

University of Naples Federico II

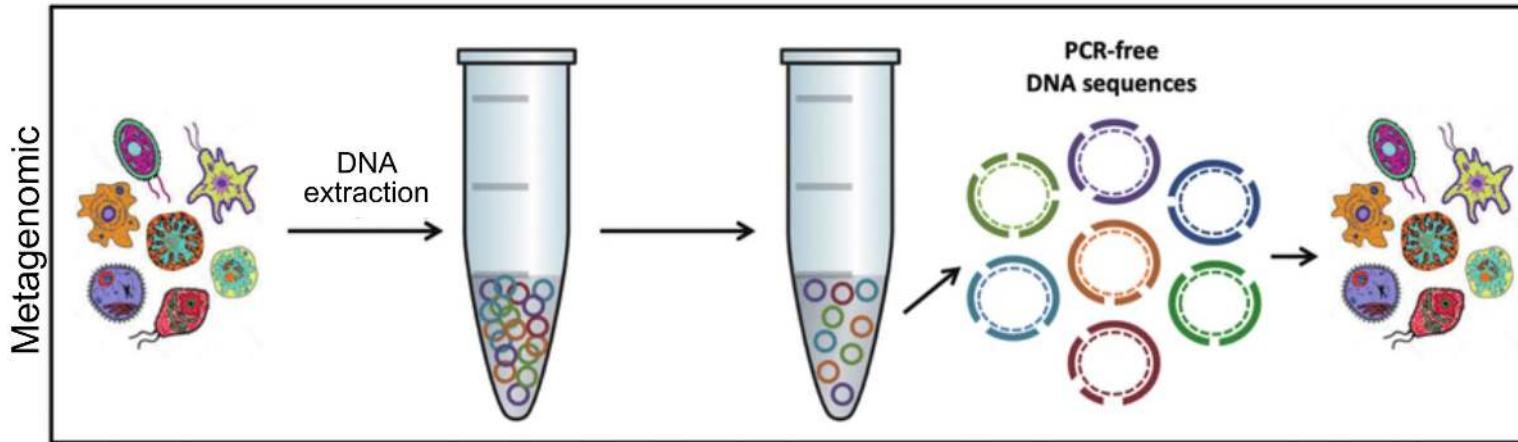
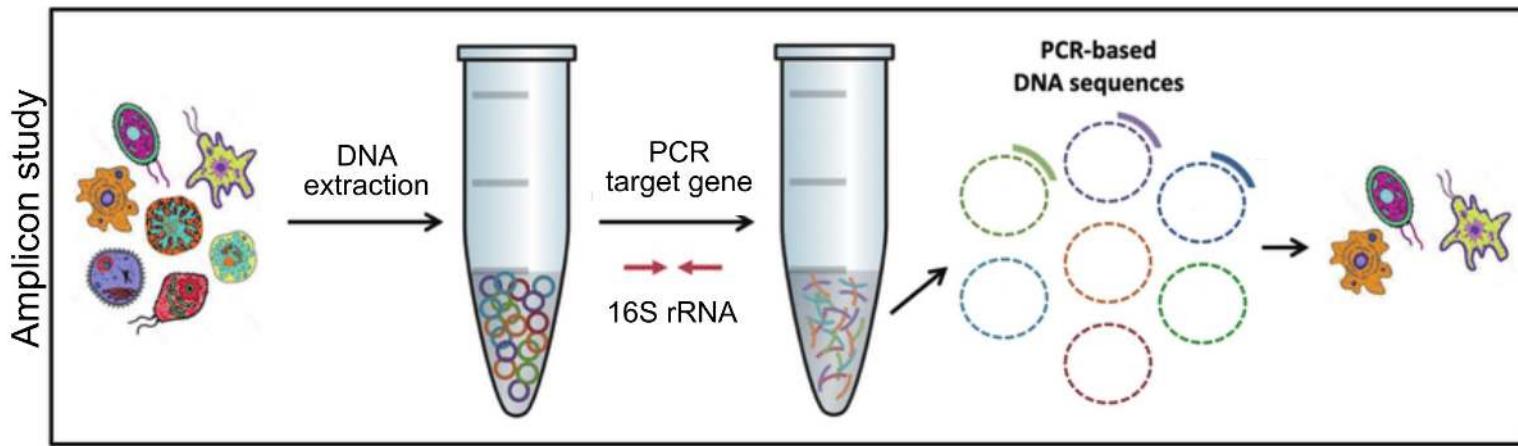
Environmental Metagenomic

aa 2020-2021

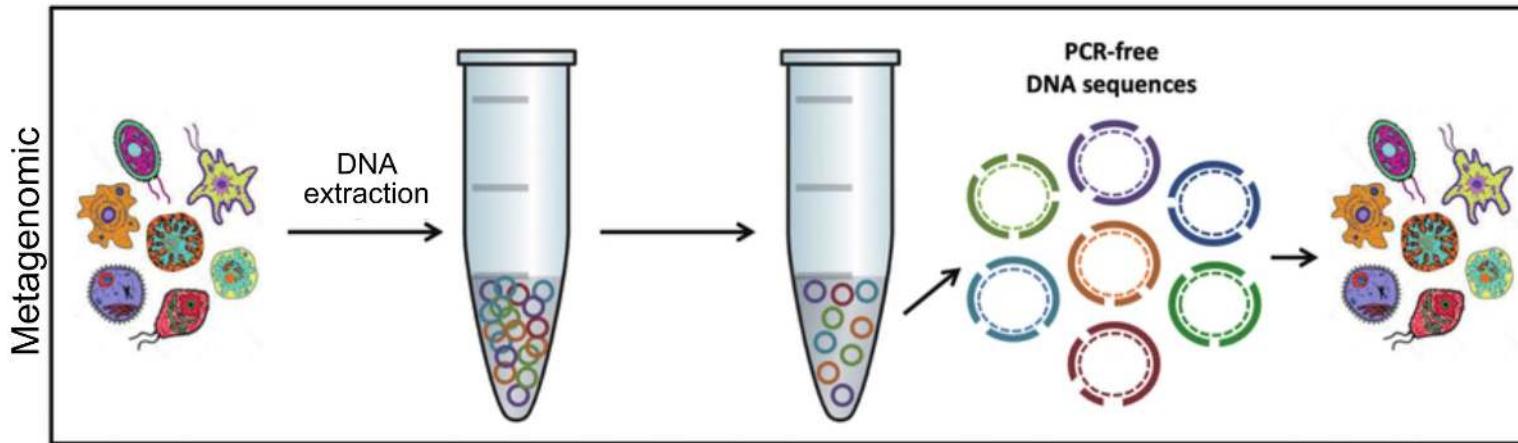
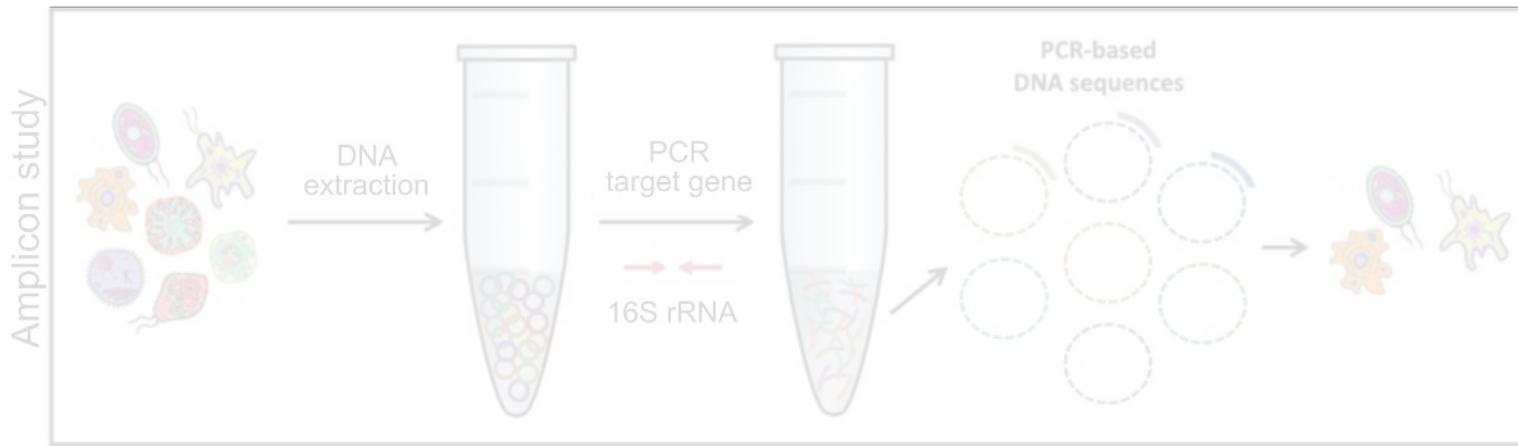
INTRODUCTION TO SHOTGUN METAGENOMIC

Donato Giovannelli

Amplicon vs Metagenomic



Amplicon vs Metagenomic



Microbial Diversity: metagenomic data

Brief history of Metagenomic

Norman R. Pace

propose the idea of cloning DNA directly from environmental sample to analyze 16S rRNA diversity



1985

1998

2002

2003

2004

2005

2007



Mya Breitbart

used environmental shotgun sequencing to show the diversity of virus in seawater



Jo Handelsman is the first to use the term METAGENOMIC referring to the analysis of community genomes



A pilot **GOS** project in the **Sargasso Sea** shows unprecedented bacterial diversity in seawater



Craig Venter leads the *Global Ocean Sampling Expedition* (GOS) to collect metagenomic samples throughout the journey

Robert Edwards published sequences generated using pyrosequencing techniques

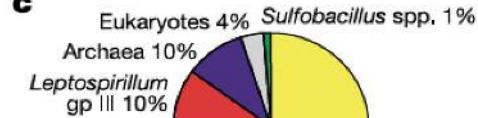


Stephan Schuster publishes the first sequences generated using high-throughput sequencing.

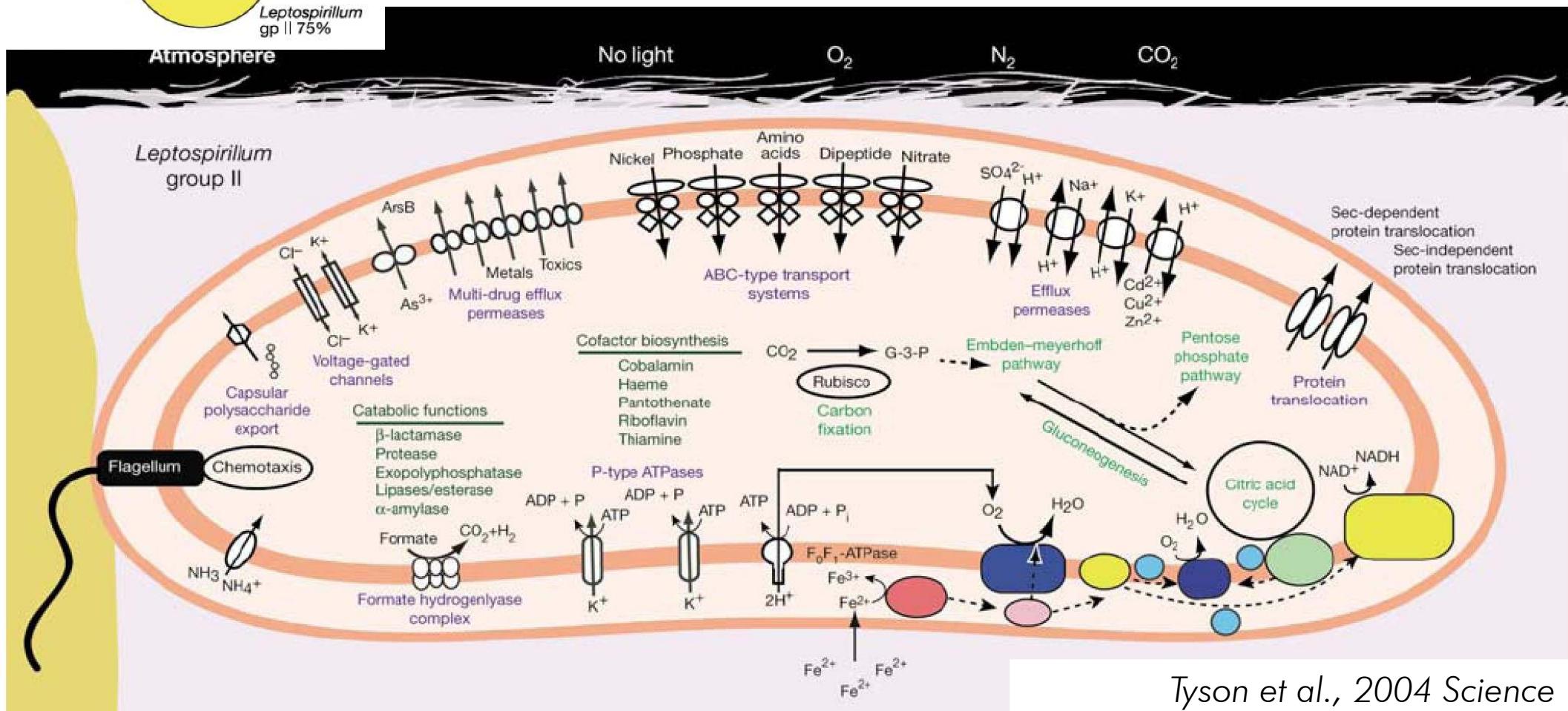
The history of metagenomics: An incomplete summary

a post by Matthew Schechter

<https://merenlab.org/2020/07/27/history-of-metagenomics/>

c

First binned genomes from an environmental metagenome were reported in 2004



The list of -omics is longer than you may expect

	Google pages
<i>Genomic</i>	2,780,000
<i>Transcriptomic</i>	221,000
<i>Proteomic</i>	857,000
<i>Secretomic</i>	18,700
<i>Metabolomic</i>	29,000
<i>Metallomeic</i>	630
<i>Regulomic</i>	696
<i>Lipidomic</i>	4,620
...	

Metagenomic

Metagenomic reveals the functional (as well as taxonomic) diversity of a community

Being based on the sequencing of DNA, it reveals the metabolic potential, and does not provide any information on the active metabolisms

Metagenomic approaches, intended as the sequencing of the community DNA, can be applied to a sample as it is (directly after extraction of the bulk community DNA) or combined with a number of other techniques (like SIP, enrichments, etc..)

*The **depth of sequencing** (how many reads are obtained for each sample) and the **read length** (how many base pairs is each read), both depending on the sequencing technology used, very important parameters that can severely affect the final result*

Metagenomic basic approaches

Metagenomic analysis can happen at three different levels:

Metagenomic basic approaches

Metagenomic analysis can happen at three different levels:

Read-based analysis: it applies to the analysis of unassembled reads. Potentially provides a large amount (ecologically speaking) of coarse level information about the metabolic potential at the community level. Put another way, we can broadly assign a lot of function without much information about their coupling in different organisms

Metagenomic basic approaches

Metagenomic analysis can happen at three different levels:

Read-based analysis: it applies to the analysis of unassembled reads. Potentially provides a large amount (ecologically speaking) of coarse level information about the metabolic potential at the community level. Put another way, we can broadly assign a lot of function without much information about their coupling in different organisms

Contig-based analysis: it analyses the assembled reads to provide genetic context to the identified functions and might provide also taxonomic information. The amount of information obtained is more detailed but represent a smaller subset of the community (i.e. only the reads that assemble)

Metagenomic basic approaches

Metagenomic analysis can happen at three different levels:

Read-based analysis: it applies to the analysis of unassembled reads. Potentially provides a large amount (ecologically speaking) of coarse level information about the metabolic potential at the community level. Put another way, we can broadly assign a lot of function without much information about their coupling in different organisms

Contig-based analysis: it analyses the assembled reads to provide genetic context to the identified functions and might provide also taxonomic information. The amount of information obtained is more detailed but represent a smaller subset of the community (i.e. only the reads that assemble)

Genome resolved metagenomic: it analyses near-complete genomes assembled from metagenomes (MAGs). It provides highly detailed informations about a small subset of the community (only the contigs that bin into high quality genomes)

Metagenomic jargon

Assembly: linking together reads that are partially overlapping

Contigs: sequences of variable length obtained from assembling the reads

Binning: Clustering together contigs of variable length (bins) based on their identification as coming from a close group of highly related strains

MAGs: Metagenomes assembled genomes

Single-copy marker gene: ubiquitous conserved genes that are present with a single copy in known genomes

Annotation: Assigning taxonomy or function to a sequence

Read-mapping: counting the number of reads that recruit (align) to a specific sequence (MAG, contig or gene)

Technical difficulties and challenges when using (metagen)omic

Study question / hypothesis

- Initial sample → Sample size and replicates
- Isolation and purification of starting material → Quality and quantity. Extraction biases
- Technology/ Reads length and quantity / number of spots and resolution of gel → Depth of sequencing and coverage, Sequencing errors, Accuracy and read quality
- Annotation and assembly → Quality of reference database, Annotation accuracy, Assembly quality, Chimeric sequences/contigs
- Taxonomic binning → Lateral gene transfer events
- Genomic reconstruction from metagenome → Multi strain genomes

Next Gen sequences are on average short (75-450 bp). This is less than the average gene (\sim 1500 bp). What are the implication for our purpose?

Illumina

50 to 350 bp Paired End

Pacbio

up to 1000 bp

Nanopore

up to 10 kbp

Does size matter?

Next Gen sequences are on average short (75-450 bp). This is less than the average gene (~1500 bp). What are the implication for our purpose?

Illumina

50 to 350 bp Paired End

Pacbio

up to 1000 bp

Nanopore

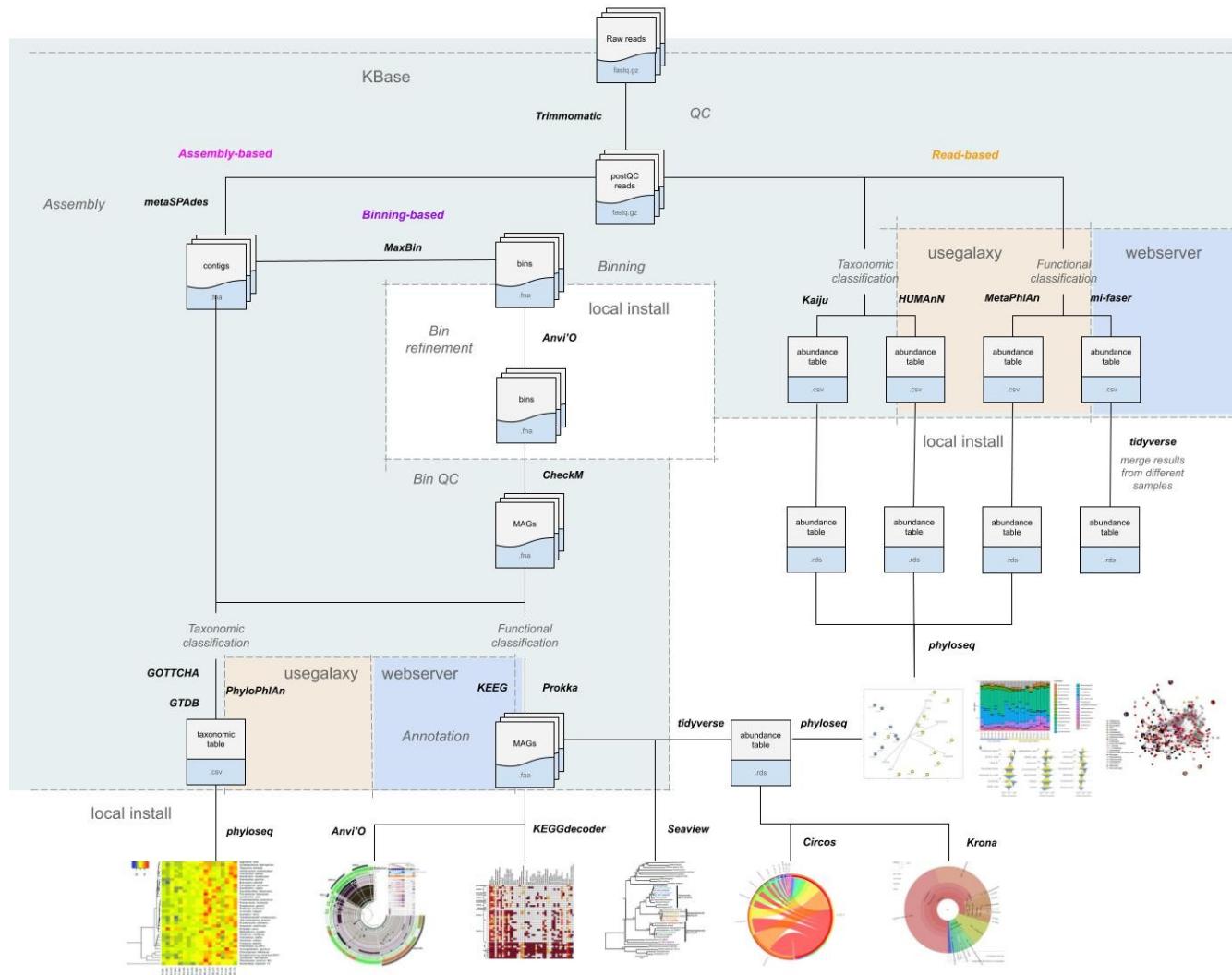
up to 10 kbp

Does size matters?

