

Introduction to metabarcoding

Antonino Malacrino

EEOB, The Ohio State University

December, 16 2020

Outline

Introduction

DNA extraction

PCR

Metabarcoding workflow

Library preparation

Illumina sequencing technology

Sequencing run

Secondary data processing

Data analysis

More tools

Introduction

Which is our question?

- ▶ Who is in there?
- ▶ What they can do?
- ▶ What are they doing?
- ▶ Are they really doing it?



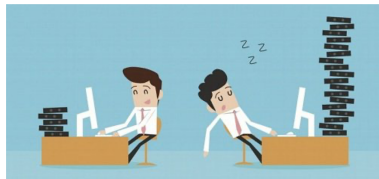
Which is our question?

- ▶ Who is in there?
- ▶ What they can do?
- ▶ What are they doing?
- ▶ Are they really doing it?



Which is our question?

- ▶ Who is in there?
- ▶ What they can do?
- ▶ What are they doing?
- ▶ Are they really doing it?



Which is our question?

- ▶ Who is in there?
- ▶ What they can do?
- ▶ What are they doing?
- ▶ Are they really doing it?

Molecular barcoding



Molecular barcoding



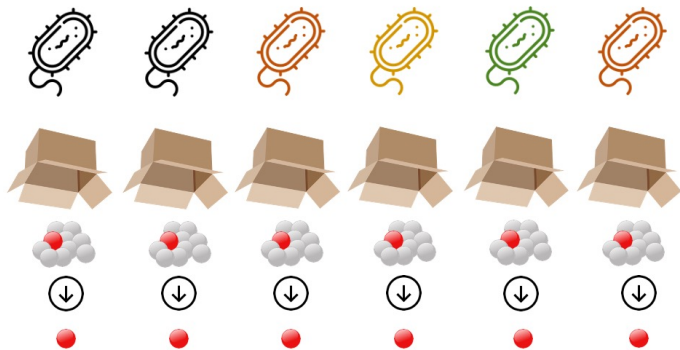
Molecular barcoding



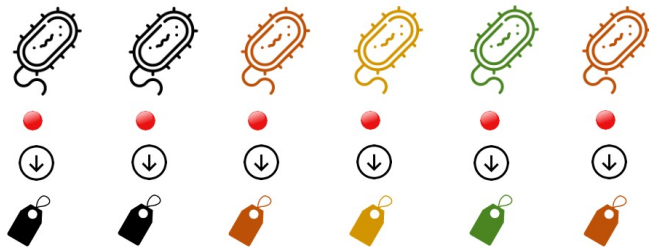
Molecular barcoding



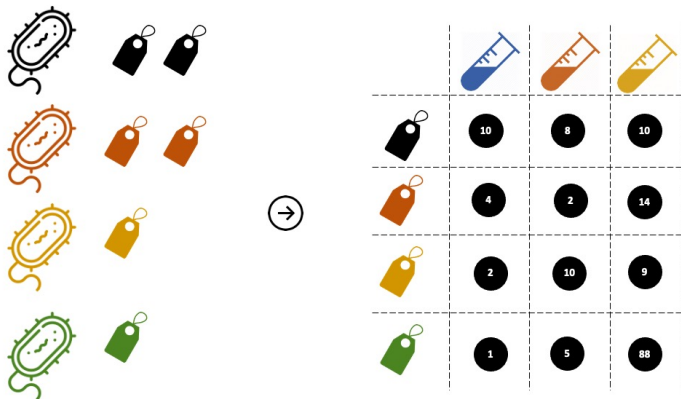
Metabarcoding



Metabarcoding



Metabarcoding



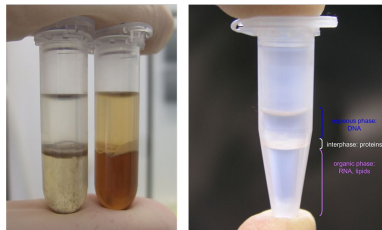
General workflow

1. Sample collection
2. DNA extraction
3. PCR amplification of our target
4. Adaptor ligation
5. Sequencing
6. Data analysis

DNA extraction

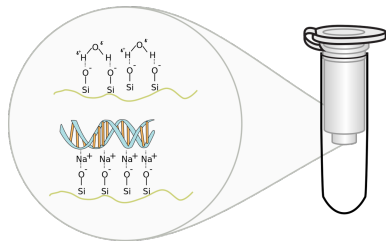
Different methods

- ▶ Organic solvents
- ▶ Spin columns
- ▶ SPRI beads
- ▶ Many other (CTAB, ...)



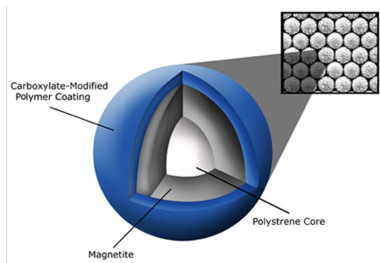
Different methods

- ▶ Organic solvents
- ▶ Spin columns
- ▶ SPRI beads
- ▶ Many other (CTAB, ...)



Different methods

- ▶ Organic solvents
- ▶ Spin columns
- ▶ SPRI beads
- ▶ Many other (CTAB, ...)

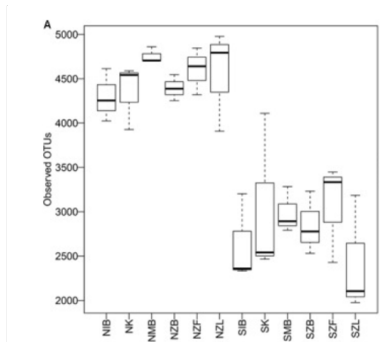


Different methods

- ▶ Organic solvents
- ▶ Spin columns
- ▶ SPRI beads
- ▶ Many other (CTAB, ...)

Technical considerations

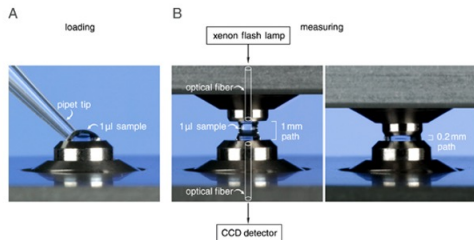
- ▶ DNA extraction introduces a bias in the final dataset
- ▶ A recent investigation on 322 studies shows they used 72 different methods. 14 did not report such info!



Technical considerations

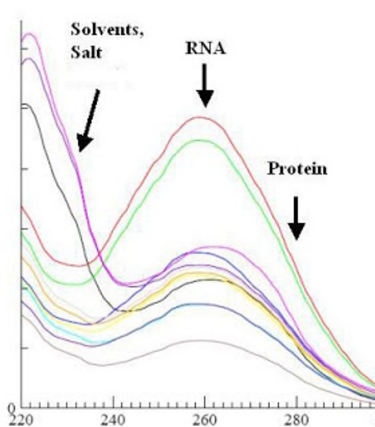
- ▶ DNA extraction introduces a bias in the final dataset
- ▶ A recent investigation on 322 studies shows they used 72 different methods. 14 did not report such info!

