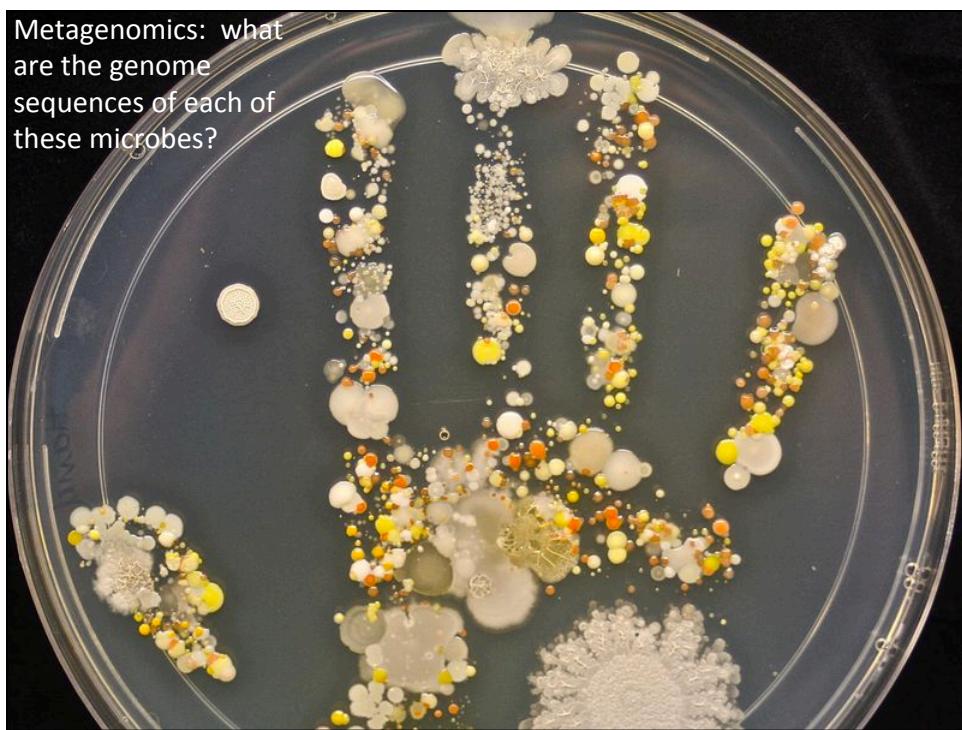


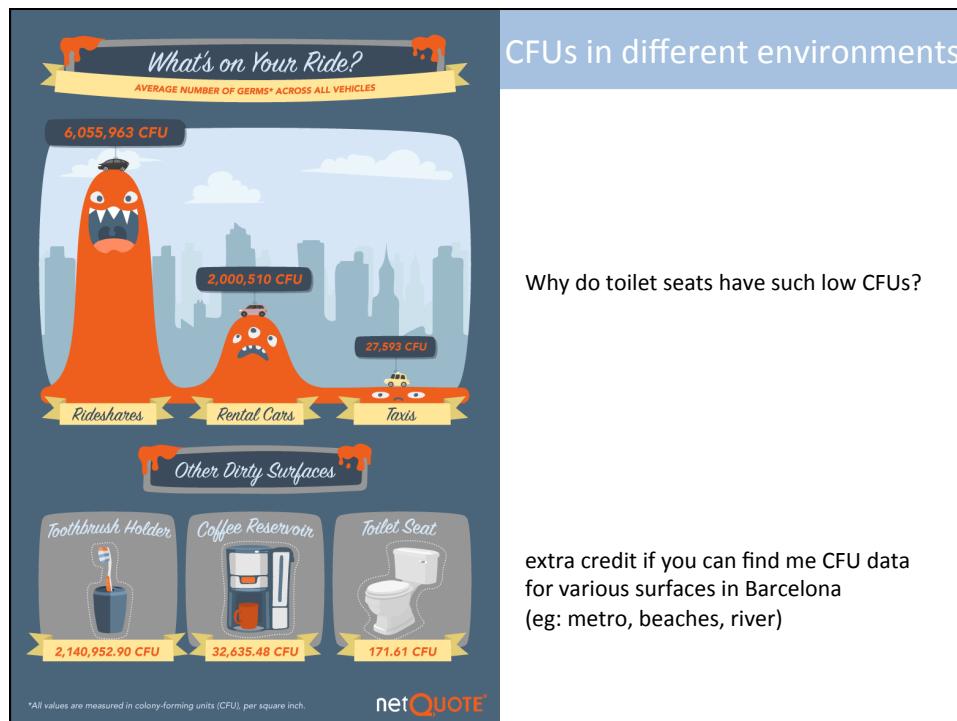
