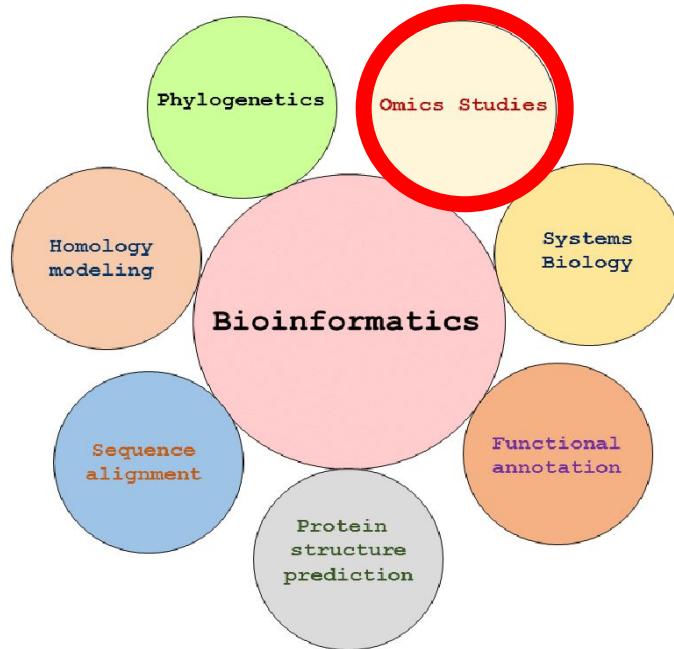


Metagenomics and amplicon sequencing

What is bioinformatics?

- Combines biology and computer science to **interpret biological data**

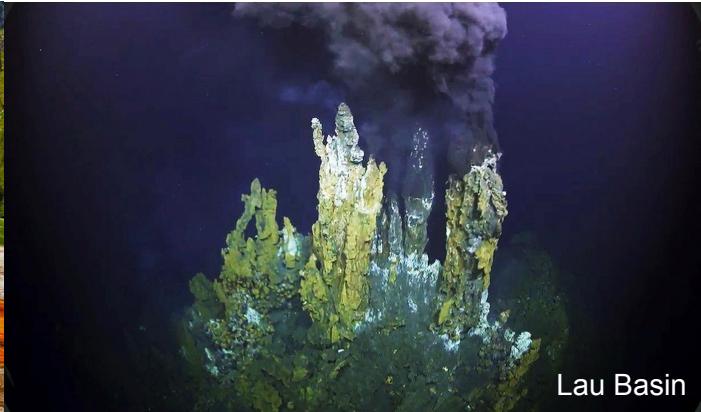


Field trip!

1. Pick a field destination/ system of your choice
2. Sample for some bacteria



Yellowstone National Park



Lau Basin



Prieto-Barajas et al., 2018

c)

Field trip!

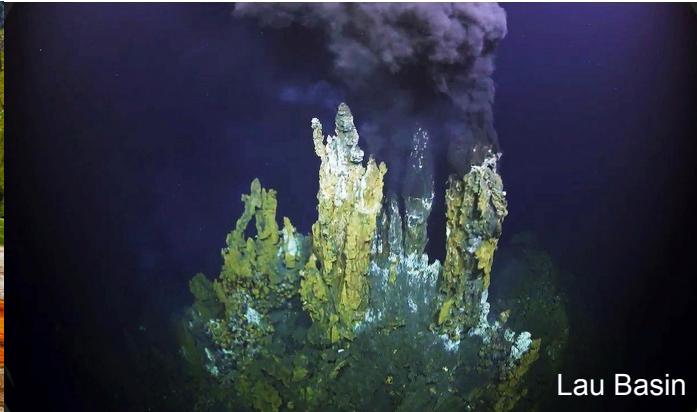
1. Pick a field destination/ system of your choice
2. Sample for some bacteria

Questions you may ask yourself:

Who is there? And what could they be doing?!



Yellowstone National Park



Lau Basin



Prieto-Barajas et al., 2018

c)

WHO is there?



WHAT could they be doing?

WHO is there?

- We often attack this question
with **amplicon sequencing**

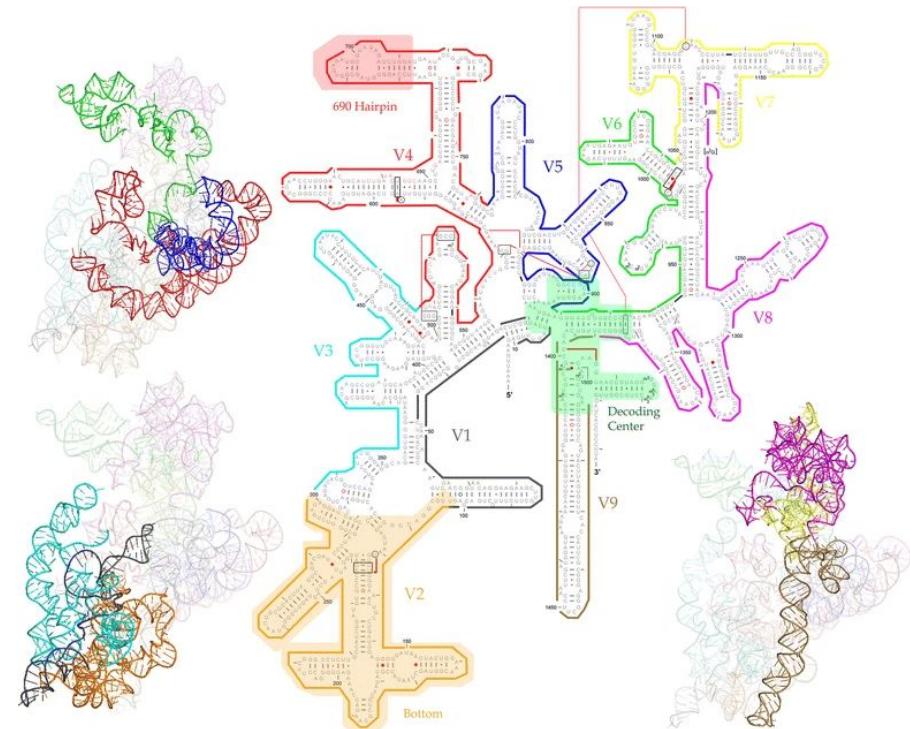
WHAT could they be doing?



Amplicon

16S rRNA

- Working with an amplified gene target
- Usually using SSU rRNA
 - 16S rRNA - prokaryotes
 - 18S rRNA - eukaryotes
 - Very conserved target among different domains of life
- Benefits:
 - Highly targeted
 - Faster
 - Less expensive
- Drawbacks:
 - Not looking at whole genomes



Yang et al., 2016

Linear map of the 9 hypervariable regions of 16S rRNA

- Primers target specific regions
- Can sequence whole sequence (1500bp) - Sanger sequencing
 - Downsides: more computationally expensive
- Next generation sequencing sequences small chunks/regions
- 16S rRNA preferred regions: V4, V5, V6

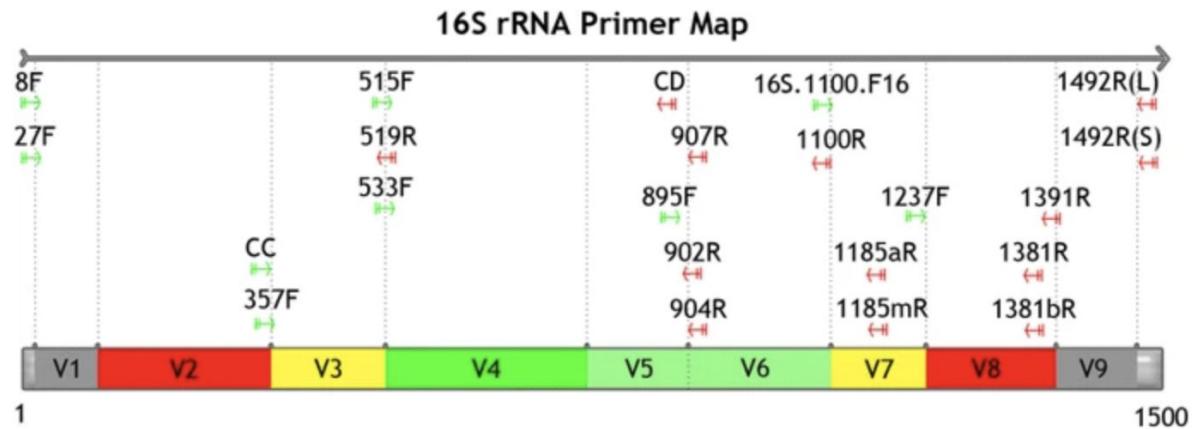
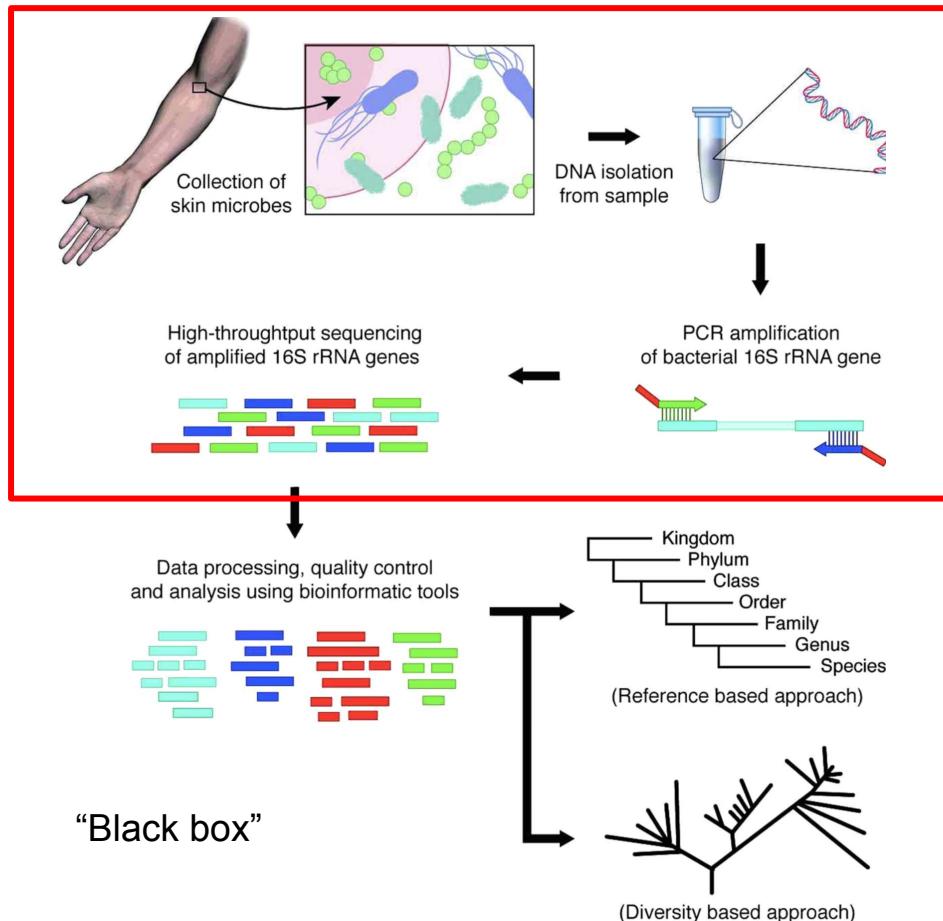


Illustration of different variable regions. Red regions (V2, V8) have a poor phylogenetic resolution at the phylum level. Green regions (V4, V5, V6) are associated with the shortest geodesic distance, which suggests that they may be the best choice for phylogeny-related analyses and the phylogenetic analysis of novel bacterial phyla. The figure refers to the primer map from Lutzonilab (<http://lutzonilab.org/16s-ribosomal-dna/>). Use of this information was approved by the original authors of the website