QC - Nanodrop



QC - Nanodrop

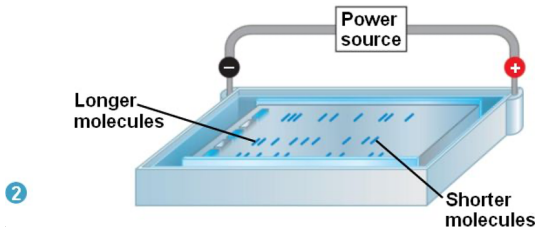
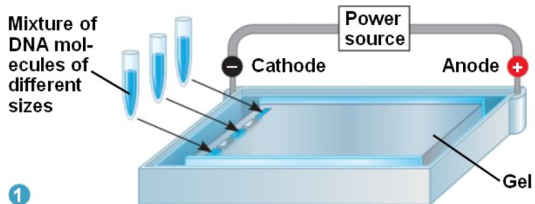
- ▶ DNA concentration ng/ul
- ▶ Ratio 260/230
- ▶ Ratio 260/280



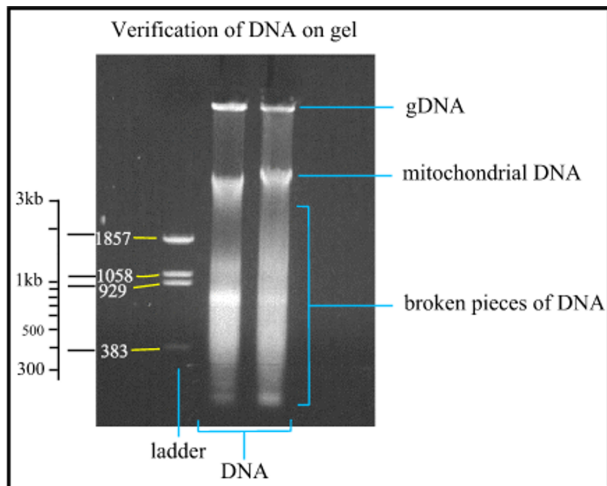
QC - Qubit



QC - Gel Electrophoresis



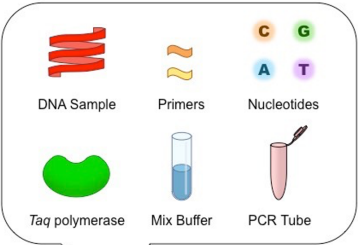
QC - Gel Electrophoresis



PCR

PCR

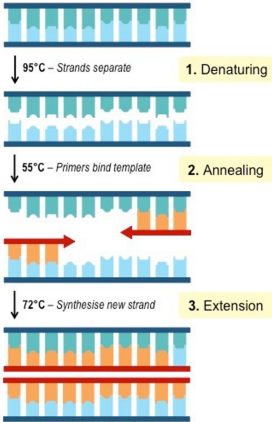
PCR Components



Thermal Cycler



PCR Process (ONE Cycle)

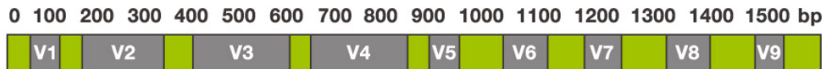


Thermocycler



Technical considerations

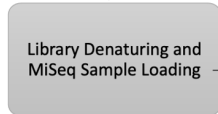
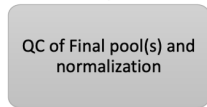
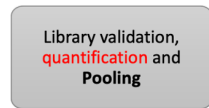
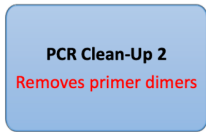
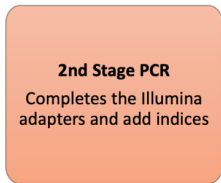
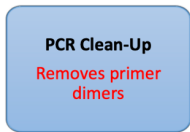
- ▶ Target gene
 - ▶ Who is our target? Bacteria, fungi, insects, fish, specific genus
 - ▶ Is the resolution optimal for our question?
 - ▶ Are PCR primers available or we have to design them?
 - ▶ Is the taxonomy database available or we have to build a custom one?
- ▶ PCR bias
 - ▶ Use a Hi-Fi polymerase
 - ▶ Use optimal annealing temperature
 - ▶ Do not exaggerate with PCR cycles
 - ▶ Run multiple PCRs on the same sample



CONSERVED REGIONS: unspecific applications

VARIABLE REGIONS: group or species-specific applications

Metabarcoding workflow



Control samples

- ▶ **Negative control 1.** Run molecular biology grade water throughout the pipeline. Pool it with the other samples even if you do not see amplification.
- ▶ **Negative control 2.** This is the negative control from your first PCR. If you see a band, discard the entire batch of samples and start again. If no band is observed, sequence anyway.
- ▶ **Mock community.** Pre-built or custom.

Control samples

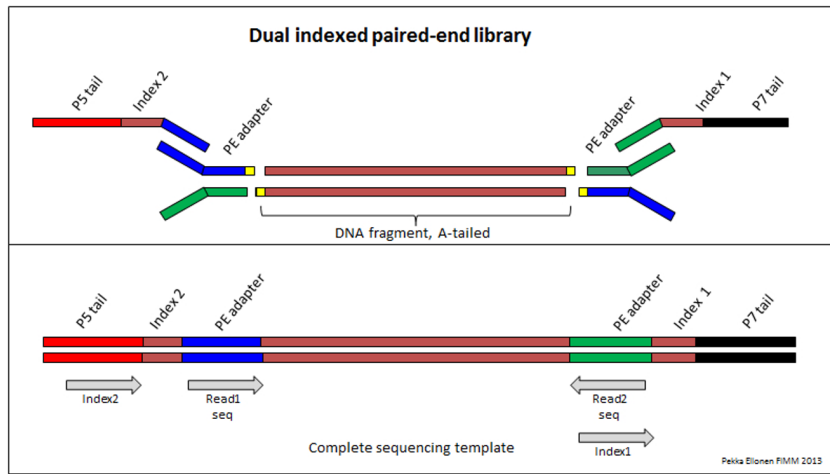
- ▶ **Negative control 1.** Run molecular biology grade water throughout the pipeline. Pool it with the other samples even if you do not see amplification.
- ▶ **Negative control 2.** This is the negative control from your first PCR. If you see a band, discard the entire batch of samples and start again. If no band is observed, sequence anyway.
- ▶ **Mock community.** Pre-built or custom.

Control samples

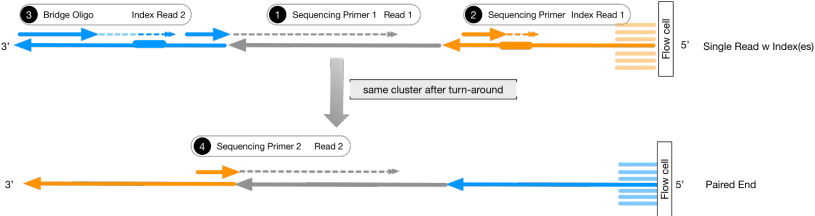
- ▶ **Negative control 1.** Run molecular biology grade water throughout the pipeline. Pool it with the other samples even if you do not see amplification.
- ▶ **Negative control 2.** This is the negative control from your first PCR. If you see a band, discard the entire batch of samples and start again. If no band is observed, sequence anyway.
- ▶ **Mock community.** Pre-built or custom.