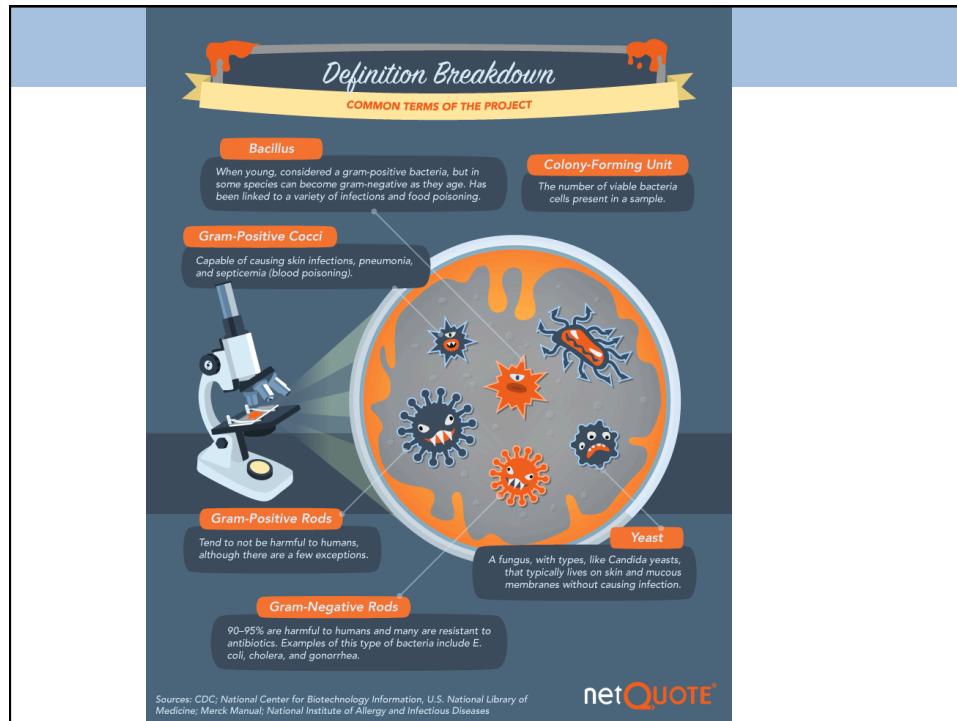
METAGENOMICS: Why should we care?