YES! In sequencing it makes a big difference. Generally speaking, the longer the reads, higher the amount of information it contains. You can deal with short reads increasing coverage (i.e. number of total transcripts per sample), but gene calling, taxonomic and functional annotation and contigs assembly is a computational challenge.

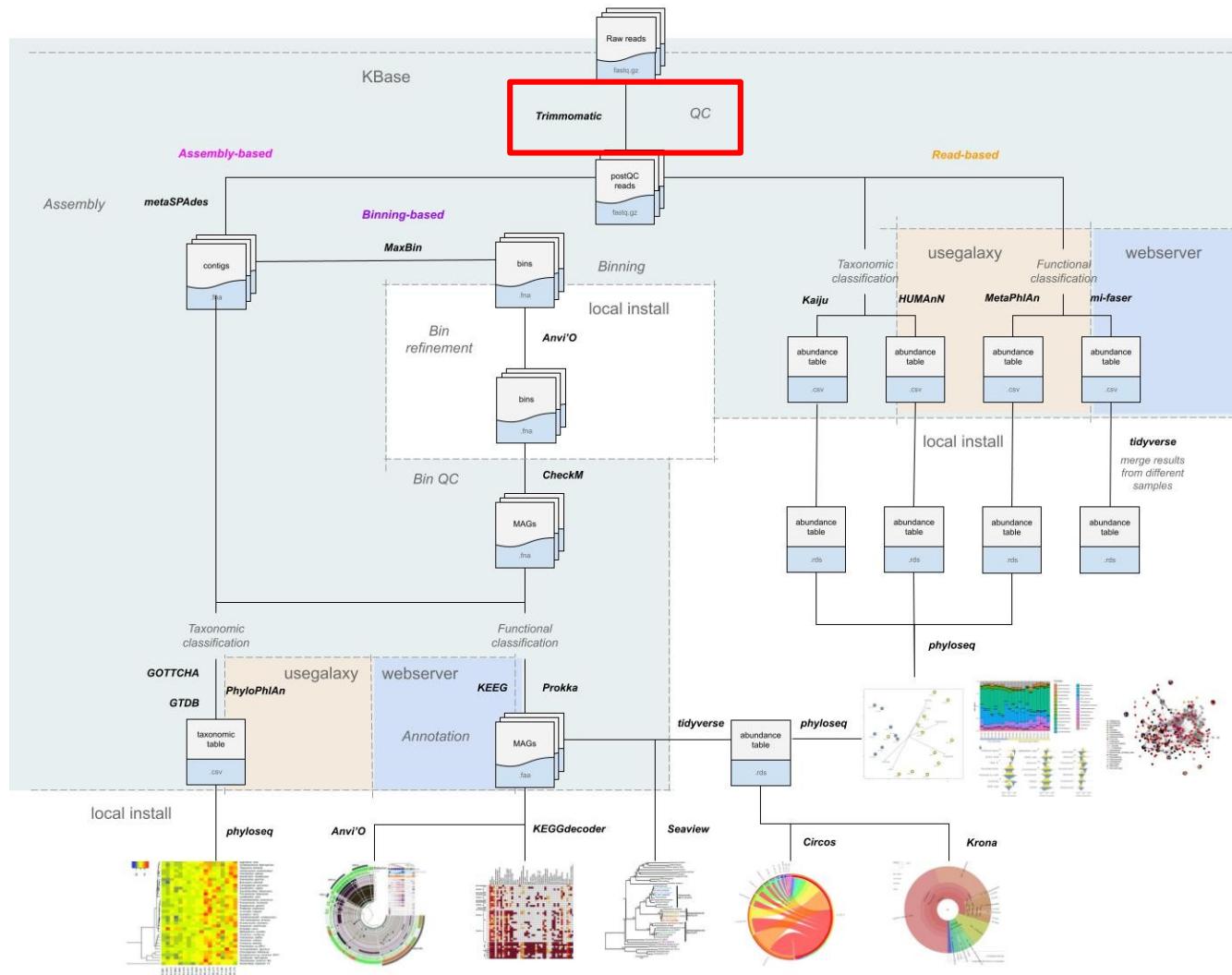
Shotgun Metagenomic

generic workflow - Giovannelli Lab 2021



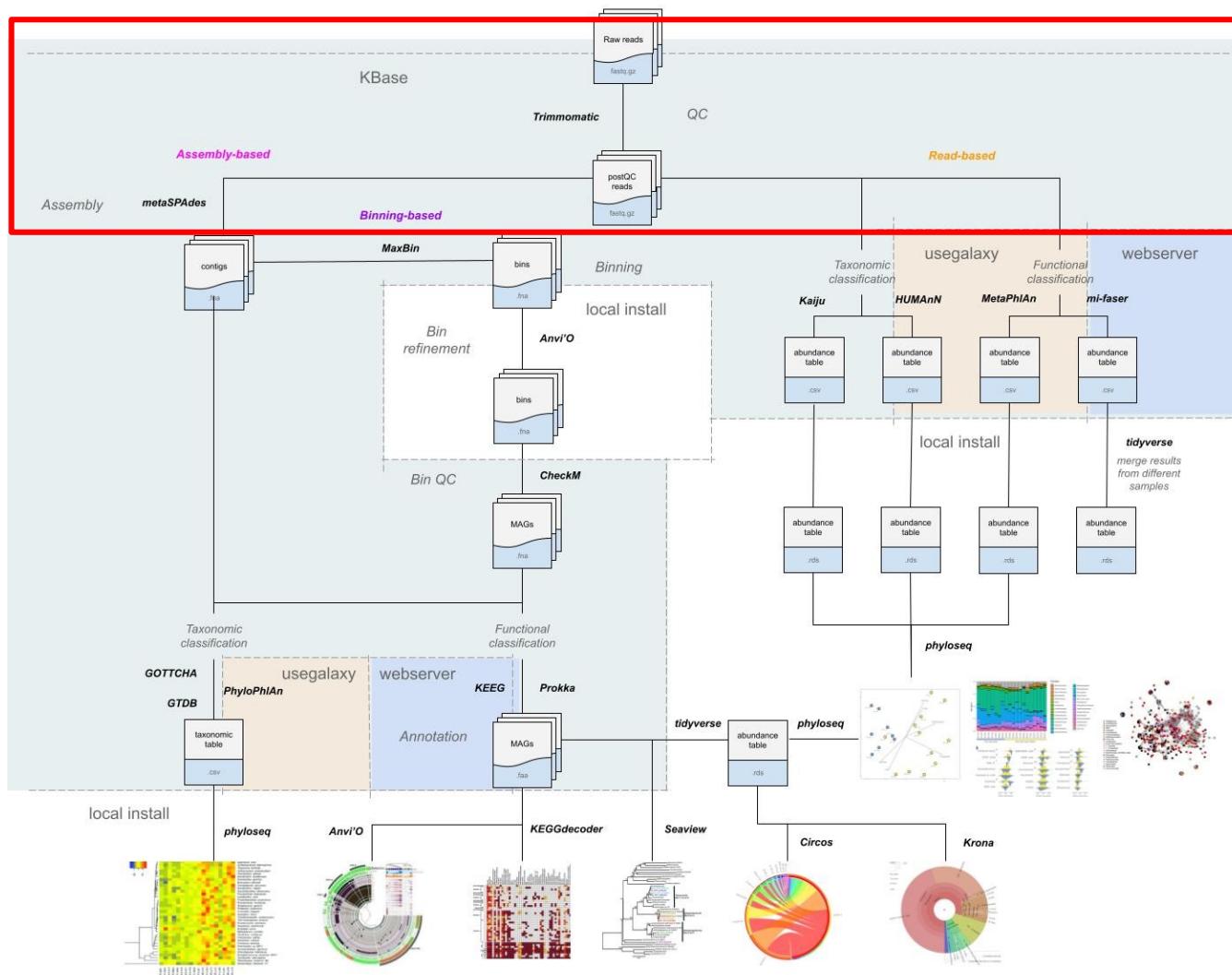
Shotgun Metagenomic

generic workflow - Giovannelli Lab 2021



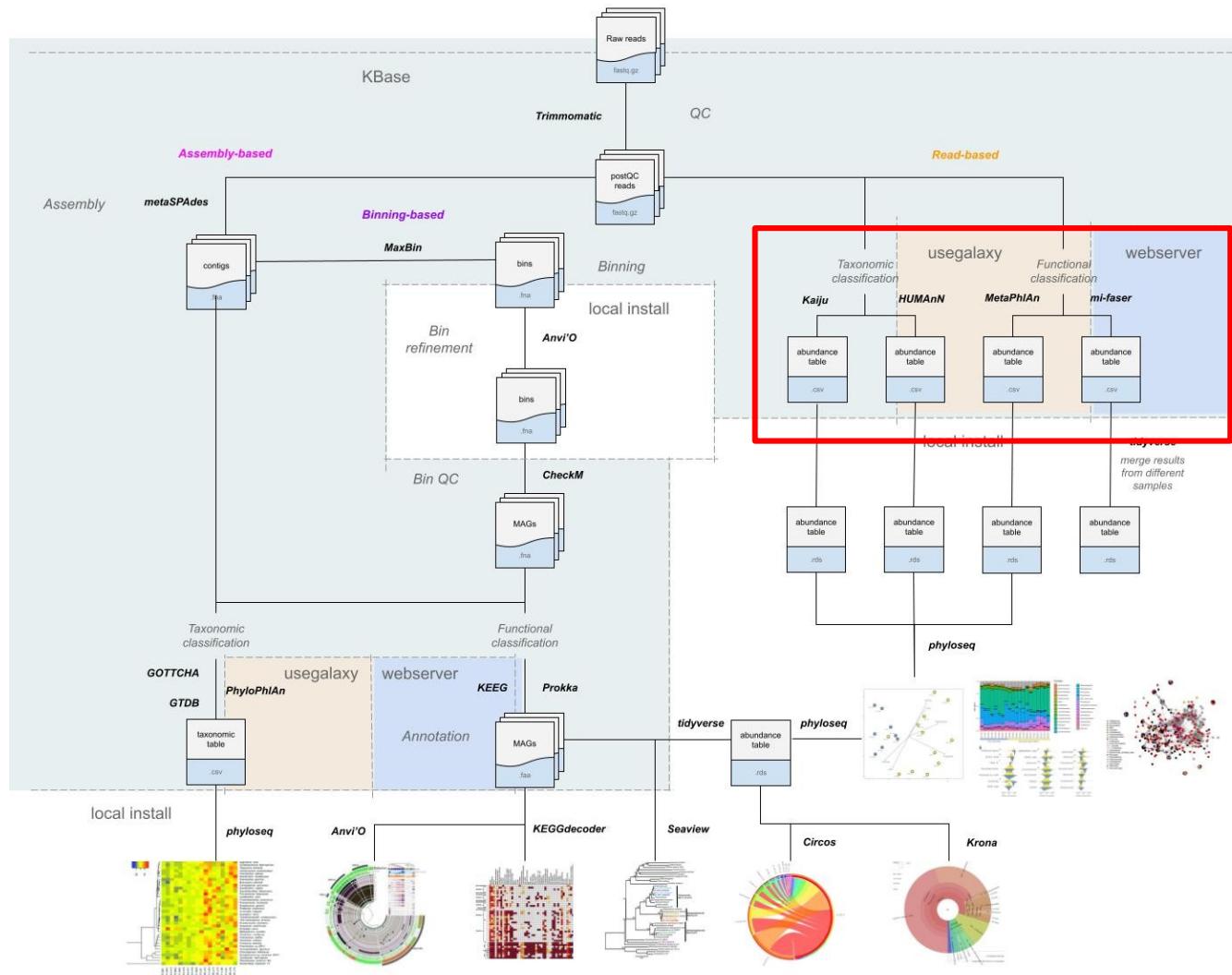
Shotgun Metagenomic

generic workflow - Giovannelli Lab 2021



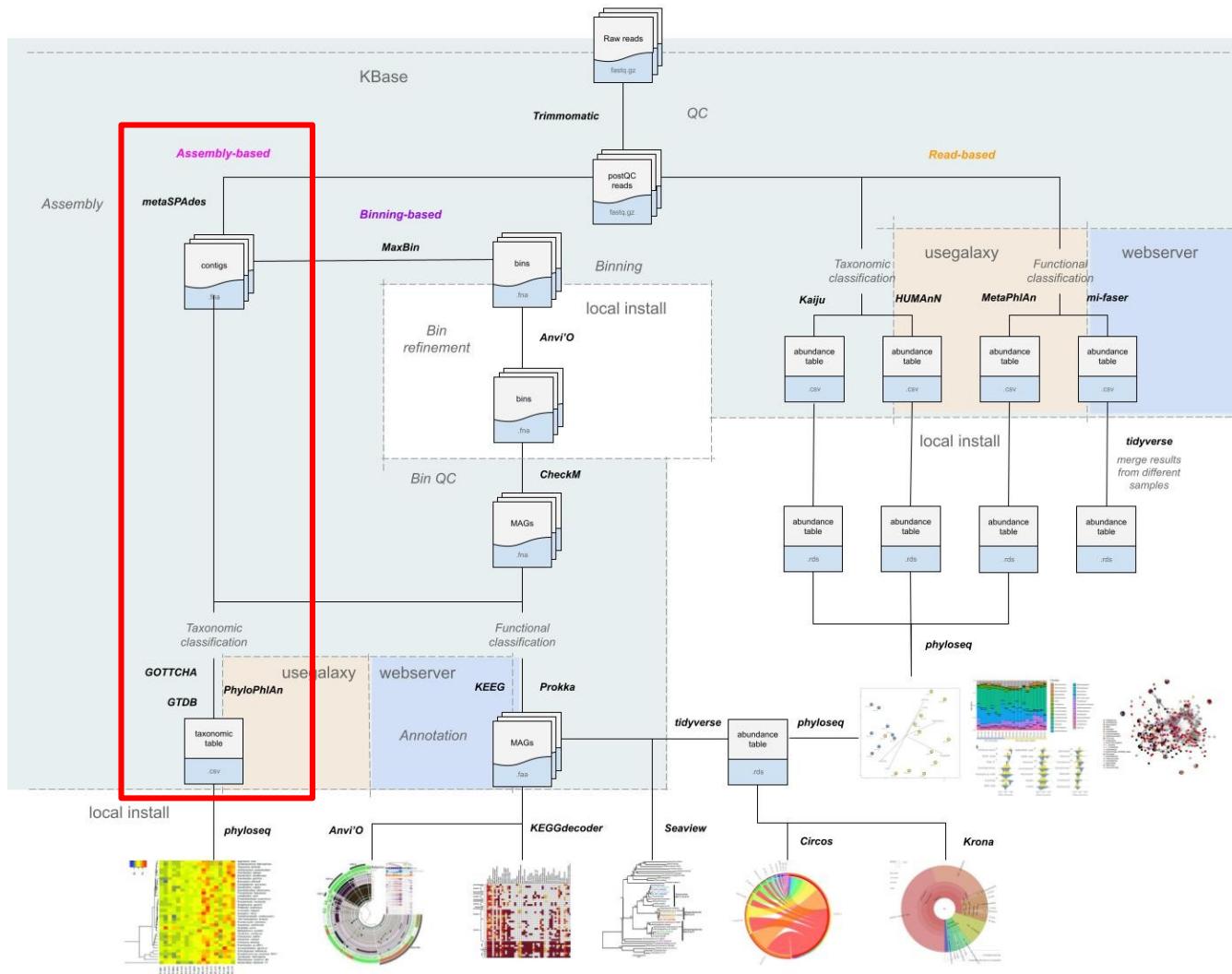
Shotgun Metagenomic

generic workflow - Giovannelli Lab 2021



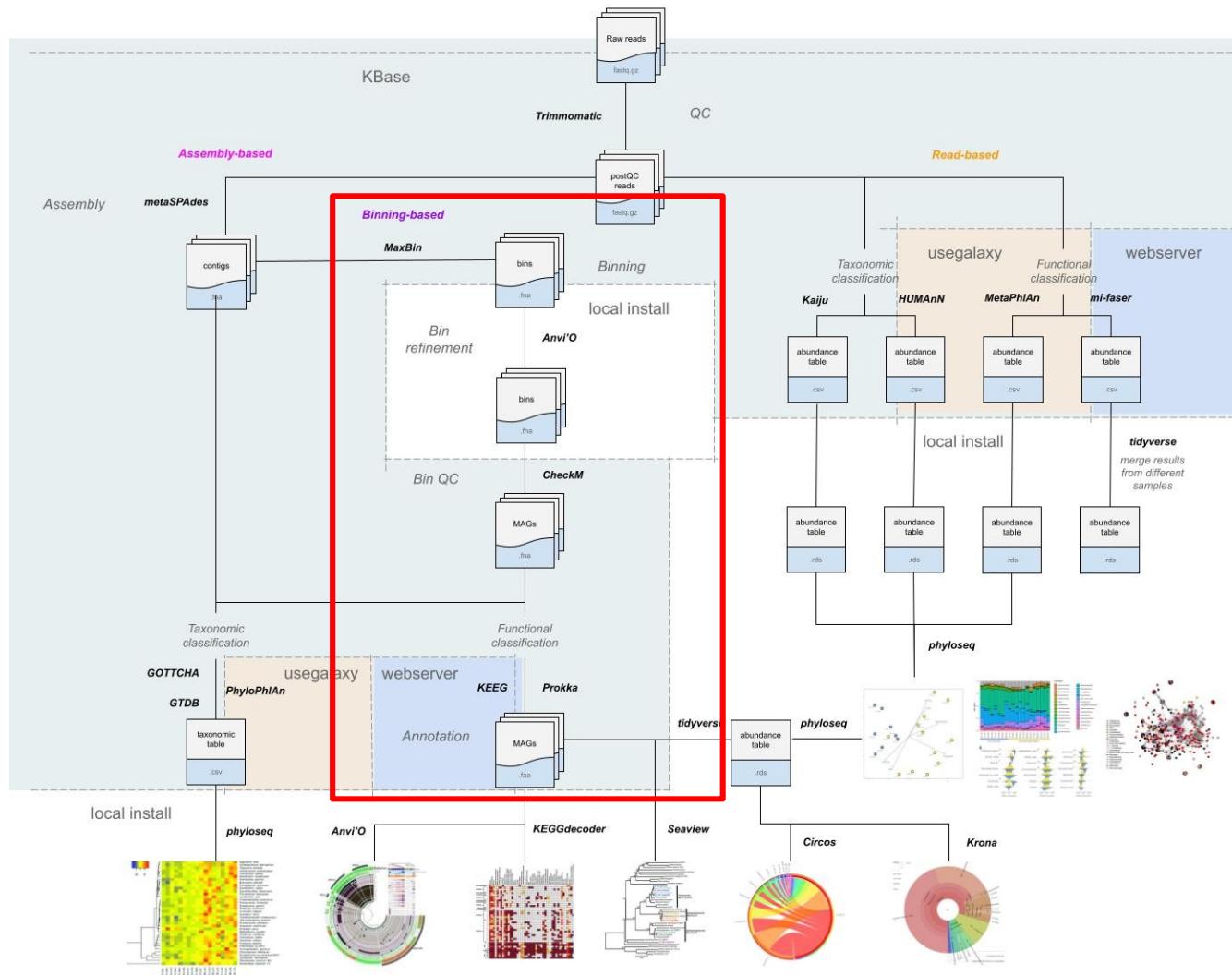
Shotgun Metagenomic

generic workflow - Giovannelli Lab 2021



Shotgun Metagenomic

generic workflow - Giovannelli Lab 2021



Read Assembly

Reference **TCCTAGAGATCCGCC**TCTTAGCGGAATAATACAGCCGAATCTTAGCGGAATTGCCAGCACAG

Reads **CCTAGAGATCCG**
GAGATCCGCCTC
ATCCGCCTCTTA
GCCTCTTAGCGG
CTTAGCGG**ATAT**
TAGCGG**ATATAA**
TATAA**ATACAGCC**
ACAGCCGAATCT
CCGAATCTTAGC
GAATCTTAGCGG
TCTTAGCGG**AAAT**
AGCGGAATTGCC
GGAATTGCCAGC
AATTGCCAGCAC
TTGCCAGCACAG