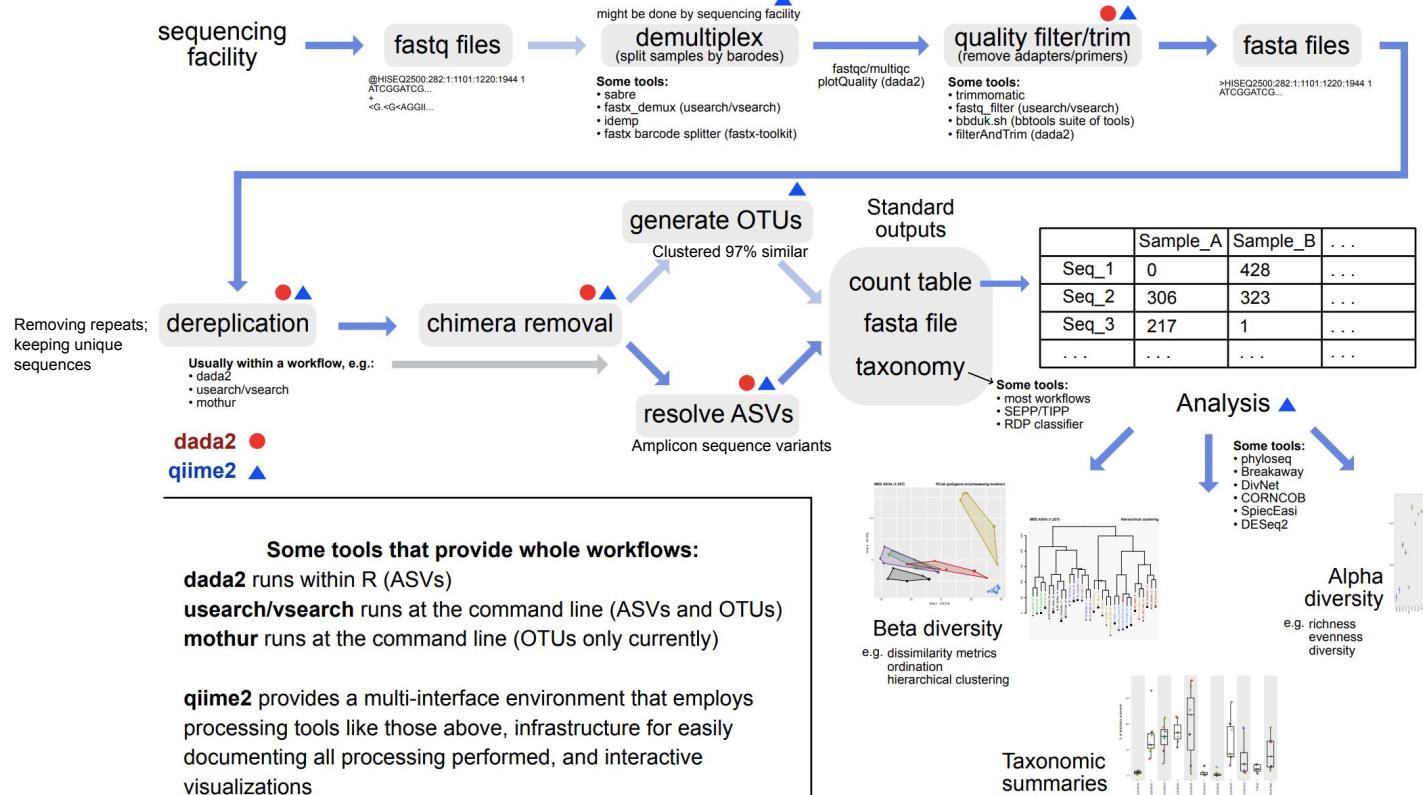
Amplicon sequencing: Lab-based workflow



Amplicon sequencing: Bioinformatic workflow

Overview of generic* amplicon workflow

*This is generic; specific workflows can vary on the order of steps here and how they are done.



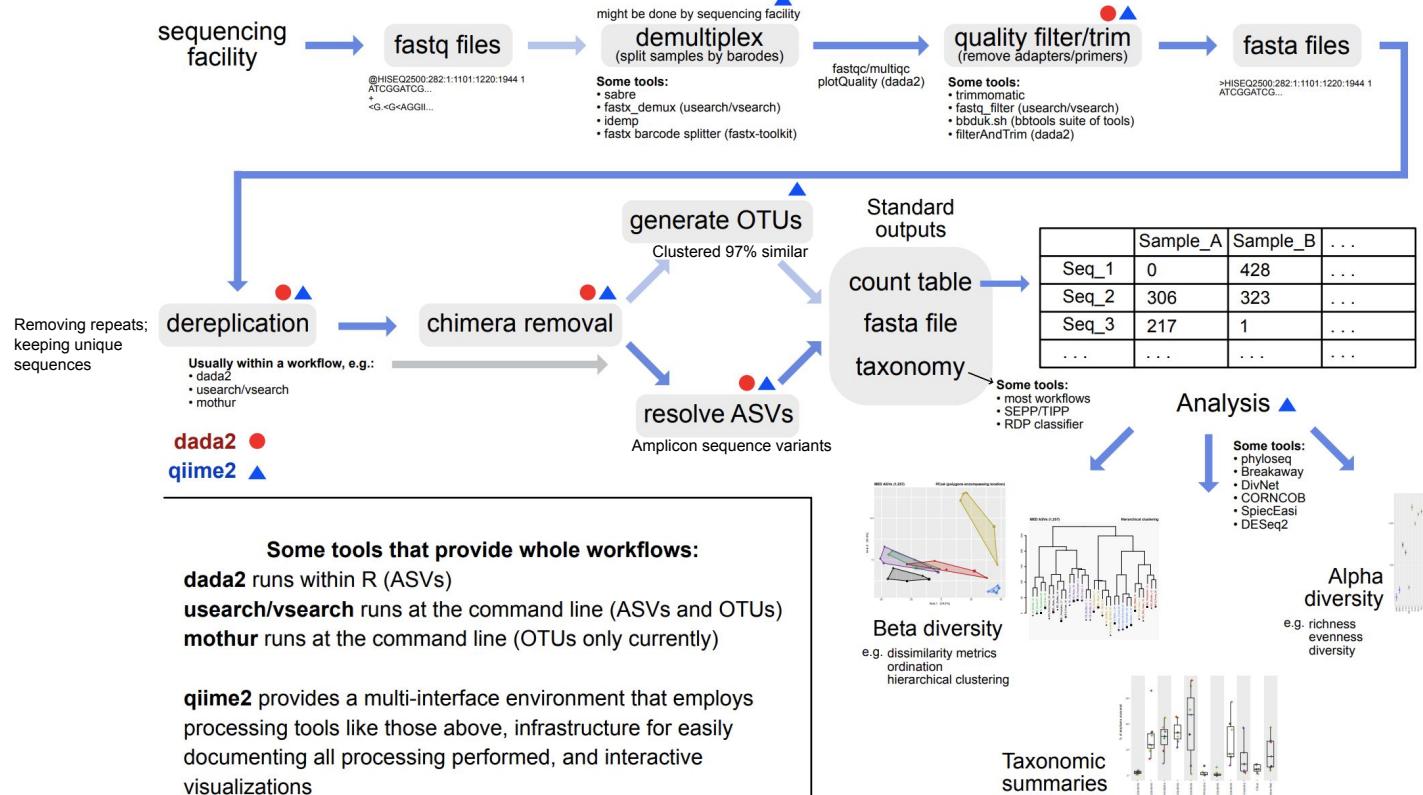
Y103.fasta

@Y103_0 HISSEQ-301:HNWGKBCXX-2:1101:6397:2187 1:N:0:CCAGTT orig_bc=AAAAAAAAAAAA new_bc=AAAAAAAAAAAA bc_diffs=0
TACGTAGGGTGCAGCGTITGCCGAATTACTGGCGTAAAGAGCTCGTAGGCGGTTTGTCGCTCGTCTGTGAAATCCCGAGCTCAACTGCGGCTTGCAAGCGATACGG
GCAAACCTGAGTACTGCAGGGAGACTGGAATTCTGGTGTAGCGGTGAAATGCGCAGATATCAGGAGGAACACCGTGGCGAAGGCAGGCTCTGGCAGTACTGACGC
TGAGGAGCGAAAGCGTGGGTAGCGAACAGG
+
|||||||||||||H||||||G|||||||||||||
|||||||||||||H||||||G|||||||||||||
@Y103_1 HISSEQ-301:HNWGKBCXX-2:1101:11507:2246 1:N:0:CCAGTT orig_bc=AAAAAAAAAAAA new_bc=AAAAAAAAAAAA bc_diffs=0
TACGTAGGGTGCAGCGTITGCCGAATTACTGGCGTAAAGAGCTCGTAGGCGGTTTGTCGCTCGTCTGTGAAATCCCGAGCTCAACTGCGGCTTGCAAGCGATACGG
GCTCTGCTGAGTGCAGGAGAGTAAGCGGAATTCTAGTGTAGCGGTGAAATGCGTAGATATTAGGAAGAACACCACTGGCGAAGGCAGGCTACTGGACGTAACTGACG
TTGAGGCTCGAAAGCGTGGGAGCGAACAGG
+
|||||||||||||H||||||G|||||||||||||
|||||||||||||H||||||G|||||||||||||
@Y103_2 HISSEQ-301:HNWGKBCXX-2:1101:18481:2208 1:N:0:CCAGTT orig_bc=AAAAAAAAAAAA new_bc=AAAAAAAAAAAA bc_diffs=0
TACGTAGGGTGCAGCGTITGCCGAATTACTGGCGTAAAGAGCTCGTAGGCGGTTTGTCGCTCGTCTGTGAAATCCCGAGCTCAACTGCGGCTTGCAAGCGATACGG
GCAGACTCGAGTACTGCAGGGAGACTGGAATTCTGGTGTAGCGGTGAAATGCGCAGATATCAGGAGGAACACCGTGGCGAAGGCAGGCTCTGGCAGTACTGACGC
TGAGGAGCGAAAGCGTGGGTAGCGAACAGG
+
|||||||||||||H||||||G|||||||||||||
|||||||||||||H||||||G|||||||||||||
@Y103_3 HISSEQ-301:HNWGKBCXX-2:1101:5935:2268 1:N:0:CCAGTT orig_bc=AAAAAAAAAAAA new_bc=AAAAAAAAAAAA bc_diffs=0
TACGAAGGGGCTAGCGTTGTCGATTACTGGCGTAAAGCGCAGCTAGGCGGATTGGTCAGTTAGAGGTGAAATCTGGAGGCTCAACTCCAGAACTGCCTTAACTG
CCAGTCTCGAGTCGGAGAGGGTAGTGTGAACTCTAGTGTAGAGGTGAAATTCTGTAGATATTAGGAAGAACACCACTGGCGAAGGCAGGCTACTGGTCGGTACTGACGC
TGAGGTGCGAAAGCGTGGGAGCGAACAGG
+
|||||||||||||H||||||G|||||||||||||
|||||||||||||H||||||G|||||||||||||
@Y103_4 HISSEQ-301:HNWGKBCXX-2:1101:9217:2438 1:N:0:CCAGTT orig_bc=AAAAAAAAAAAA new_bc=AAAAAAAAAAAA bc_diffs=0
TACGAAGGGTCAAGCGTTAACCGGAATTACTGGCGTAAAGCGCGCTAGGCGTTTGTCGTTAAAGTTGATGTGAAAGCCCCGGCTCAACTGGGACTGCATCCAAACT
GGCAAGCTAGAGTATGGCAGGGGTGTGGAATTCTGTAGCGGTGAAATGCGTAGATATAGGAAGAACACCACTGGCGAAGGCAGCACCTGGCTAAACTGACA
CTGAGGTGCGAAAGCGTGGGAGCGAACAGG
+
|||||||||||||H||||||G|||||||||||||
|||||||||||||H||||||G|||||||||||||
@Y103_5 HISSEQ-301:HNWGKBCXX-2:1101:5325:2570 1:N:0:CCAGTT orig_bc=AAAAAAAAAAAA new_bc=AAAAAAAAAAAA bc_diffs=0
TACGTAGGGTGCAGCGTITGCCGAATTACTGGCGTAAAGAGCTCGTAGGCGGTTTGTCGCTCGTCTGTGAAATCCCGAGCTCAACTGCGGCTTGCAAGCGATACGG
GCAAACCTGAGTACTGCAGGGAGACTGGAATTCTGGTGTAGCGGTGAAATGCGCAGATATCAGGAGGAACACCGTGGCGAAGGCAGGCTCTGGCAGTACTGACGC
TGAGGAGCGAAAGCGTGGGTAGCGAACAGG
+
|||||||||||||H||||||G|||||||||||||
|||||||||||||H||||||G|||||||||||||
@Y103_6 HISSEQ-301:HNWGKBCXX-2:1101:17617:2520 1:N:0:CCAGTT orig_bc=AAAAAAAAAAAA new_bc=AAAAAAAAAAAA bc_diffs=0
TACGAAGGGTCAAGCGTTAACCGGAATTACTGGCGTAAAGCGCGCTAGGCGTTTGTCGTTAAAGTTGATGTGAAAGCCCCGGCTCAACTGGGACTGCATCCAAACT
GGCGAGCTAGAGTACCGTAGAGGGTAGTGTGAAATTCTGTAGCGGTGAAATGCGTAGATATAGGAAGAACACCACTGGCGAAGGCAGCACCTGGACTGACA
CTGAGGTGCGAAAGCGTGGGAGCGAACAGG
+
|||||||||||||H||||||G|||||||||||||
|||||||||||||H||||||G|||||||||||||
@Y103_7 HISSEQ-301:HNWGKBCXX-2:1101:19548:2731 1:N:0:CCAGTT orig_bc=AAAAAAAAAAAA new_bc=AAAAAAAAAAAA bc_diffs=0
TACGTAGGGTGCAGCGTITGCCGAATTACTGGCGTAAAGAGCTCGTAGGCGGTTTGTCGCTCGTCTGTGAAATCCCGAGCTCAACTGCGGCTTGCAAGCGATACGG
GCAAACCTGAGTACTGCAGGGAGACTGGAATTCTGGTGTAGCGGTGAAATGCGCAGATATCAGGAGGAACACCGTGGCGAAGGCAGCACCTGGACTGACA
TGAGGAGCGAAAGCGTGGGTAGCGAACAGG
+