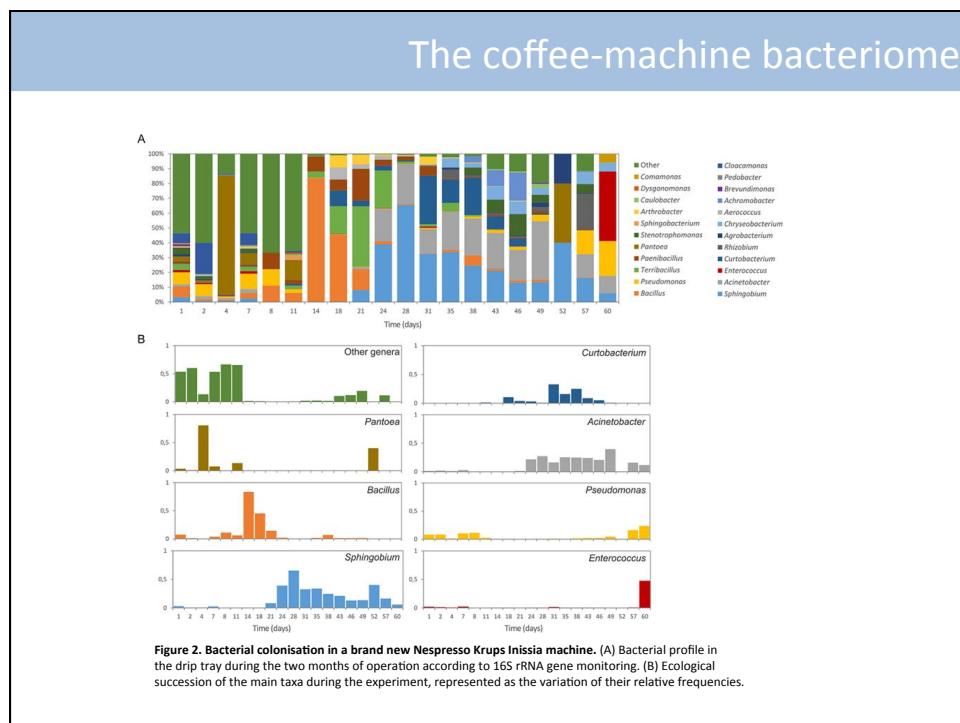
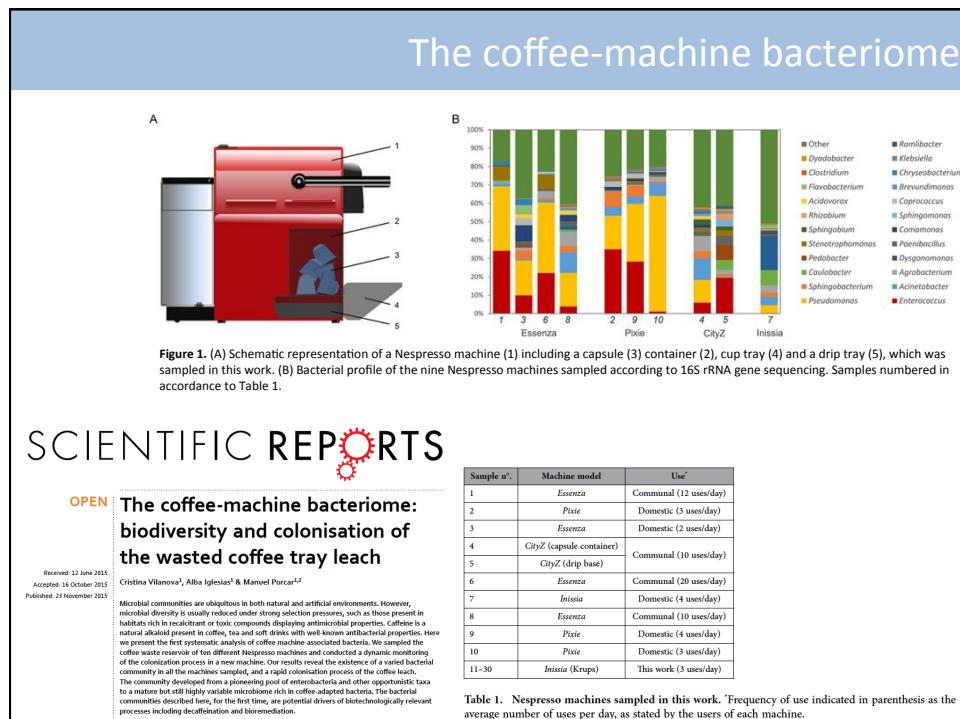
- (1) Almost all **antibiotics** and **anticancer drugs** are naturally occurring metabolites produced by microbes. Without finding new microbes, we won't find new drugs.
- (2) Microbes in our bodies cause disease, **prevent disease**, and interact with our diet and genome to modify their effect on health. Without knowing the makeup of each person's internal microbiome, our ability to diagnose and cure diseases will always be limited.
- (3) Microbes are found inside of many tumors, and **treatment of tumors** with antibiotics (which do not affect human cells) can either increase or decrease the growth rate of the tumor: microbes play strange and poorly understood roles in many diseases that we wouldn't normally think of as being microbial in nature.
- (4) The best way to **understand evolution** is to obtain sequences for thousands of genomes.
- (5) Microbes found in contaminated waters (eg: pig shit in farms, heavy-metal contamination from old factories) can often eat toxic compounds, and turn them into non-toxic compounds, providing a possible solution to pollution. **Microbes as cleaners**.