Contigs **CCTAGAGATCCGCC**TCTTAGCGG
CTTAGCGG**ATATAA****ATACAGCC**GAATCTTAGCGG
TCTTAGCGG**AAATTGCCAGC**
AGCGGAATTGCCAGCACAG

Reference **TCCTAGAGATCCGCC**TCTTAGCGGAATAATACAGCCGAATCTTAGCGGAATTGCCAGCACAG

Reads **CCTAGAGATCCG**
GAGATCCGCCTC
ATCCGCCTTA
GCCT**CTTAGCGG**
CTTAGCGGATAT
TAGCGGATATAA
TATAATACAGCC
ACAGCCGAATCT
CCGAATCTTAGC
GAAT**CTTAGCGG**
CTTAGCGGAAT
AGCGGAATTGCC
GGAATTGCCAGC
AATTGCCAGCAC
TTGCCAGCACAG

Contigs **CCTAGAGATCCGCC**TCTTAGCGG
CTTAGCGGATATAATACAGCCGAATCTTAGCGG
CTTAGCGGAATTGCCAGCACAG

Velvet

Spades

IDBA-UD

Unicycler

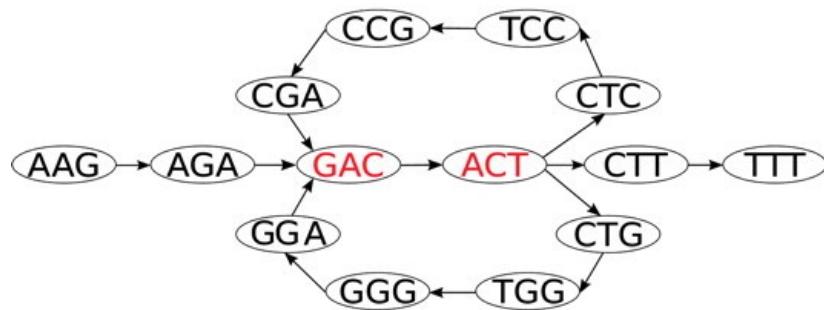
AbySS

MEGAHIT

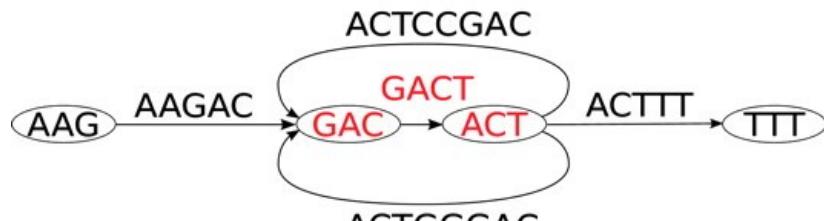
Trinity

MaSuRCA

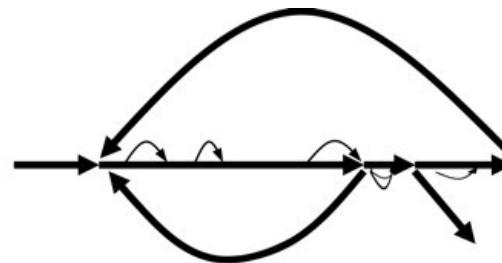
AA**GACT**CC**GACT**GG**GACT**TT



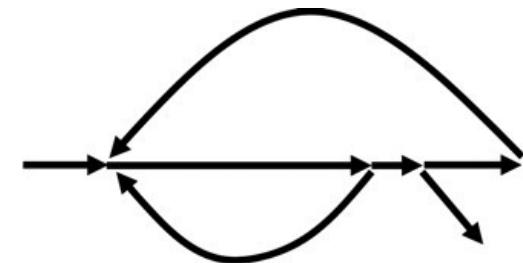
A de Bruijn graph of a sequence



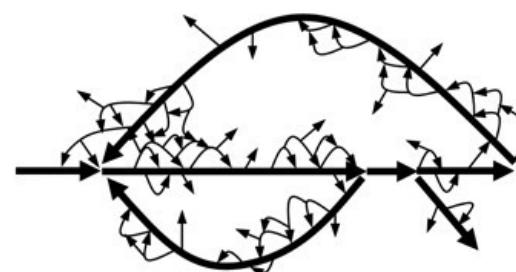
B condensed de Bruijn graph



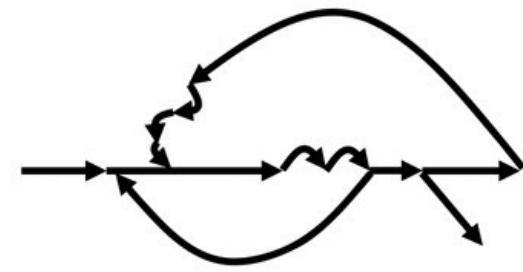
C de Bruijn graph of a genome



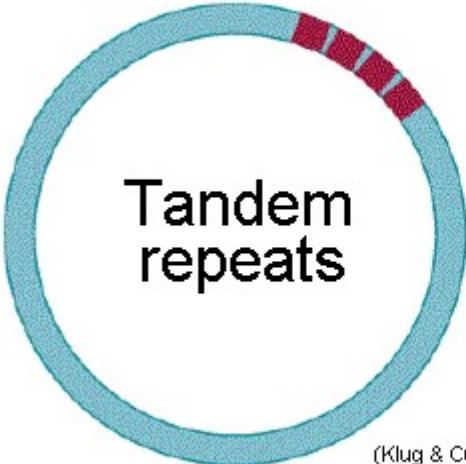
E repeat graph of a genome



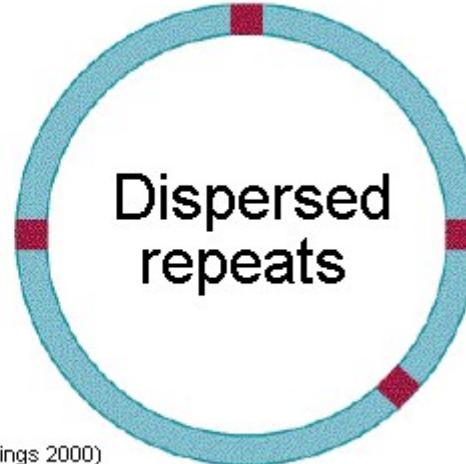
D de Bruijn graph of a set of reads



F repeat graph on a set of reads



Tandem
repeats



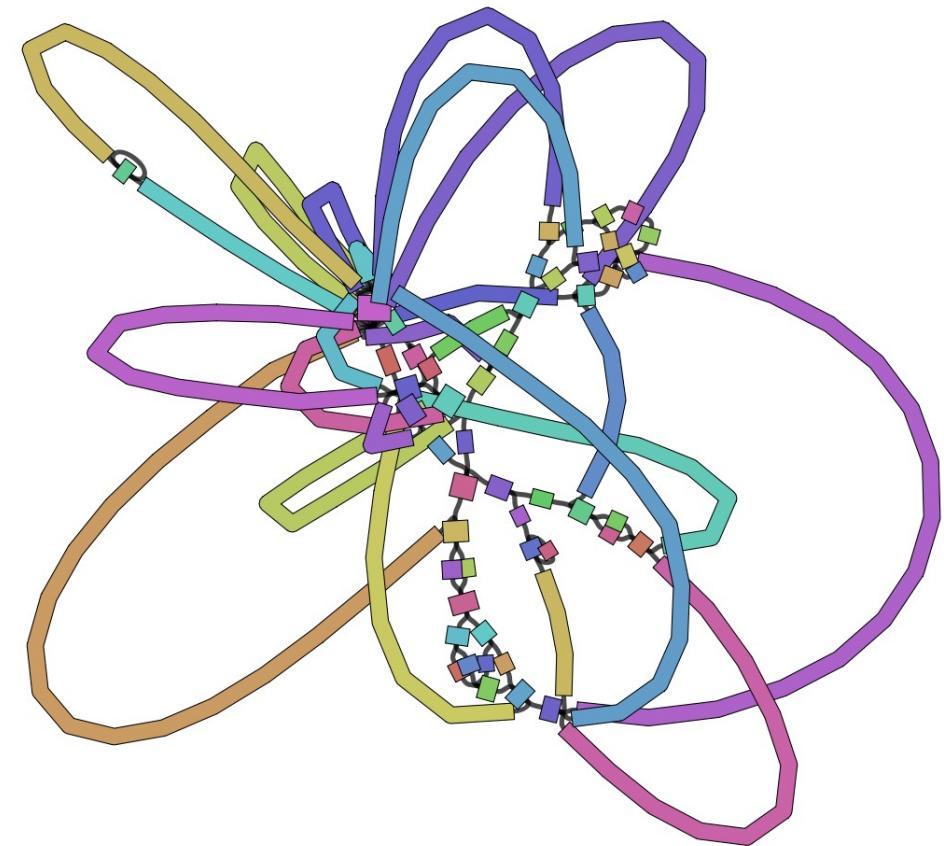
Dispersed
repeats

(Klug & Cummings 2000)

Tandem
repeats

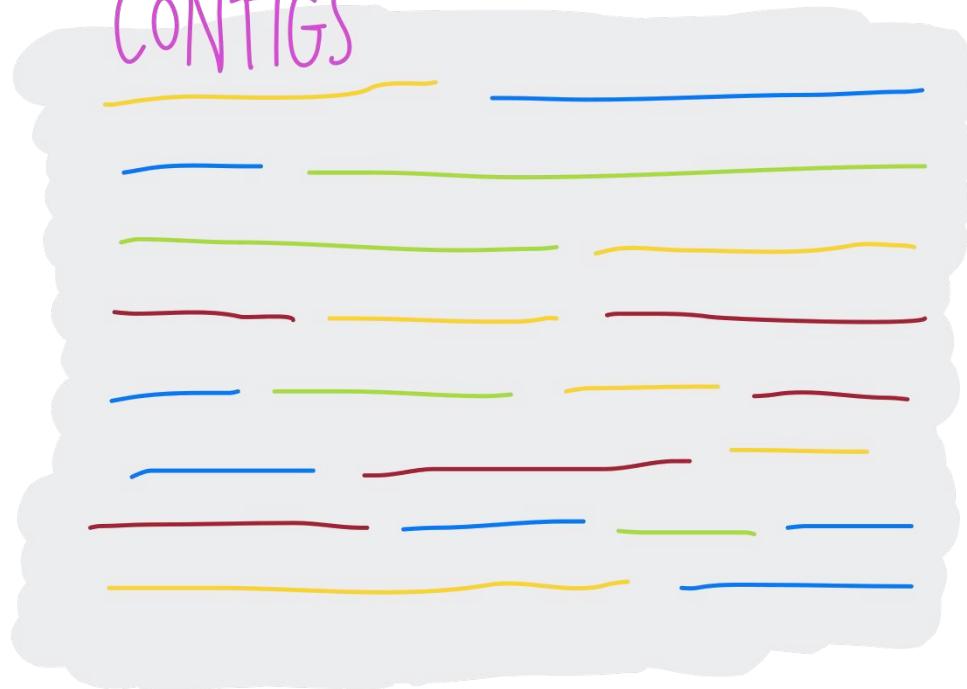
Dispersed
repeats

(Klug & Cummings 2000)



Contig binning

CONTIGS



SEQUENCE COMPOSITION

CONTIGS



MAGs



Peter A. Noble¹

Robert W. Citek²

Oladele A. Ogunseitan³

¹Belle W. Baruch Institute for
Marine Biology and Coastal
Research, University of South
Carolina, Columbia, SC, USA

²Department of Soil and
Environmental Science, University
of California at Riverside,
Riverside, CA, USA

³Department of Environmental
Analysis and Design, University
of California at Irvine,
Irvine, CA, USA

Tetranucleotide frequencies in microbial genomes

A computational strategy for determining the variability of long DNA sequences in microbial genomes is described. Composite portraits of bacterial genomes were obtained by computing tetranucleotide frequencies of sections of genomic DNA, converting the frequencies to color images and arranging the images according to their genetic position. The resulting images revealed that the tetranucleotide frequencies of genomic DNA sequences are highly conserved. Sections that were visibly different from those of the rest of the genome contained ribosomal RNA, bacteriophage, or undefined coding regions and had corresponding differences in the variances of tetranucleotide frequencies and GC content. Comparison of nine completely sequenced bacterial genomes showed that there was a nonlinear relationship between variances of the tetranucleotide frequencies and GC content, with the highest variances occurring in DNA sequences with low GC contents (less than 0.30 mol). High variances were also observed in DNA sequences having high GC contents (greater than 0.60 mol), but to a much lesser extent than DNA sequences having low GC contents. Differences in the tetranucleotide frequencies may be due to the mechanisms of intercellular genetic exchange and/or processes involved in maintaining intracellular genetic stability. Identification of sections that were different from those of the rest of the genome may provide information on the evolution and plasticity of bacterial genomes.

Tetranucleotide frequency

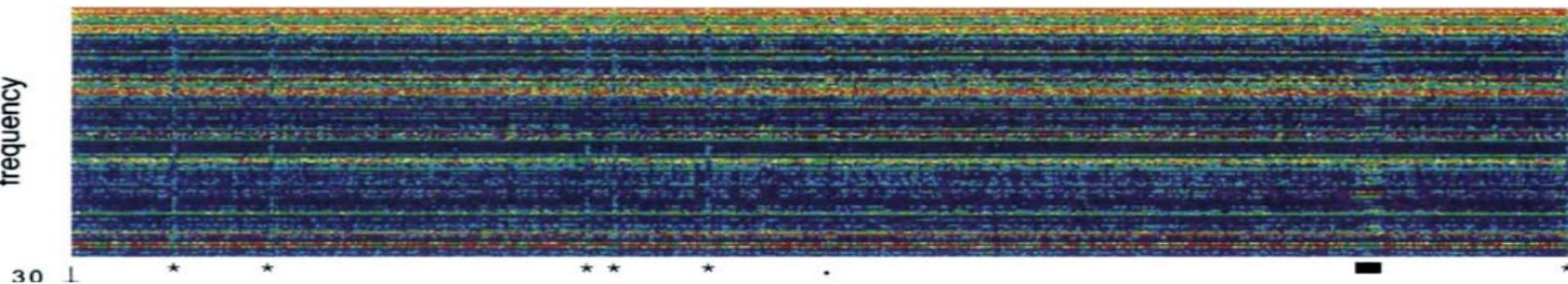
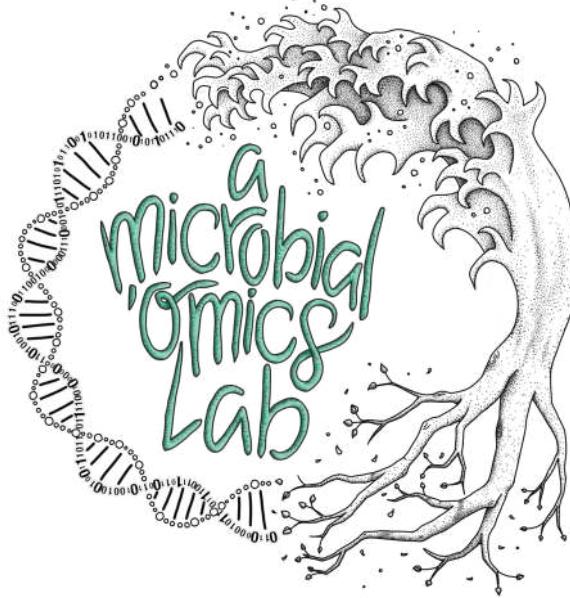


Figure 1. Fingerprints, variances of tetranucleotide frequencies, and GC values of sections of the *Haemophilus influenzae* Rd genome are consecutively ordered from the *NotI* restriction site [9]. Each column of the color image represents the fingerprint obtained from the analysis of one DNA sequence (*i.e.*, a 3000 bp section). Each row represents the frequency of a specific tetranucleotide and its complement. Tetranucleotides are arranged alphabetically on the *y*-axis. Each tetranucleotide is represented by a box, whose color is determined by its frequency, ranging from purple (low) to red (high). A star (*) identifies sections containing ribosomal RNA. The black bar identifies the location of the cryptic Mu-like bacteriophage. The variance and GC values were computed from the analysis of one section.