Amplicon sequencing: Bioinformatic workflow

Overview of generic* amplicon workflow

*This is generic; specific workflows can vary on the order of steps here and how they are done.

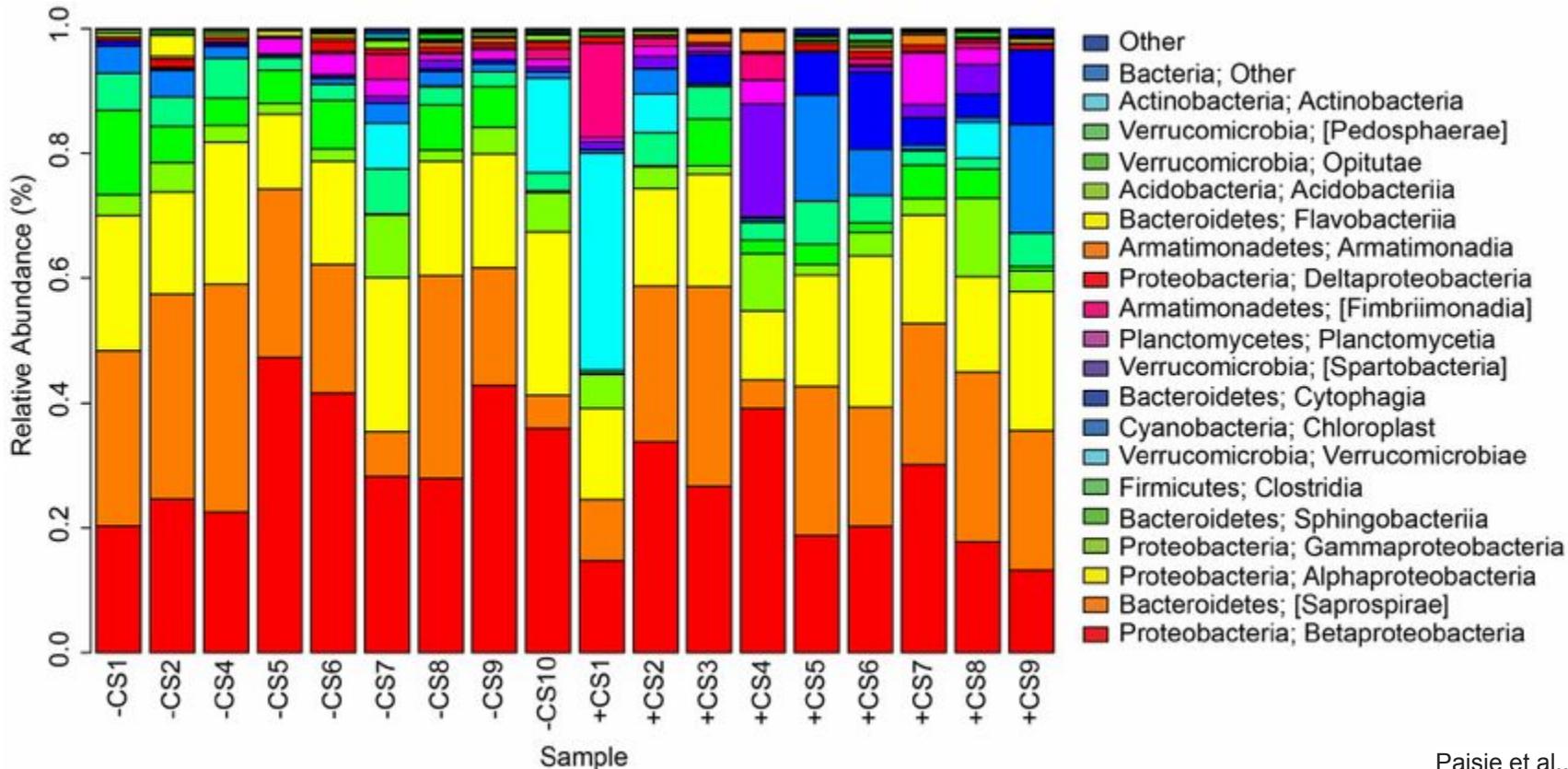


Reference databases

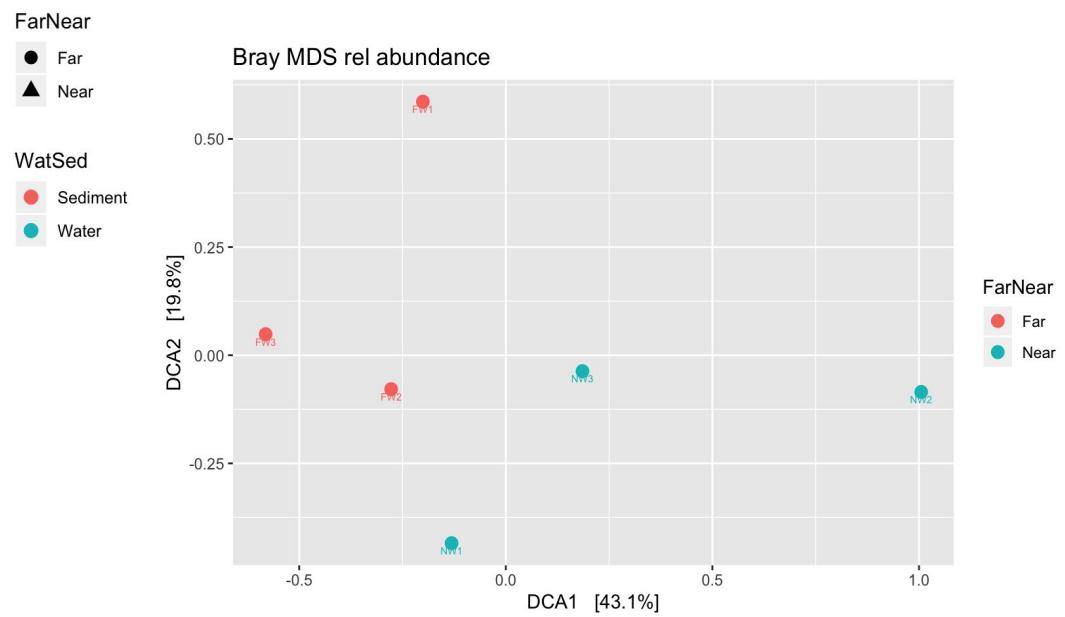
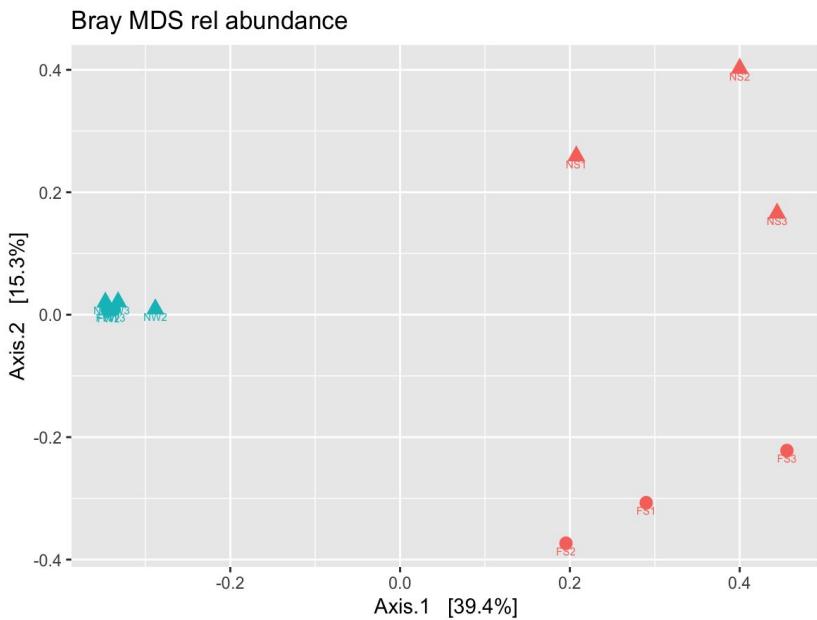
It's important to know and understand what reference database you are using, especially for taxonomy assignment.

<u>Database</u>	<u>Target Amplicon</u>
SILVA	16S & 18S rRNA
RDP	16S rRNA & 28S rRNA (fungi)
PR2	18S rRNA
Greengenes	16S rRNA (<i>no longer maintained</i>)
UNITE	Fungal ITS

Example figures



16S NMDS Plot



Questions?

WHO is there?



- We often attack this question
with **amplicon sequencing**

WHAT could they be doing?



WHO is there?



- We often attack this question with **amplicon sequencing**

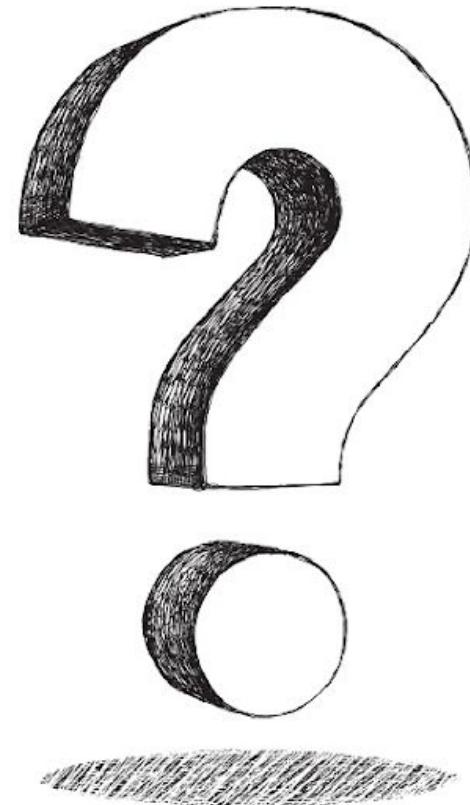


WHAT could they be doing?