Metagenomics: How do we do it?

What microbes can we find where?





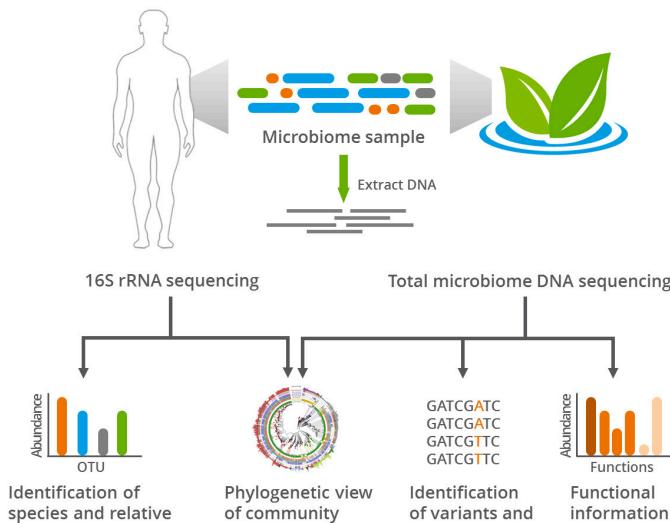


Group Activity

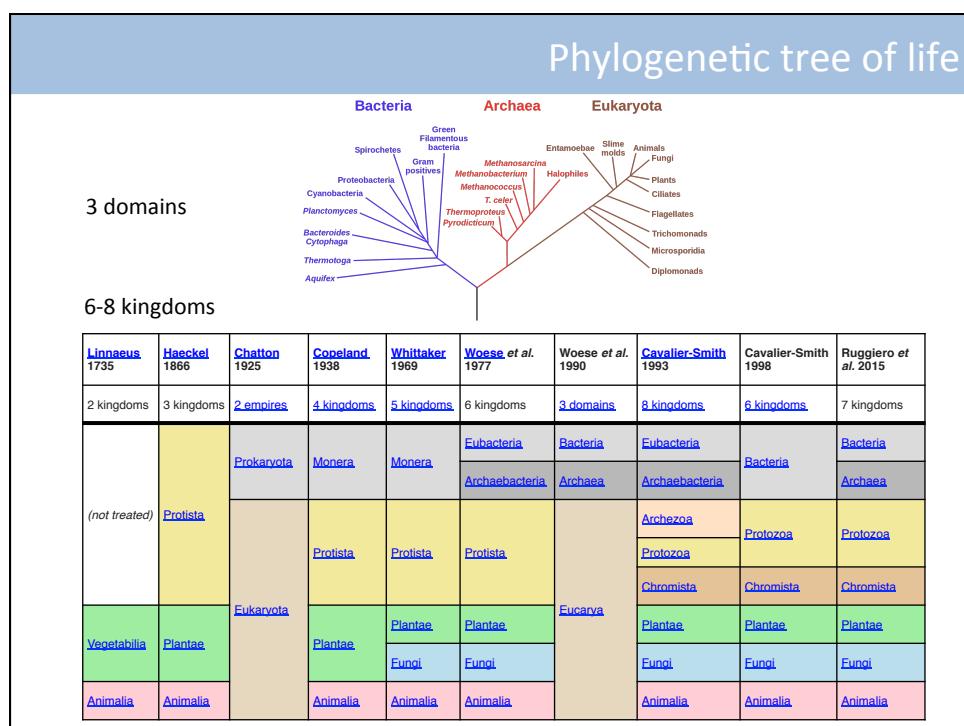
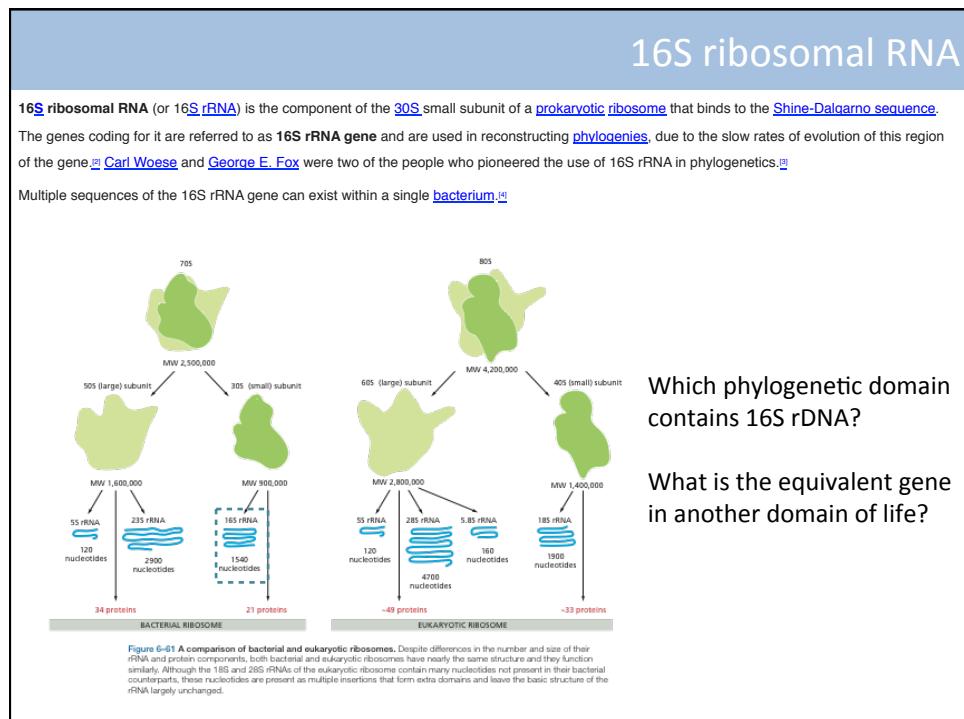
Make groups of four. As a group, take ten minutes to write down answers to each of the following questions:

1. Write down a **definition of metagenomics**
2. How would you **count the number of species** in 1L of sea water from Somorrostro beach?
3. How would you determine **how many viable individual bacterial cells** are in 1L of sea water from Somorrostro beach? What are some benefits and problems of your method? (Why is it good? What might it miss?)
4. How would you determine **how many unique bacterial species** are there in 1L of sea water from Somorrostro beach?
- Without sequencing? - With sequencing?
hints for sequencing: different species have different alleles (versions) of each gene. some genes, but not all genes, are unique to each species.

Sequencing to identify species



OTU = Operational Taxonomic Unit, a group of very similar 16S sequences



16S ribosomal RNA

What does the S in 16S stand for?

What are the units of this measure? nucleotides? seconds? amino acids?

