# Metagenomics Lesson 1
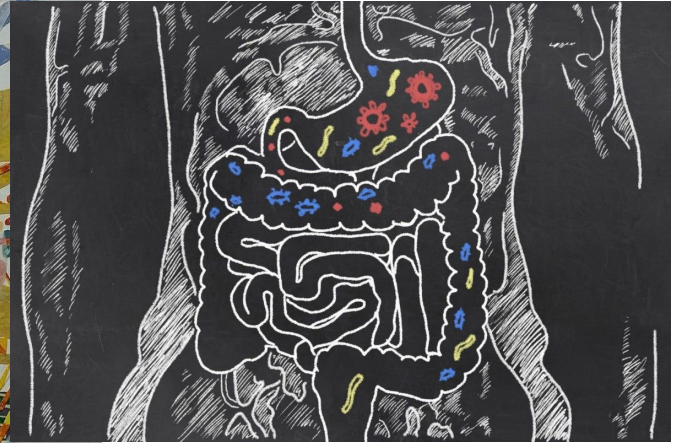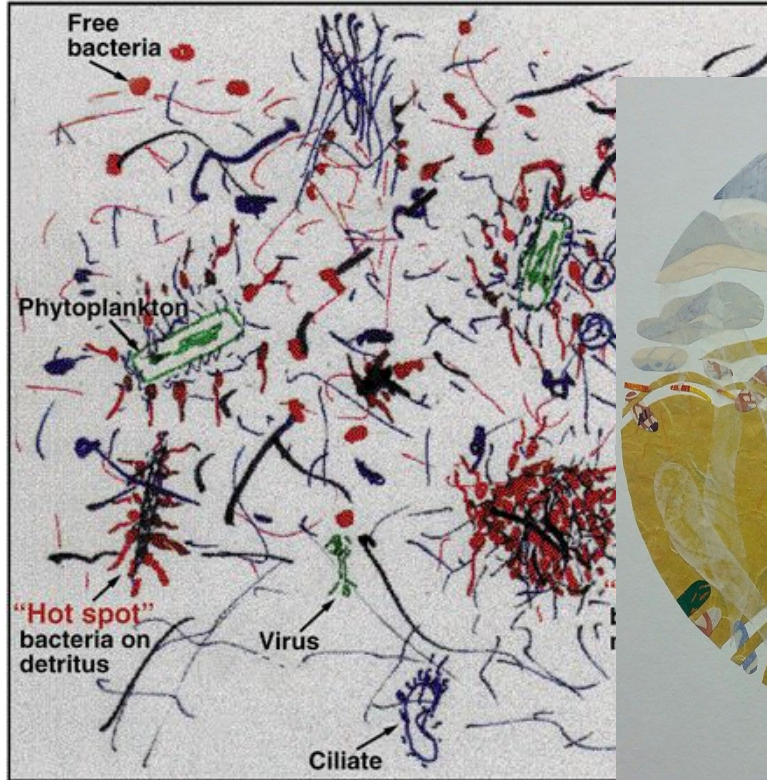