Library preparation

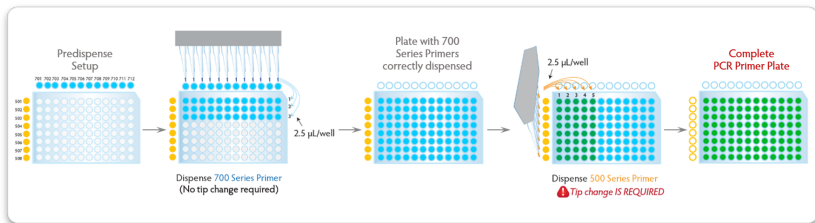
Multiplexing



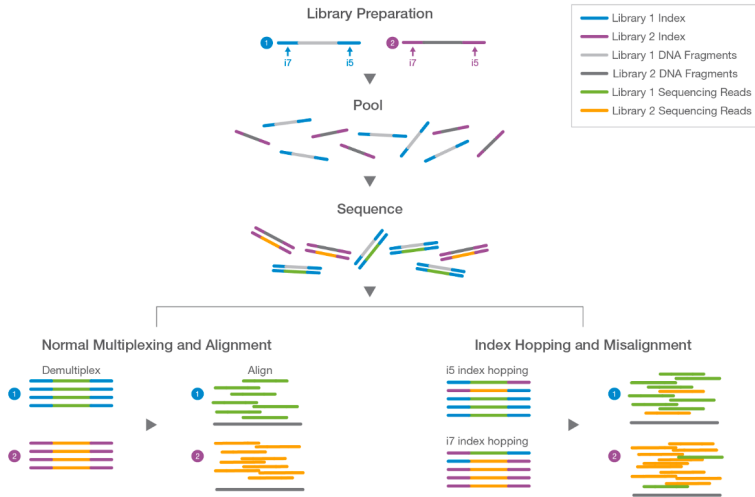
Multiplexing



Multiplexing



Index hopping



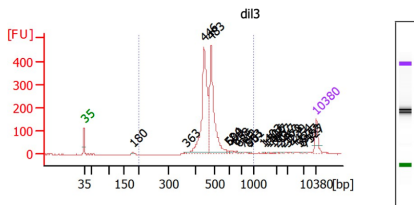
Equimolar pooling

Goal: guarantee that all samples are sequenced at the same depth

1. Qubit
2. Calculate nM
3. Pool according to sample concentration

Final quality control

- ▶ Bioanalyzer / Tapestation
- ▶ qPCR
- ▶ Qubit



Overall Results for sample 3 : dil3

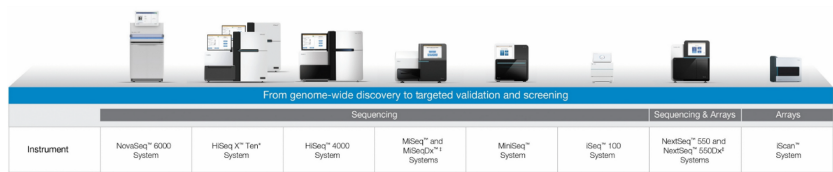
Number of peaks found: 23
 Noise: 1.2
 Corr. Area 1: 2,193.1

Region table for sample 3 : dil3

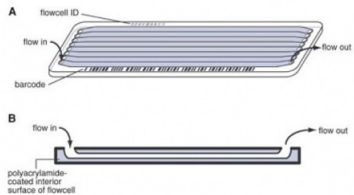
From [bp]	To [bp]	Corr. Area	% of Total	Average Size [bp]	Size distribution in CV [%]	Conc. [pg/ μ l]	Molarit y [pmol/l]	Co lo r
200	1,000	2,193.194	477	12.8		1,496.61	4,823.2	■