<https://merenlab.org/momics/>



GTTTGGCATGATTAAGGGAGTTCTTTGTGCTTC

GTTTGGCATGATTAAGGGAGTTCTTTGTGCTTC

k=2

GTTTGGCATGATTAAGGGAGTTCTTTGTGCTTC

AA | AC | AG | AT | CA | CC | CG | CT | GA | GC | GG | GT | TA | TC | TG | TT

k=2

GTTTGGCATGATTAAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

k=2

GT TTGGCATGATTAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0

k=2

C TTT TTGGCATGATTAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	1

k=2

G-TT-GGCATGATTAAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	2

k=2

GT~~T~~**TT**GGCATGATTAAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	3

k=2

GTTT **TGG** CATGATTAAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
0	0	0	0	0	0	0	0	0	0	0	1	0	0	1	3

k=2

GTTT-GG-CATGATTAAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
0	0	0	0	0	0	0	0	0	0	1	1	0	0	1	3

k=2

GTTTTCGCATGATTAAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
0	0	0	0	0	0	0	0	0	1	1	1	0	0	1	3

k=2

GTTTGCA**CAT**-GATTAAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
0	0	0	0	1	0	0	0	0	1	1	1	0	0	1	3

k=2

GTTTTGGC ATG GATTAAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
0	0	0	1	1	0	0	0	0	1	1	1	0	0	1	3

k=2

GTTTTGGCA.TGATTAAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
0	0	0	1	1	0	0	0	0	1	1	1	0	0	2	3

k=2

GTTTGCGCATGA-TAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
0	0	0	1	1	0	0	0	1	1	1	1	0	0	2	3

k=2

GTTTTGGCATTGATTAAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
0	0	0	2	1	0	0	0	1	1	1	1	0	0	2	3

k=2

GTTTGGCATGATTAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
0	0	0	2	1	0	0	0	1	1	1	1	0	0	2	4

k=2

GTTTGGCATGATTAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
1	0	2	2	1	0	0	2	2	2	2	3	1	2	4	10

k=2

GTTTGCGATGATTAAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
1	0	2	2	1	0	0	2	2	2	2	3	1	2	4	10

GAAGCACAAAAGAAACTCCTTAATCATGCCAAAAAC

9

k=2

GTTTGCGATGATTAAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
1	0	2	2	1	0	0	2	2	2	2	3	1	2	4	10

9

GAAGCACAAAAAGAAAACTCCTTAATCATGCCAAAAAC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
10	3	2	2	4	2	0	2	2	2	0	0	1	2	1	1

k=2

GTTTGGCATGATTAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
1	0	2	2	1	0	0	2	2	2	2	3	1	2	4	10

9

GAAGCACAAAAGAAAACTCCTTAATCATGCCAAAAAC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
10	3	2	2	4	2	0	2	2	2	0	0	1	2	1	1

GTTTGGCATGATTAAGGGAGTTCTTTGTGCTTC
GAAGCACAAAAGAAAACTCCTTAATCATGCCAAAAAC

k=2

GTTTGGCATGATTAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
1	0	2	2	1	0	0	2	2	2	2	3	1	2	4	10



GAAGCACAAAAAGAAAACTCCTTAATCATGCCAAAAAC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
10	3	2	2	4	2	0	2	2	2	0	0	1	2	1	1

GTTTGGCATGATTAAGGGAGTTCTTTGTGCTTC
GAAGCACAAAAAGAAAACTCCTTAATCATGCCAAAAAC

AA	AC	AG	GA	CA	CC	CG	GC	AT	TA

k=2

GTTTGGCATGATTAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
1	0	2	2	1	0	0	2	2	2	2	3	1	2	4	10



GAAGCACAAAAAGAAAACTCCTTAATCATGCCAAAAAC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
10	3	2	2	4	2	0	2	2	2	0	0	1	2	1	1

GTTTGGCATGATTAAGGGAGTTCTTTGTGCTTC
GAAGCACAAAAAGAAAACTCCTTAATCATGCCAAAAAC

AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
11	3	4	4	5	2	0	2	2	1

k=2

GTTTGGCATGATTAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
1	0	2	2	1	0	0	2	2	2	2	3	1	2	4	10



GAAGCACAAAAAGAAAACTCCTTAATCATGCCAAAAAC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
10	3	2	2	4	2	0	2	2	2	0	0	1	2	1	1

GTTTGGCATGATTAAGGGAGTTCTTTGTGCTTC
GAAGCACAAAAAGAAAACTCCTTAATCATGCCAAAAAC

AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
11	3	4	4	5	2	0	2	2	1

k=2

GTTTGGCATGATTAAGGGAGTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
1	0	2	2	1	0	0	2	2	2	2	3	1	2	4	10



GAAGCACAAAAAGAAAACTCCTTAATCATGCCAAAAAC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
10	3	2	2	4	2	0	2	2	2	0	0	1	2	1	1

GTTTGGCATGATTAAGGGAGTTCTTTGTGCTTC
GAAGCACAAAAAGAAAACTCCTTAATCATGCCAAAAAC

AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
11	3	4	4	5	2	0	2	2	1

→ PALINDROMES :)

k=2

GTTTGCGATGATTAAAGGGAGTTCTTTGTGCTTC

	AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
X	11	3	4	4	5	2	0	2	2	1

k=2

GTTTGCGATGATTAAAGGGAGTTCTTTGTGCTTC

	AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
X	11	3	4	4	5	2	0	2	2	1
Y										
Z										
L										
K										
M										

k=2

ACTTCCGCAGTCGGGCATTACGCGTTGTGGAATGA

	AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
X	11	3	4	4	5	2	0	2	2	1
Y	4	5	2	4	5	4	4	3	2	1
Z										
L										
K										
M										

k=2

AC TT GCGC AG TCGCGC ATTAC GCGT AGT GGAATAA

	AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
X	11	3	4	4	5	2	0	2	2	1
Y	4	5	2	4	5	4	4	3	2	1
Z	4	5	3	2	4	1	5	5	2	3
L										
K										
M										

k=2

GGAGCGTTTATTAGTACCGGTTTGAAGTTAAC

	AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
X	11	3	4	4	5	2	0	2	2	1
Y	4	5	2	4	5	4	4	3	2	1
Z	4	5	3	2	4	1	5	5	2	3
L	11	6	3	2	2	3	2	1	1	4
K										
M										

k=2

GCCGCGAGCGGCCCGGCCGGCTTCGGCGCCGCAC

	AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
X	11	3	4	4	5	2	0	2	2	1
Y	4	5	2	4	5	4	4	3	2	1
Z	4	5	3	2	4	1	5	5	2	3
L	11	6	3	2	2	3	2	1	1	4
K	1	1	2	2	1	8	9	10	0	0
M										

k=2

GGGCCTGCGCCGGTCCAGTCACCCGGCTGCGACCT

	AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
X	11	3	4	4	5	2	0	2	2	1
Y	4	5	2	4	5	4	4	3	2	1
Z	4	5	3	2	4	1	5	5	2	3
L	11	6	3	2	2	3	2	1	1	4
K	1	1	2	2	1	8	9	10	0	0
M	0	4	4	3	4	10	4	5	0	0

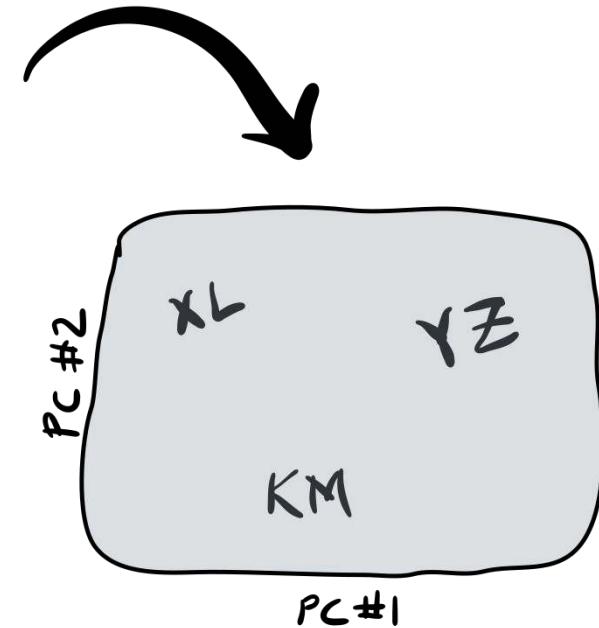
k=2

	AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
X	11	3	4	4	5	2	0	2	2	1
Y	4	5	2	4	5	4	4	3	2	1
Z	4	5	3	2	4	1	5	5	2	3
L	11	6	3	2	2	3	2	1	1	4
K	1	1	2	2	1	8	9	10	0	0
M	0	4	4	3	4	10	4	5	0	0

k=2

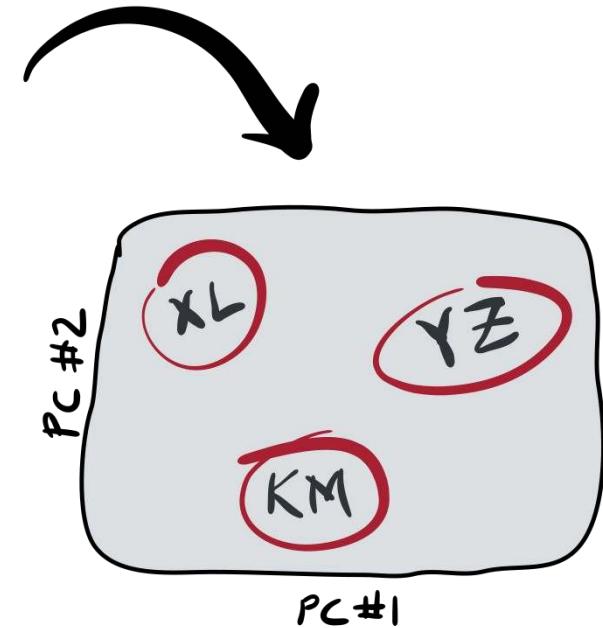
	AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
X	11	3	4	4	5	2	0	2	2	1
Y	4	5	2	4	5	4	4	3	2	1
Z	4	5	3	2	4	1	5	5	2	3
L	11	6	3	2	2	3	2	1	1	4
K	1	1	2	2	1	8	9	10	0	0
M	0	4	4	3	4	10	4	5	0	0

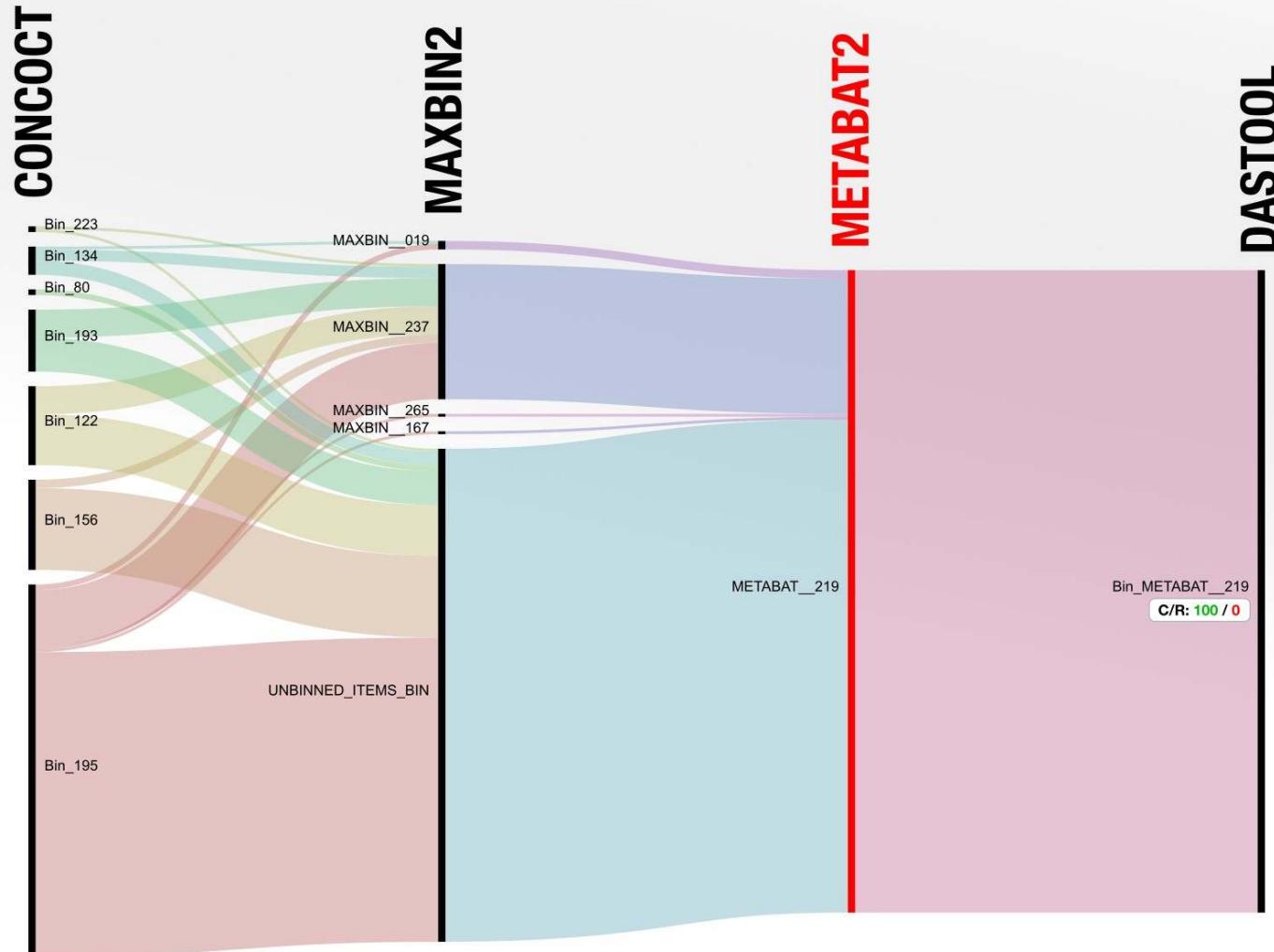
k=2



	AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
X	11	3	4	4	5	2	0	2	2	1
Y	4	5	2	4	5	4	4	3	2	1
Z	4	5	3	2	4	1	5	5	2	3
L	11	6	3	2	2	3	2	1	1	4
K	1	1	2	2	1	8	9	10	0	0
M	0	4	4	3	4	10	4	5	0	0

k=2





./anvi-script-gen-alluvial --algorithm DASTOOL --bin Bin_METABAT__219

Checking Bin quality: CheckM

CheckM uses single-copy marker genes to estimate a bin completeness and contamination

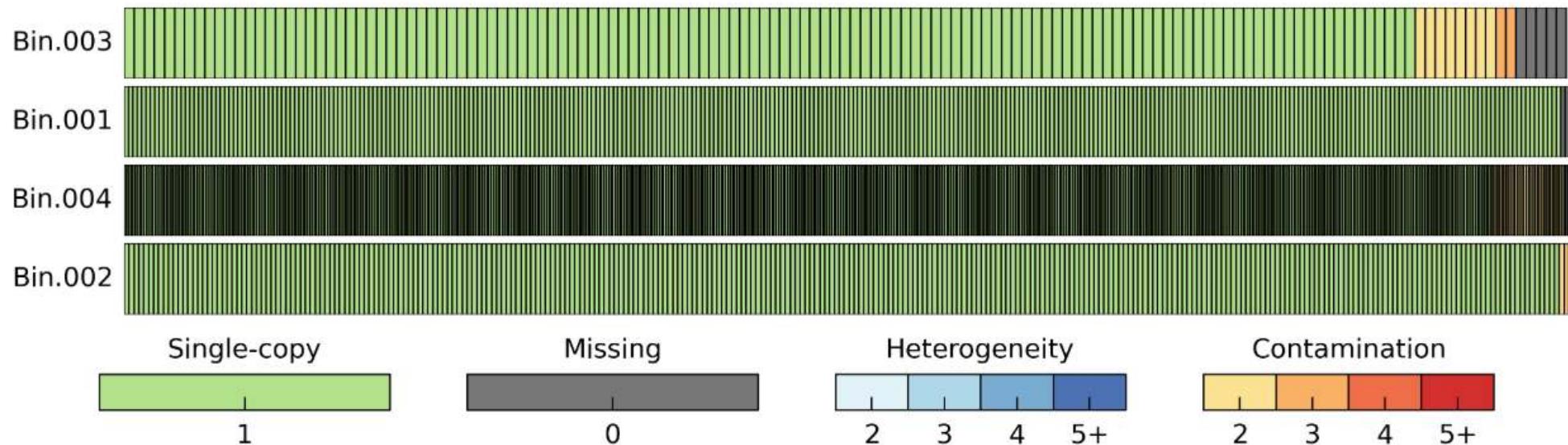
Evaluating the quality of obtained MAGs is very important in order to minimize the likelihood of extracting chimeric information

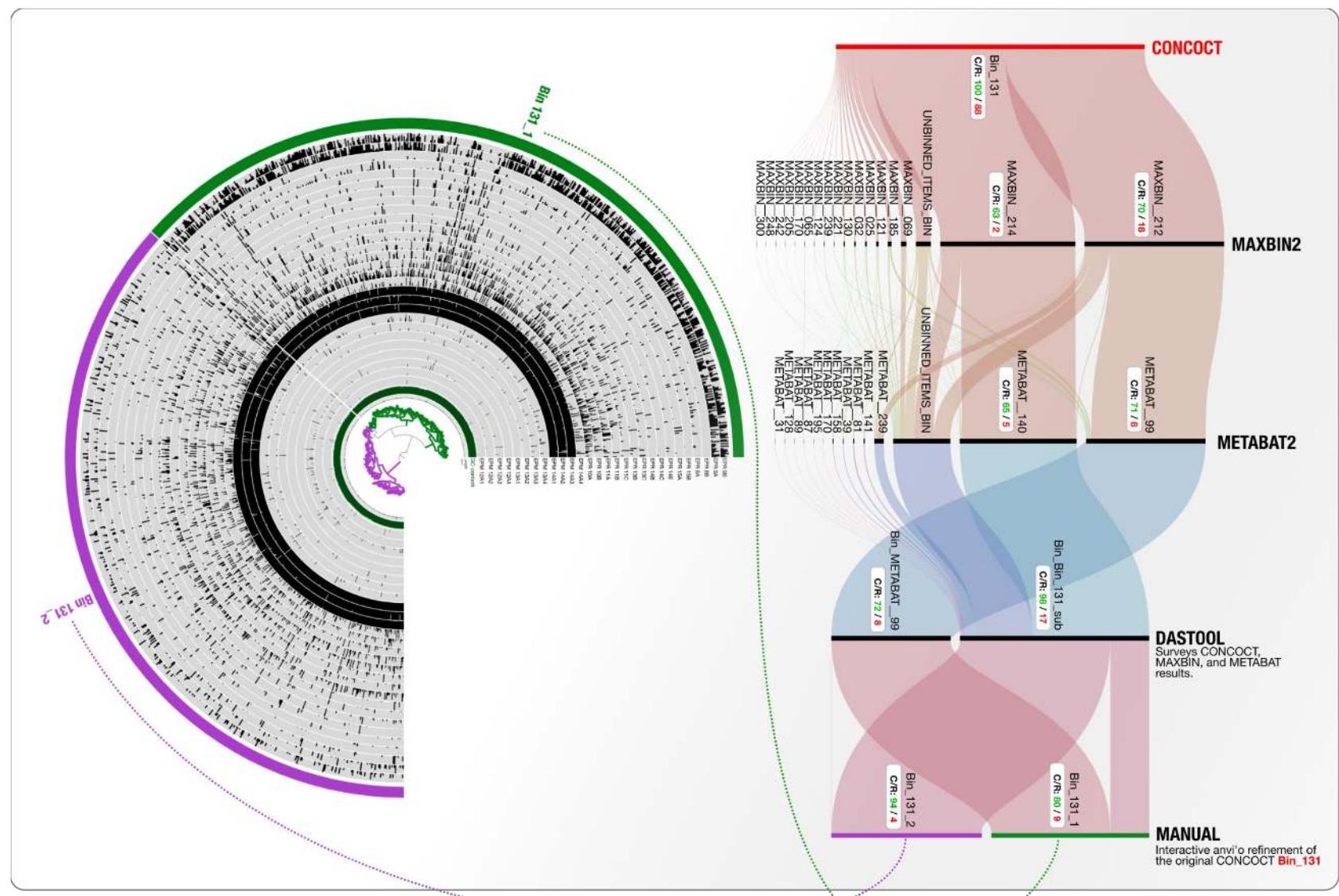
A manual inspection and refinement of the MAGs usually can improve the quality of the obtained bins

Checking Bin quality: CheckM

CheckM uses single-copy marker genes to estimate a bin completeness and contamination

[CheckM PLOT](#) | [CheckM Table](#)





Annotation: assigning functions

Annotation: assigning functions

Annotation of sequences (both for functional and taxonomic assignment) is usually done through sequence similarity

Basically, if a sequence is similar (within a certain pre-selected threshold) to a “known” sequence it gets annotated with the same function (or taxonomic assignment)

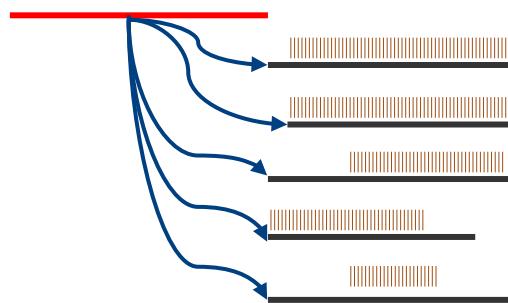
In practice, annotation is complex and can be approached in many different ways, and even similarity-based annotation can have many approaches

Several problems are derived from the functional and taxonomic annotation of sequences

- Find related or similar sequences by mapping letters of two sequences, with some spacers (indels),

```
76 GGMLKPIEGGTYEVNEAMVEDLKIGVQGPHASNLLGGILSNEIAKEIGKRAFIVDPVVVDE 135
      ||||||| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
61 GGMLKPIEGGTYEVNEAMVEDLKIGFEGPHAXNLGGILSNEIAKKLGKRAFIVDPVVVDX 120
```

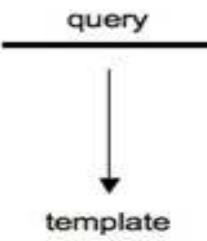
- Parse similarities, determine “best hit”



- Examples of pairwise search tools - Basic Local Alignment Search Tool or BLAST, LAST, etc.

A.

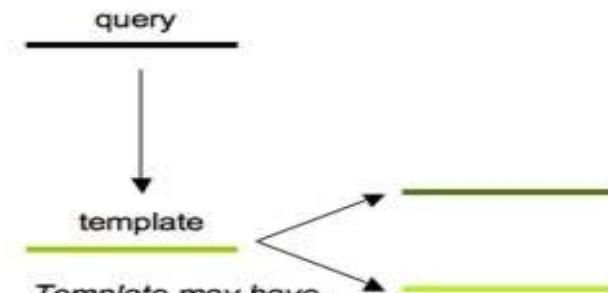
Paralogs problem



Template is a paralog, more likely have diverged functionally

B.