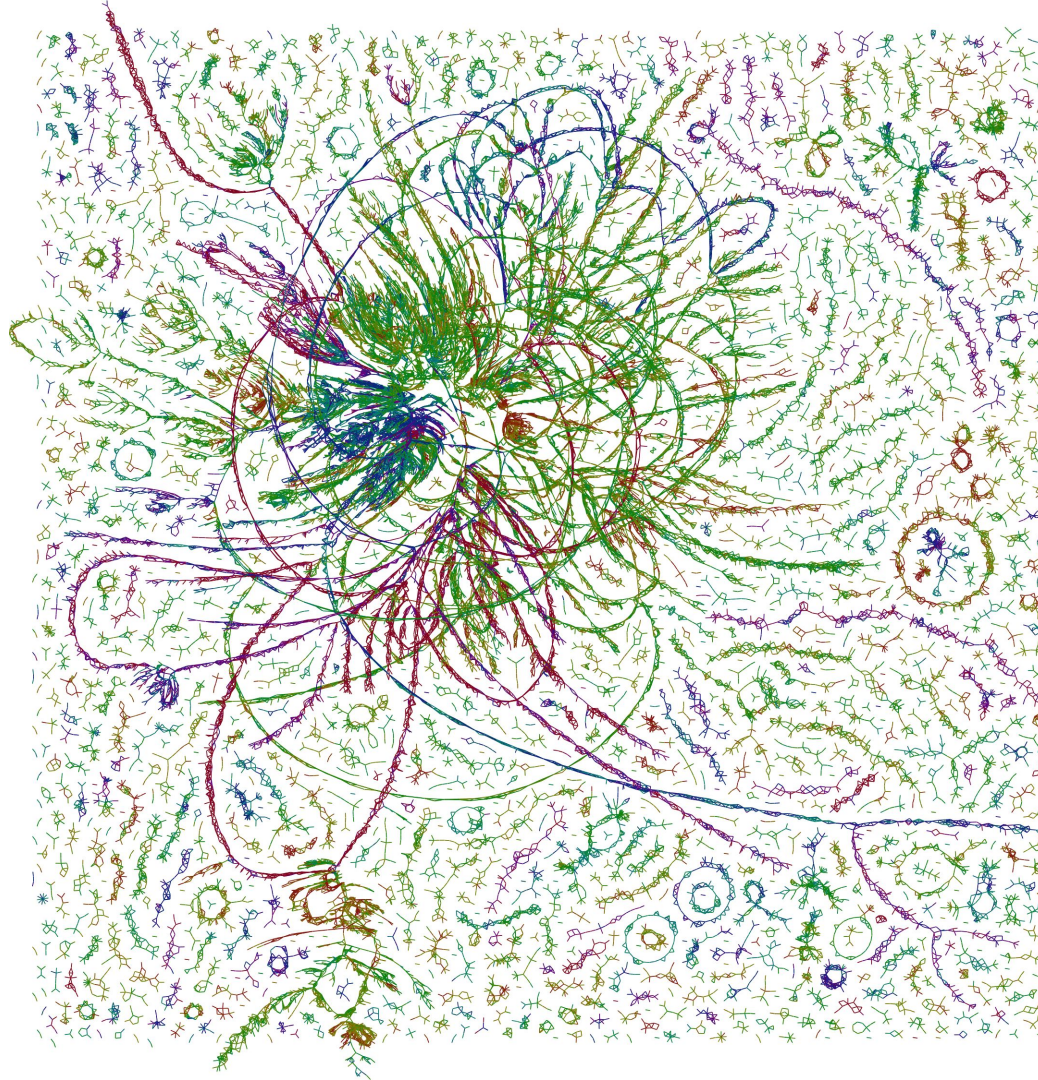# Metagenomics Lesson 1