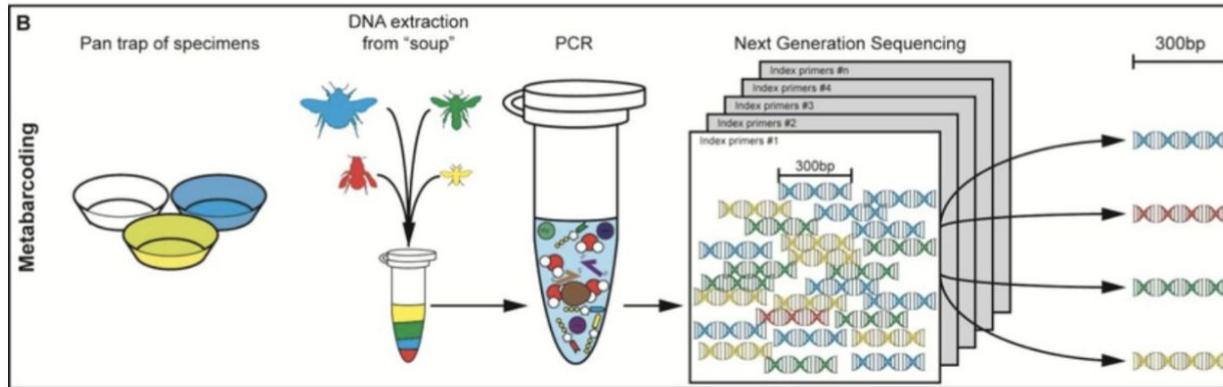
- We often attack this question with **metagenomic sequencing**

Metagenomic sequencing

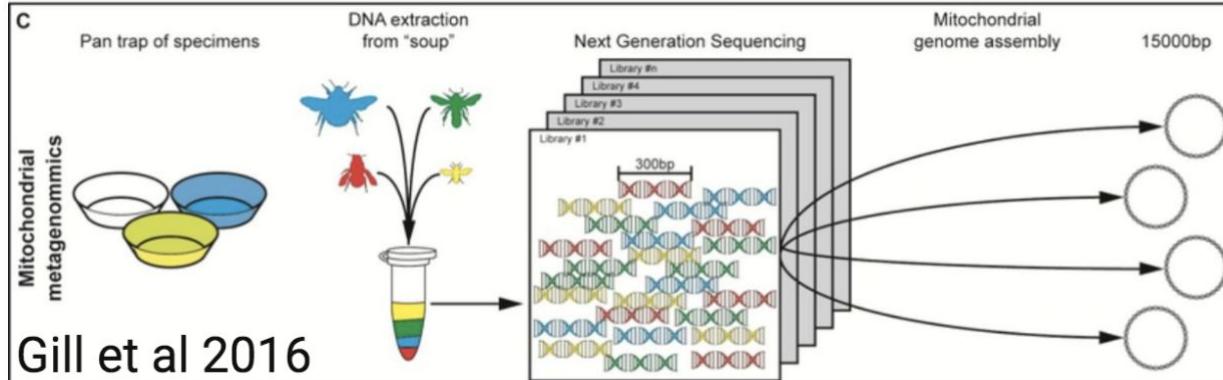
- Also referred to as shotgun sequencing
- Looking at whole genomes
- Benefits:
 - LOTS of data
- Downside:
 - Expensive
 - Longer sequencing process



Amplicon sequencing vs. metagenomic “shotgun” sequencing



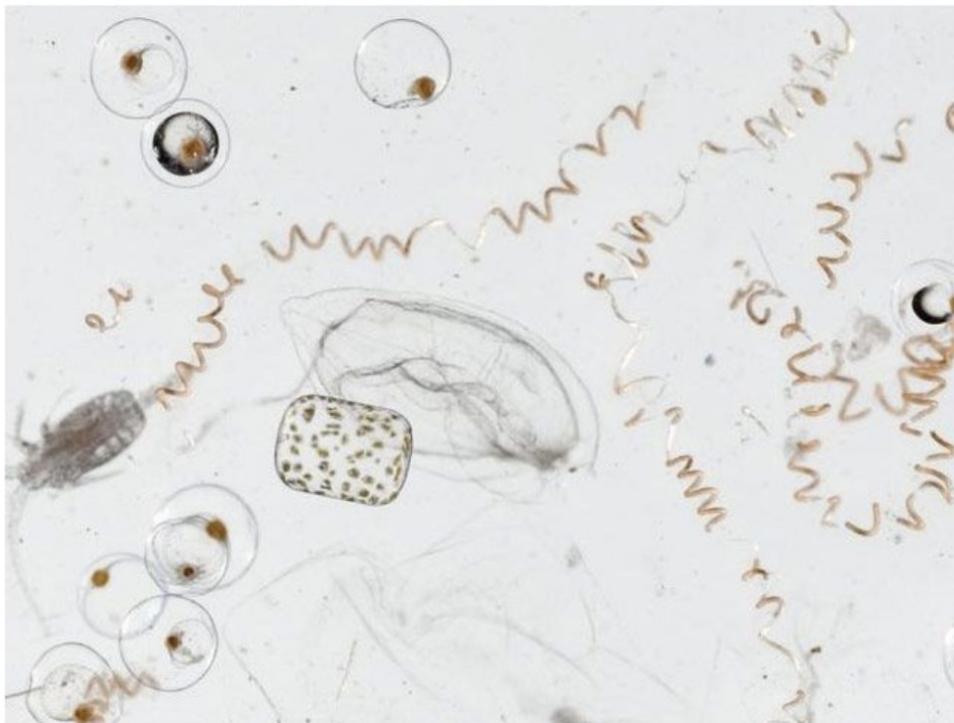
“Amplicon sequencing” or “metabarcoding” is sequencing a specific target region from many genomes (e.g. 16S rRNA gene, *nifH* gene)



Gill et al 2016

“Shotgun Metagenomics” is (incomplete) sequencing of a mixture of genomes using an untargeted approach

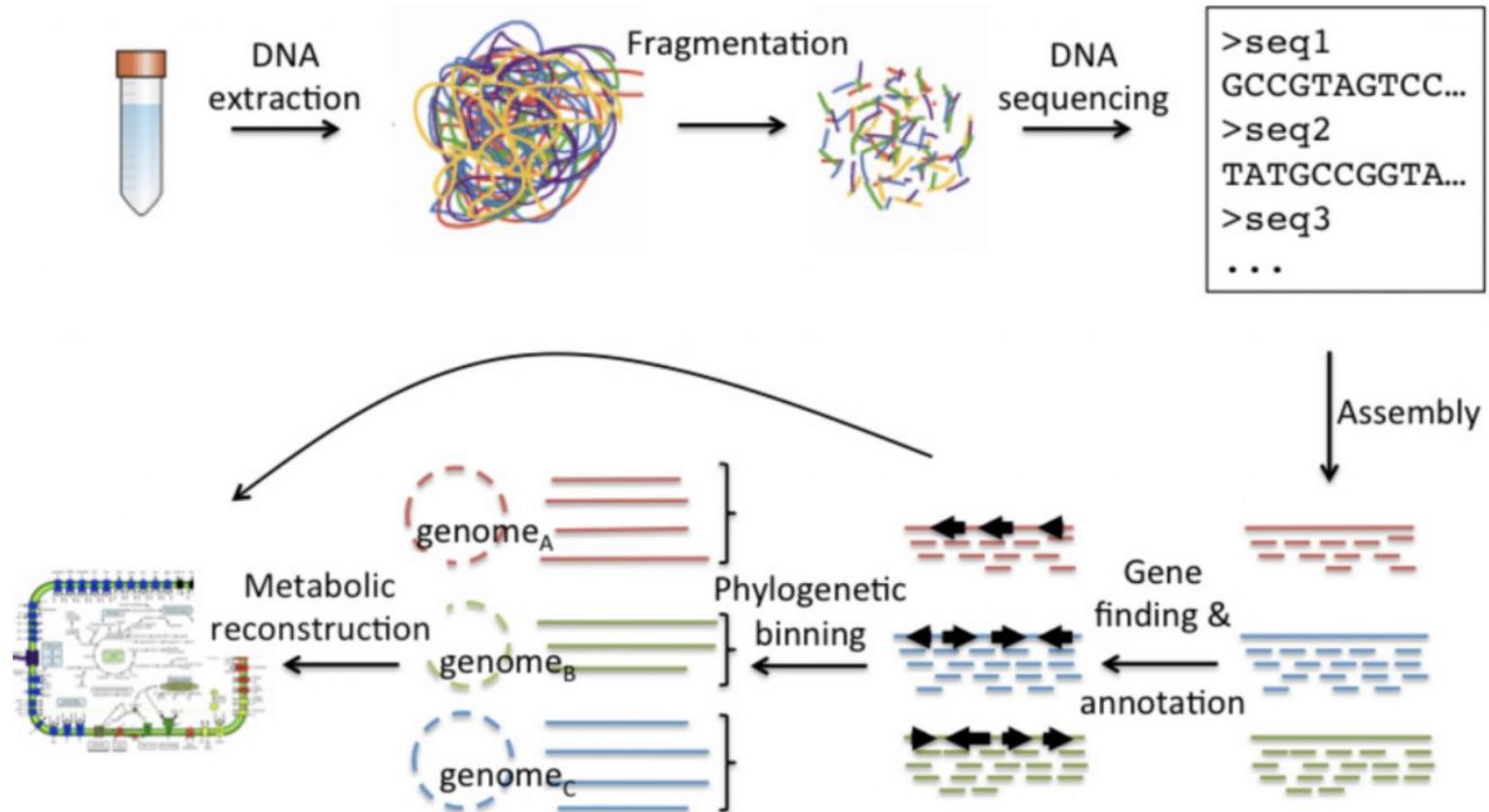
Incomplete because one drop of seawater contains about...



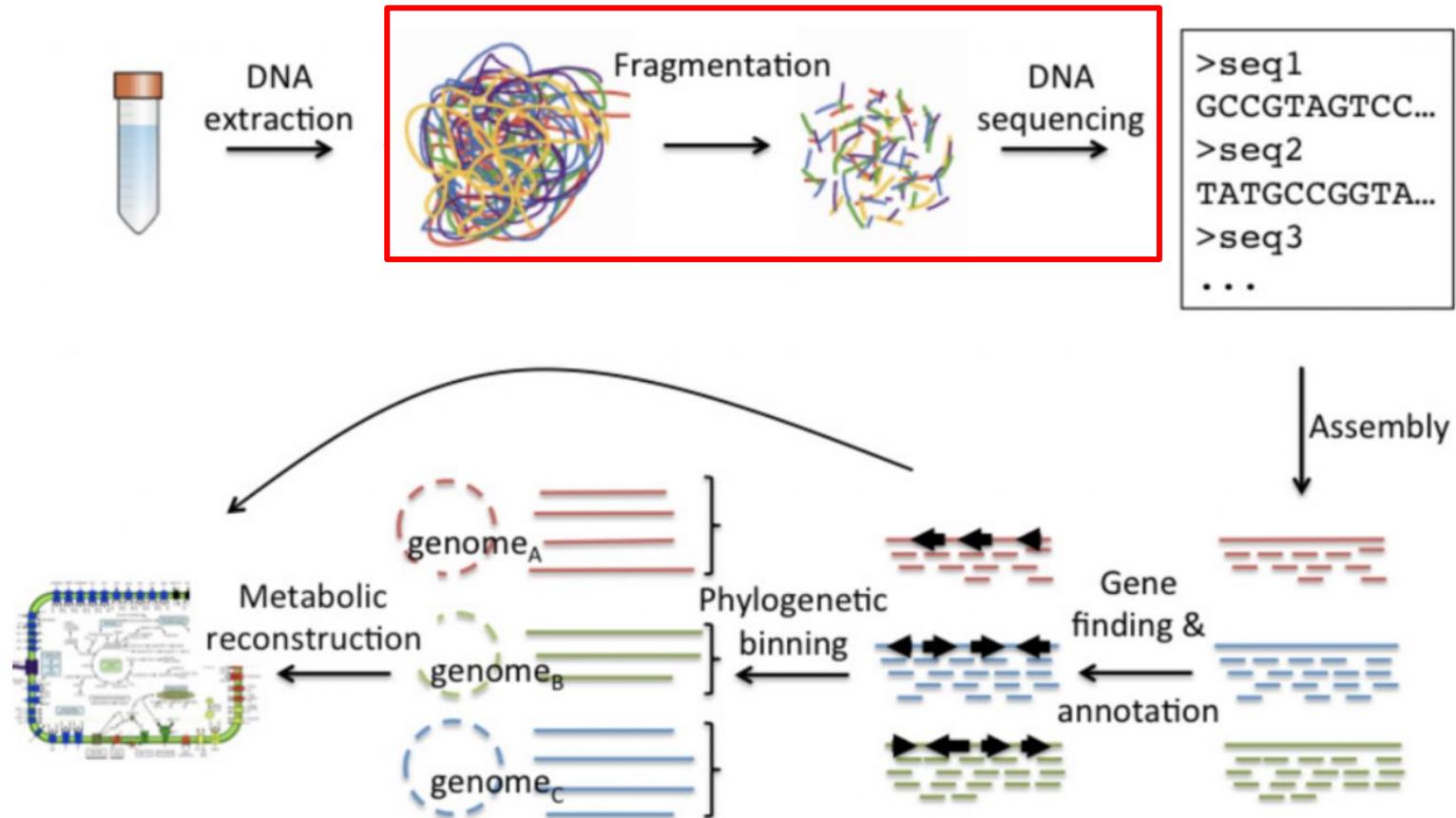
$$10^6 \text{ bacteria/mL} * 3 \times 10^6 \text{ bp/bacteria} = \\ 3 \times 10^{12} \text{ bp/mL}$$