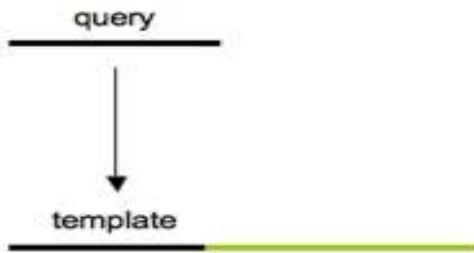
Moonlighting problem



Template may have more than one function

C.

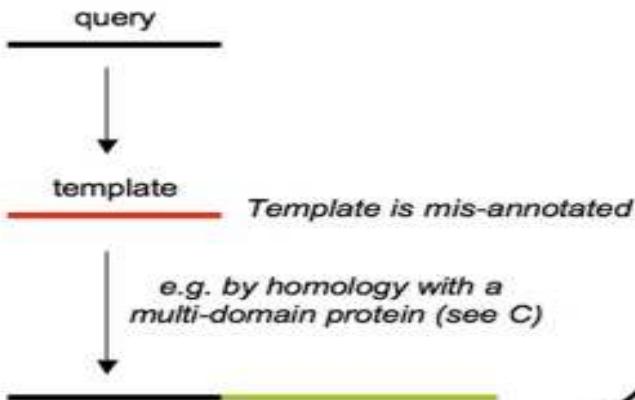
Multi-domain proteins problem



Template annotation may be based on a non-matching domain

D.

Database mis-annotations problem



Tully, USC

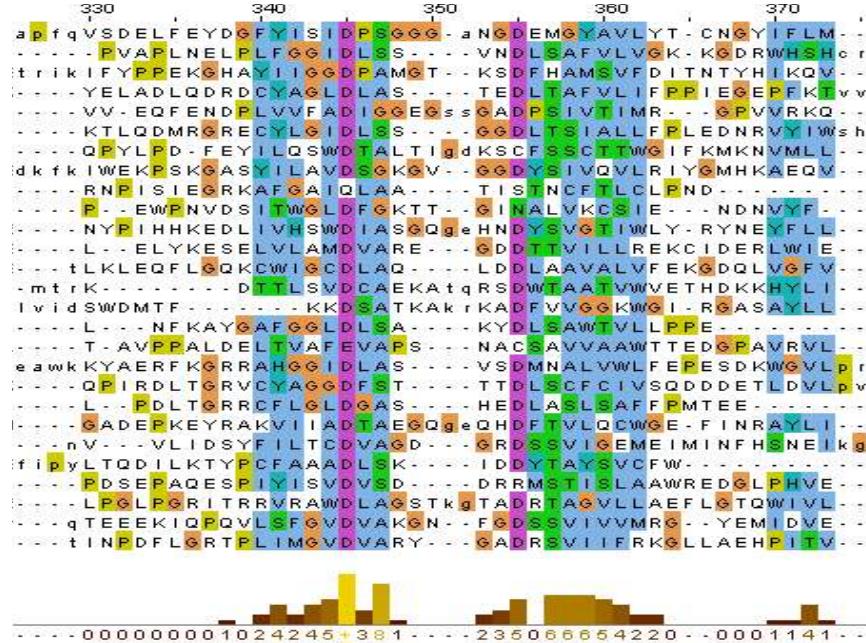
- “Feature prediction” - location of genes and other features, such as noncoding (e.g., rRNA, tRNA), protein coding sequences (CDS), and more (e.g., regulatory sites, repeats, frameshifts)
- Next, an attempt to “name” and interpret function/role by comparing against functional “databases”
 - NCBI Guidelines on “names” (~20 rules) e.g.,
 - “concise name, not a description or phrase”
 - “The protein name should not contain specific characteristics of the protein, (e.g., subcellular location, domain structure...”)
- Archives of accumulated biological data and knowledge
 - Genome, gene sequences, mutations
 - Gene regulation, expression, splice variants
 - Protein sequence, post-translational modifications
 - Protein tertiary structure, localization, networks
 - Enzyme kinetics, metabolites, metabolic networks
 - E.g., nr, swissprot, pdb, Pfam



Tully, USC

Other approach to annotation

- Alternative to Blast
- Multiple sequence alignment (MSA) of known sequences – detect “regions of similarity” – build consensus
- Use structural and mechanistic information (catalytic sites)
- Generate profiles using Hidden Markov Models (HMMs) or Position-Specific Scoring Models (PSSMs)
- More sensitive than pairwise – detect distant relationships
- Example of profile/HMM search tools: RPS-BLAST, Hmmer

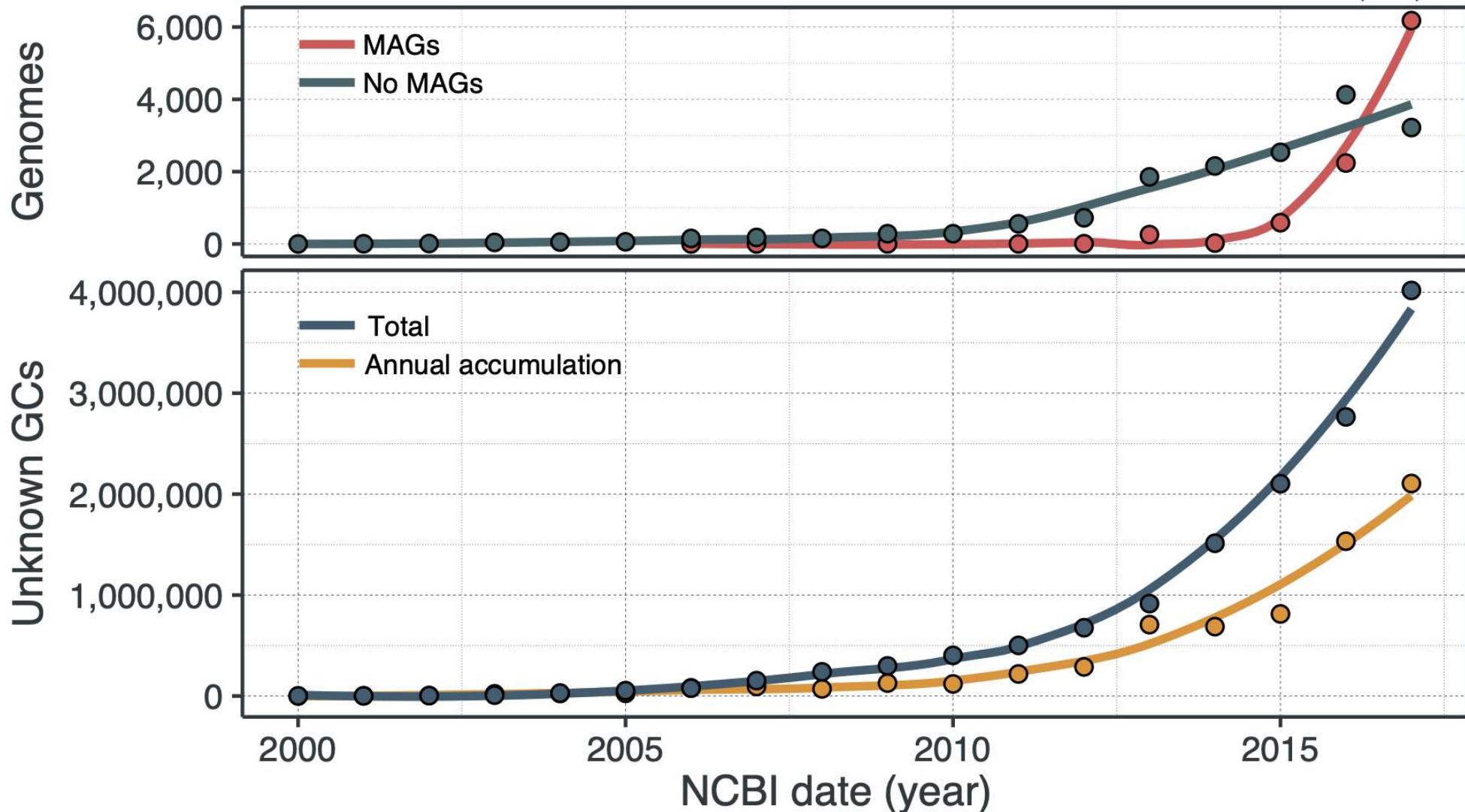


The ONLY way to confirm the function of a coding gene is through biochemical analyses and functional assays. No matter of similar two sequences are there might still be a functional difference.

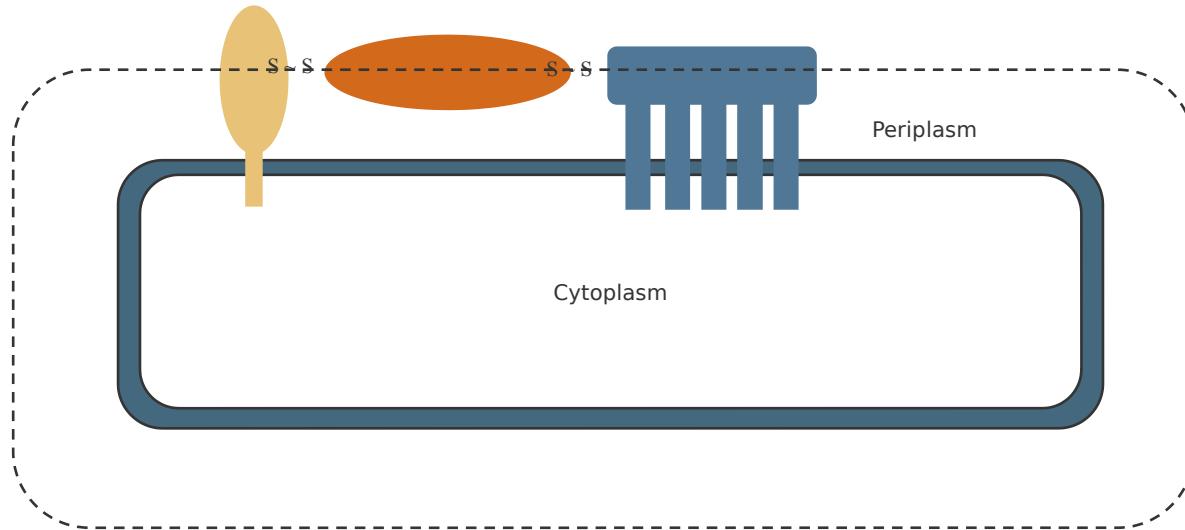
A large portion of the identified genes codes for hypothetical proteins (proteins with no known function). In the E.coli genome these account for ~30% of the total coding genes. In metagenome this number can be higher.

If the gene has never been seen before is usually called and ORFan

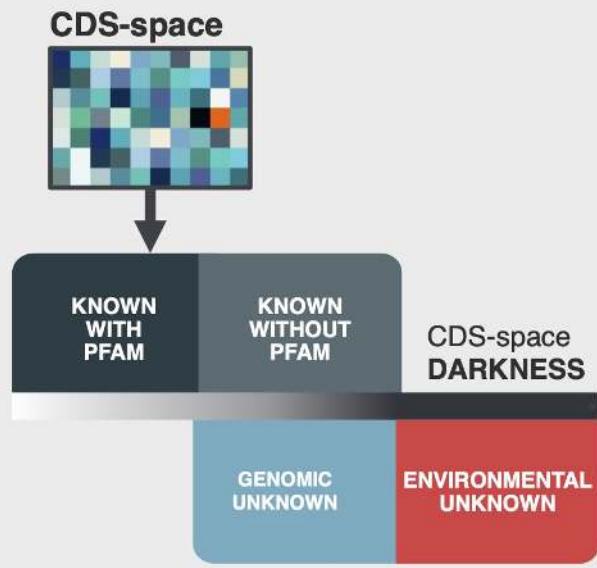
Genomes: 25,717 (GTDB r86)
Unknown GCs: 4,017,947



- Gene Context
- Sub-cellular localization – PSORTb
- Topological features
- Prediction of binding residues – ex. CAZy: Carbohydrate Binding Motifs



CONCEPTUAL FRAMEWORK



Known with Pfam

GC annotated to contain one or more Pfam entries but excluding DUFs

Known without Pfam

GC that have a known function but lack a Pfam annotation

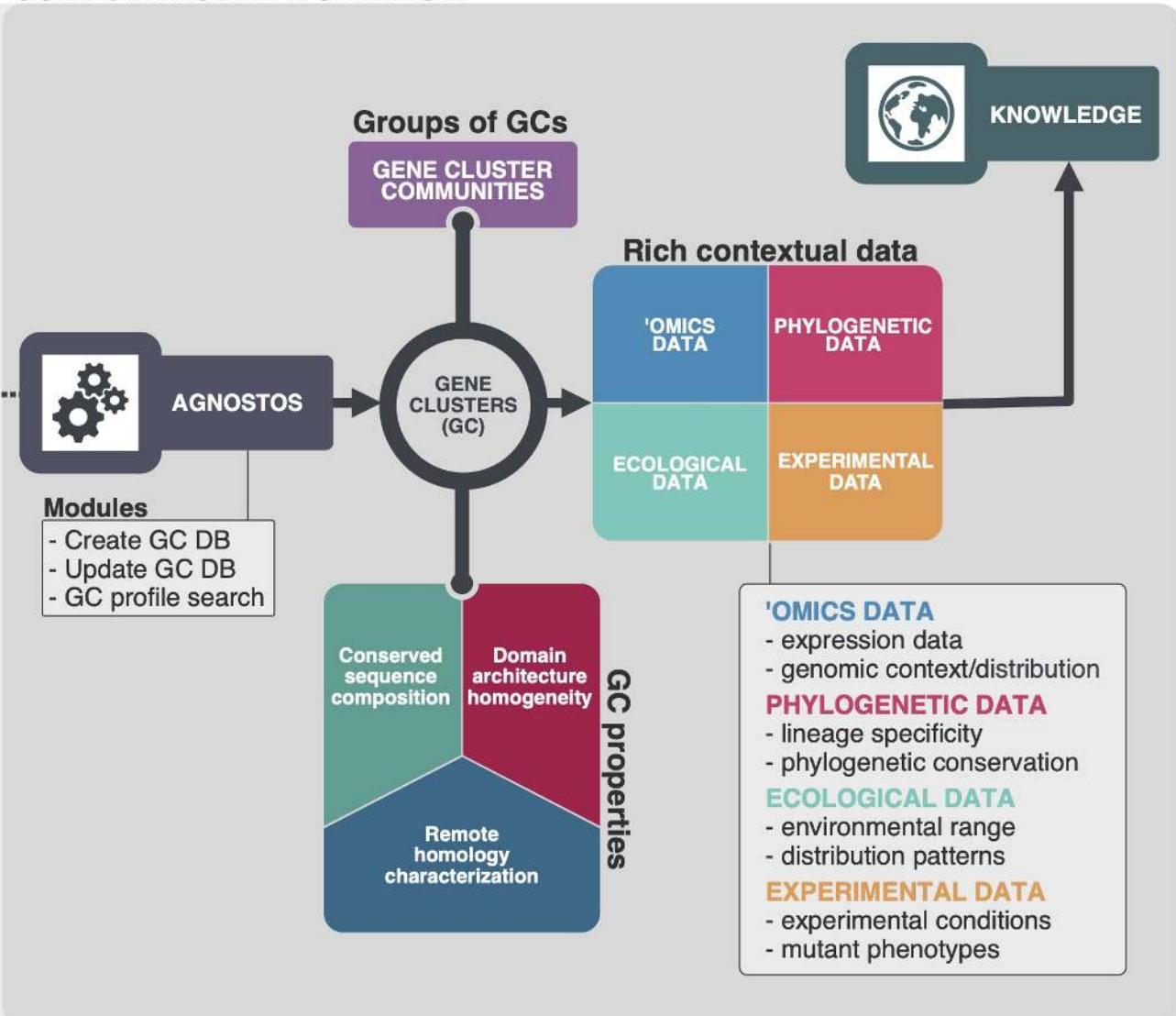
Genomic unknown

GC of unknown function (DUFs are included here) and found in sequenced or draft genomes

Environmental unknown

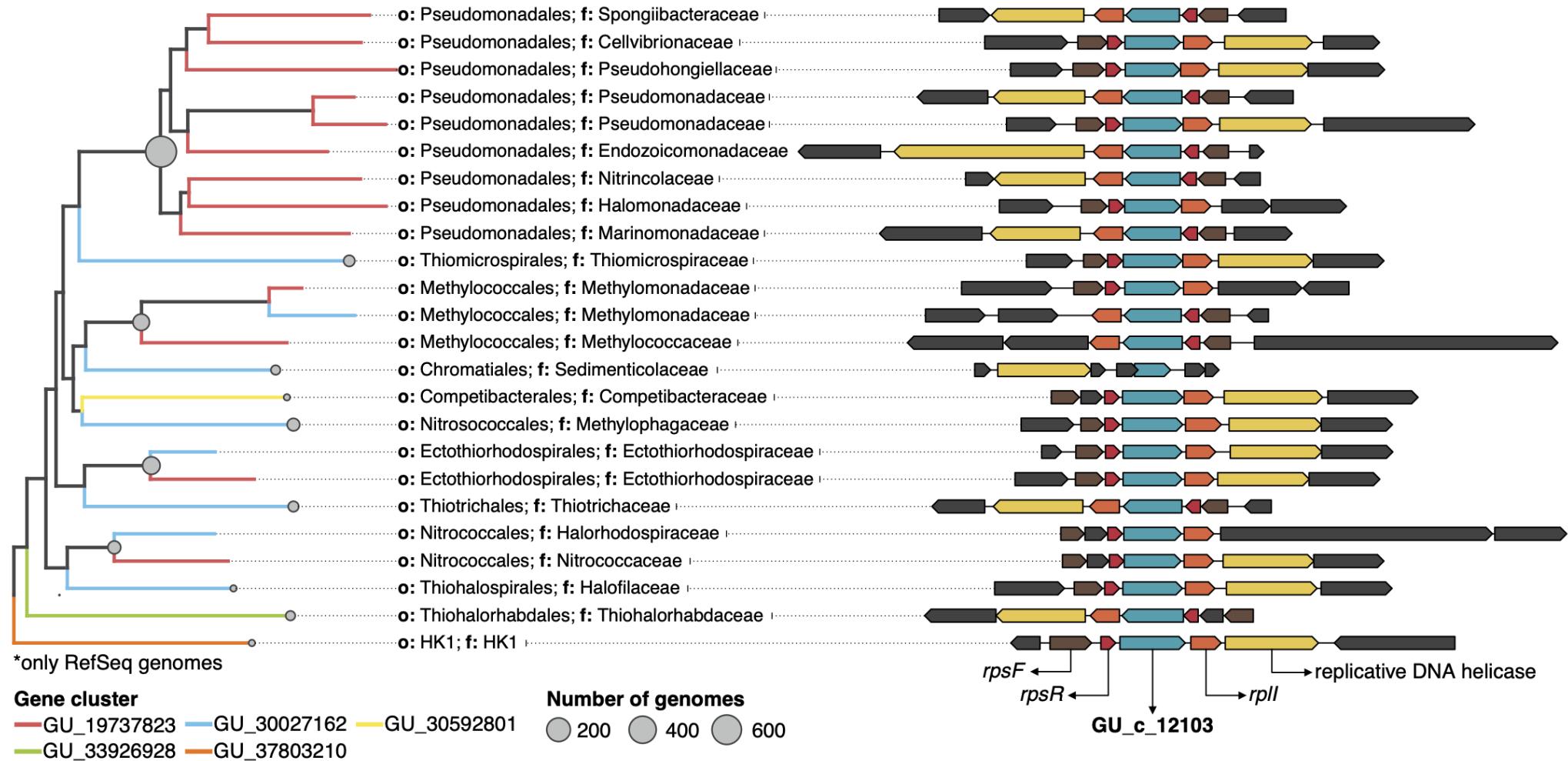
GC of unknown function not detected in sequenced or draft genomes, but only in environmental metagenomes or MAGs