**What is Metagenomics?**

**What kinds of questions can metagenomics be used to answer?**

**Is metagenomics right for me?**
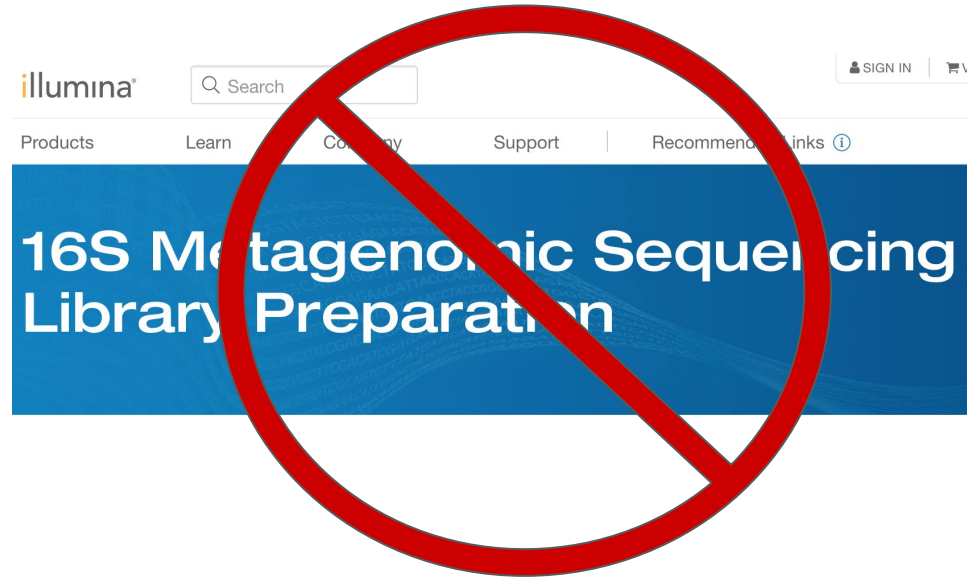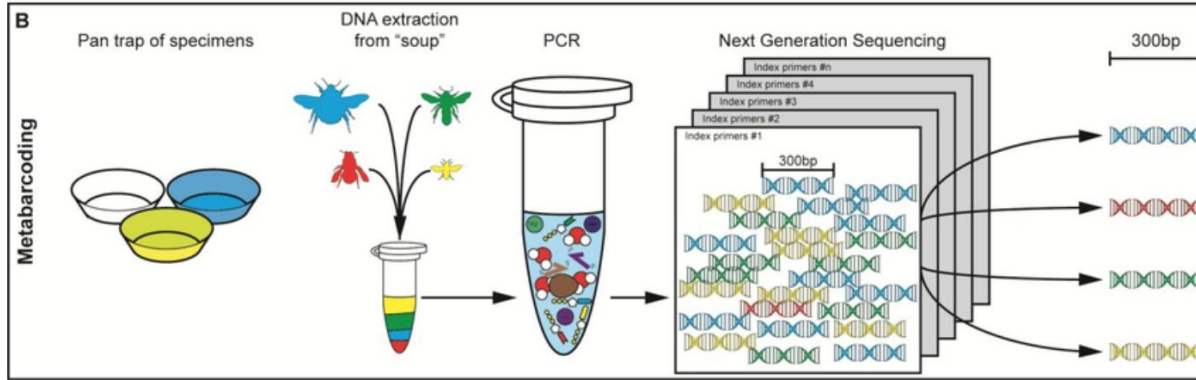
**Yay! You got data! Now what?**

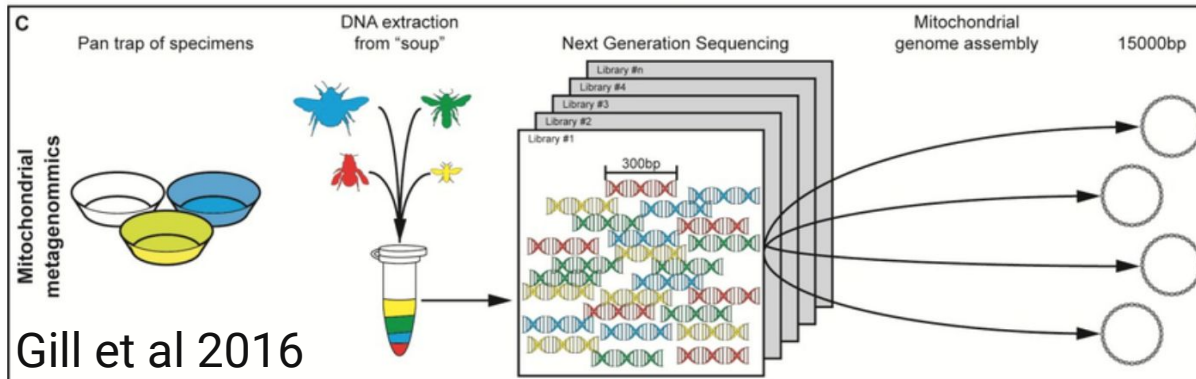# What is Metagenomics?