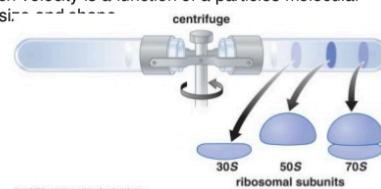
16S ribosomal RNA

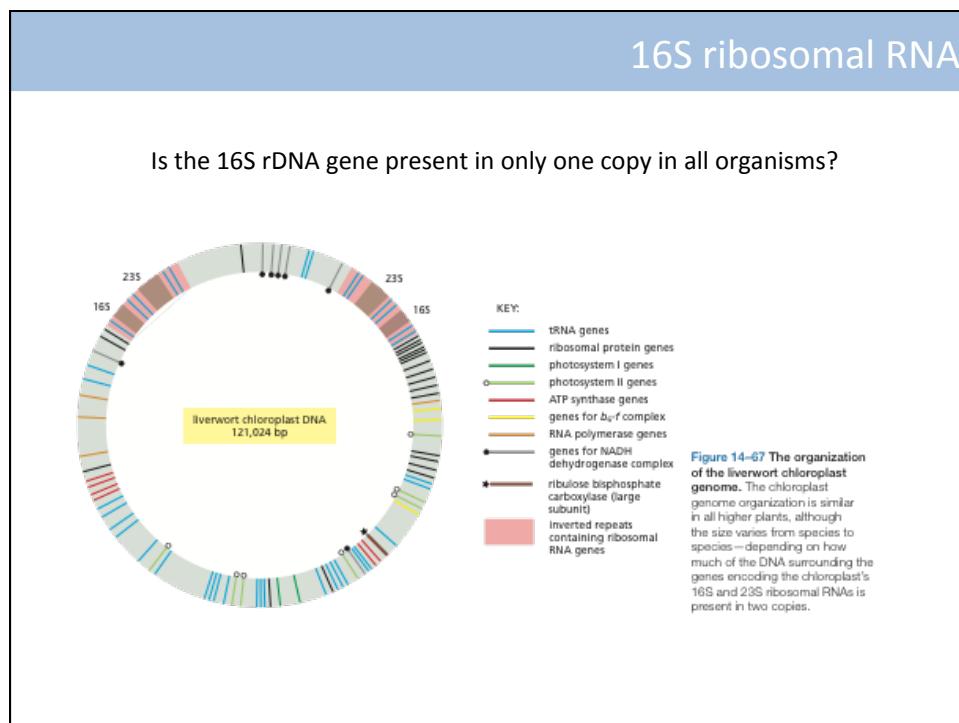
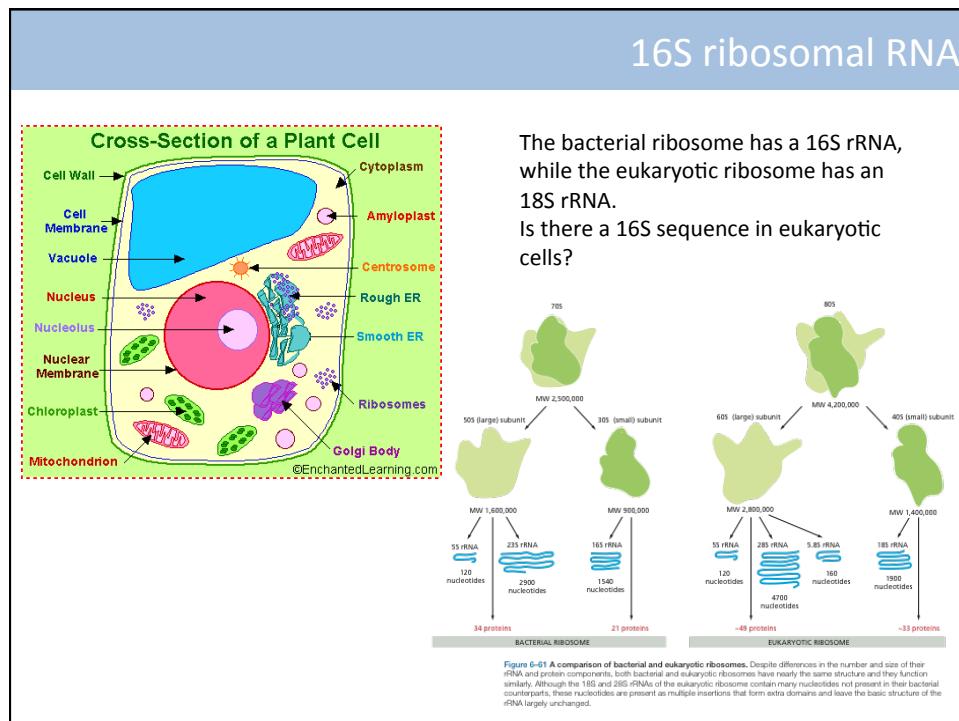
What does the S in 16S stand for?

What are the units of this measure? nucleotides? seconds? amino acids?

Svedberg Unit (S)

- The large and small subunit of ribosome are named according to the velocity of sedimentation when subjected to centrifugal force.
- The unit used to measure sedimentation velocity is **Svedberg (S)**.
- The larger value faster is the sedimentation velocity, hence larger the molecule.
- Named after inventor of ultracentrifuge Theodor Svedberg.
- The sedimentation velocity is a function of a particles molecular weight, volume, st^{-1}





How does 16S sequencing actually work?

In groups of four, design a PCR & sequencing experiment to determine which species are on the plate.