COMPUTATIONAL WORKFLOW



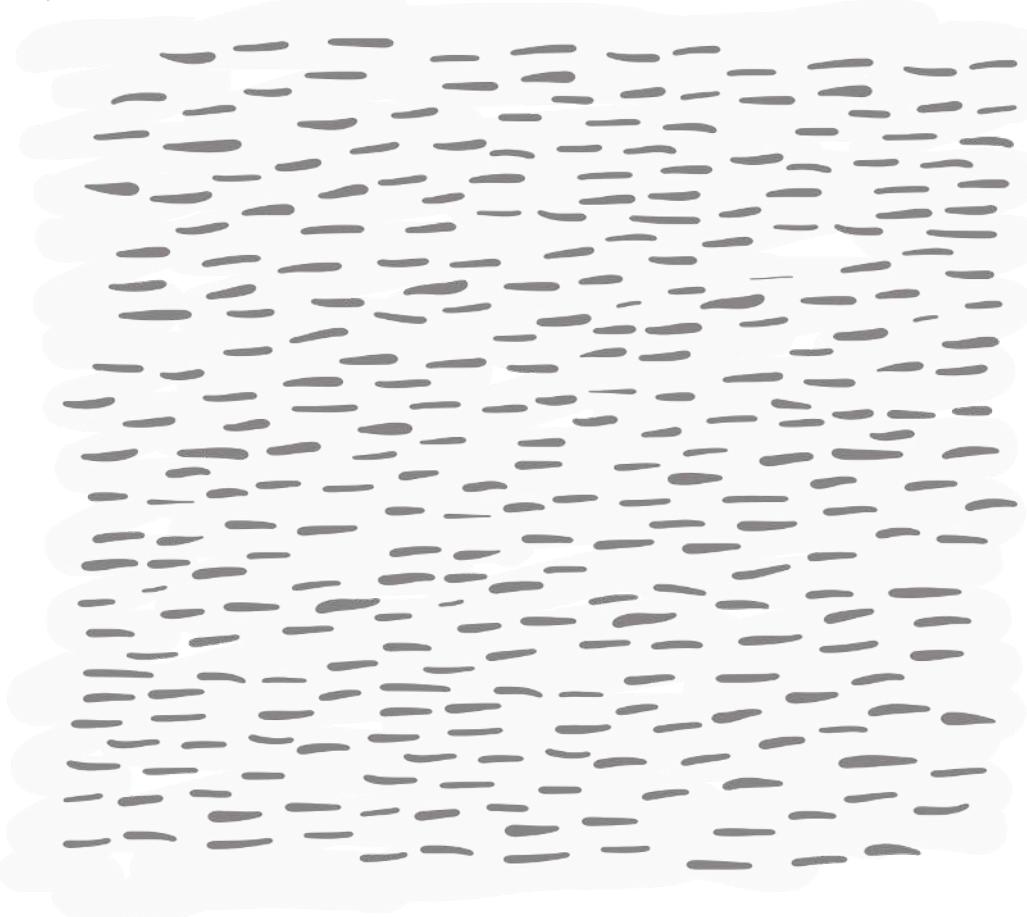
GTDB r86 | Bacteria; Proteobacteria; Gammaproteobacteria*

Genomic architecture surrounding GCC::GU_c_21103 genes

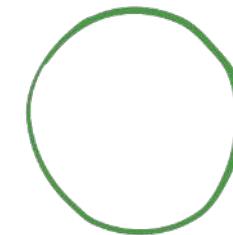
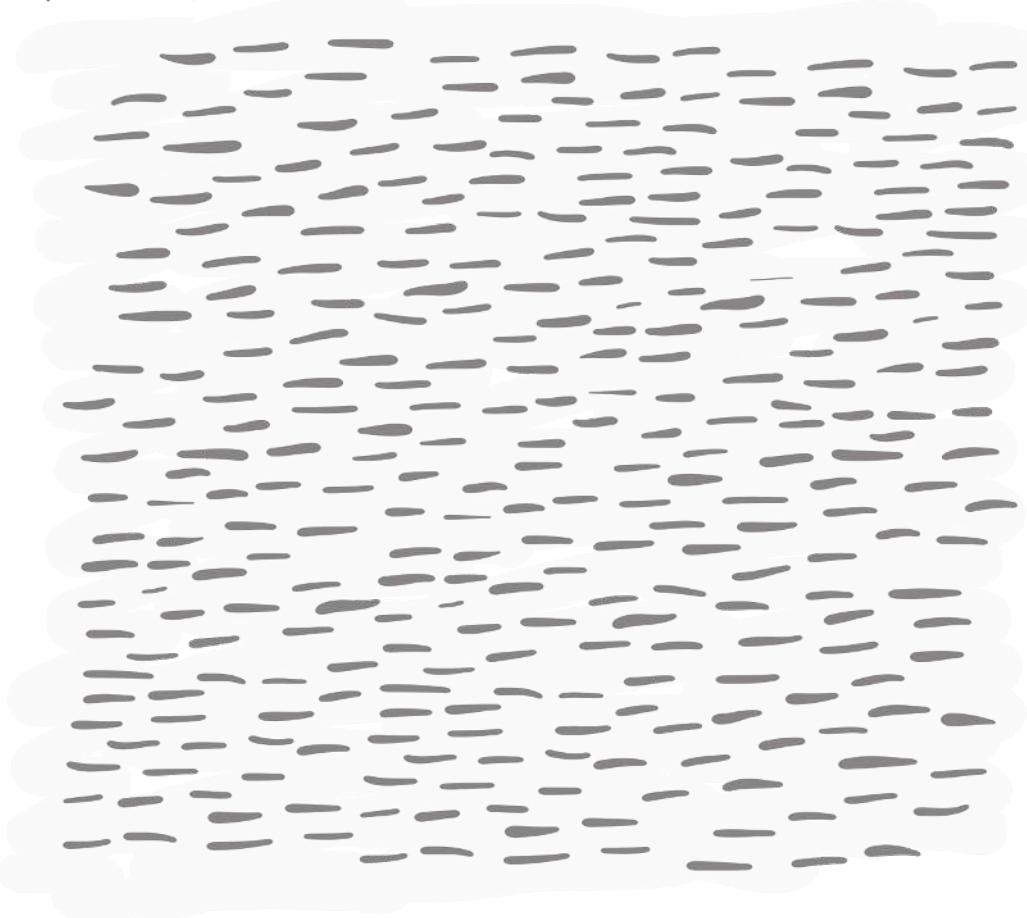


Read mapping: quantify functions

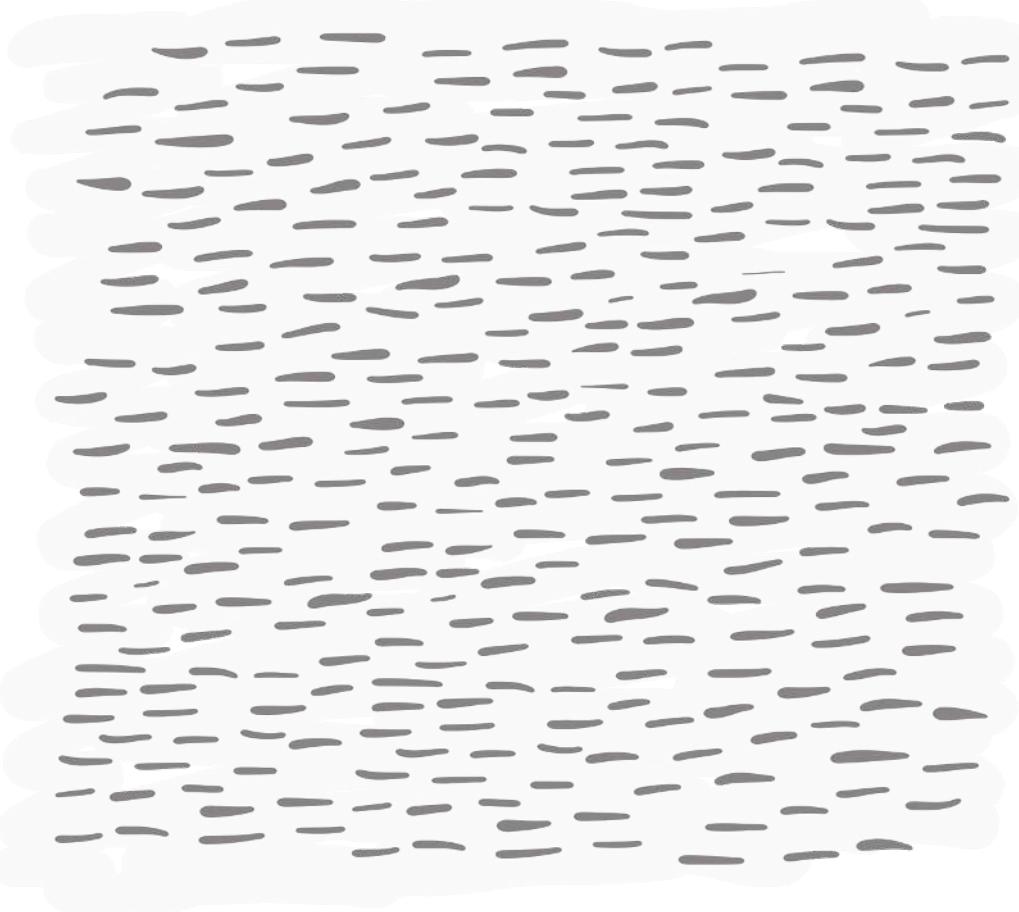
METAGENOMIC SHORT READS



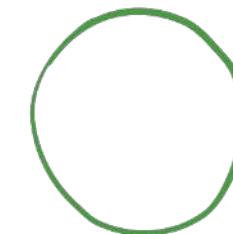
METAGENOMIC SHORT READS



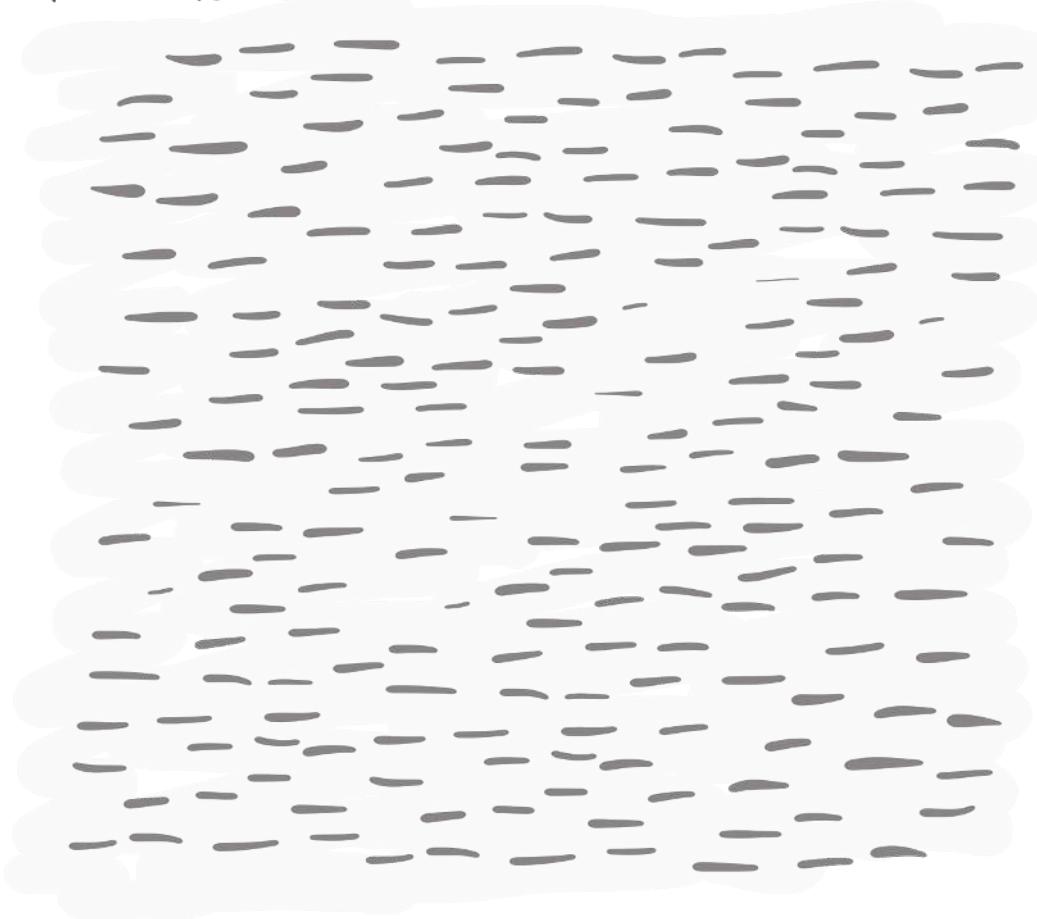
METAGENOMIC SHORT READS



READ
RECRUITMENT →



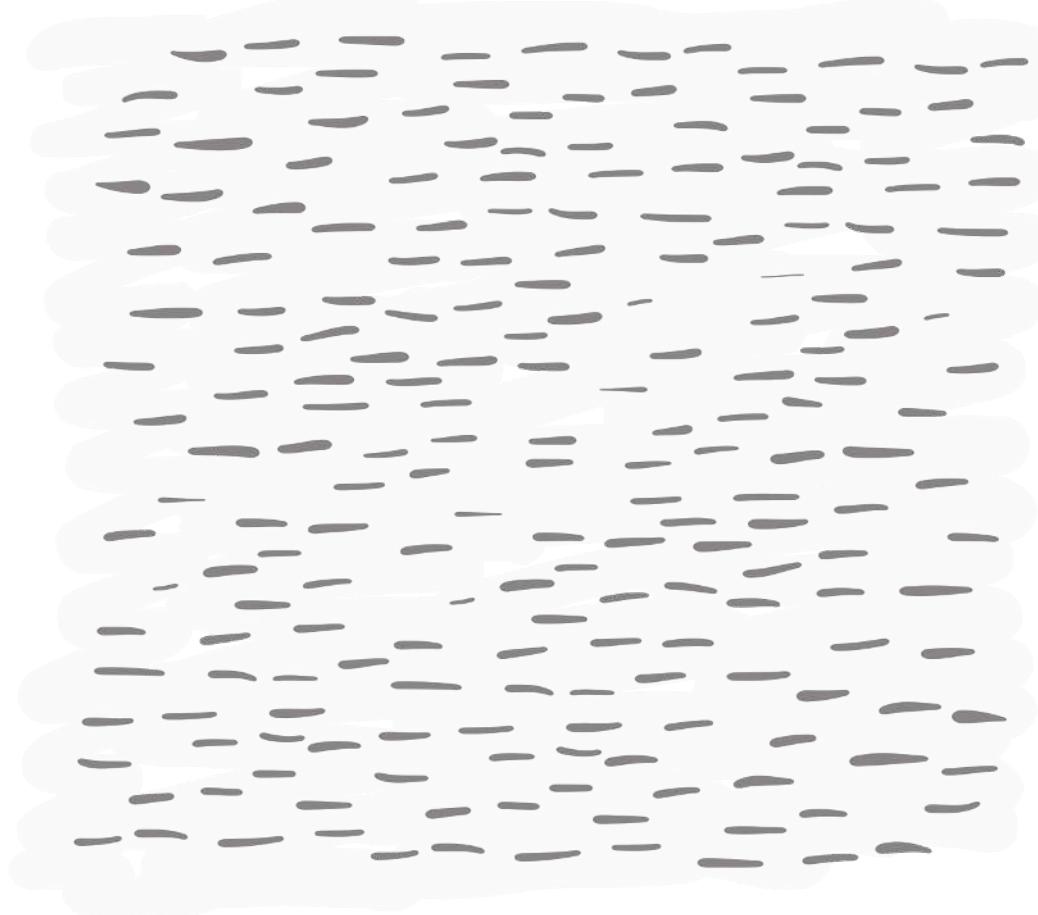
METAGENOMIC SHORT READS



READ
RECRUITMENT →



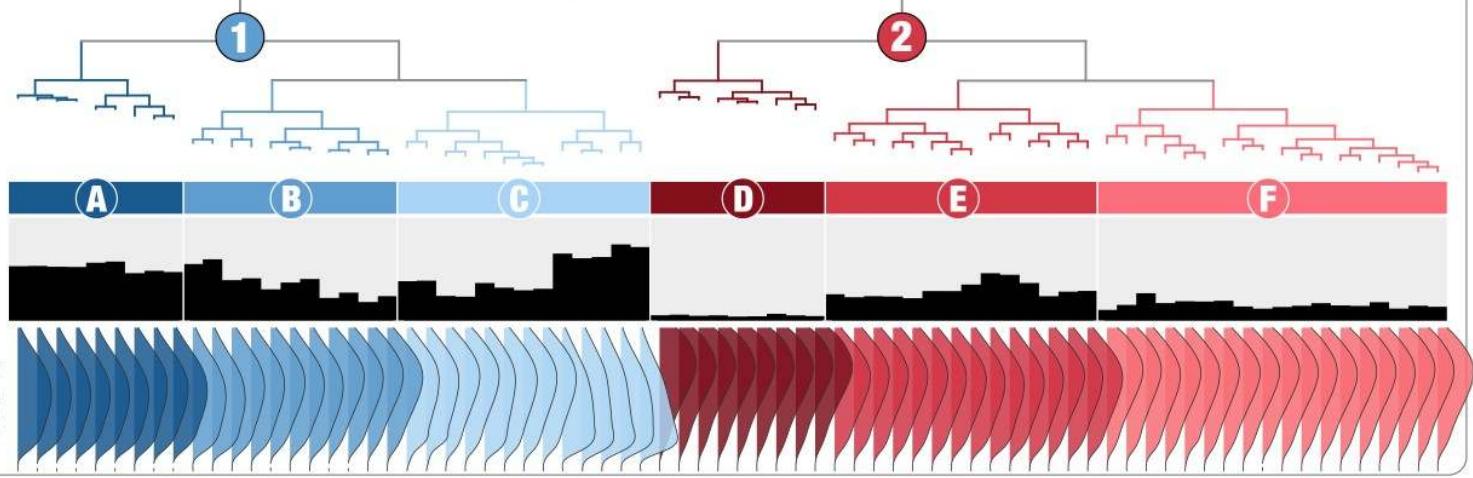
METAGENOMIC SHORT READS



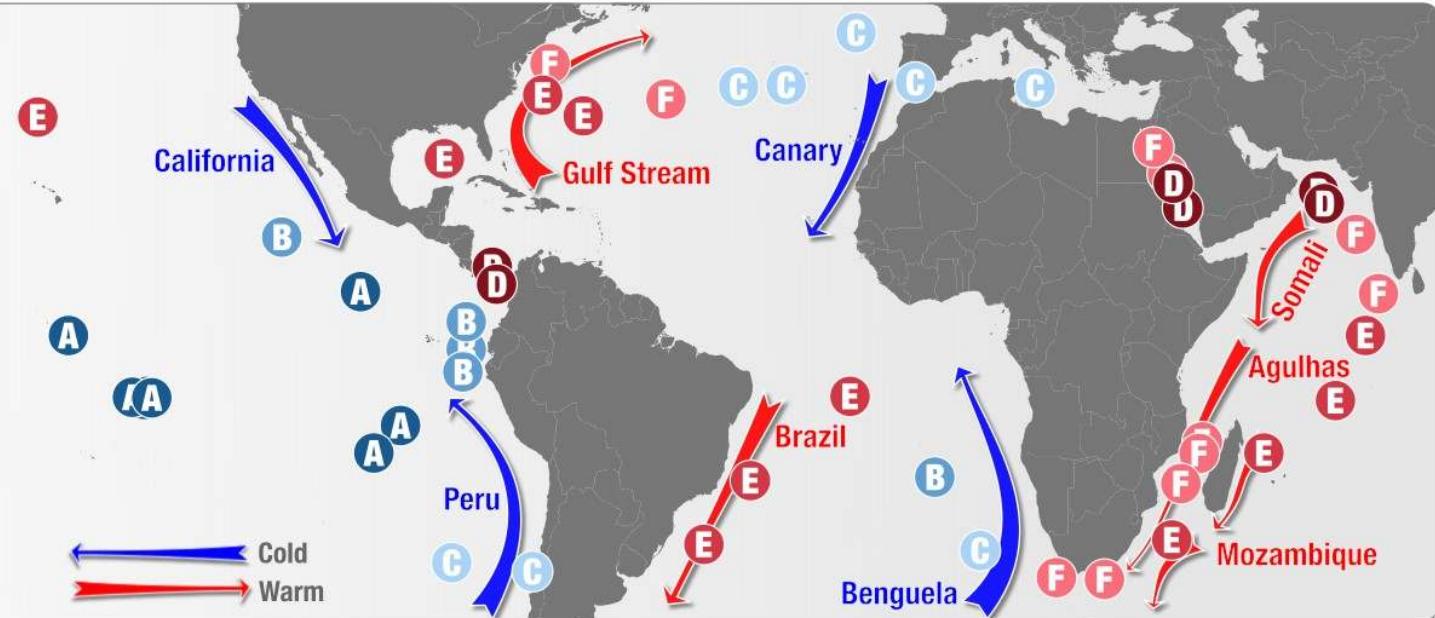
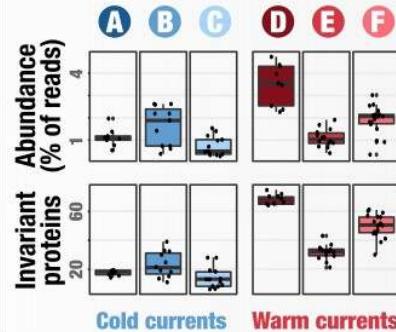
Results

