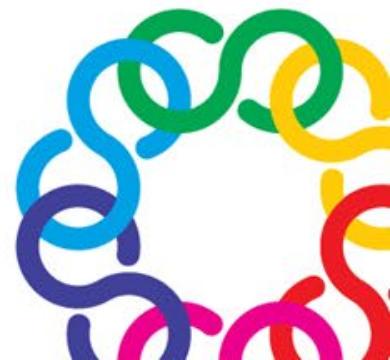
Press the "Submit" button to upload the sequences and begin the classification.

If you provided an e-mail address you will receive an e-mail when your job finishes or fails including a link to the results. You will also be able to download the results as a CSV file.

Only sequences encoding the catalytic subunits of hydrogenases will be classified, i.e. those binding the [NiFe]-centre (NiFe-hydrogenases), [FeFe]-centre (FeFe-hydrogenases), or [Fe]-centre (Fe-hydrogenases). Electron-transfer subunits, accessory proteins, and maturation factors cannot be classified by this service.

Limits

<https://services.birc.au.dk/hyddb/>



Gene annotation

- The **PSORT** family - prediction of protein localization sites in cells.
- Useful for making cell schematics!



Updates | Documentation | Resources | Contact

Submit a Sequence to PSORTb version 3.0.2

Based on a study last performed in 2010, PSORTb v3.0.2 is the most precise bacterial localization prediction tool available. PSORTb v3.0.2 has a number of [improvements](#) over PSORTb v2.0.4. Version 2 of PSORTb is maintained [here](#).

You can currently submit one or more Gram-positive or Gram-negative bacterial sequences or archaeal sequences in FASTA format (?). Copy and paste your FASTA-formatted sequences into the textbox below or select a file containing your sequences to upload from your computer. Web display mode is limited to the analysis of approximately 100 proteins. For larger analyses, either enter your email address in the form below (results of up to 5000 per submission returned by email) or for even larger analyses we can help you or you can download the standalone version.

See also:

- [Updates](#)
- [Precomputed genome results](#)
- [Limitations of PSORTb v.3.0](#)
- [PSORTb User's Guide](#)
- [Docker PSORTb web service](#) (what is [docker](#)?)
- [Download standalone PSORTb](#)
- [Docker standalone PSORTb](#) (what is [docker](#)?)

<http://psort.org/>

Choose an organism type (?): Required

Choose Gram stain (?): Required

Advanced Gram stain options (?): Required

Output format (?):

Show results (?):

Email address:

Copy and paste your FASTA sequences below

or upload from file: no file selected
(uploads limited to 50KB, approximately 100 proteins, in Web display mode, enter an email address to use email mode if you need to analyze more proteins)



Gene annotation

SignalP 6.0 (<https://services.healthtech.dtu.dk/service.php?SignalP>)

- Predicts signal peptides on amino acid sequences across kingdoms of life.
 - Uses Protein Language Models

sp_Q45071_XYND_BACSU_Arabinoxylan_arabinofuranohydrolase_OS_Bacillus_subtilis_Prediction: Lipoprotein signal peptide (Sec/SPII)

Submit data

Sequence submission: paste the sequence(s) and/or upload a local file

Info Protein sequences should be not less than 10 amino acids. The maximum number of proteins is 5000.
Info The long output format might timeout for more than 100 entries.

Mirror Use SignalP 6.0 on BioLib if this server is heavily loaded.

Enter protein sequence(s) in fasta format...

For example proteins [Click here](#)

For example proteins [click here](#)

Format directly from your local computer

Organism

Organism

- Eukarya
- Other
- "Eukarya" only predicts

Sec/SPI SPs.

[Submit](#) [Clear fields](#)

Output format

Output format.

- Long output
- Short output (no figures)

Model mode

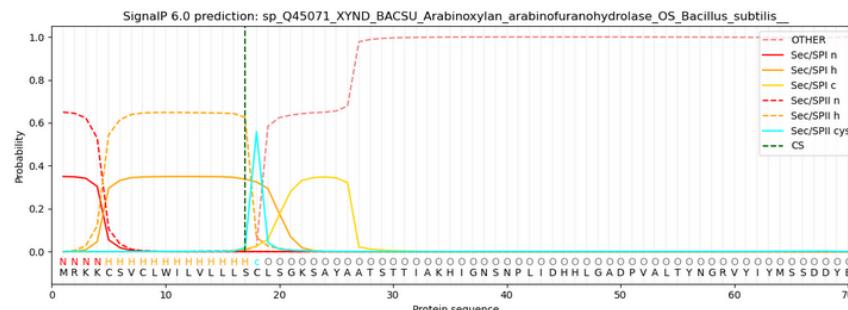
Mode
Easy

Fast
Slow

i The slow mode takes 6x longer to compute. Use when accurate region borders are needed.