in
this
house
we

don't use "metagenomics"
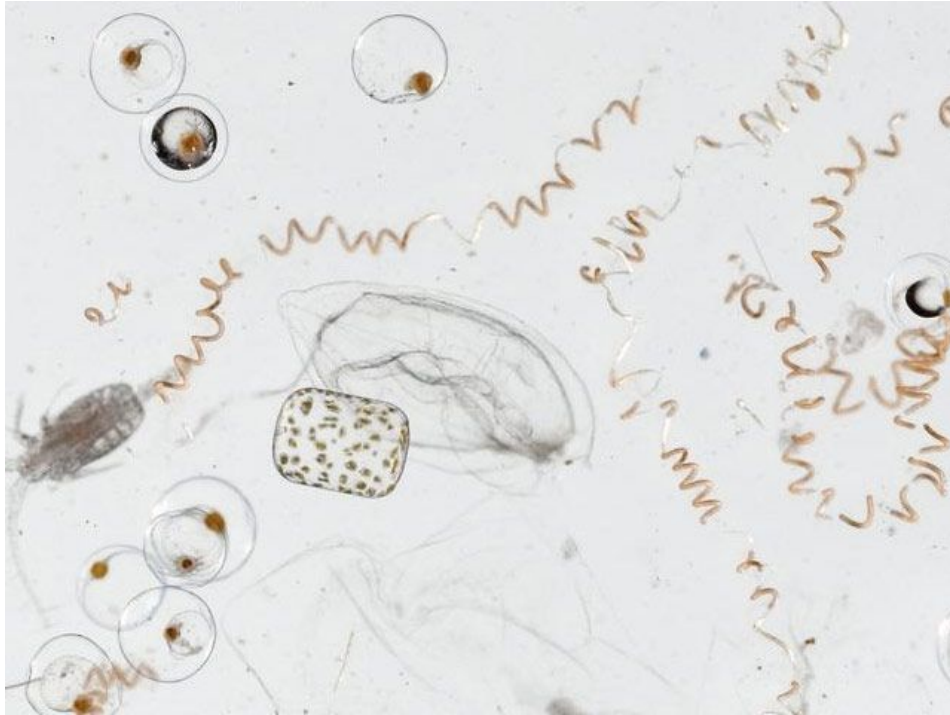to mean "amplicon sequencing"

# What is Metagenomics?



Gill et al 2016

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

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

*Incomplete* because one drop of seawater contains about...



$10^6$ bacteria/mL * $3\times10^6$ bp/bacteria = $3\times10^{12}$ bp/mL

+ $10^3$ euks/mL * $3\times10^8$ bp/euk = $3\times10^{11}$ bp/mL

= 3.3 Tbp/mL

= 210 MiSeq runs/mL

= 0.5 NovaSeq run/mL

0.5 Genbank/mL = 0.2 SRA / L

# One scoop of soil contains about...



$$10^{10} \text{ bacteria/g} * 4\times10^6 \text{ bp/bacteria} =$$
$$4\times10^{16} \text{ bp/g}$$

$$+ \quad 10^5 \text{ euks/g} * 3\times10^8 \text{ bp/euk} =$$
$$3\times10^{13} \text{ bp/g}$$

$$= 40 \text{ Pbp/g}$$

$$= 2.6\text{M MiSeq runs/g}$$

$$= 6{,}000 \text{ NovaSeq runs/g}$$

$$6{,}000 \text{ Genbank/g} = 2.8 \text{ SRA / g}$$

# One pinch of stool contains about...



$10^{11}$ bacteria/g * $4\times10^6$ bp/bacteria = $4\times10^{17}$ bp/g

+ $10^6$ euks/g * $3\times10^8$ bp/euk = $3\times10^{14}$ bp/g

+ $10^8$ colonocytes/g * $3\times10^9$ bp/cell = $3\times10^{17}$ bp/g

= 700 Pbp/g

= 45M MiSeq runs/g

= 100k NovaSeq runs/g

100k Genbank/g = 40 SRA / g

What kinds of questions can metagenomics be used to answer?

Metagenome Sequencing

Quality Control

**Who is there?**

Taxonomic Diversity
Phylogenetic Diversity ← Marker Gene Analysis

Taxonomic Diversity
Phylogenetic Diversity ← Binning
Novel Taxa

Genome Diversity
Novel Genomes ← Assembly

**What are they doing?**

Gene Prediction → Gene Diversity
Novel Genes

Functional Annotation → Protein Family Diversity
Functional Diversity

**Comparative Metagenomics**

Intercommunity Similarity
Metadata Correlations
Biomarker Detection

Sharpton (2014)

# What kinds of questions can metagenomics be used to answer?

**Who is there?**
**(Taxonomy & Molecular Evolution)**

- Is this gene present in this sample?
- How many homologs of this gene appear in this sample?
- Which genomes encode this gene?
- Is this pathogen present in this environment?
- How closely related is this uncultured strain to this cultured representative?
- How many ecotypes of this bacterium appear in this environment?
- How the h*ll many prokaryotic Phyla are there in the world??