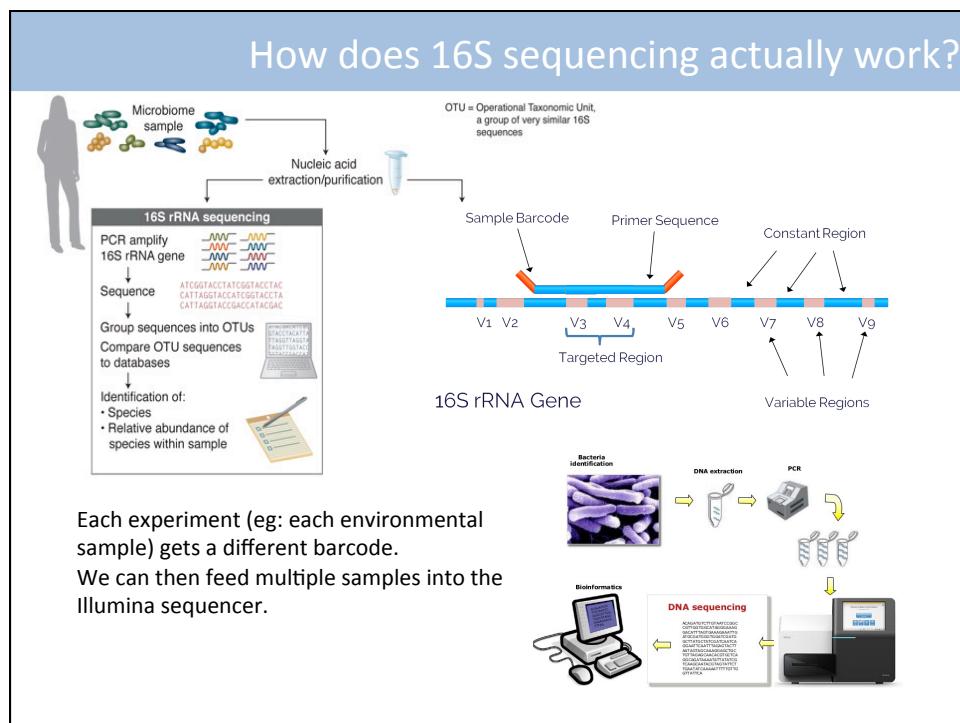
16S rRNA Gene

The diagram illustrates the PCR process on the 16S rRNA gene. It shows a template DNA strand with primers 1 and 2 annealed to the constant region. The PCR cycle results in an exponential increase in the number of DNA molecules. The final product is a double-stranded DNA molecule where each strand has a primer at one end and a sequence from the targeted region (V1-V9) in the middle. The variable regions (V1-V9) are highlighted in pink, while the constant region is blue.

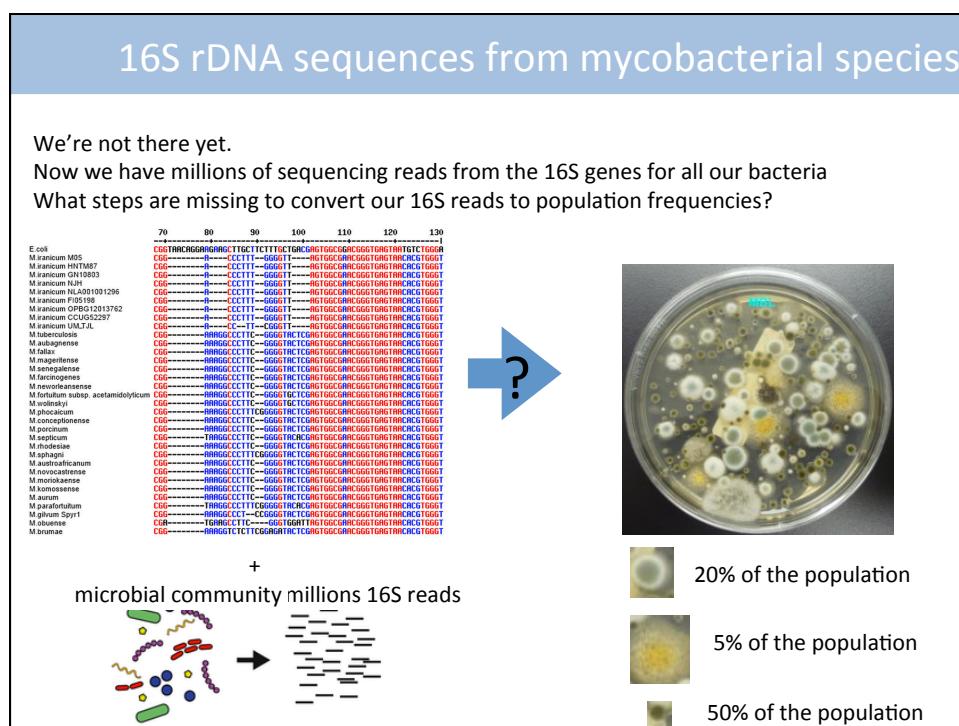
OTU Legend:

- Other
- Dyadobacter
- Clostridium
- Flavobacterium
- Acidovorax
- Rhizobium
- Sphingobium
- Stenotrophomonas
- Pedobacter
- Caulobacter
- Sphingobacterium
- Pseudomonas
- Ramlibacter
- Klebsiella
- Chryseobacterium
- Brevundimonas
- Coprococcus
- Sphingomonas
- Comamonas
- Paenibacillus
- Dysgonomonas
- Agrobacterium
- Acinetobacter
- Enterococcus