Protein type	Other	Signal Peptide (Sec/SPI)	Lipoprotein signal peptide (Sec/SPII)	TAT signal peptide (Tat/SPI)	TAT Lipoprotein signal peptide (Tat/SPII)	Pilin-like signal peptide (Sec/SPIII)
Likelihood	0.0012	0.3394	0.6584	0.0004	0.0003	0.0003

Download: [PNG](#) / [EPS](#) / [Tabular](#)



Gene annotation

BRENDA: <https://www.brenda-enzymes.org/>

- Enzyme database
- Has reaction and pathway schematics
 - In-house pathways
 - KEGG pathways

Information on EC 3.2.1.23 - beta-galactosidase

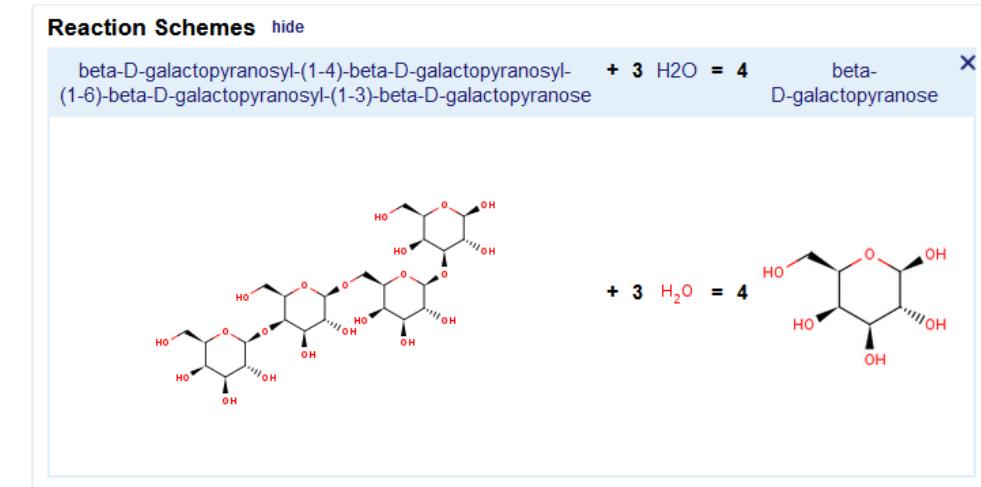
for references in articles please use BRENDA:EC3.2.1.23

EC Tree

- └ 3 Hydrolases
 - └ 3.2 Glycosylases
 - └ 3.2.1 Glycosidases, i.e. enzymes that hydrolyse O- and S-glycosyl compounds
 - └  3.2.1.23 beta-galactosidase

IUBMB Comments

Some enzymes in this group hydrolyse alpha-L-arabinosides; some animal enzymes also hydrolyse beta-D-fucosides and beta-D-glucosides; cf. EC 3.2.1.108 lactase.



[△ top](#) [print](#) [hide 7 entries](#) [Go to Pathway Search](#)

PATHWAY SOURCE ▲▼	PATHWAYS ▲▼
BRENDA	metabolism of disaccharids
KEGG	Galactose metabolism , Glycosaminoglycan degradation , Glycosphingolipid biosynthesis - ganglio series , Other glycan degradation , Sphingolipid metabolism
-	Sphingolipid metabolism



Summary of online resources

Resources to help interpret your data:

- KEGG: <https://www.genome.jp/kegg/pathway.html>
- BioCyc: <https://biocyc.org/>
- MetaCyc: <https://metacyc.org/>
- HydDB: <https://services.birc.au.dk/hyddb/>
- PSORT: <https://psort.hgc.jp/>
- SignalP: <https://services.healthtech.dtu.dk/service.php?SignalP>
- BRENDA: <https://www.brenda-enzymes.org/>
- ANNOTREE: <http://annotree.uwaterloo.ca/annotree/>



Task: Gene annotation

Tasks:

- View KEGG annotation in KEGG website

KEGG: <https://www.genome.jp/kegg/pathway.html>

- View functional distribution in ANNOTREE

KEGG PATHWAY Database
Wiring diagrams of molecular interactions, reactions and relations

Menu PATHWAY BRITE MODULE KO GENES LIGAND NETWORK DISEASE DRUG DBGET

Select prefix Enter keywords

[New pathway maps | Update history]

Pathway Maps

KEGG PATHWAY is a collection of manually drawn pathway maps representing our knowledge on the molecular interaction, reaction and relation networks for:

1. Metabolism
2. Genetic Information Processing
3. Environmental Information Processing
4. Cellular Processes
5. Organismal Systems
6. Human Diseases
7. Drug Development

KEGG PATHWAY is the reference database for pathway mapping in [KEGG Mapper](#).