$$+ \quad 10^3 \text{ euks/mL} * 3 \times 10^8 \text{ bp/euk} = \\ 3 \times 10^{11} \text{ bp/mL} \\ = 3.3 \text{ Tbp/mL}$$

Metagenomic workflow



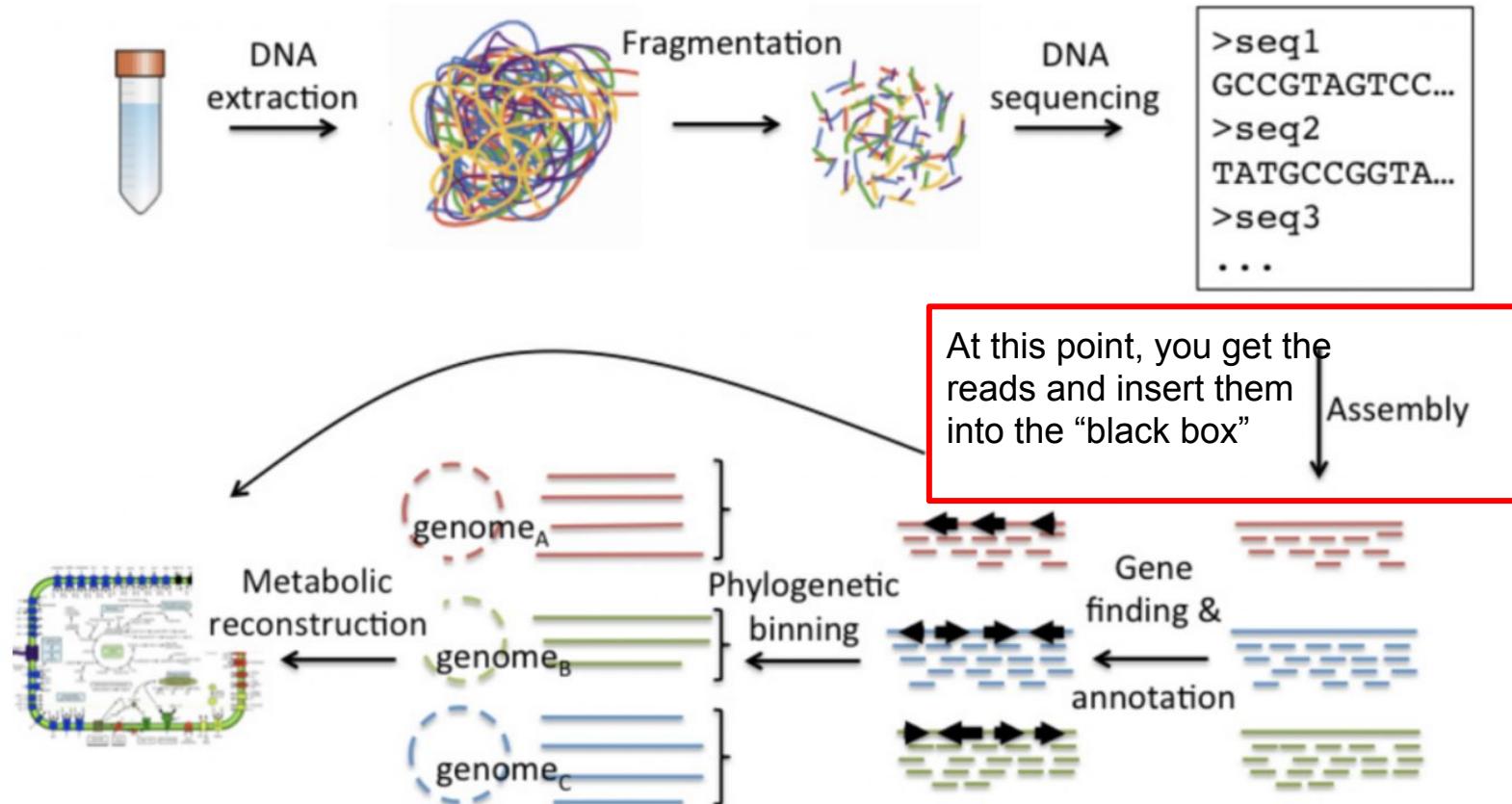
Metagenomic workflow

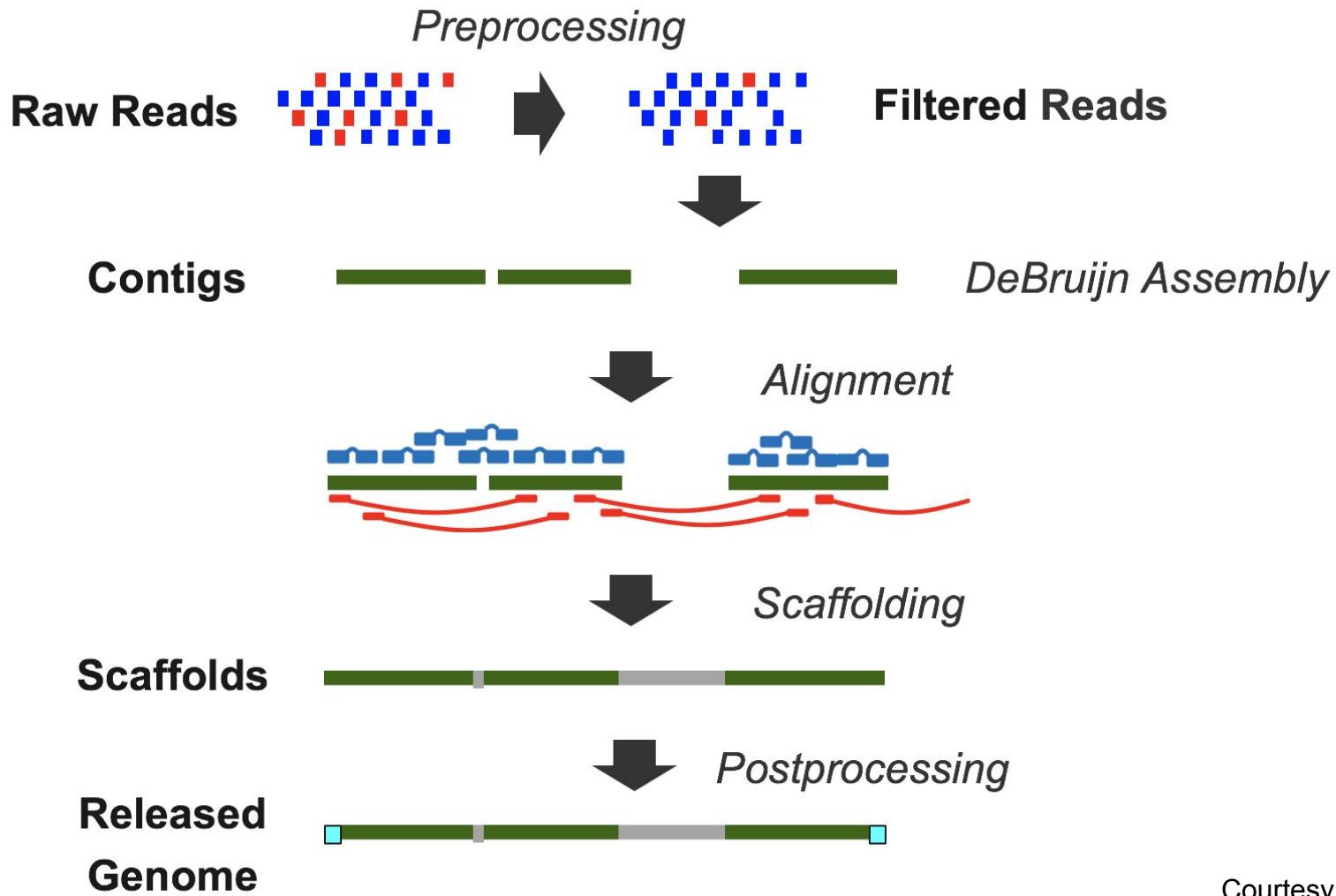






Metagenomic workflow





Courtesy: Brian Bushnell

**“It was the best of times, it was the worst of times,
it was the age of wisdom, it was the age of
foolishness, it was the epoch of belief, it was the
epoch of incredulity, ... “**

Dickens, Charles. *A Tale of Two Cities*. 1859. London: Chapman Hall

itwas the best of times it was the worst of times it was the age of wisdom it was the age of foolishness...

Generate random 'reads'



How do we assemble?



fincreduli geoffolis Itwasthebe Itwasthebe geofwisdom itwastheep epochofinc timesitwas stheepocho nessitwast wastheageo theepochof stheepocho hofincredu estoftimes eoffoolish lishnessit hofbeliefi pochofincr itwasthewo twastheage toftimesit domitwasth ochofbelie eepochofbe eepochofbe astheworst chofincred theageofwi iefitwasth ssitwasthe astheepoch efitwasthe wisdomitwa ageoffooli twasthewor ochofbelie sdomitwast sitwasthea eepochofbe ffoolishne eofwisdomi hebestofti stheageoff twastheepo eworstofti stoftimesi theepochof esitwasthe heepochofi theepochof sdomitwast astheworst rstoftimes worstoftim stheepocho geoffolis ffoolishne timesitwas lishnessit stheageoff eworstofti orstoftime fwisdomitw wastheageo heageofwis incredulit ishnessitw twastheepo wasthewors astheepoch heworstoft ofbeliefit wastheageo heepochofi pochofincr heageofwis stheageofw fincreduli astheageof wisdomitwa wastheageo astheepoch olishnessi astheepoch itwastheep twastheage wisdomitwa fbeliefitw bestoftime epochofbel theepochof sthebestof lishnessit hofbeliefi Itwasthebe ishnessitw sitwasthew ageofwisdo twastheage esitwasthe twastheage shnessitwa fincreduli fbeliefitw theepochof mesitwasth domitwasth ochofbelie heageofwis oftmesitw stheepocho bestoftime twastheage foolishnes ftimesitwa thebestoft itwastheag theepochof itwasthewo ofbeliefit bestoftime mitwasthea imesitwast timesitwas orstoftime estoftimes twasthebes stoftimesi sdomitwast wisdomitwa theworstof astheworst sitwasthew theageoffo eepochofbe

De Bruijn solution...