- A count table with the abundance of specific functions
- A count table with the abundance of specific MAGs
- The reconstruction of specific metabolic pathways
- Functional genes to be used for phylogenetic analyses
- Cluster of genes to be analyzed for their variability, adaptations or structural variations

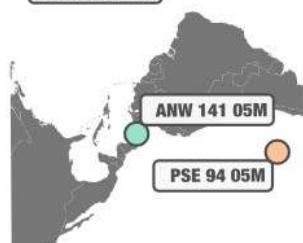
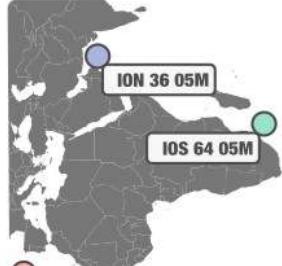
A Clustering of metagenomes based on SAAVs profiles



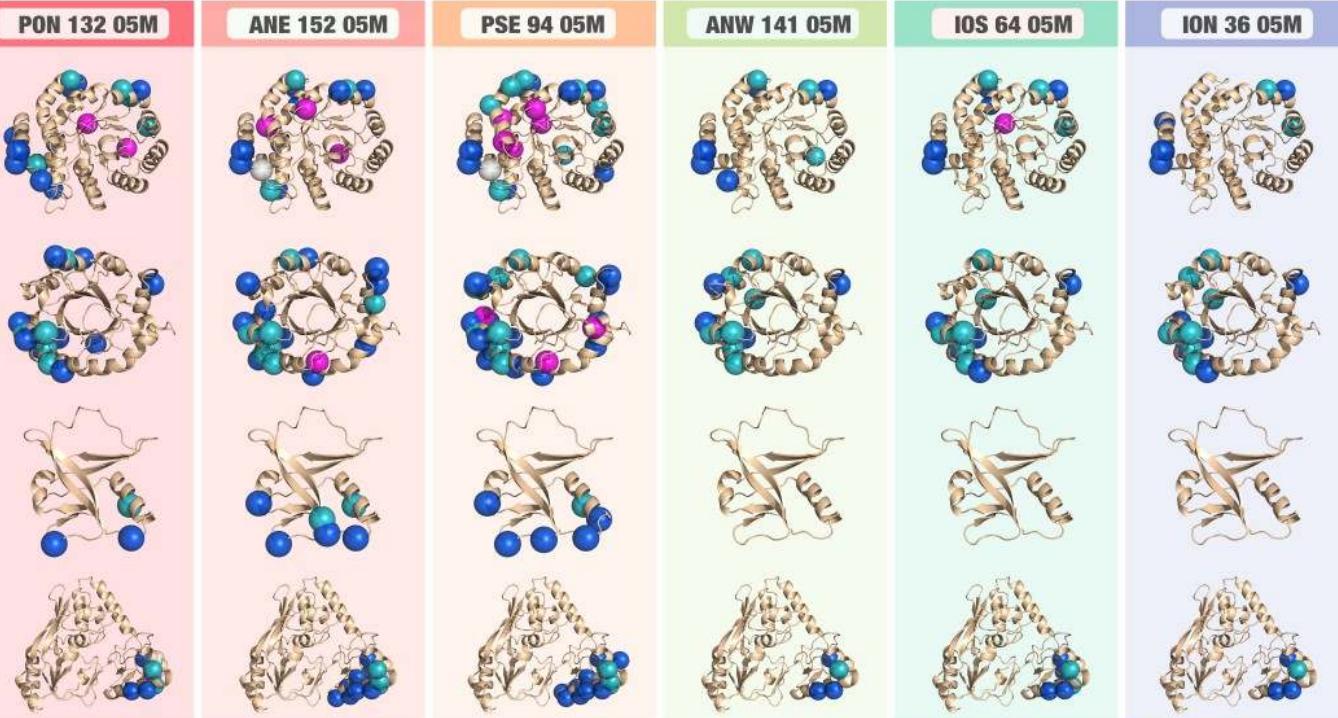
B Biogeography of 1a.3.V proteotypes (surface samples)



C Example proteins with SAAVs



Gene 1313
Dihydrodipicolinate synthetase family



Solvent accessibility: ● Exposed ● Intermediate ● Buried ● NA



KBase - <https://www.kbase.us/>

The KBase platform integrates several key features:

- Data**: Represented by a folder icon, the Data section shows a list of recent genome assemblies and pangenomes for *Shewanella* strains, such as *Shewanella_tree*, *Shewanella_oneidensis_MR1_NCB1*, and *Shewanella_amazonensis_SB2B_NCB1*.
- Apps**: Represented by a circular arrow icon, the Apps section lists various analytical tools including Annotation, Assembly, Communities, Comparative Genomics, Expression, Metabolic Modeling, Reads, Sequence, Uncategorized, Upload, and Util.
- Analysis Steps**: Represented by a stack of cubes icon, the Analysis section shows a detailed phylogenetic tree for the *Shewanella* unknown strain, comparing it against other strains like *Shewanella amazonensis* and *Shewanella oneidensis*.
- Sharing**: Represented by a share icon, the Sharing section includes a "share" button in the top right corner of the main interface.
- Commentary**: Represented by a speech bubble icon, the Commentary section is shown as a narrative text box describing the comparative genomics analysis.
- Visuals**: Represented by a chart icon, the Visuals section is shown as a phylogenetic tree visualization.
- Custom Scripts**: Represented by a code editor icon, the Custom Scripts section is shown as a "Change layout" button on the phylogenetic tree interface.

Shewanella Comparative Genomics

In this Narrative, we continue the study of an unknown strain of *Shewanella* isolated from a site contaminated with iron and manganese. In the previous Narrative, we assembled and annotated the genome of the strain. In this Narrative, we use comparative genomics tools to compare the new unknown strain to other *Shewanella* strains available in KBase's public genome collection. This analysis allows us to assess the closest related strains to elucidate similarities and possibly unique characteristics of our newly discovered strain.

Species Insert Genome Into Species Tree

Add one or more genomes to the KBase species tree.

Input Objects

Genome: *Shewanella_unknown*

Parameters

Neighbor public genome count: 1 ≤ 10 ≤ 200

Output Objects

Output Tree: *Shewanella_unknown_tree*

Shewanella_unknown_tree

v1 - KBaseTrees.Tree-1.0

Change layout

Phylogenetic tree showing relationships between various *Shewanella* strains. Nodes are labeled with strain names and IDs, and support values are indicated at the nodes.

- Root node (0.837)
- Node 1 (0.771) leading to *Shewanella amazonensis* 5029 (v01g.26557)
- Node 2 (1.000) leading to *Shewanella* sp. ANA-3 (v01g.3779)
- Node 3 (0.709) leading to *Shewanella* sp. MR-4 (v01g.2626)
- Node 4 (0.709) leading to *Shewanella* sp. MR-7 (v01g.2627)
- Node 5 (1.000) leading to *Shewanella oneidensis* MR-1 (v01g.372)
- Node 6 (1.000) leading to *Shewanella violacea* (v01g.24813)
- Node 7 (1.000) leading to *Shewanella baltica* OS189 (v01g.1305)
- Node 8 (0.771) leading to *Shewanella baltica* OS195 (v01g.1285)
- Node 9 (0.771) leading to *Shewanella baltica* OS155 (v01g.26354)

Practical: read-based metagenomic analysis

Practical: read-based metagenomic analysis



- Google Scholar (<https://scholar.google.com/>) – To find the relevant paper



- NCBI SRA (<https://www.ncbi.nlm.nih.gov/sra>) – To retrieve the accession number and sequences



- Kaiju (<http://kaiju.binf.ku.dk/>) – To taxonomically annotate the reads



- mi-faser (<https://services.bromberglab.org/mifaser/>) – To functionally annotate the reads



- KEGG Mapper (https://www.genome.jp/kegg/tool/map_pathway2.html) – To explore the metabolic pathways



- Excel or Google Sheet – To manipulate the results for plotting and statistical testing



- RawGraphs (<https://app.rawgraphs.io>) – To plot the results for visualization



- Powerpoint or Google Slides – To present the results

Web server - Submit job

Please use the form to upload the sequencing file(s).

Once uploading is completed, press the Submit button at the bottom of the page.
Only submit one data set at a time.

Job Name

You can give a custom name to your submission.

e-mail

Receive a notification after your submission has been processed. [?]

File with sequencing reads *

Nucleotide sequences must be in compressed FASTA or FASTQ format [?]

Select file

File name:

Start upload

Progress:

Upload a second file for paired-end sequencing

Options

Reference Database

- RefSeq Genomes - proteins from completely assembled RefSeq genomes: Bacteria, Archaea, Viruses
- proGenomes - proteins from the representative genomes in proGenomes: Bacteria, Archaea, Viruses
- NCBI BLAST nr - non-redundant protein database: Bacteria, Archaea, Viruses
- NCBI BLAST nr+euk - as above, but also including fungi and microbial eukaryotes.

SEG filter

- Filter low complexity protein query sequences

Run mode

- MEM - for maximum exact matches.
- Greedy - allows mismatches.

Minimum match length

11

Only applicable for Greedy mode:

Minimum match score

75

Allowed mismatches

5

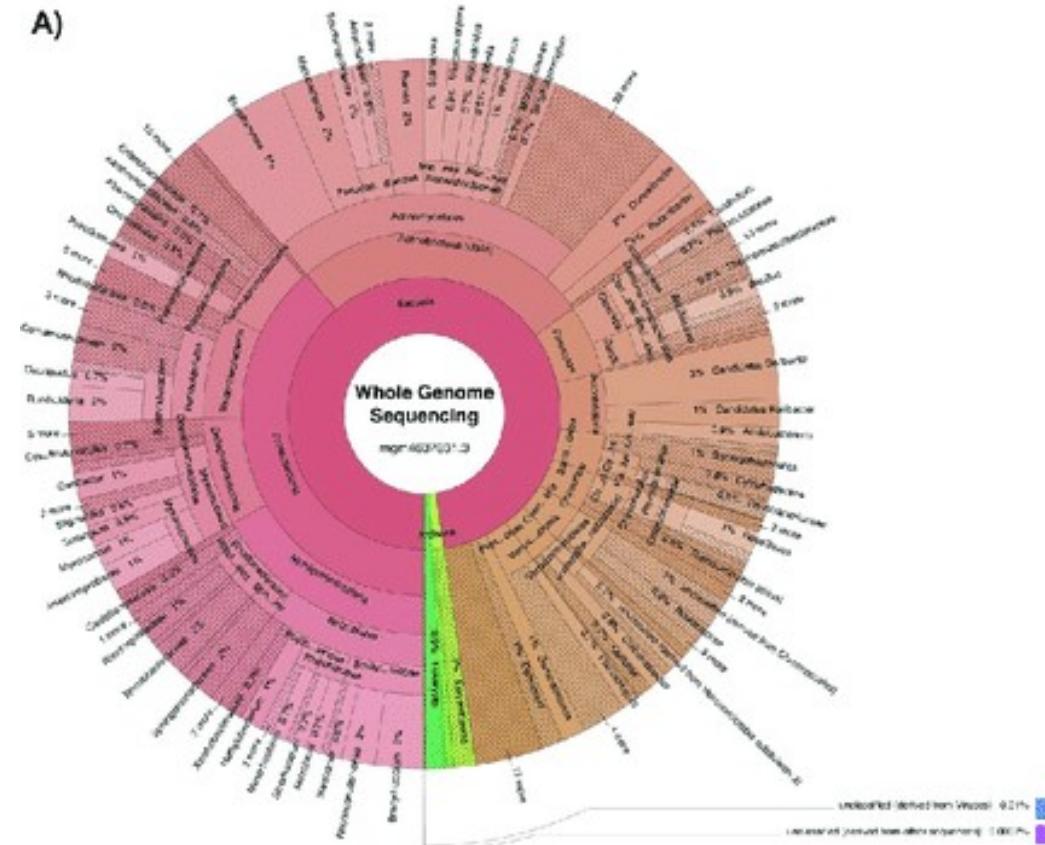
Submit

<http://kaiju.binf.ku.dk/server>

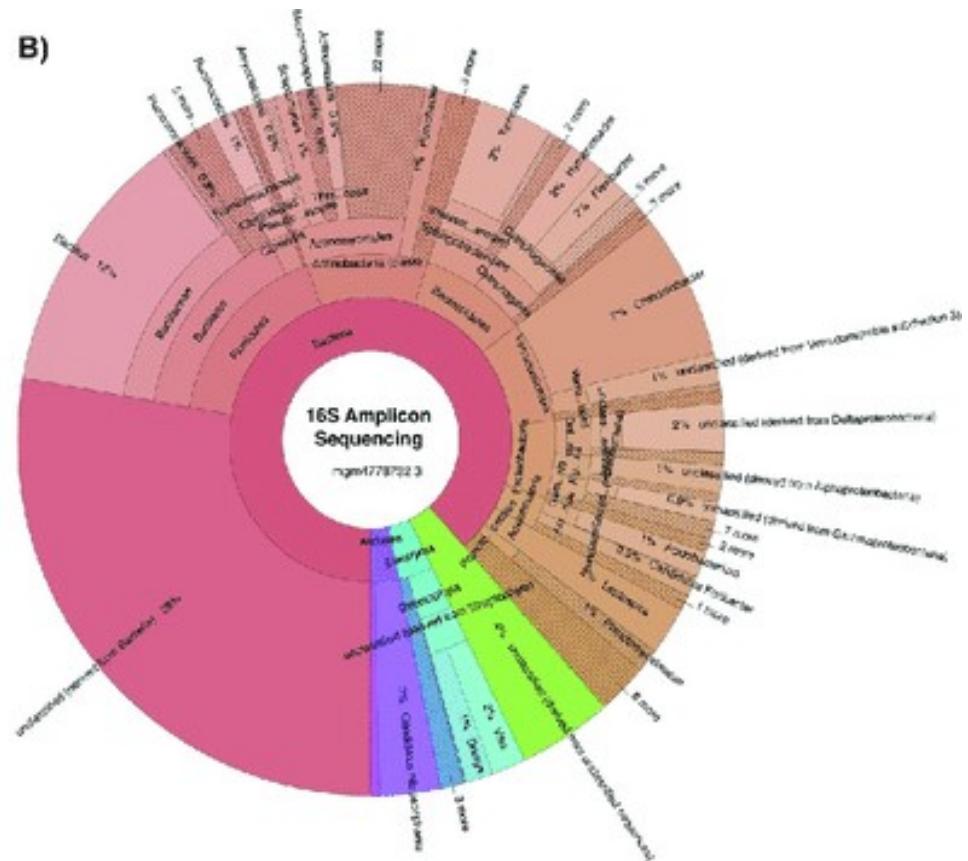


Krona

A)



B)



mi-faser, microbiome - functional annotation of sequencing reads

is a super-fast (< 20min/10GB of reads) and accurate (> 90% precision) method for annotation of molecular functionality encoded in sequencing read data without the need for assembly or gene finding.

Uploads

[File](#) [SRA](#) [Url](#)

Add a file or drag it here.

- up to 100 files at the same time
- **compressed** upload possible
- fasta/fastq format required

[+ Add file](#)

Note: Registered users can resume uploads within 3 hours.

[Start upload](#)[Pause upload](#)

Submission

Details of current submission.

Upload:

0%

Database

GS-21-all

reference db (2021)
proteins 15,524 / E.C.s 2,716

Options

Auto Submit

ON

Submit once upload complete

Paired-end mode

OFF

Submit paired-end files

Quality control

OFF

Quality control for fastq files

Read mapping

OFF

Save raw reads per E.C.

- The default reference database has been updated [03/2021] -

⌚ <https://services.bromberglab.org/mifaser/submit>

mi-faser, microbiome - functional annotation of sequencing reads

is a super-fast (< 20min/10GB of reads) and accurate (> 90% precision) method for annotation of molecular functionality encoded in sequencing read data without the need for assembly or gene finding.

Uploads

[File](#) [SRA](#) [Url](#)

Add a file or drag it here.

- up to 100 files at the same time
- **compressed** upload possible
- fasta/fastq format required

[+ Add file](#)

Note: Registered users can resume uploads within 3 hours.

[Start upload](#)[Pause upload](#)

Submission

Details of current submission.

[Upload:](#)

0%

Database

[GS-21-all](#)

reference db (2021)
proteins 15,524 / E.C.s 2,716

Options

[Auto Submit](#)[ON](#)

Submit once upload complete

[Paired-end mode](#)[OFF](#)

Submit paired-end files

[Quality control](#)[OFF](#)

Quality control for fastq files

[Read mapping](#)[OFF](#)

Save raw reads per E.C.