Pathway Identifiers

Each pathway map is identified by the combination of 2-4 letter prefix code and 5 digit number (see [KEGG Identifier](#)). The prefix has the following meaning:

map	manually drawn reference pathway
ko	reference pathway highlighting KOs
ec	reference metabolic pathway highlighting EC numbers
rn	reference metabolic pathway highlighting reactions
<org>	organism-specific pathway generated by converting KOs to gene identifiers

and the numbers starting with the following:

011	global map (lines linked to KOs)
012	overview map (lines linked to KOs)
010	chemical structure map (no KO expansion)
07	drug structure map (no KO expansion)
other	regular map (boxes linked to KOs)

are used for different types of maps.

1. Metabolism

1.0 Global and overview maps

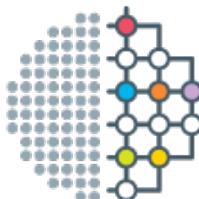
01100 Metabolic pathways
01110 Biosynthesis of secondary metabolites
01120 Microbial metabolism in diverse environments
01130 Biosynthesis of antibiotics
01200 Carbon metabolism
01210 2-oxocarboxylic acid metabolism





Delving into plastisphere metagenomes: The joys of the unassembled metagenomic analyses

Jessica Wallbank



NeSI
New Zealand eScience
Infrastructure


AIM²
Aotearoa Impacts &
Mitigation of Microplastics

Outline of my project

- Multi-omics approach
 - Amplicon data
 - Metagenomic data
- 7 substrate types
(5 plastics, glass and seawater)
- 3 timepoints
(3, 9 and 12 months)
- 106 samples in total

Amplicon data revealed potential
'plastic degrading microbes' were
present in very low abundances.

Aims of my project

- Compare taxonomic assignments from metagenomic data to amplicon data
- Investigate the broad-scale functional potential within microbial communities
- Determine whether core functions were conserved regardless of substrate/timepoint
- Identify potential genes conferring plastic degradation

Methods

Sample
collection and
DNA extractions

Quality checking, read trimming
and adapter removal



usadellab/
Trimmomatic

Shotgun sequencing



Taxonomic assignment

Sample collection and DNA extractions

Quality checking, read trimming and adapter removal



Taxonomic assignment using Metaxa2

Shotgun sequencing



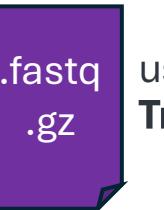
Metaxa2

- Unassembled reads
- Identifies and extracts partial rRNA sequences
 - small subunit (SSU; 16S/18S) rRNA gene
 - large subunit (LSU; 23S/28S) rRNA gene
- Assigns taxonomy to archaeal, bacterial, eukaryotic, mitochondrial or chloroplast origin

Assembled but not binned

Sample
collection and
DNA extractions

Quality checking, read trimming
and adapter removal



usadellab/
Trimmomatic

Shotgun sequencing



Assembly
metaSPAdes

Taxonomic
assignment
using
Metaxa2

Assembled but not binned

Sample collection and DNA extractions

Quality checking, read trimming and adapter removal



.fastq
.gz

usadellab/
Trimmomatic

Shotgun sequencing



Taxonomic assignment using Metaxa2

Assembly
metaSPAdes

hyattpd/**Prodigal**
Prediction of open reading frames (ORFs)



.daa

Annotating sequences by aligning contigs with reference databases

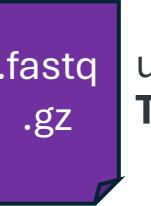
Assembled but not binned

Sample collection and DNA extractions



Shotgun sequencing

Quality checking, read trimming and adapter removal



usadellab/
Trimmomatic

Assembly
metaSPAdes

Taxonomic assignment using Metaxa2

hyattpd/**Prodigal**
Prediction of open reading frames (ORFs)