Courtesy: Brian Bushnell

Convert reads into “Kmers”

Kmer: a substring of defined length

Reads:	theageofwi	s the best of	as the age of	worst of tim	ime sit was
Kmers : (k=3)	the	sth	ast	wor	ime
he a	the	sth	ors	mes	
eag	heb	the	rst	esi	
age	ebe	hea	sto	sit	
geo	bes	eag	t of	itw	
eof	est	age	oft	twa	
of w	sto	geo	fti	was	
fwi	to f	eof	tim	ast	

.....etc for all reads in the dataset

Using DNA sequences, it would look more like this...

sequence

ATGGAAGTCGCAGAATC

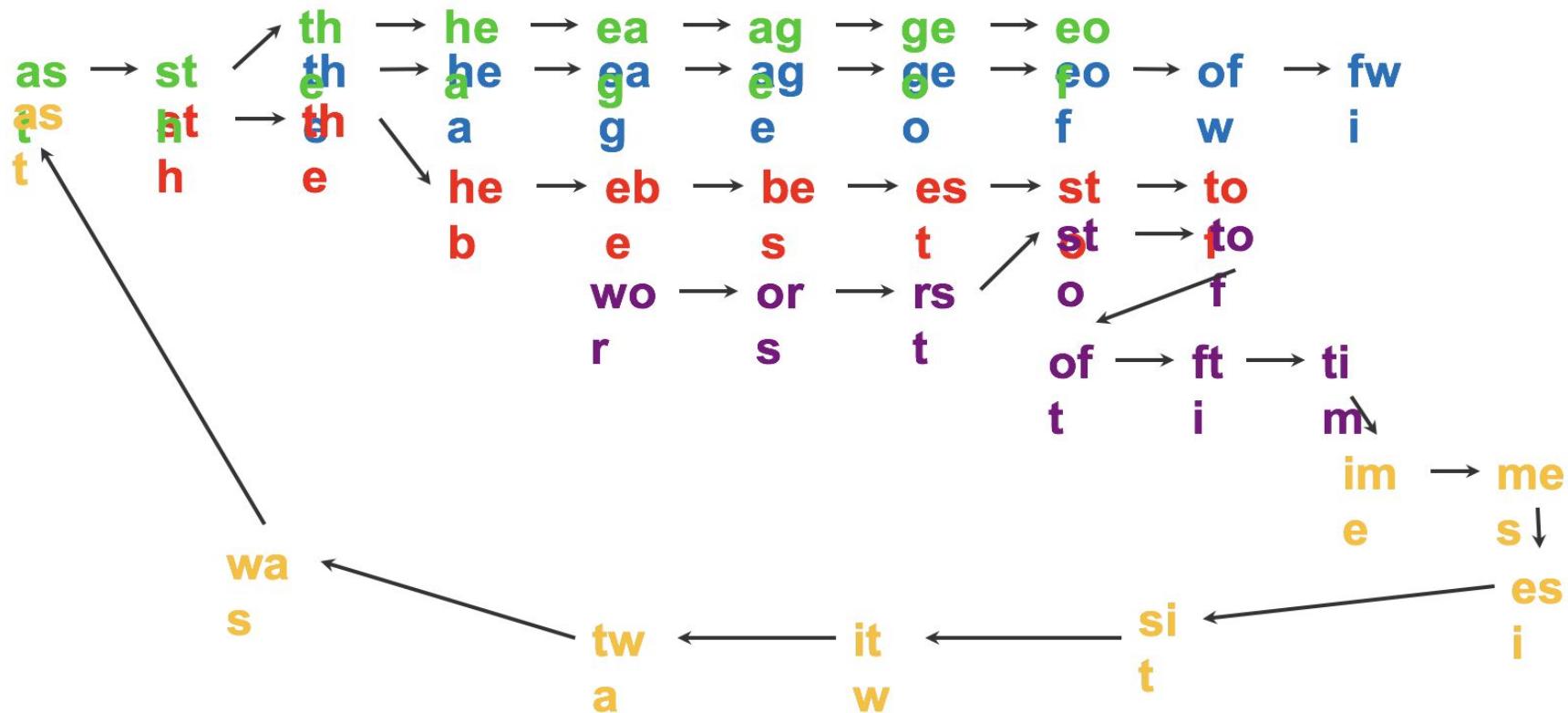
7mers

ATGGAAG
TGGAAAGT
GGAAGTC
GAAGTCG
AAGTCGC
AGTCGCG
GTCGCGG
TCGCGGG
CGCGGGAA
GCGGAAT
CGGAATC

Modified from:

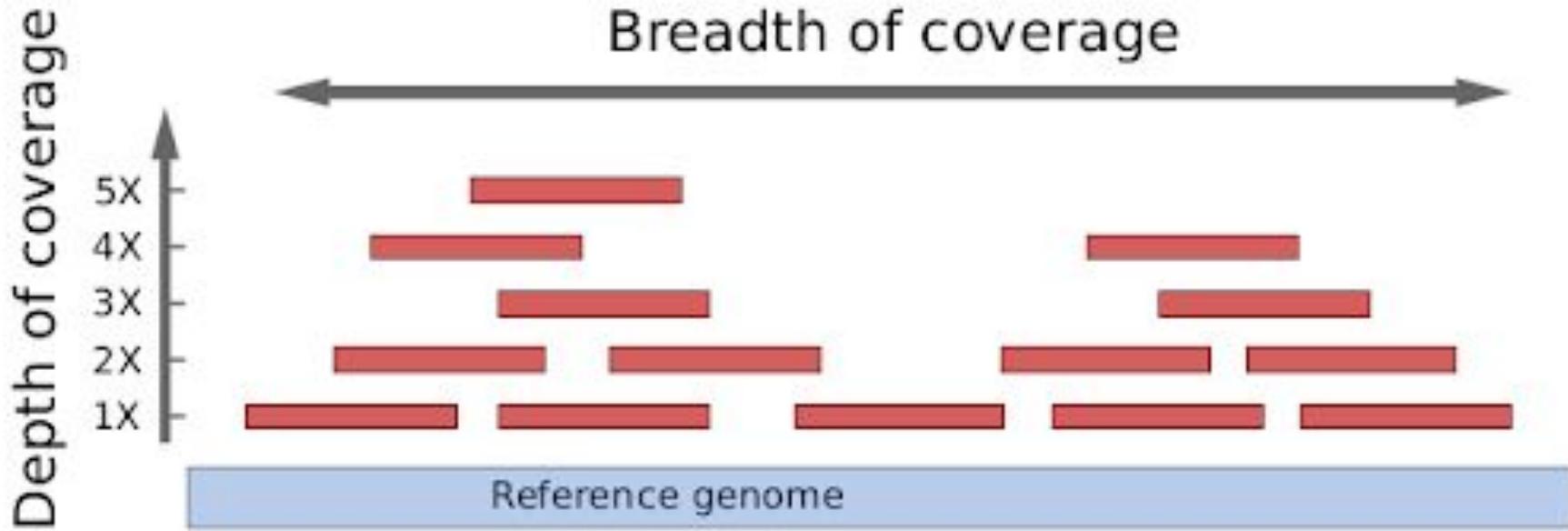
<https://homolog.us/Tutorials/book4/p2.1.html>

Build a De-Bruijn graph from the kmers

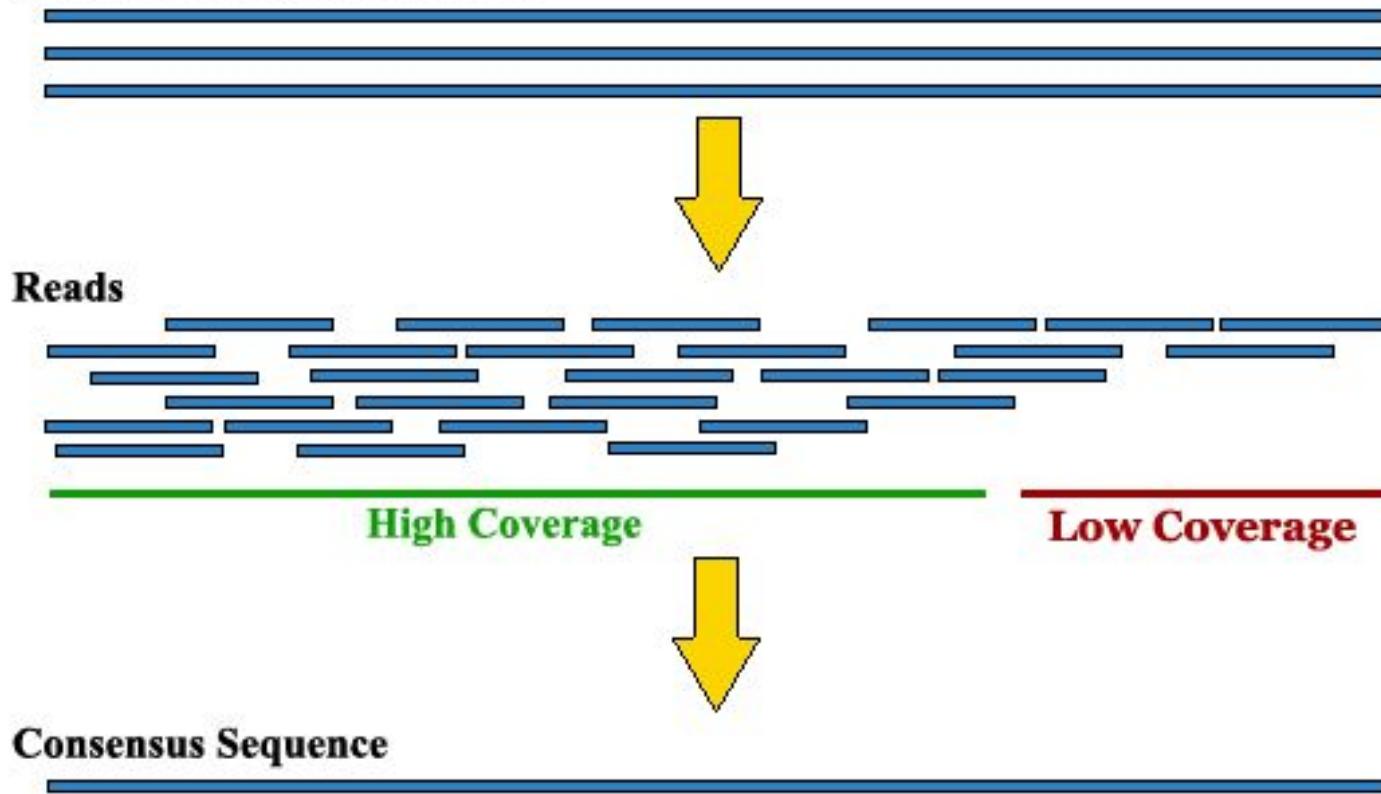


**‘It was the best of times, it was the worst of times,
it was the age of wisdom, it was the age of
foolishness, it was the epoch of belief, it was the
epoch of incredulity, ... “**