- The default reference database has been updated [03/2021] -

⌚ <https://services.bromberglab.org/mifaser/submit>

Submission uYEHN9Ls2DI9KGNJ

Summary statistics

Number of unique functions found:	746
Reads that map unambiguously to a function:	310231 (96%)
Reads that map to multiple functions ():	11489 (4%)

 Annotations Multiple Mappings Read Mappings Quality Control

Submission details

label CV88_metag

reads 9537740

created 2020-04-17 09:39:42

runtime 00h : 10m : 50s 

qa control 00h : 02m : 04s

file name CV88_metag.fastq.gz

file size 1988MB

database GS+

user donato.giovannelli@gmail.com

your jobs CV88_metag.fastq.gz

Annotated functions

 regexSearch 

10

Select all

Deselect all

Select Search

Show all

Filter selected

 Download filtered (csv)

Enzyme E.C.#	Annotation	Read Count	Color 
2.7.7.6	DNA-directed RNA polymerase	14011	
5.99.1.3	DNA topoisomerase (ATP-hydrolysing)	7867	
3.6.3.14	H ⁺ -transporting two-sector ATPase	7675	
3.6.4.13	RNA helicase	6223	
4.2.1.11	phosphopyruvate hydratase	6018	
1.7.99.4	nitrate reductase	4935	
3.6.4.12	DNA helicase	4161	
2.7.2.3	phosphoglycerate kinase	4032	
6.3.4.4	adenylosuccinate synthase	3782	

Submission uYEHN9Ls2DI9KGNJ

Summary statistics

Number of unique functions found:

746

Reads that map unambiguously to a function:

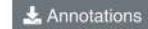
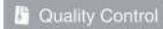
310231

(96%)

Reads that map to multiple functions ():

11489

(4%)

 Annotations Multiple Mappings Read Mappings Quality Control

Submission details

label CV88_metag

reads 9537740

created 2020-04-17 09:39:42

runtime 00h : 10m : 50s 

qa control 00h : 02m : 04s

file name CV88_metag.fastq.gz

file size 1988MB

database GS+

user donato.giovannelli@gmail.com

your jobs CV88_metag.fastq.gz

Annotated functions

 regexSearch 

10

Select all

Deselect all

Select Search

Show all

Filter selected

 Download filtered (csv)

Enzyme E.C.#	Annotation	Read Count	Color 
2.7.7.6	DNA-directed RNA polymerase	14011	
5.99.1.3	DNA topoisomerase (ATP-hydrolysing)	7867	
3.6.3.14	H ⁺ -transporting two-sector ATPase	7675	
3.6.4.13	RNA helicase	6223	
4.2.1.11	phosphopyruvate hydratase	6018	
1.7.99.4	nitrate reductase	4935	
3.6.4.12	DNA helicase	4161	
2.7.2.3	phosphoglycerate kinase	4032	
6.3.4.4	adenylosuccinate synthase	3782	

[About KEGG Mapper](#)[Reconstruct Pathway
\(and Brite, Module\)](#)[Search Pathway
\(and Brite, Module,
Network, Disease\)](#)[Search&Color Pathway
\(and Brite, Module\)](#)[Color Pathway](#)
[Join Brite](#)[Convert ID](#)
[Annotate Sequence](#)[BlastKOALA](#)
[Map Taxonomy](#)
[KEGG](#)**Target databases:** Pathway, Brite hierarchy, Module**Search mode:** Reference Organism-specifichsa

Enter: org, ko, ec, rn, hsadd

Optional use of outside ID: Select**Enter objects one per line followed by bgcolor, fgcolor:**

Examples:

 Select**Or upload file:** Choose File No file chosen**If necessary, change default bgcolor:** pink Include aliases Use uncolored diagrams Search pathways containing all the objects (AND search) Exec Clear

KEGG Mapper Search Result

[Pathway \(130\)](#)[Brite \(0\)](#)[Module \(212\)](#)[Sort by the pathway list](#)[Show matched objects](#)

map01100 Metabolic pathways (473)

map01110 Biosynthesis of secondary metabolites (188)

map01120 Microbial metabolism in diverse environments (141)

map00520 Amino sugar and nucleotide sugar metabolism (49)

map00230 Purine metabolism (35)

map00541 O-Antigen nucleotide sugar biosynthesis (33)

map00270 Cysteine and methionine metabolism (27)

map00010 Glycolysis / Gluconeogenesis (27)

map00620 Pyruvate metabolism (26)

map00720 Carbon fixation pathways in prokaryotes (25)

map00240 Pyrimidine metabolism (24)

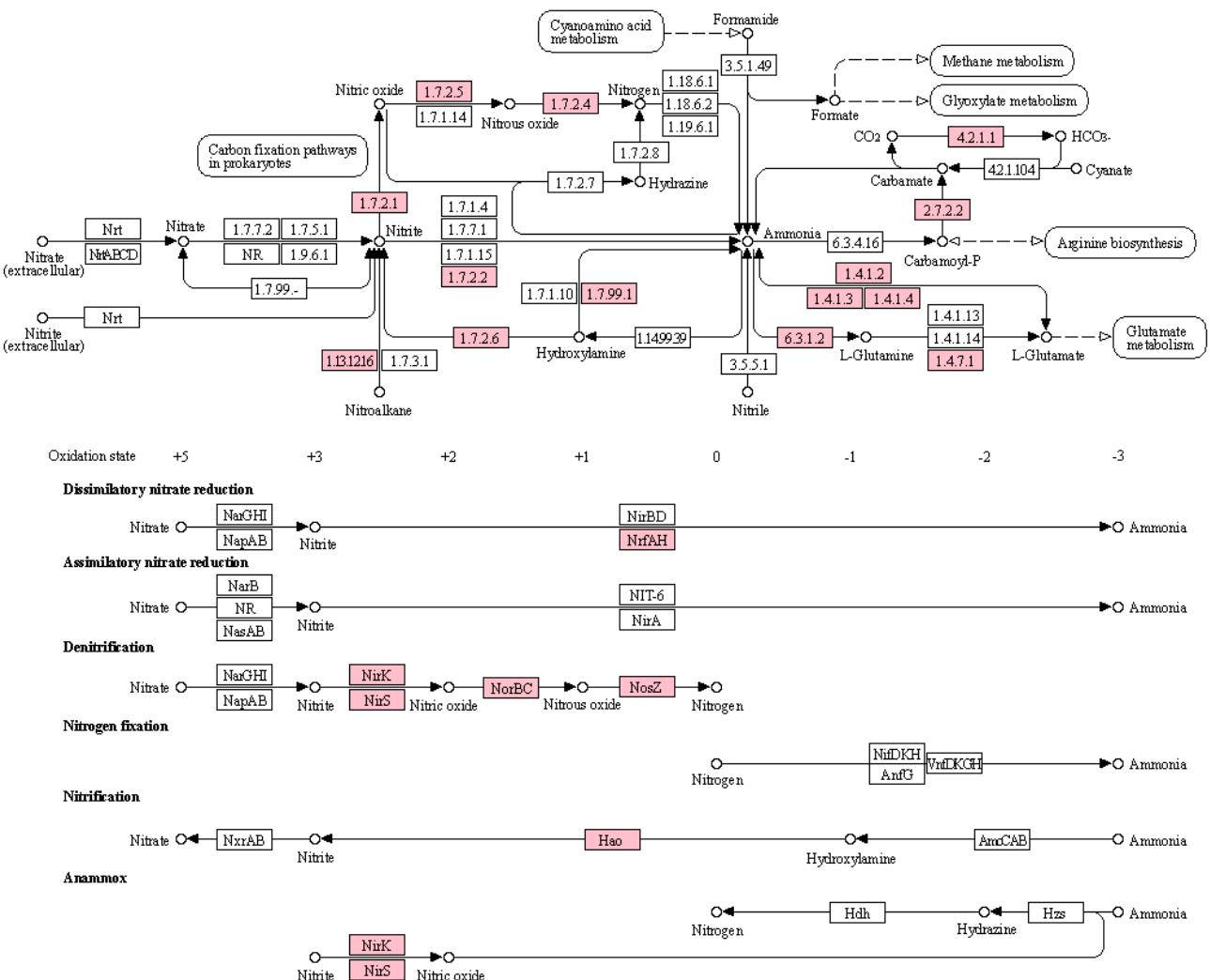
map00250 Alanine, aspartate and glutamate metabolism (23)

map00630 Glyoxylate and dicarboxylate metabolism (23)

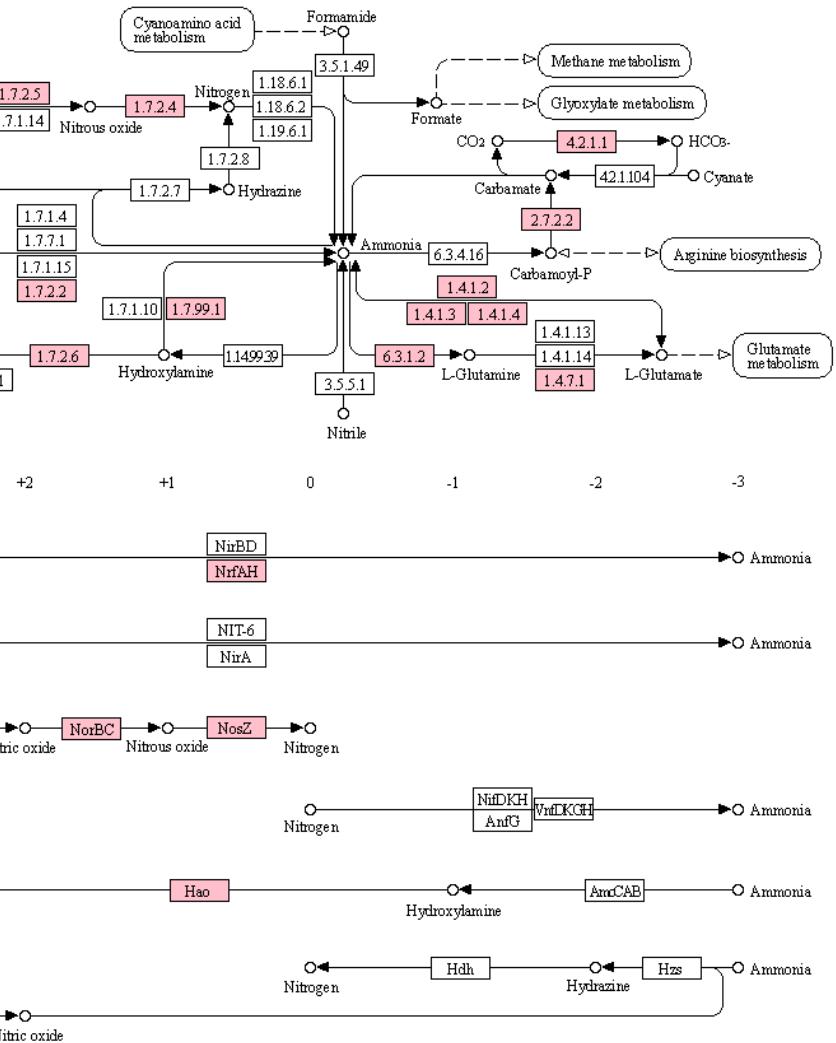
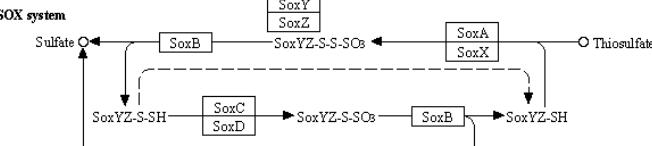
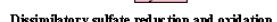
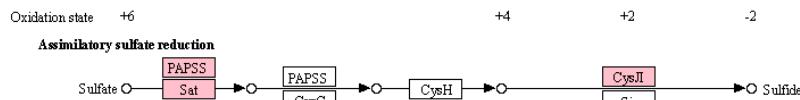
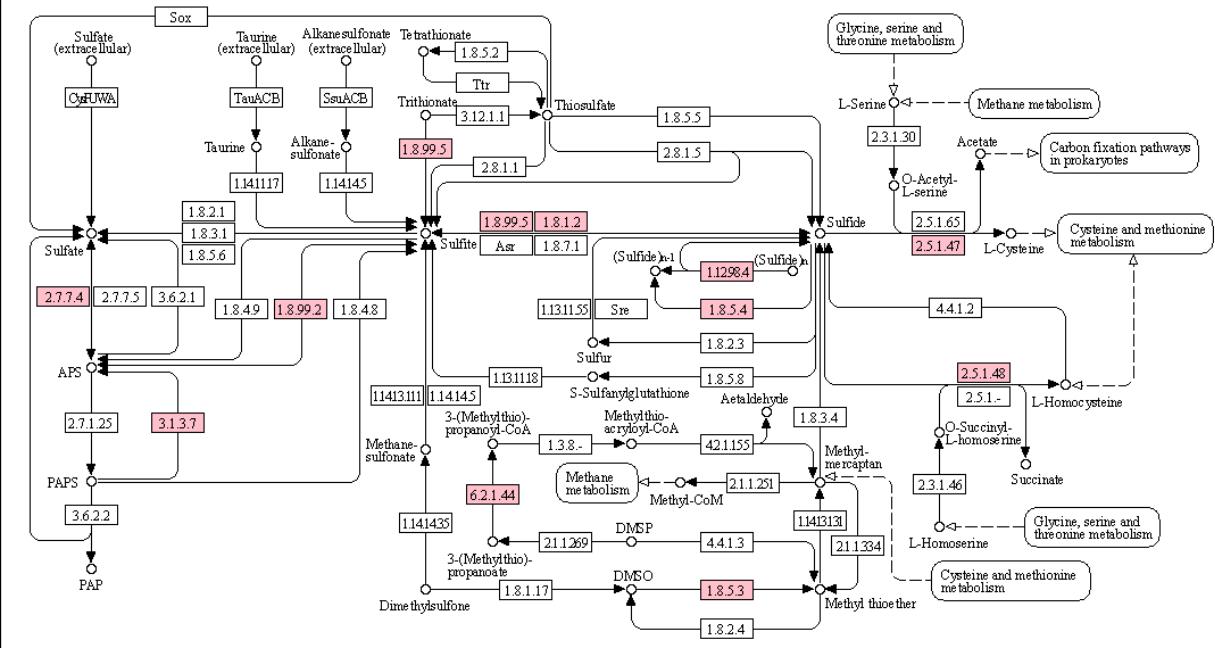
map00680 Methane metabolism (21)

map00260 Glycine, serine and threonine metabolism (20)

NITROGEN METABOLISM

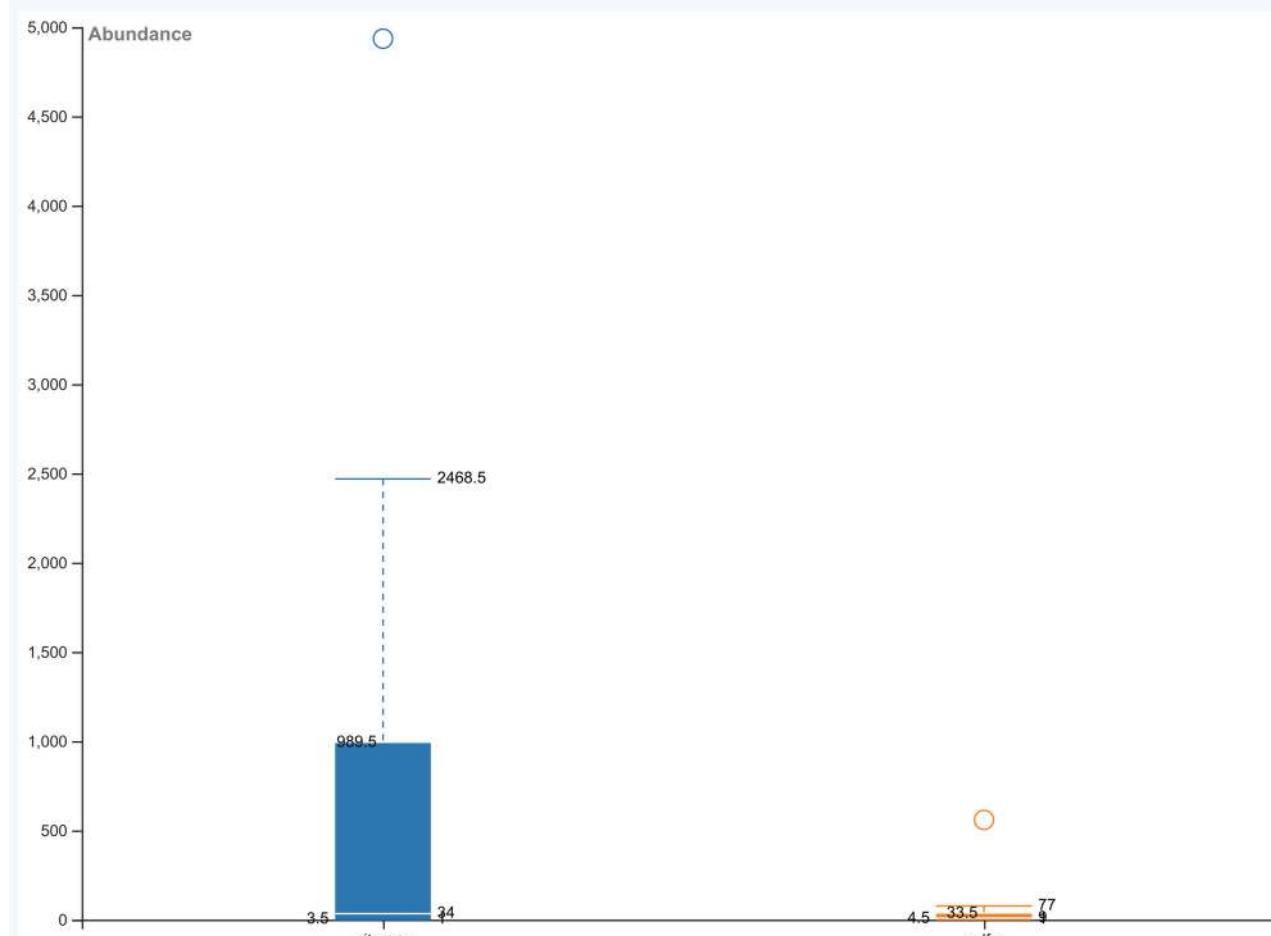


SULFUR METABOLISM



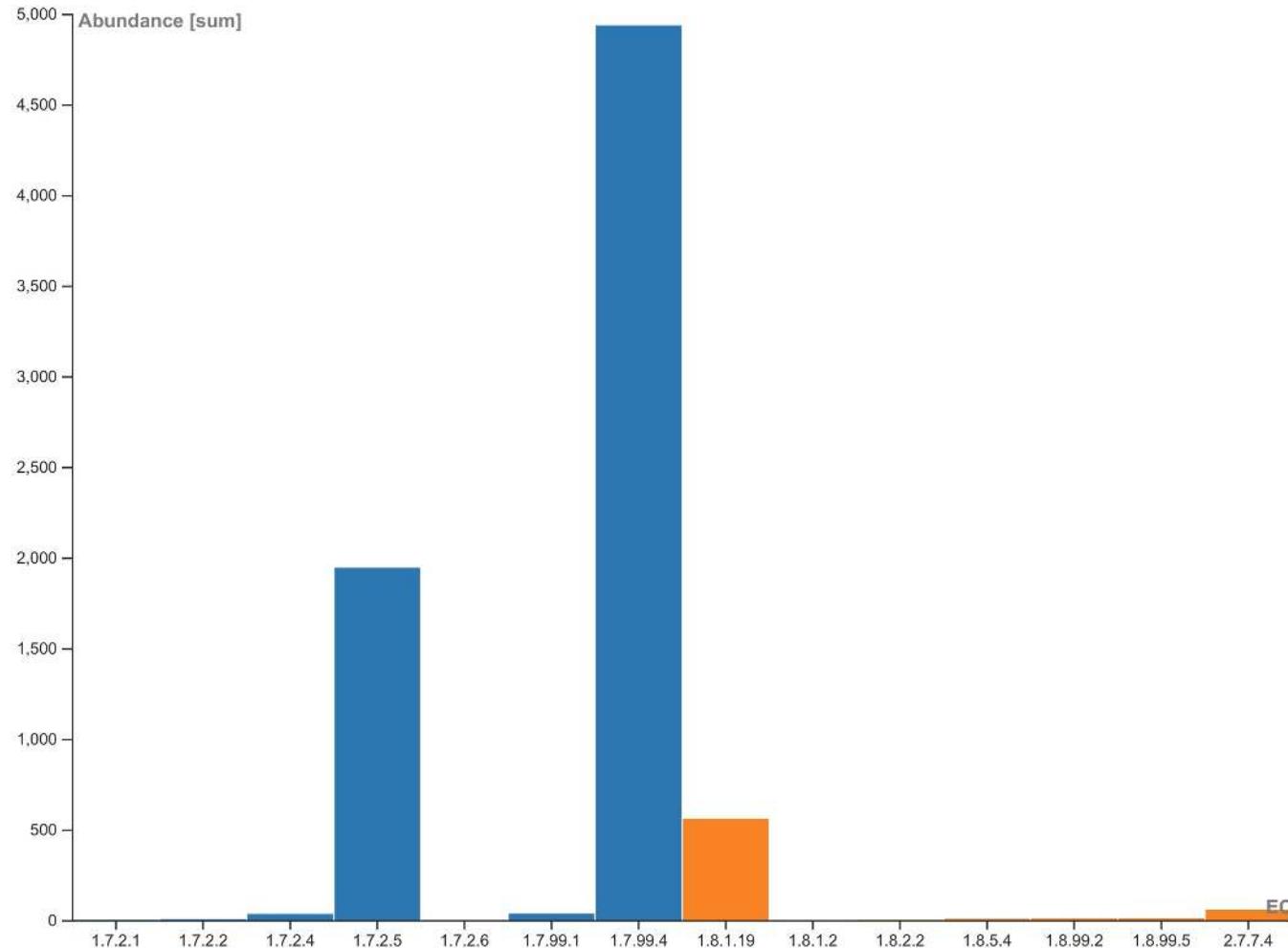


RawGraphs (<https://app.rawgraphs.io>)





RawGraphs (<https://app.rawgraphs.io>)





RawGraphs (<https://app.rawgraphs.io>)