Illumina sequencing technology

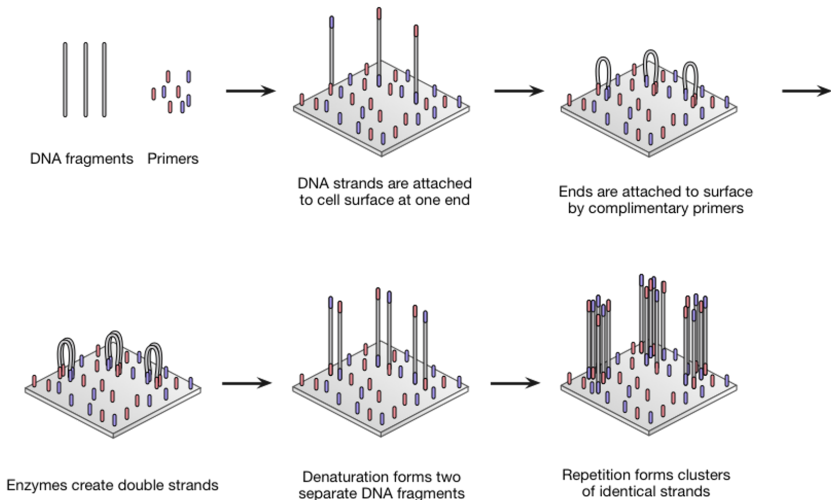
Illumina sequencing platforms



Flow cell

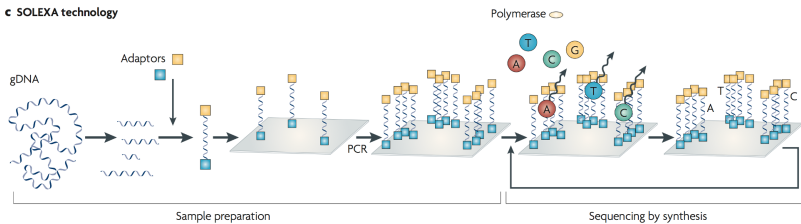


Illumina sequencing technology

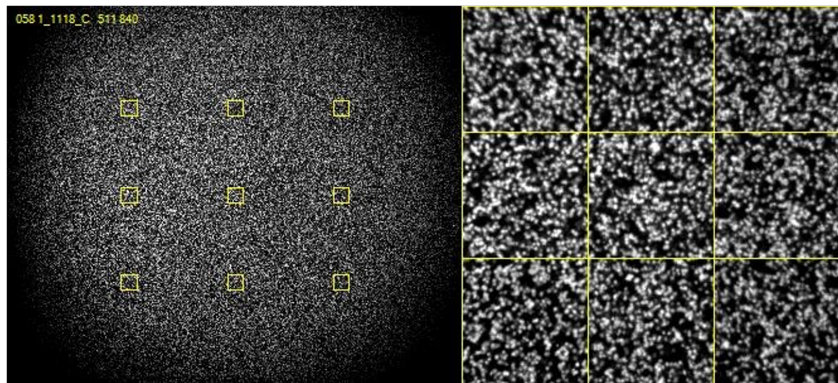


Illumina sequencing technology

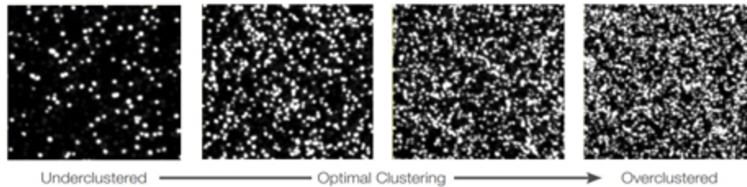
c SOLEXA technology



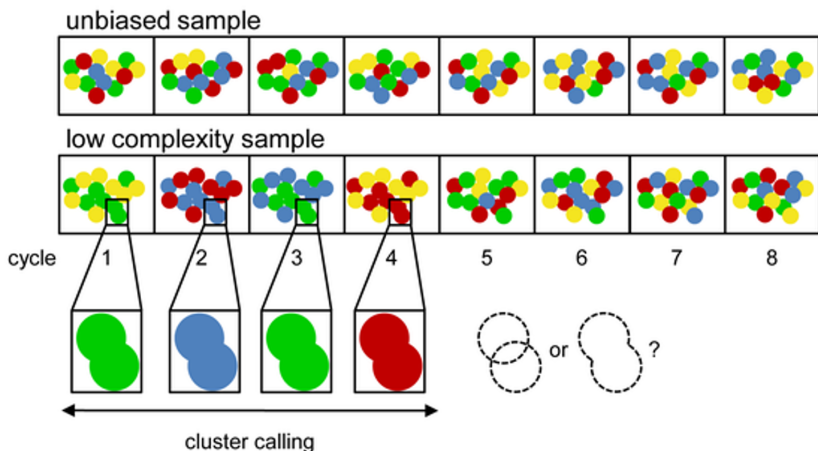
Illumina sequencing technology



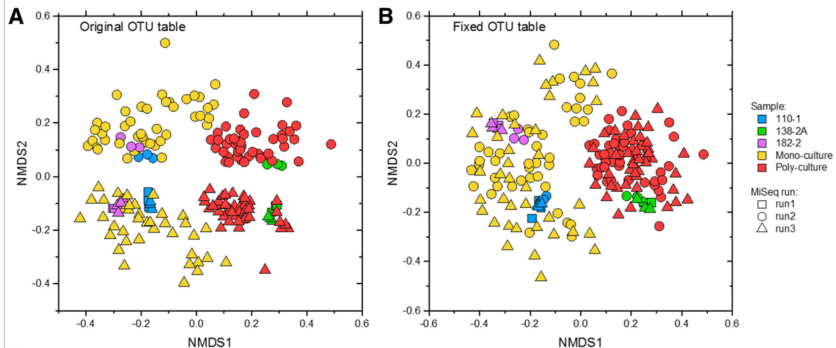
Illumina sequencing technology



Illumina sequencing technology



Technical considerations



Sequencing run

Sequencing Analysis Viewer

Run Folder: Q:1170214_K00150_0166_AHHJJHBBXX

Browse

Refresh

Analysis **Imaging** Summary Indexing

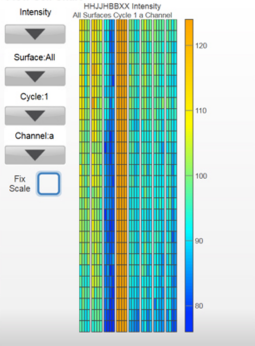
Status

Extracted: 166

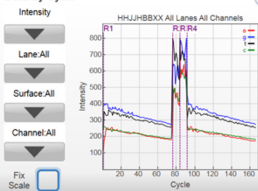
Called: 166

Scored: 166

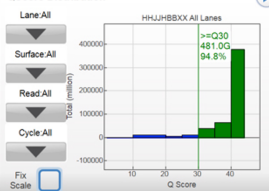
Flow Cell Chart



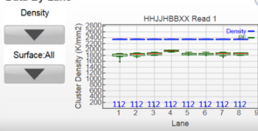
Data By Cycle



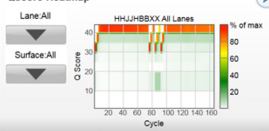
QScore Distribution



Data By Lane

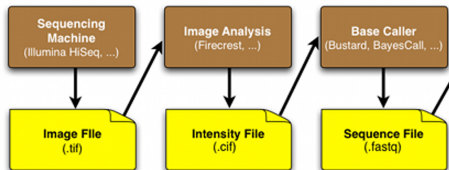


QScore Heatmap



Primary data analysis

- ▶ Primary analysis
 - ▶ From image to base calling
 - ▶ Cluster detection (4th cycle)
 - ▶ Cluster intensity correction
 - ▶ Base calling
 - ▶ Clusters are filtered (CPF, 25th cycle)
 - ▶ Q scores are assigned to each base
- ▶ CASAVA
 - ▶ Demultiplex
 - ▶ Create fastq files



Primary data analysis

- ▶ Primary analysis
 - ▶ From image to base calling
 - ▶ Cluster detection (4th cycle)
 - ▶ Cluster intensity correction
 - ▶ Base calling
 - ▶ Clusters are filtered (CPF, 25th cycle)
 - ▶ Q scores are assigned to each base
- ▶ CASAVA
 - ▶ Demultiplex
 - ▶ Create fastq files

Secondary data processing

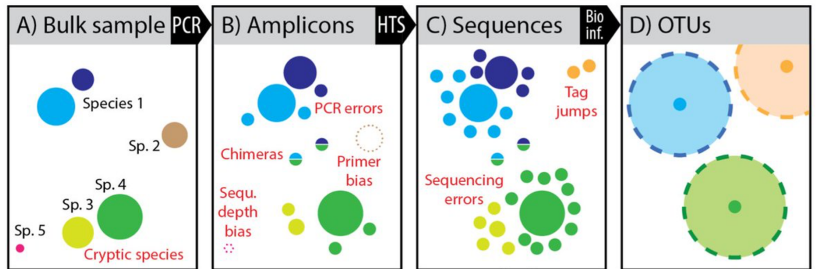
Workflow

- ▶ De-multiplexing
- ▶ Reads pre-processing
- ▶ Dereplication
- ▶ Clustering of variants
- ▶ Filtering of artefacts
- ▶ Alignment to references
- ▶ Taxonomy annotation
- ▶ Downstream analysis

First choice!

- ▶ OTUs (Operational Taxonomic Unit)
- ▶ AVSs (Amplicon Sequence Variant)

Remember: bias is everywhere!



Which is our goal?

Metadata

	sampleid	Treatment	Years
1			
2	AM-16S-1	maize-mono	1999
3	AM-16S-2	maize-mono	1999
4	AM-16S-3	maize-mono	1999
5	AM-16S-4	push-pull	1999
6	AM-16S-5	push-pull	1999

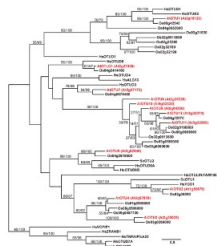
Taxonomy

ASV148428	Bacteria(100);Proteobacteria(100);Gammaproteobacte...
ASV212114	Bacteria(100);Cyanobacteria(100);Cyanobacteria(100);...
ASV9620	Bacteria(100);Proteobacteria(100);Alphaproteobacteri...
ASV147186	Bacteria(100);Proteobacteria(100);Betaproteobacteri...
ASV89359	Bacteria(100);Proteobacteria(100);Alphaproteobacteri...
ASV1061	Bacteria(100);Proteobacteria(100);Gammaproteobacte...
ASV328581	Bacteria(100);Bacteroidetes(100);Bacteroidia(100);Bac...
ASV86104	Bacteria(100);Proteobacteria(100);Alphaproteobacteri...

OTU table

	AM-16S-1	AM-16S-2	AM-16S-3
AM-16S-1	10	8	10
AM-16S-2	4	7	14
AM-16S-3	2	10	9
AM-16S-4	1	5	8

Phylogenetic tree



*.fastq files

```
Identifier ● @SRR566546.970 HWUSI-EAS1673_11067_FC7070M:4:1:2299:1109 length=50
Sequence ● TTGCCTGCCTATCATTTTAGTGCCTGTGAGGTGGAGATGTGAGGATCAGT
'+' sign ● +
Quality scores ● hhhhhhhhhghghghhhhhfhhhhfffffe'ee['X]b[d[ed'[Y[~Y
Identifier ● @SRR566546.971 HWUSI-EAS1673_11067_FC7070M:4:1:2374:1108 length=50
Sequence ● GATTTGTATGAAAGTATACAACATAAACTGCAGGTGGATCAGAGTAAGTC
'+' sign ● +
Quality scores ● hhhhgfhhcghghggfcffdhfehhhhcehdchhdhahehffffde'bVd
```

Demultiplexing

A
Library Preparation



B
Pool



C
Sequence



Sequence Output to Data File

CATTTCGACGGATCG
AACTGAGTCCGATA
AACTGATCGGATCC
CATTTCGTGGCAGTC
AACTGAACCTGATG
AACTGAGATTACAA
CATTTCGCAGTTCATT
CATTTCGAACTTCGA

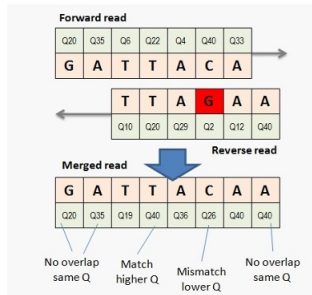
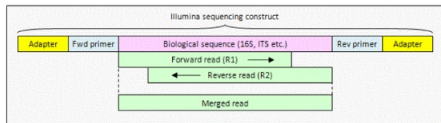
D
Demultiplex

1
CATTTCGACGGATCG
CATTTCGTGGCAGTC
CATTTCGCAGTTCATT
CATTTCGAACTTCGA

2
AACTGAGTCCGATA
AACTGATCGGATCC
AACTGAACCTGATG
AACTGAGATTACAA

- Library 1 Barcode
- Library 2 Barcode
- Sequencing Reads
- DNA Fragments
- Reference Genome

Merge PE reads



Remove suspicious reads

```
>GQY1XT001A6MUA
AATGGTACCCGTCAATTCATTGATCCTTGCGGTTTCGTTTACGGCGTGGACTACCAGTCGCACTCCAGTCATACAGTTTCCAATG
>GQY1XT001BTRWS
AATGGTACCCGTCAATTCCTTTGATCCTTGCGGGCCGTTTACGGCGTGGACTACCAGTCGCACTCGAGCTGCACAGTTTCCAAGCAGTTCCGGGGTTGGG
>GQY1XT001AK4J0
TCTAGCCGCACAGTTTCAAAGCACTCCCAGGGTT
>GQY1XT001BBPBR
AATGGTACCCGTCAATTCATTGACCGTTGCCCCCGTTTACTGTGCGGACTACCAGTCGCACTCAAGGCCCCAGTTTCAACGG
>GQY1XT001BDDE9
AATGGTACCCGTCAATTCCTTTAATCTTGCGGGTTCGTTTACGGCGTGGACTACCAGTCGCACTCCAGTTACACAGTTTCCAGAG
>GQY1XT001CIUF3
AATGGTACCCGTCAATTCCTTTGATCCTTGCGGGCCTTTACGGCGTGGACTACCAGGCGCCCTCCAGCCCGGCAGTTTCCAGTGCAGTCCCGGGTT
>GQY1XT001BKRP5
AATGGTACCCGTCAATTCATTAATCTCTCCCCCTTTCCCCCCCCCCCCCTTTCCCCCCCCCCCCCTTTCCCCCCCCCCC
>GQY1XT001B44ZE
AATGGTACCCGTCAATTCATTAACCTTGCGGGGTTTTACCGCGTGGACTACCAGGCGCCCTCAAGAAGAACAGTTTGAACGCAGCTATGGGTT
>GQY1XT001CIW3P
AATGGTACCCGTCAATTCATTGACCGTTGCCTCTCGTTTACTGCGTGGACTACCAGTCGCACTCAAGGCCCCCA
>GQY1XT001A731D
AATGGTACCCGTCAATTCATTAACGTTGCCCCGTTACTGCGTGGACTACCAGGGCAATCAAGACTGCCA
```