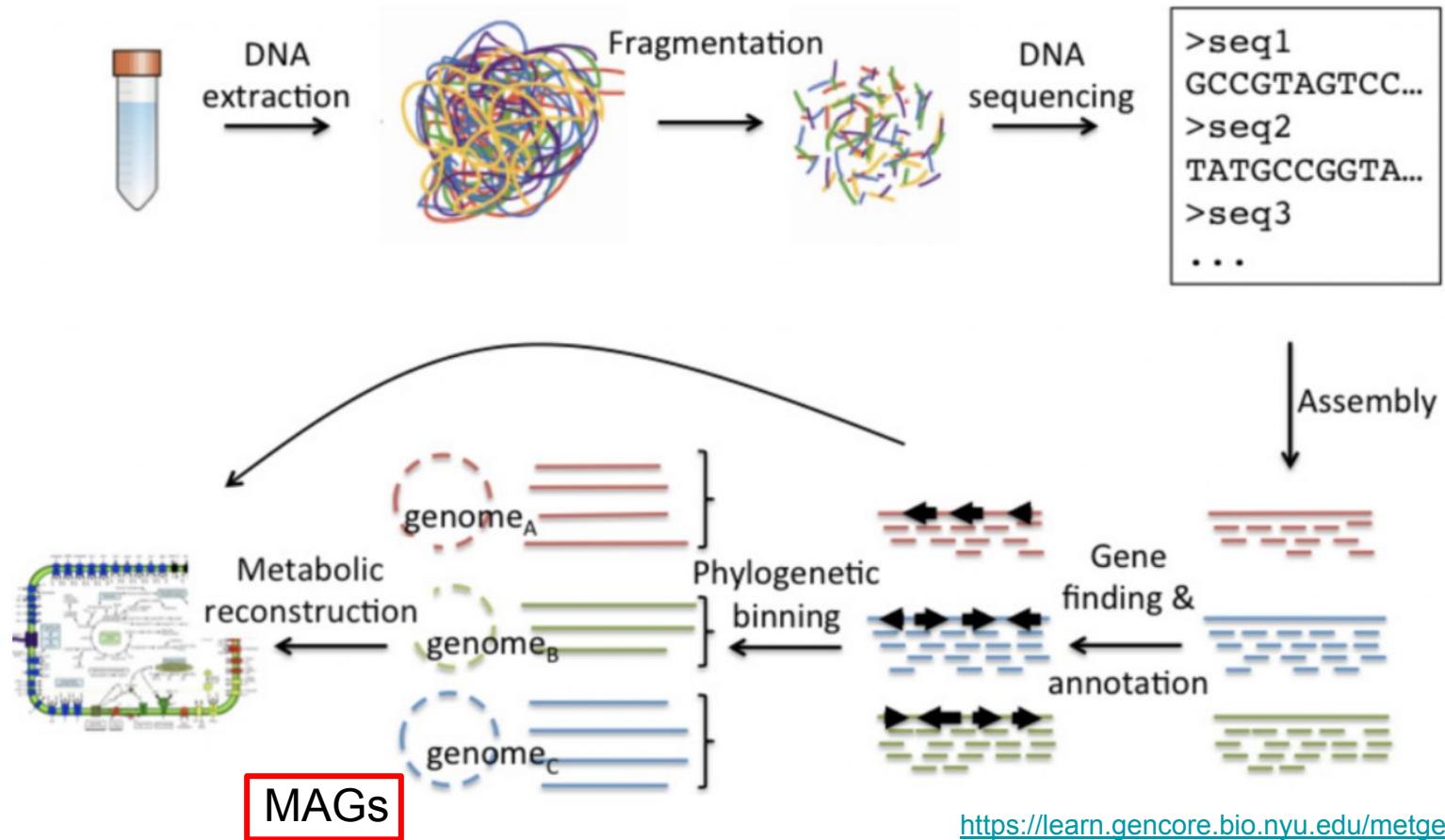
Dickens, Charles. *A Tale of Two Cities*. 1859. London: Chapman Hall



Multiple Copies of a Genome

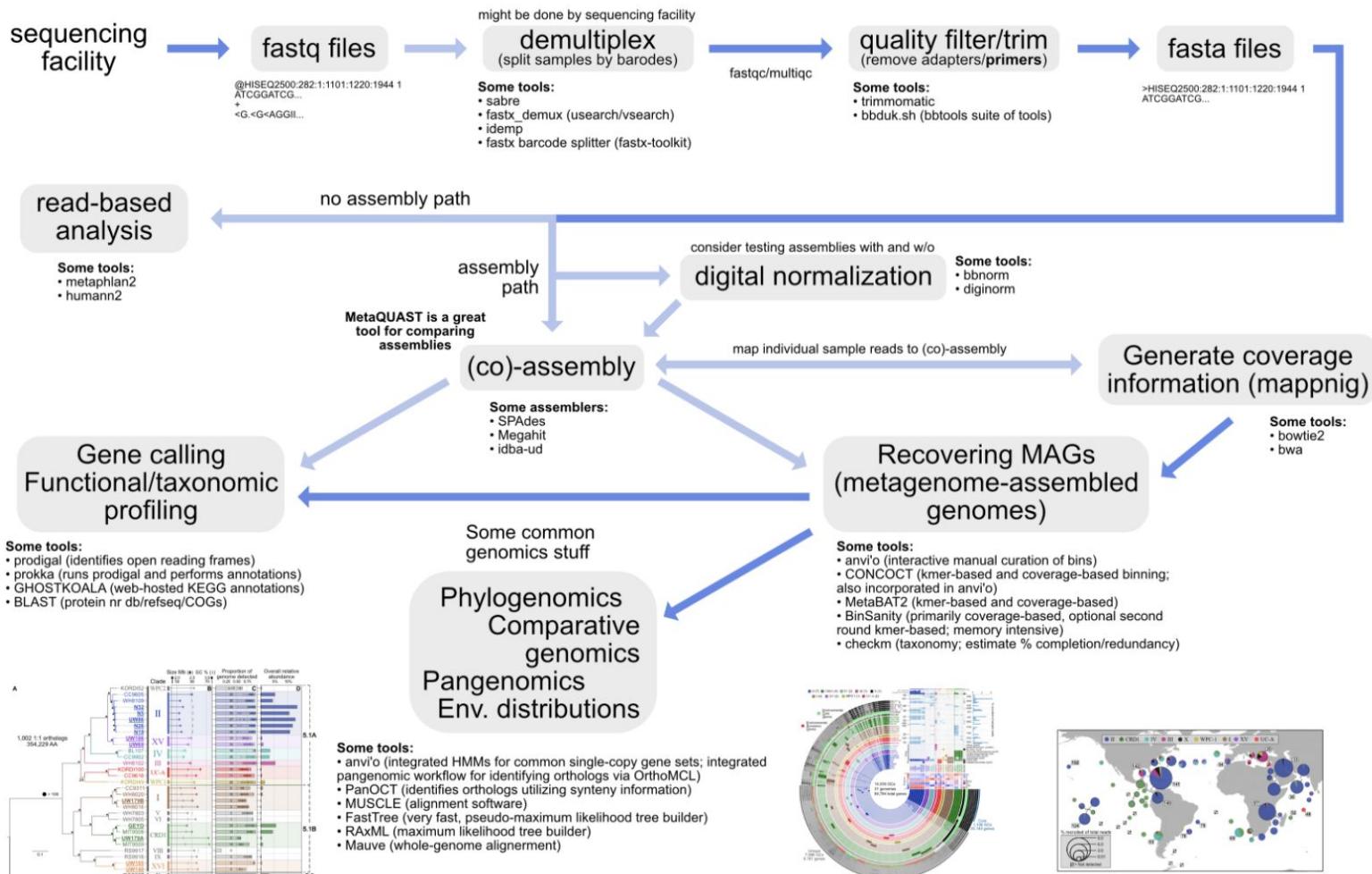


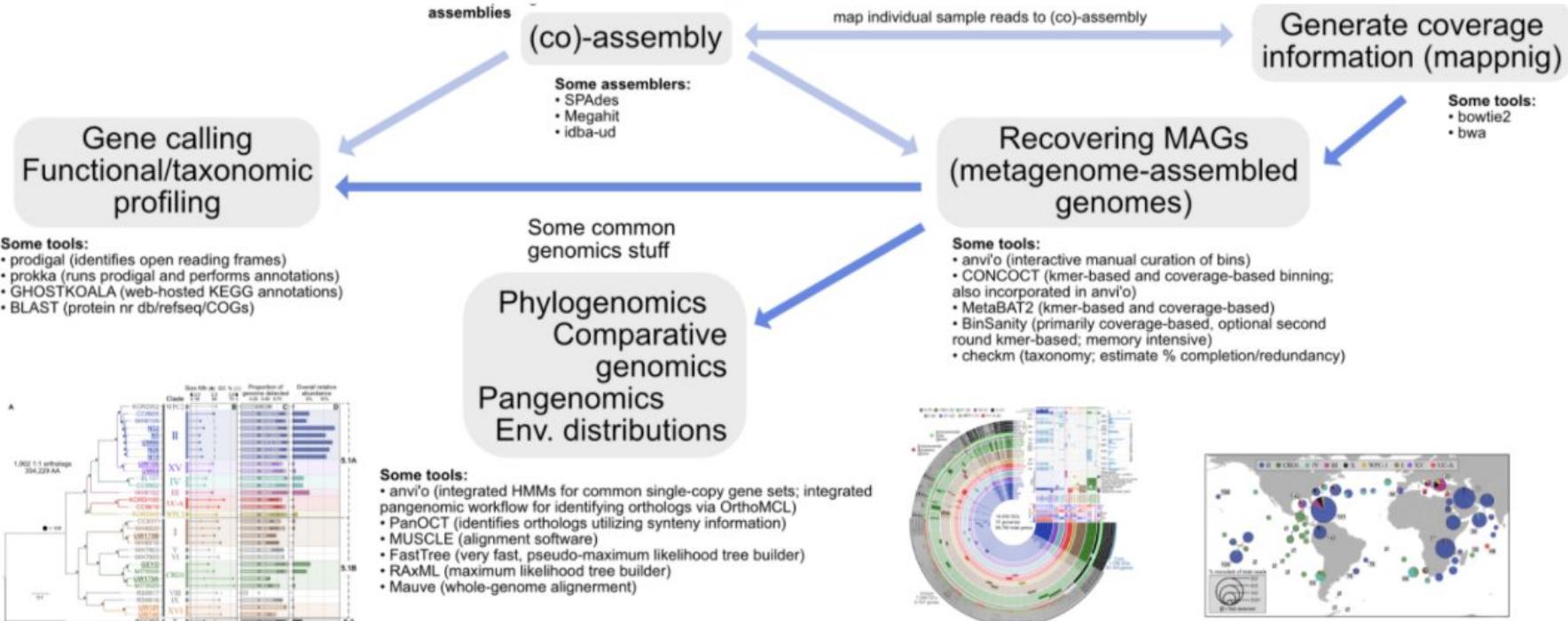
Metagenomic workflow



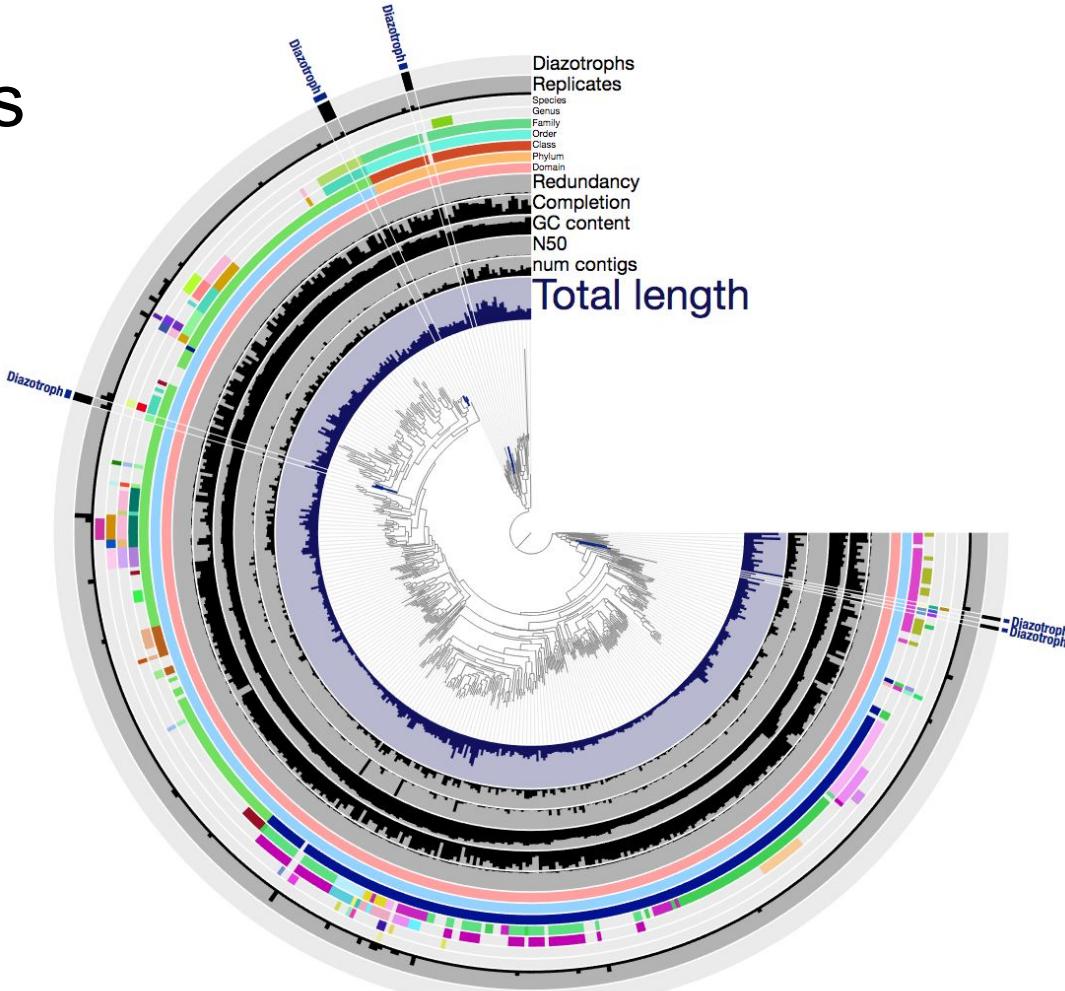
Overview of generic metagenomics workflow

When working with your own data you should never follow any pipeline blindly. There can be critical differences based on your data.