khan academy



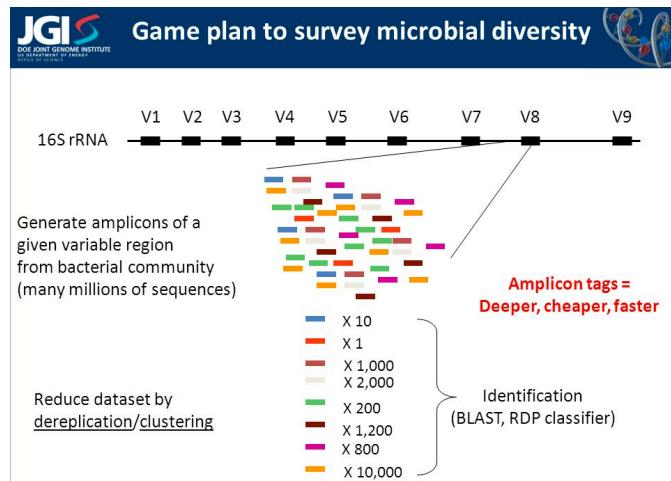
16S rDNA sequences from mycobacterial species												
Which species can you uniquely identify from these data? which can you not?												
	70	80	90	100	110	120	130					
E.coli	C	G	G	T	A	R	C	A	G	G	A	R
M. iranicum M05	C	G	G	T	A	R	C	T	C	T	T	G
M. iranicum HNTM87	C	G	G	T	A	R	C	T	C	T	T	G
M. iranicum GN10803	C	G	G	T	A	R	C	T	C	T	T	G
M. iranicum NJH	C	G	G	T	A	R	C	T	C	T	T	G
M. iranicum NL001001296	C	G	G	T	A	R	C	T	C	T	T	G
M. iranicum FI05198	C	G	G	T	A	R	C	T	C	T	T	G
M. iranicum OPBG12013762	C	G	G	T	A	R	C	T	C	T	T	G
M. iranicum CCUG5297	C	G	G	T	A	R	C	T	C	T	T	G
M. iranicum UMJL	C	G	G	T	A	R	C	T	C	T	T	G
M. tuberculosis	C	G	G	T	A	R	C	T	C	T	T	G
M. aubagnense	C	G	G	T	A	R	C	T	C	T	T	G
M. fallax	C	G	G	T	A	R	C	T	C	T	T	G
M. magiterrense	C	G	G	T	A	R	C	T	C	T	T	G
M. senegalense	C	G	G	T	A	R	C	T	C	T	T	G
M. farcinogenes	C	G	G	T	A	R	C	T	C	T	T	G
M. neworleansense	C	G	G	T	A	R	C	T	C	T	T	G
M. fortuitum subsp. acetamidolyticum	C	G	G	T	A	R	C	T	C	T	T	G
M. wolinskii	C	G	G	T	A	R	C	T	C	T	T	G
M. phocaicum	C	G	G	T	A	R	C	T	C	T	T	G
M. conceptionense	C	G	G	T	A	R	C	T	C	T	T	G
M. porcinum	C	G	G	T	A	R	C	T	C	T	T	G
M. septicum	C	G	G	T	A	R	C	T	C	T	T	G
M. rhodesiae	C	G	G	T	A	R	C	T	C	T	T	G
M. sphagni	C	G	G	T	A	R	C	T	C	T	T	G
M. austroafricanum	C	G	G	T	A	R	C	T	C	T	T	G
M. novocastrense	C	G	G	T	A	R	C	T	C	T	T	G
M. moriokaense	C	G	G	T	A	R	C	T	C	T	T	G
M. komossense	C	G	G	T	A	R	C	T	C	T	T	G
M. aurum	C	G	G	T	A	R	C	T	C	T	T	G
M. parafortuitum	C	G	G	T	A	R	C	T	C	T	T	G
M. galvini Spyrl	C	G	G	T	A	R	C	T	C	T	T	G
M. obuense	C	G	G	T	A	R	C	T	C	T	T	G
M. brumae	C	G	G	T	A	R	C	T	C	T	T	G



16S rDNA sequences from mycobacterial species

Two options:

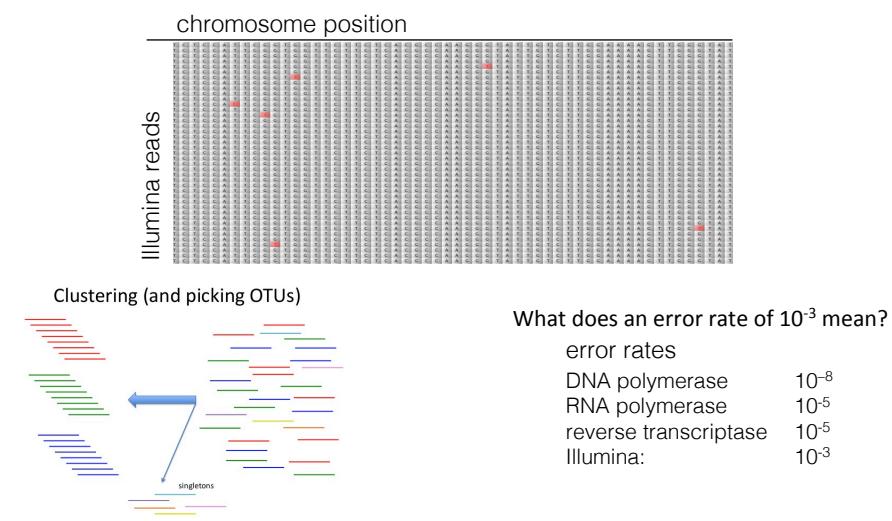
- (1) map (align) Illumina reads to known 16S sequences from known species
- (2) cluster similar sequence reads, and assume that each cluster represents one species (closely related group of organisms)



Problem: Illumina sequencing is not error free

Why are sequencing (and PCR amplification) errors a problem when counting species using 16S read counts? How do we solve this problem?

hint: will individual errors be more or less common than individual species?