Trimming same length

```
>GQY1XT001A6MUA
AATGGTACCCGTCAATTCATTTGATCTTGCGGTTGTTTTACGGCGTGGACTACCAGTCGCACTCCAGTCATACAGTTTCCAATG
>GQY1XT001BTRWS
AATGGTACCCGTCAATTCCTTTGATCTTGCGGGCCGTTTACGGCGTGGACTACCAGTCGCACTCGAGCTGCACAGTTTCCAAGCAGTTCCGGGGTTGGG
>GQY1XT001AK4J0
TCTAGCCGCACAGTTTCAAAGCACTCCCAGGGTT
>GQY1XT001BBPBR
AATGGTACCCGTCAATTCATTTGACGTTGCCCCCGTTTTACTGTGCGGACTACCAGTCGCACTCAAGGCCCCAGTTTCAACGG
>GQY1XT001BDDE9
AATGGTACCCGTCAATTCCTTTAATCTTGCGGGTCTTTACGGCGTGGACTACCAGTCGCACTCCAGTTACAGTTTCCAGAG
>GQY1XT001CIUF3
AATGGTACCCGTCAATTCCTTTGATCTTGCGGGCCTTTACGGCGTGGACTACCAGGCCCTCCAGCCGGAGTTTCCAGTGCAGTCCCGGGTT
>GQY1XT001BKRP5
AATGGTACCCGTCAATTCATTTAATCTCTCCCCCTTTCCCGCCCCCCCCCTTTCCCGCCCCCCCCCTTTCCCGCCCCCCC
>GQY1XT001B44ZE
AATGGTACCCGTCAATTCATTTAACCTTGCGGGGTTTTACGGCGTGGACTACCAGGCCCTCAAGAAGAAAGTTTTGAACGCAGCTATGGGT
>GQY1XT001CIW3P
AATGGTACCCGTCAATTCATTTGACGTTGCCCTCTGTTTTACTGCGTGGACTACCAGTCGCACTCAAGGCCCCCA
>GQY1XT001A731D
AATGGTACCCGTCAATTCATTTAACGTTGCCCCGTTACTGCGTGGACTACCAGGGCAATCAAGACTGCCCA
```

Dereplication

>GQY1XT001A6MUA

AATGGTACCCGTCAATTCATTTGATCTTGCGGTTTCGTTTACGGCGTGGACTACCAGTCGCACTCCAGTCATA

>GQY1XT001BTRWS

AATGGTACCCGTCAATTCCTTTGATCTTGCGGGCCGTTTACGGCGTGGACTACCAGTCGCACTCGAGCTGCA

>GQY1XT001BBPBR

AATGGTACCCGTCAATTCATTTGATCTTGCGGTTTCGTTTACGGCGTGGACTACCAGTCGCACTCCAGTCATA

>GQY1XT001BDDE9

AATGGTACCCGTCAATTCATTTGATCTTGCGGTTTCGTTTACGGCGTGGACTACCAGTCGCACTCCAGTCATA

>GQY1XT001CIUF3

AATGGTACCCGTCAATTCCTTTGATCTTGCGGGCCGTTTACGGCGTGGACTACCAGTCGCACTCGAGCTGCA

>GQY1XT001B44ZE

AATGGTACCCGTCAATTCATTTGATCTTGCGGTTTCGTTTACGGCGTGGACTACCAGTCGCACTCCAGTCATA

>GQY1XT001CIW3P

AATGGTACCCGTCAATTCCTTTGATCTTGCGGGCCGTTTACGGCGTGGACTACCAGTCGCACTCGAGCTGCA

>GQY1XT001A731D

AATGGTACCCGTCAATTCATTTGATCTTGCGGTTTCGTTTACGGCGTGGACTACCAGTCGCACTCCAGTCATA

Dereplication

```
>GQY1XT001A6MUA DEPTH = 5  
AATGGTACCCGTC AATTCATTTGATCTTGCGGTTTCGTTTACGGCGTGGACTACCAGTCGCACTCCAGTCATA  
>GQY1XT001BTRWS DEPTH = 3  
AATGGTACCCGTC AATTCCTTTGATCTTGCGGGCCGTTTACGGCGTGGACTACCAGTCGCACTCGAGCTGCA
```

```
>GQY1XT001BBPBR  
AATGGTACCCGTC AATTCATTTGATCTTGCGGTTTCGTTTACGGCGTGGACTACCAGTCGCACTCCAGTCATA  
>GQY1XT001BDDE9  
AATGGTACCCGTC AATTCATTTGATCTTGCGGTTTCGTTTACGGCGTGGACTACCAGTCGCACTCCAGTCATA  
>GQY1XT001CIUF3  
AATGGTACCCGTC AATTCCTTTGATCTTGCGGGCCGTTTACGGCGTGGACTACCAGTCGCACTCGAGCTGCA  
>GQY1XT001B44ZE  
AATGGTACCCGTC AATTCATTTGATCTTGCGGTTTCGTTTACGGCGTGGACTACCAGTCGCACTCCAGTCATA  
>GQY1XT001CIW3P  
AATGGTACCCGTC AATTCCTTTGATCTTGCGGGCCGTTTACGGCGTGGACTACCAGTCGCACTCGAGCTGCA  
>GQY1XT001A731D  
AATGGTACCCGTC AATTCATTTGATCTTGCGGTTTCGTTTACGGCGTGGACTACCAGTCGCACTCCAGTCATA
```

Cluster variants

```
>*S16-0000006
TACGTTTATCGCGTT-AGCTTCGCCAAGCACAGCATCCTGCGCTTAGCCAACGTACATCGTTTAGGGTGTGGACTAA
>#S16-0000046
TACGTTTATCGCGTTTAGCTTCGCCAAGCACAGCATCCTGCGCTTAGCCAACGTACATCGTTTAGGGTGTGGACTAA
>#S16-0000241
TACGTTTATCGCGTT-AGCTTCGCCAAGCACAGCATCCTGCGCTTAGCCAACGTACATCGT-TAGGGTGTGGACTAA
>#S16-0000375
TACGTTTATCGCATT-AGCTTCGCCAAGCACAGCATCCTGCGCTTAGCCAACGTACATCGTTTAGG-TGTGGACTAA
>*S16-0000001
GGCACTTAAAGCGTTAGCTACGGCGCAGAAACCACGGGTGG-CCCCCACACCTAGTGCCCAACGTTTACAGCGTGGT
>#S16-0000209
GGCACTTAAAGCGTTAGCTACGGCGCAGAAACCACGGGTGGTCCCCCACACCTAGTGCCCAACGTTTACAGCGTGGG
>#S16-0000667
GGCACTTAAAGCGTTAGCTACGGCGCAGAAACCACGGGTGG-CCCCCACACCTAGTGC-CAACGTTTACAGCGTGGT
>*S16-0000004
TCGACTTAACGCGTTAGCTCCGGAAGCCACGCCTCAAGG-GCACAACCTCCAAGTCGACATCGTTTACGGCGTGGAT
>#S16-0000625
TCGACTTAACGCGTTAGCTCCGGAAGCCACGCCTCAAGG-GCACAACCTCCAAGTCGACATCGT-TACGGCGTGGAT
>#S16-0000673
TCGACTTAACGCGTTAGCTCCGGAAGCCACGCCTCAAGGGGCACAACCTCCAAGTCGACATCGTTTACGGCGTGGAT
```