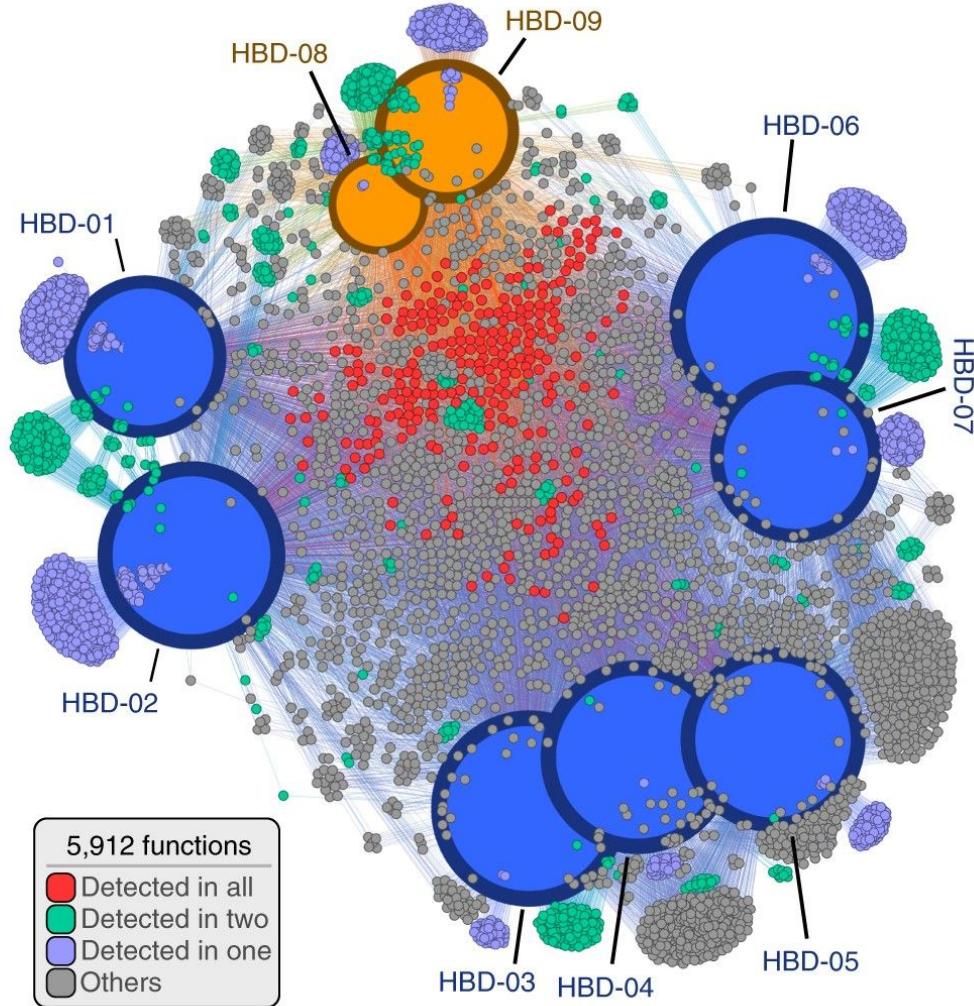




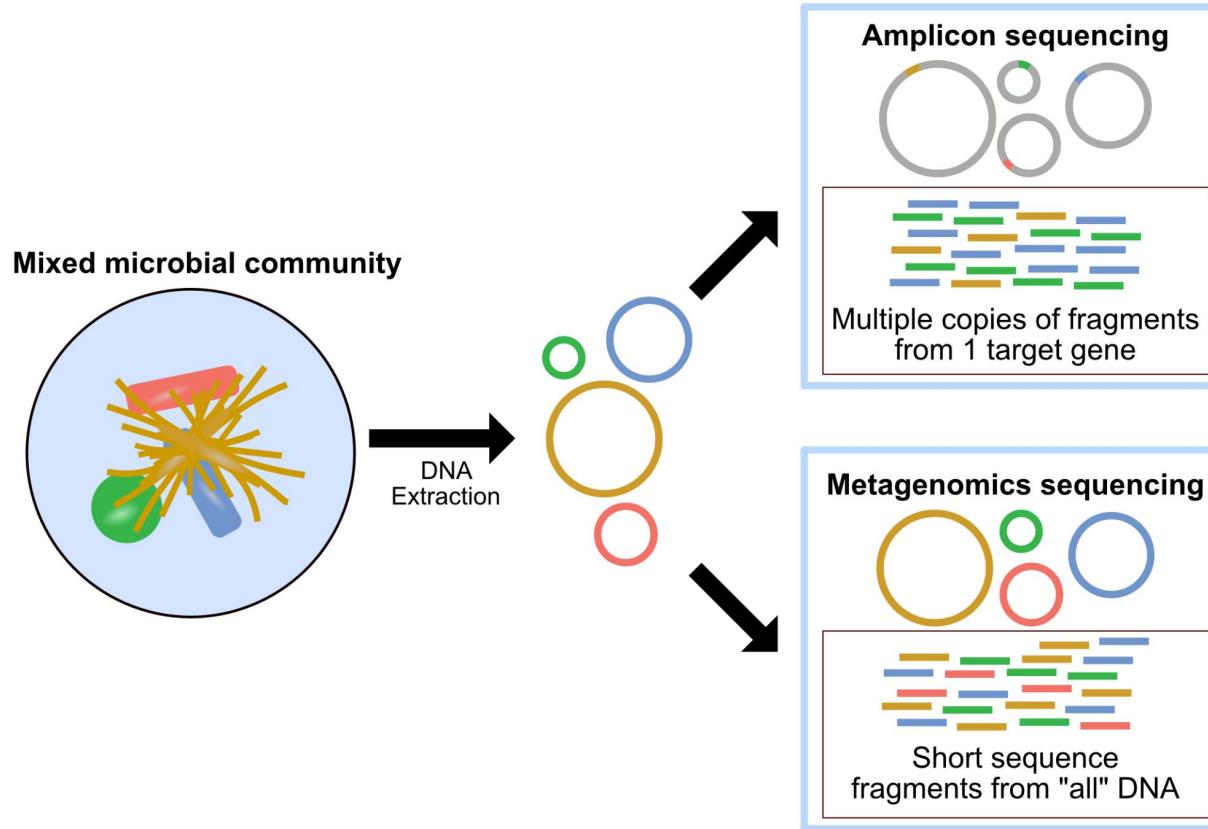
Example figures



b Functional network



Recap: Metagenomic vs amplicon sequencing



Metagenomic vs genomic sequencing

Metagenomic Sequencing

- **Purpose:** Metagenomic sequencing is used to analyze the **collective genome of a microbial community** directly from an environmental sample, without the need to culture organisms in the lab.

Genomic Sequencing

- **Purpose:** Genomic sequencing is focused on determining the complete DNA sequence of a **single organism's genome**.

Metagenomic vs genomic sequencing

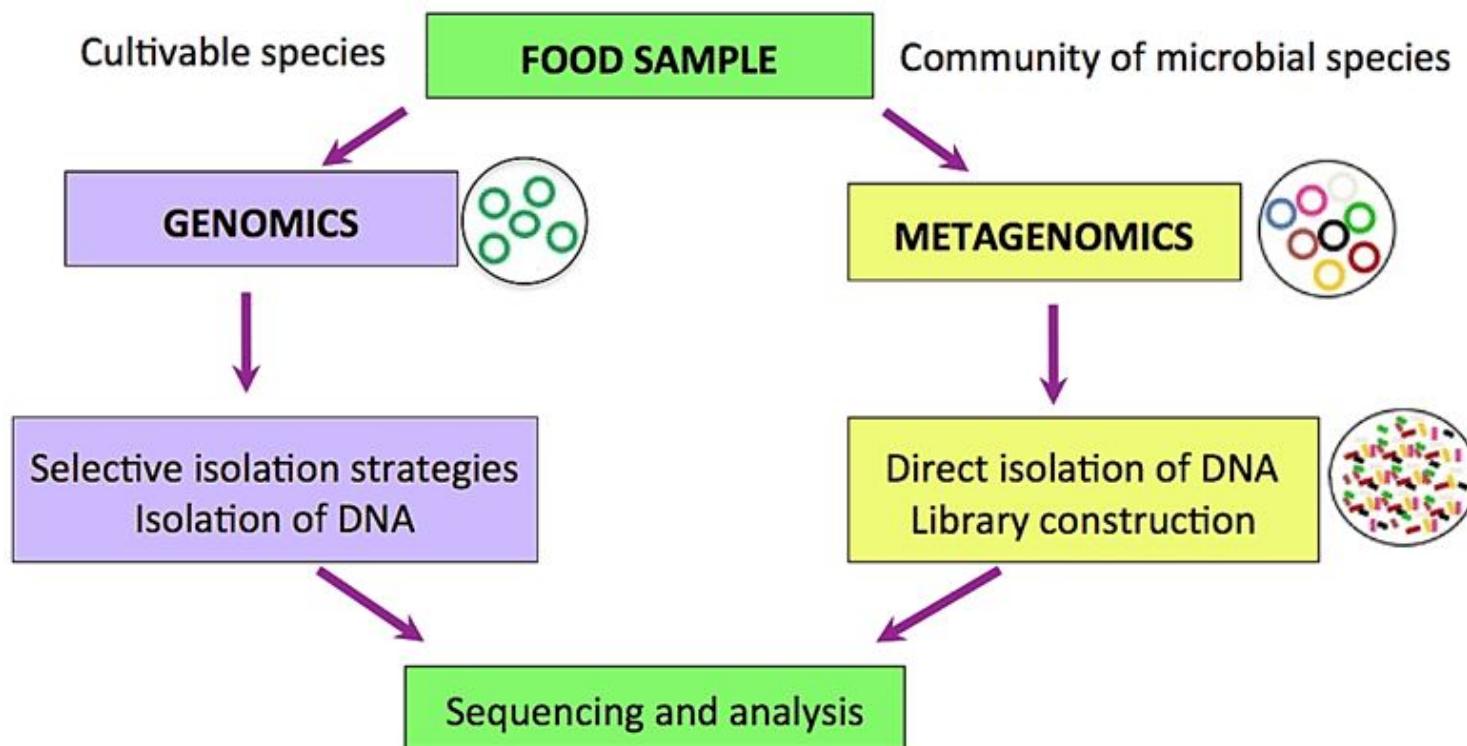
Metagenomic Sequencing

- **Purpose:** Metagenomic sequencing is used to analyze the **collective genome of a microbial community** directly from an environmental sample, without the need to culture organisms in the lab.
- Extract DNA from an entire sample with multiple community members

Genomic Sequencing

- **Purpose:** Genomic sequencing is focused on determining the complete DNA sequence of a **single organism's genome**.
- Usually requires isolating a single organism from a community of organisms and then extracting this DNA

Metagenomic vs genomic sequencing



Questions?