**What are they doing?**
**(Community Ecology & Function)**

- What proteins do symbionts encode to mediate relationships with their host?
- Which genes/pathways/genomes co-occur in this environment?
- What antibiotic resistance genes does this community encode?
- How many different carbon fixation pathways exist in hydrothermal vents?
- Are there novel CRISPR-Cas systems yet to be discovered?

# Is metagenomics right for me?

**You might try amplicon sequencing if...**

- You need to detect rare genes or species
- You're working with eukaryotes
- You have many (1000s) samples to run

**You might try Quantitative PCR if...**

- You only care about presence/absence
- You want to quantify how many copies of a gene/species is present in a sample

**You might try isolate genomics if…**

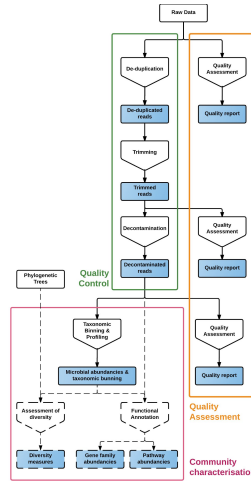- You can isolate your organism of interest
- You're working with eukaryotes

**You might try single-cell genomics if…**

- Your interest is population genetics
- Your interest is in novel taxa
- Your interest is horizontal gene transfer and pangenomes

**You might try meta-/transcriptomics if…**

- You want levels of gene expression
- You're working with eukaryotes

# MGnify

Submit, analyse, discover and compare microbiome data

Overview | Submit data | Text search | Sequence search | Browse data | Genomes | API | About | Help | Login

# Future Metagenomics Lessons

Taxonomic Classification

Assembly

Binning

For more info go to:
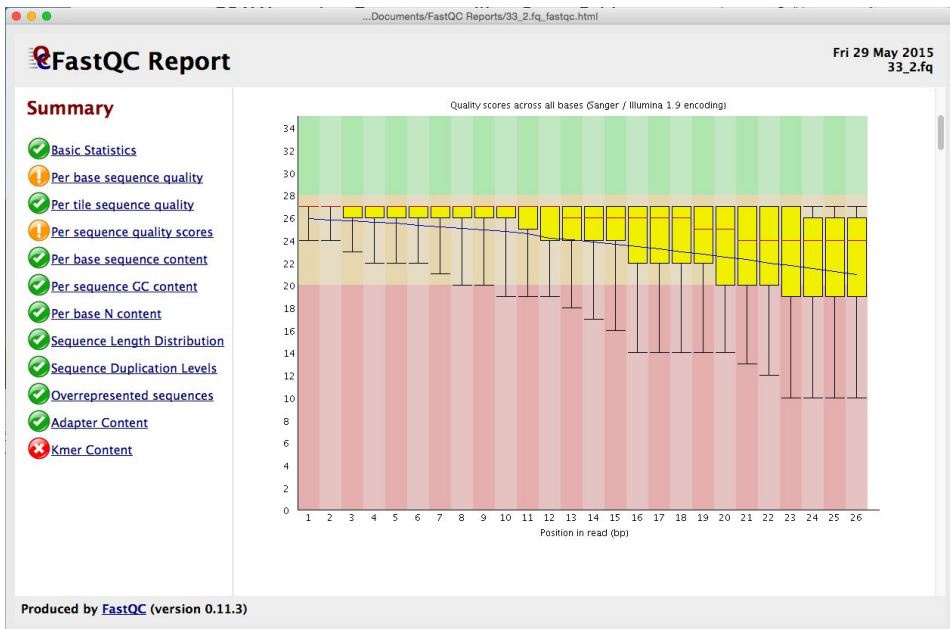https://github.com/biovcnet/topic-metagenomics

# Other BVCN Topics

#amplicons

#functionalannotation

#transcriptomics

#networkscience

#population-genetics-and-comparative-genomics

# Demo #1 by Alexis Marshall on Quality Control



https://www.youtube.com/watch?v=7jRTyfdIXLo