Cluster variants

```
>*S16-0000006 DEPTH + 3
TACGTTTATCGCGTT-AGCTTCGCCAAGCACAGCATCCTGCGCTTAGCCAACGTACATCGTTTAGGGTGTGGACTAA
>*S16-0000001 DEPTH + 2
GGCACTTAAAGCGTTAGCTACGGCGCAGAAACCACGGGTGG-CCCCCACACCTAGTGCCCAACGTTTACAGCGTGGT
>*S16-0000004 DEPTH + 2
TCGACTTAACGCGTTAGCTCCGGAAGCCACGCCTCAAGG-GCACAACCTCCAAGTCGACATCGTTTACGGCGTGGAT
```

```
>#S16-0000046
TACGTTTATCGCGTTTAGCTTCGCCAAGCACAGCATCCTGCGCTTAGCCAACGTACATCGTTTAGGGTGTGGACTAA
>#S16-0000241
TACGTTTATCGCGTT-AGCTTCGCCAAGCACAGCATCCTGCGCTTAGCCAACGTACATCGT-TAGGGTGTGGACTAA
>#S16-0000375
TACGTTTATCGCAATT-AGCTTCGCCAAGCACAGCATCCTGCGCTTAGCCAACGTACATCGTTTAGG-TGTGGACTAA
>#S16-0000209
GGCACTTAAAGCGTTAGCTACGGCGCAGAAACCACGGGTGG-TCCCCCACACCTAGTGCCCAACGTTTACAGCGTGGG
>#S16-0000667
GGCACTTAAAGCGTTAGCTACGGCGCAGAAACCACGGGTGG-CCCCCACACCTAGTGC-CAACGTTTACAGCGTGGT
>#S16-0000625
TCGACTTAACGCGTTAGCTCCGGAAGCCACGCCTCAAGG-GCACAACCTCCAAGTCGACATCGT-TACGGCGTGGAT
>#S16-0000673
TCGACTTAACGCGTTAGCTCCGGAAGCCACGCCTCAAGGGCACAACCTCCAAGTCGACATCGTTTACGGCGTGGAT
```

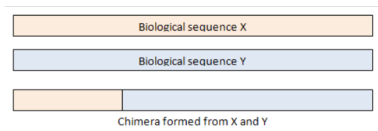
www.sixthstep.com

Filtering artefacts

- ▶ PCR errors
 - ▶ Most *Taq* polymerases introduce point mutations (error) at a rate of 1 every 1000 bases
 - ▶ **Solution:** use Hi-Fi polymerases with lower error rates (\$\$\$)
- ▶ Chimeras
 - ▶ Chimeras are sequences formed by two or more biological sequences joined together
 - ▶ **Solution:** reduce number of PCR cycles and increase annealing temperature

Filtering artefacts

- ▶ PCR errors
 - ▶ Most *Taq* polymerases introduce point mutations (error) at a rate of 1 every 1000 bases
 - ▶ **Solution:** use Hi-Fi polymerases with lower error rates (\$\$\$)
- ▶ Chimeras
 - ▶ Chimeras are sequences formed by two or more biological sequences joined together
 - ▶ **Solution:** reduce number of PCR cycles and increase annealing temperature



Filtering artefacts

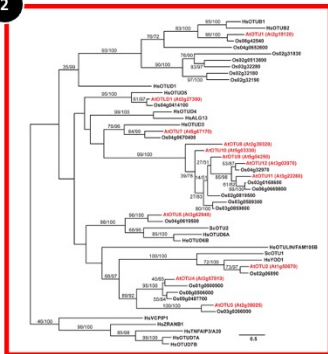
```
>*16S-0000011 | depth=44 | freq=2.42
TTCAGTCGCTCCCTAGCTTTCGCACTTCAGCGTCAGTTGCCGTCCAGTGAACATCTTCATCATCGGCATT
CCTGCACATATCTACGAATTTACTCTACTCGTGCAGTTCGGTCCACCTCTCCAGCACCTAGCCAACAG
>*16S-0000076 | depth=33 | freq=1.82
TTC AATGTTTGCTCCCCACGCTTTCGAGCCTCAGCGTCAGTTAC AAGCCAGAGAGCCGCTTTCGCCACCGGT
GTTCTCCATATATCTACGCATTTACCCGCTACACATGGAATTCCACTCTCCCTCTTGCACTCAAGTTAAA
>*16S-0000052 | depth=32 | freq=1.76
TTCACGATACCCGCACCTTCGAGCTTAAGCGTCAGTGGCGCTCCCGTCAGCTGCCTTCGCAATCGGAGTTCT
TCGTATATCTAAGCATTTACCCGCTACACGACGAATTCGGCAACGTTGTGCGTACTCAAGGAAACCAGTA
>*16S-0000141 | depth=15 | freq=0.83
TTC AACGTTTCGCTCCCTGGCTTTCGCGCCTCAGCGTCAGTTTTCGTCCAGAAAGTCGCCTTCGCCACTGGT
GTTCTTCCTAATATCTACGCATTTACCCGCTACACTAGGAATTCCACTTTCCTCTCCGATACTT
>#16S-0000058 | depth=12 | freq=0.66
TTCAGTCGCTCCCTAGCTTTCGCACTTCAGCGTCAGTTGCCGTCCAGTGAACATCTTCATCATCGGCATT
GTTCTCCATATATCTACGCATTTACCCGCTACACATGGAATTCCACTCTCCCTCTTGCACTCAAGTTAAA
>*16S-0000098 | depth=10 | freq=0.55
TTTAGTCCTGTTGCTCCCCACGCTTTCGCTCCTCAGCGTCAGTAACGGCCAGAGACCCGCCCTTCGCCACC
GGTGTCTTCTTGATATCTCGGCATTTACCCGCTACACAGGAGTTCCAGCTCCCT
>#16S-0000295 | depth=2 | freq=0.11
TTCACGATACCCAGCTTTCGAGCATCAGCGTCAGTTGCGCTACAGTAAGCTGCCTTCGCAATCGGAGTTCT
TCGTGATATCTAAGCATTTACCCGCTACACAGGAATTCGGCTAGTTTCGGCGCACTCAAGCCCCCAGTT
>#16S-0000021 | depth=1 | freq=0.06
TTC AACGTTTCGCTCCCTGGCTTTCGCGCCTCAGCGTCAGTTTTCGTCCAGAAAGTCGCCTTCGCCACTGGT
```

Chimera

Contaminations








Phylogenetic tree!

2



OTU table!

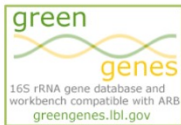
3

			
	10	8	10
	4	2	14
	2	10	9
	1	5	8

Align to reference

```
>*16S-0000002 | depth=42 | freq=2.31
TTCAACCTTGC GGTCG TACTCCCCAGGCGGAGTGC TTAATGCGT TAGCTGCGGCACTAAACCCCGGAAAGGGTCTAACACCTAGCACTCATCGTT
TACGGCGTGGACTACCAGGGTATCTAATCCTGTTGCTFCCCACGCTTTCGAGCCTCAGCGTCAGTTACAAGCCAGAGAGCCGCTTTCGCCACCG
GTGTTCTCCATATATCTAGCATTTCACCGCTACACATGGAATCCACTCTCCCCTCTTGCAC TCAAGTTAAACAGTTTCCAAAGCGTACTATG
GTTAAGCCACAGCCTTTAACTTCAGACTTATCT
>*16S-0000019 | depth=12 | freq=0.66
TTCAGCCTTGC GGCCGTACTCCCCAGGCGGATTACTTATCGCATTGCTTCGGCACAGACAGTCTTCTGCCACACCCAGTAATCATCGTTTAC
GGCCGGACTACCAGGGTATCTAATCCTGTTGCTFCCCAGGCTTTCGCACTTCAGCGTCAGTTACCGTCCAGTGAACATCTTCATCATCGGCA
TTCTCTGCATATCTACGAATTCACCTCTACTCGTGCAGTTCGTCACCTCTCCGGTACTCCAGCCTATCAGTTTCAAAGGCAGGCTCGCGT
TGAGCCGCAGGTTTTACCCCTGACTTGAAAGG
```

VS.