Moving to whole genomes

The obvious method:

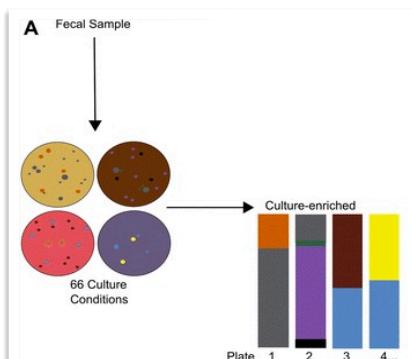


Isolate unique species



Do whole-genome sequencing of each species

hint:



What is a problem with this approach?

<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC492978>

Extracting genomes

Solution: assemble into contigs as best you can, then search databases of sequenced genomes. **But we want to find NEW stuff.**



...you only see what is in the database

Extracting genomes

Binning

Genomes from different species have different statistical properties

PhD student

Complex sample

Genomic signatures:

- GC / Codon usage
- Tetranucleotide frequency + statistical method

"Binning"

CENTER FOR MICROBIAL COMMUNITIES | AALBORG UNIVERSITY

Functional metagenomics

Maybe you don't care about the complete genome?
Just about the functions encoded in genes.

search for novel antibiotics:

```

graph LR
    A[Sewage sample] --> B[Genomic DNA extraction]
    B --> C[Genomic DNA sheared into -2kb fragments]
    C --> D[DNA fragments ligated into plasmid vector]
    D --> E[Transformation of electocompetent E.coli]
    E --> F[Functional screen on antibiotic-containing media]
    F --> G[PCR amplification of functionally-selected genes]
    G --> H[Sequencing and annotation]
    H --> I[CATCGA]
    
```

The flowchart illustrates the workflow for finding novel antibiotics. It starts with a "Sewage sample", followed by "Genomic DNA extraction", "Genomic DNA sheared into -2kb fragments", "DNA fragments ligated into plasmid vector", "Transformation of electocompetent *E.coli*", "Functional screen on antibiotic-containing media", "PCR amplification of functionally-selected genes", "Sequencing and annotation", and finally the sequence "CATCGA".

M. Sommer, DTU, Denmark (in prep)

CENTER FOR MICROBIAL COMMUNITIES | AALBORG UNIVERSITY

Activity 1

ARTICLE

OPEN

Insights into the phylogeny and coding potential of microbial dark matter

Christian Hahn¹, Patrick Schenck², Alexander Sczyrbi³, Natalia N. Ivanova⁴, Luis J. Andrade⁵, Ian Fang Cheng⁶, George Church⁷, Michael L. Smith⁸, Daniel C. Raskin⁹, Daniel C. Connon¹⁰, Daniel R. Relman¹¹, George Church¹², Stefan M. Sievert¹³, Wen-Tau Liu¹⁴, Jonathan A. Eisen¹⁵, Steven J. Hallam¹⁶, Nikos C. Kyrpides¹⁷, Kostas Konstantinidis¹⁸, Edward M. Rubin¹⁹, Philipp Hugenholtz²⁰, Sven Datta²¹

Genomic sequencing enhances our understanding of the biological world by providing blueprints for the evolutionary and functional diversity that shapes the biosphere. However, microbial genomes that are currently available are of limited diversity and coverage. To address this challenge, we have developed a strategy to target and sequence 30 uncultivated archaeal and bacterial cells from nine diverse habitats. By applying single-cell genomics to target and sequence 30 uncultivated archaeal and bacterial cells from nine diverse habitats, we are able to resolve many intra- and inter-phylum level relationships and to identify novel genes and metabolic pathways. This study greatly expands the genomic representation of the tree of life and challenges established hierarchies between the three domains of life. These include a novel archaeal clade and the rapid evolution of the tree of life. Our results also provide a framework for understanding the function of uncultivated microorganisms in their natural habitats, facilitating organism-level interpretation of ecosystem function. This study greatly expands the genomic representation of the tree of life and provides a systematic step towards a better understanding of biological evolution on our planet.

What problem in the culture, isolate, sequence method does single-cell whole-genome sequencing solve? Why?

Why did they assembly only 201 unique genomes from 3300 whole-genome sequences?

https://www.youtube.com/watch?v=9Q_lyY9qampI
<https://www.youtube.com/watch?v=fCWR-zdA2Xg>

Activity 2

Couldn't we just :

1. isolate all the cells from the environmental sample
2. extract the DNA
3. fragment the DNA and feed it into a sequencer
4. assemble the resulting reads into genomes

Why not?
What problems will this face?

The below pipeline is difficult / impossible. Why?
recovering genomes from metagenomic data

hint: imagine we tried to do this with people to assemble personal genomes. What % of reads would be unique to each individual?

Activity 3

How can we use Hi-C contact maps to solve the assembly problem in metagenomics?