Annotating sequences by aligning contigs with reference databases

Read mapping
BBMap / Bowtie

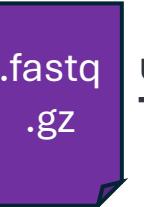
Count data

featureCounts

Unassembled

Sample collection and DNA extractions

Quality checking, read trimming and adapter removal



usadellab/
Trimmomatic

Shotgun sequencing



Taxonomic assignment using Metaxa2

Prediction of open reading frames (ORFs)

gatech-genemark/
MetaGeneMark-2

hyattpd/**Prodigal**

gaberoo/
FragGeneScan

Unassembled

Sample collection and DNA extractions



Shotgun sequencing



Quality checking, read trimming and adapter removal

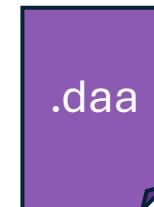


usadellab/
Trimmomatic

Taxonomic assignment using Metaxa2



Annotating sequences by aligning them with reference databases



gatech-genemark/
MetaGeneMark-2

hyattpd/**Prodigal**

gaberoo/
FragGeneScan

Prediction of open reading frames (ORFs)

Unassembled

Sample collection and DNA extractions



Quality checking, read trimming and adapter removal



usadellab/
Trimmomatic

Shotgun sequencing

Taxonomic assignment using Metaxa2

gatech-genemark/
MetaGeneMark-2

hyattpd/**Prodigal**

Prediction of open reading frames (ORFs)

gaberoo/
FragGeneScan

Annotating sequences by aligning them with reference databases



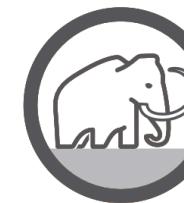
MEGANIZER



community edition

MEGAN6

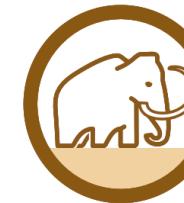
- Software for interactive metagenomics analyses
- Specialised for unassembled reads
- Command-line tools
- Taxonomic analyses – NCBI
- Functional analyses - SEED, InterPro2GO, eggNOG