Assign taxonomy

AY053482.1;tax=k:Bacteria,p:Firmicutes,c:Bacilli,o:Lactobacillales,f:Streptococcaceae,
g:Streptococcus,s:pseudopneumoniae

Sequence ID: lc|Query_210570 Length: 1429 Number of Matches: 1

Range 1: 565 to 882 [Graphics](#)

Score	Expect	Identities	Gaps	Strand
588 bits(318)	7e-172	318/318(100%)	0/318(0%)	Plus/Minus

Taxonomy!

4

ASV148428	Bacteria(100);Proteobacteria(100);Gammaproteobacte...
ASV212114	Bacteria(100);Cyanobacteria(100);Cyanobacteria(100);...
ASV9620	Bacteria(100);Proteobacteria(100);Alphaproteobacteri...
ASV147186	Bacteria(100);Proteobacteria(100);Betaproteobacteria...
ASV89359	Bacteria(100);Proteobacteria(100);Alphaproteobacteri...
ASV1061	Bacteria(100);Proteobacteria(100);Gammaproteobacte...
ASV328581	Bacteria(100);Bacteroidetes(100);Bacteroidia(100);Bac...
ASV86104	Bacteria(100);Proteobacteria(100);Alphaproteobacteri...

Data analysis

What we will use?



What do we need to start?

```
phyloseq-class experiment-level object
otu_table() OTU Table: [ 43879 taxa and 289 samples ]
sample_data() Sample Data: [ 289 samples by 9 sample variables ]
tax_table() Taxonomy Table: [ 43879 taxa by 7 taxonomic ranks ]
phy_tree() Phylogenetic Tree: [ 43879 tips and 43878 internal nodes ]
```

otu_table()

```
      16S.SOI.41 16S.SOI.48 16S.SOI.18 16S.SOI.46 16S.SOI.59 16S.SOI.34 16S.R00.57 16S.R00.43 16S.R00.58
denovo7709      1      1      0      0      0      0      0      0      0
denovo7708      0      0      1      1      1      0      0      0      0
denovo22216     0      0      0      0      0      1      0      0      0
```

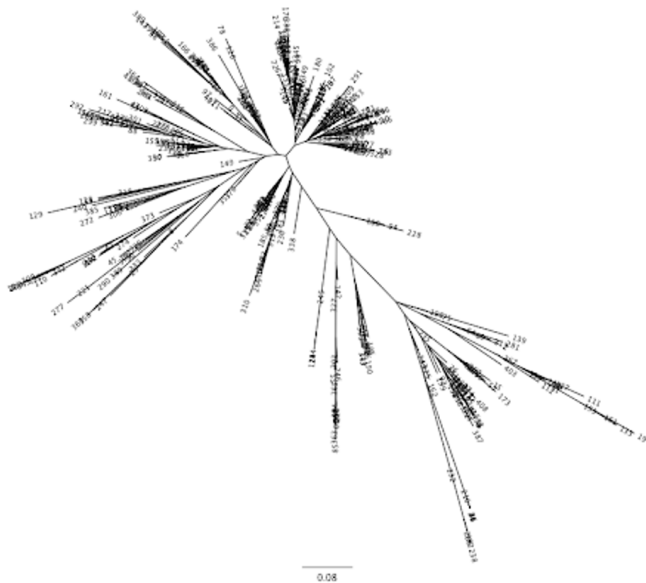
sample_data()

	A	B	C	D	E	F	G	H
1	SampleID	Community	Category	Sample_type	Genotype	Soil	Aphid	H_defensa
2	16S.APH.1	Bacterial	Experimenta	Aphid	TBR	WHS	Present	Absent
3	16S.APH.2	Bacterial	Experimenta	Aphid	TBR	WHS	Present	Absent
4	16S.APH.3	Bacterial	Experimenta	Aphid	TBR	WHS	Present	Absent
5	16S.APH.4	Bacterial	Experimenta	Aphid	TBR	WHS	Present	Absent
6	16S.APH.5	Bacterial	Experimenta	Aphid	TBR	WHS	Present	Absent
7	16S.APH.6	Bacterial	Experimenta	Aphid	TBR	WHS	Present	Present
8	16S.APH.7	Bacterial	Experimenta	Aphid	TBR	WHS	Present	Present
9	16S.APH.8	Bacterial	Experimenta	Aphid	TBR	WHS	Present	Present
10	16S.APH.9	Bacterial	Experimenta	Aphid	TBR	WHS	Present	Present
11	16S.APH.10	Bacterial	Experimenta	Aphid	TBR	WHS	Present	Present
12	16S.APH.11	Bacterial	Experimenta	Aphid	TBR	MICROB	Present	Absent
13	16S.APH.12	Bacterial	Experimenta	Aphid	TBR	MICROB	Present	Absent
14	16S.APH.13	Bacterial	Experimenta	Aphid	TBR	MICROB	Present	Absent
15	16S.APH.14	Bacterial	Experimenta	Aphid	TBR	MICROB	Present	Absent
16	16S.APH.15	Bacterial	Experimenta	Aphid	TBR	MICROB	Present	Absent
17	16S.APH.16	Bacterial	Experimenta	Aphid	TBR	MICROB	Present	Present
18	16S.APH.17	Bacterial	Experimenta	Aphid	TBR	MICROB	Present	Present
19	16S.APH.18	Bacterial	Experimenta	Aphid	TBR	MICROB	Present	Present

tax_table()

```
Rank1 Rank2
denovo7709 "D_0__Bacteria" "D_1__Proteobacteria"
denovo7708 "D_0__Bacteria" "D_1__Proteobacteria"
denovo22216 "D_0__Bacteria" "D_1__Bacteroidetes"
denovo11322 "D_0__Bacteria" "D_1__Bacteroidetes"
denovo44859 "D_0__Bacteria" "D_1__Chloroflexi"
```

phy_tree()



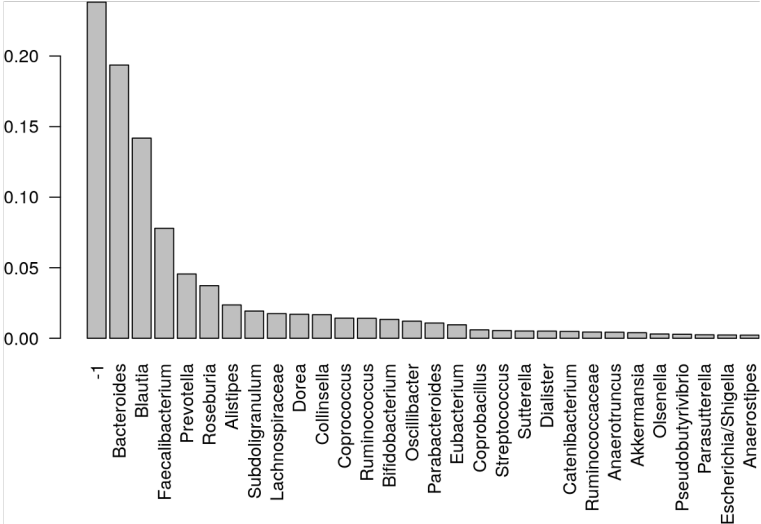
Terminology

- ▶ Quantitative versus qualitative metrics
 - ▶ qualitative metrics only account for whether an organism is present or absent
 - ▶ quantitative metrics account for abundance
- ▶ Phylogenetic versus non-phylogenetic metrics
 - ▶ non-phylogenetic metrics treat all OTUs as being equally related
 - ▶ phylogenetic metrics incorporate evolutionary relationships between the OTUs