MEGAN⁶

community edition

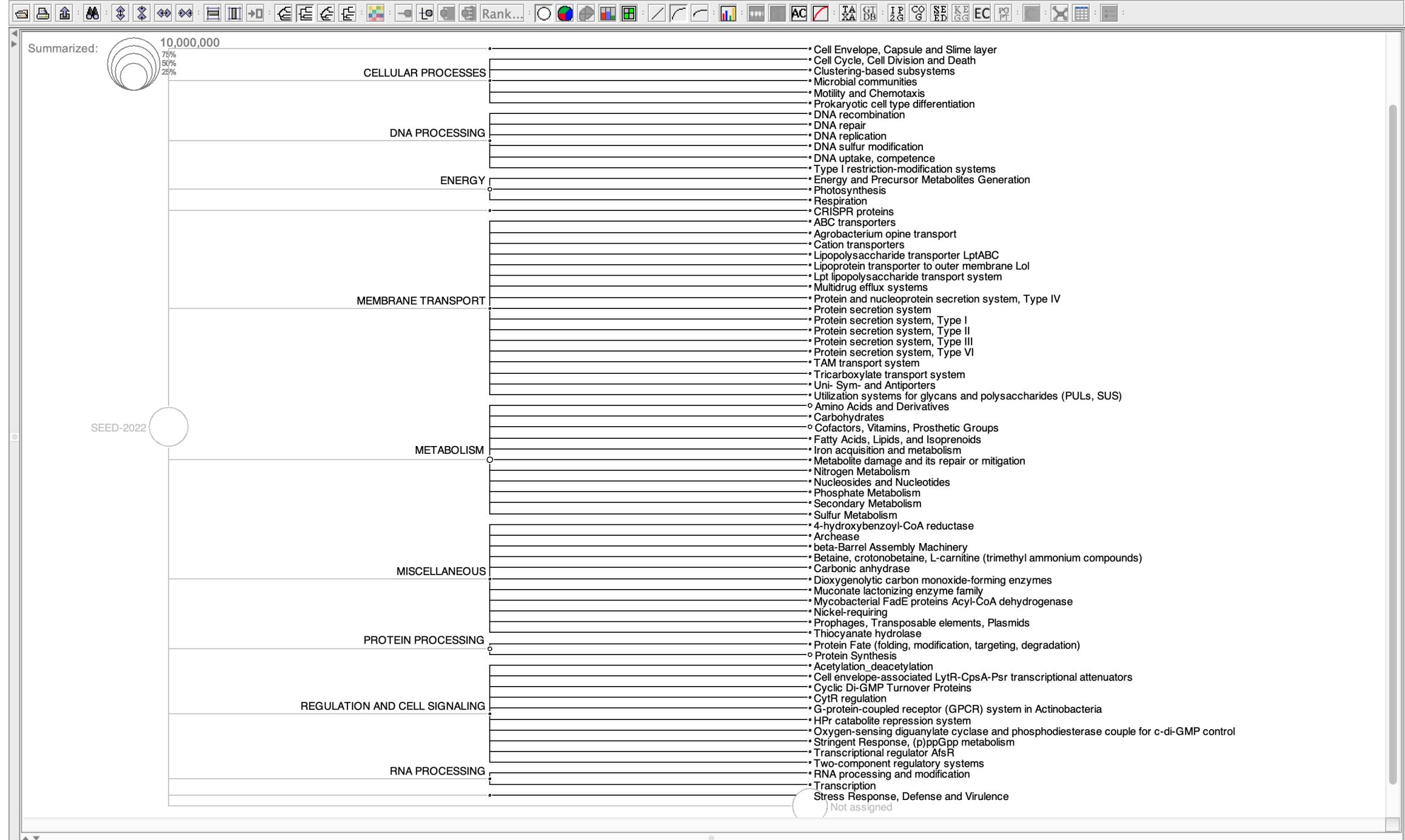
Community edition is free



MEGAN⁶

ultimate edition

Ultimate edition requires a paid subscription but includes up-to-date KEGG pathways



Visualisation of broad-scale functioning (SEED)

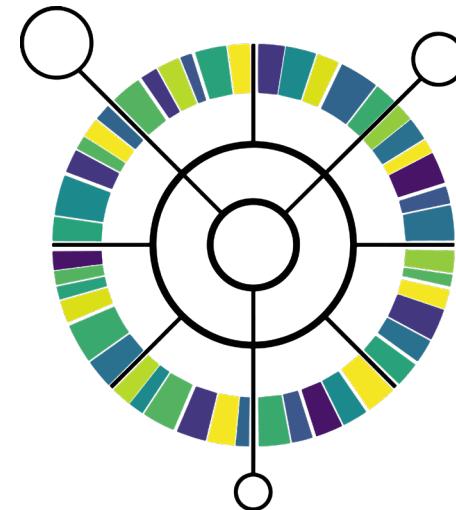
Count tables were normalised before Bray-Curtis dissimilarity matrices were calculated.

Specialised gene annotations



Item	Update
Microorganisms	753
Proteins	219
Last update	13/05/2024
Last addition	Protein 00219

CARD
(Comprehensive
Antibiotic Resistance
Database)



Custom diamond databases can be made and used with diamond.

Identification of genes associated with plastic degradation



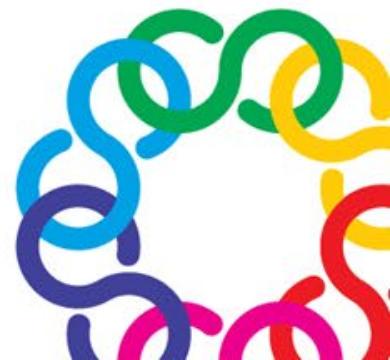
Item	Update
Microorganisms	753
Proteins	219
Last update	13/05/2024
Last addition	Protein 00219

Task: DRAM

[Go to Github MGSS webpage](#)

Tasks:

- MAG annotation with DRAM



Task: Coverage calculation

[Go to Github MGSS webpage](#)

Tasks:

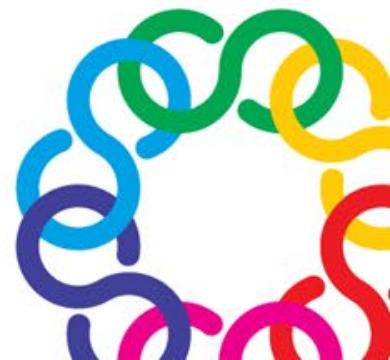
- Coverage calculation using Bowtie2



Task: Pick group challenge!

Tasks:

- Introduce group project goals
- Dividing into working groups / get a group name
- Select a goal from your project



Task: Pick group challenge!

Determine which genome(s) have the following attributes, and the genetic mechanisms used for these attributes:

1. Denitrification (Nitrate or nitrite to nitrogen)
2. Ammonia oxidation (Ammonia to nitrite or nitrate)
3. Anammox (Ammonia and nitrite to nitrogen)
4. Sulfur oxidation (SOX pathway, thiosulfate to sulfate)
5. Sulfur reduction (DSR pathway, sulfate to sulfide)
6. Photosynthetic carbon fixation
7. Non-photosynthetic carbon fixation (Reverse TCA or Wood-Ljungdahl)
8. Non-polar flagella expression due to a chromosomal deletion
9. Plasmid-encoded antibiotic resistance
10. Aerobic (versus anaerobic) metabolism