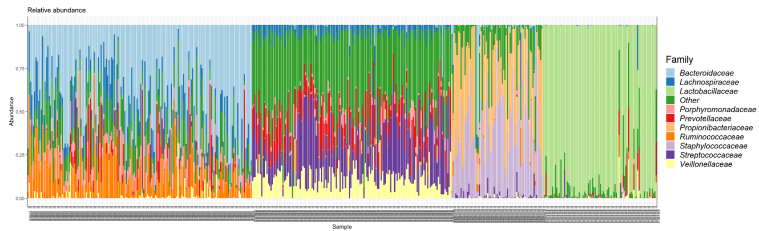
Terminology

- ▶ Quantitative versus qualitative metrics
 - ▶ qualitative metrics only account for whether an organism is present or absent
 - ▶ quantitative metrics account for abundance
- ▶ Phylogenetic versus non-phylogenetic metrics
 - ▶ non-phylogenetic metrics treat all OTUs as being equally related
 - ▶ phylogenetic metrics incorporate evolutionary relationships between the OTUs

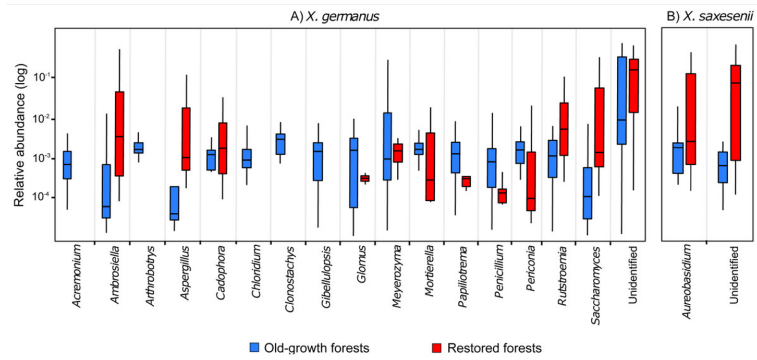
Exploratory plots



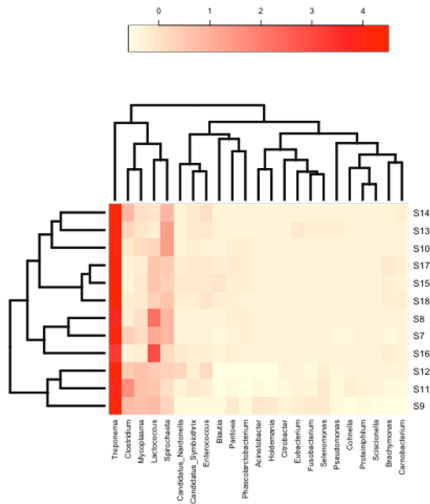
Who is in there?



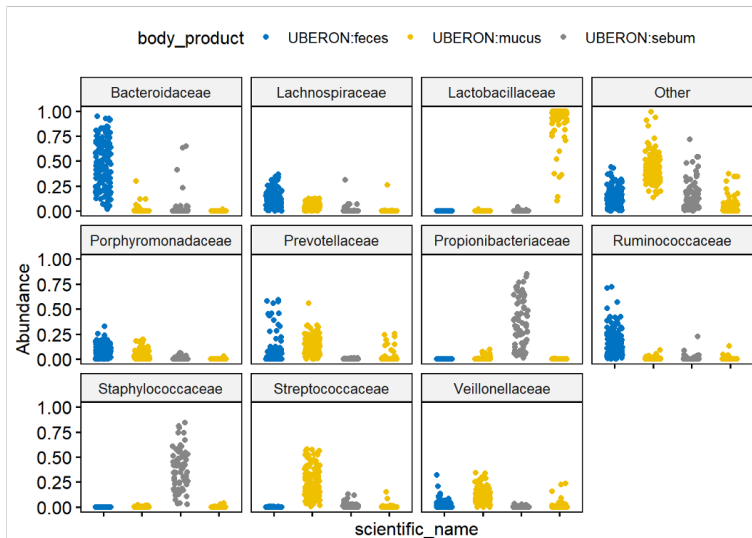
Who is in there?



Who is in there?



Who is in there?



Does community structure vary between treatments?

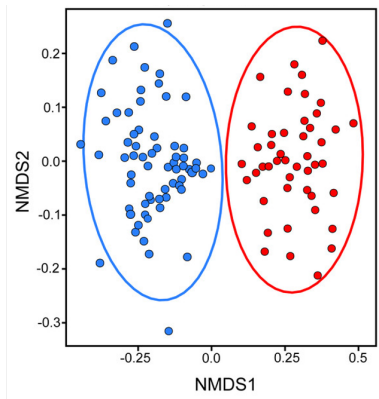
► Multivariate statistics

- Dissimilarity matrices (Bray-Curtis, Jaccard, (w)unifrac)
- Ordination plots (PCA, PCoA, NMDS, ...)
- Statistical tests (PERMANOVA, ANOSIM, ...)

0								
4.15	0							
11.02	15.01	0						
7.16	3.03	18.02	0					
43.72	47.49	32.80	50.41	0				
54.37	58.23	43.36	61.19	11.12	0			
46.34	50.20	35.34	53.16	3.78	8.03	0		
55.42	59.27	44.42	62.23	12.05	1.12	9.08	0	

Does community structure vary between treatments?

- ▶ Multivariate statistics
 - ▶ Dissimilarity matrices (Bray-Curtis, Jaccard, (w)unifrac)
 - ▶ Ordination plots (PCA, PCoA, NMDS, ...)
 - ▶ Statistical tests (PERMANOVA, ANOSIM, ...)



Does community structure vary between treatments?

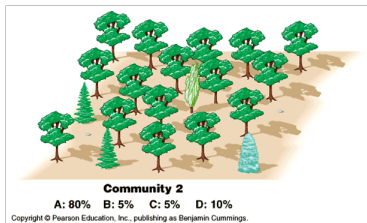
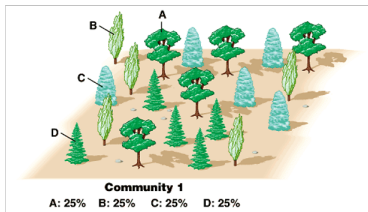
► Multivariate statistics

- Dissimilarity matrices (Bray-Curtis, Jaccard, (w)unifrac)
- Ordination plots (PCA, PCoa, NMDS, ...)
- Statistical tests (PERMANOVA, ANOSIM, ...)

```
##  
## Call:  
## adonis(formula = erie_bray ~ Station, data = sampledf)  
##  
## Permutation: free  
## Number of permutations: 999  
##  
## Terms added sequentially (first to last)  
##  
##          Df SumsOfSqs MeanSqs F.Model    R2 Pr(>F)  
## Station    2   0.6754  0.33772   2.7916 0.09531 0.003 **  
## Residuals 53   6.4118  0.12098             0.90469  
## Total     55   7.0872                1.00000  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '
```

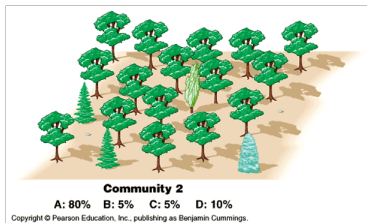
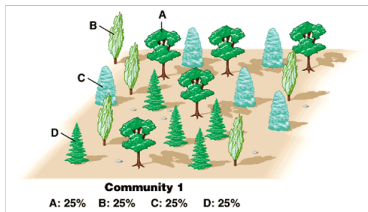
Does the community diversity vary between treatments?

- ▶ **Richness** refers to how many different types of organisms are present in a sample
- ▶ **Evenness** tells us how even or uneven the distribution of species abundances are in a given environment
- ▶ **Diversity** is a measurement of species richness combined with evenness



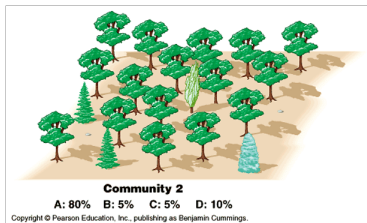
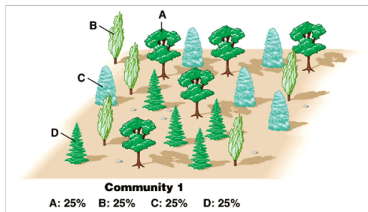
Does the community diversity vary between treatments?

- ▶ **Richness** refers to how many different types of organisms are present in a sample
- ▶ **Evenness** tells us how even or uneven the distribution of species abundances are in a given environment
- ▶ **Diversity** is a measurement of species richness combined with evenness



Does the community diversity vary between treatments?

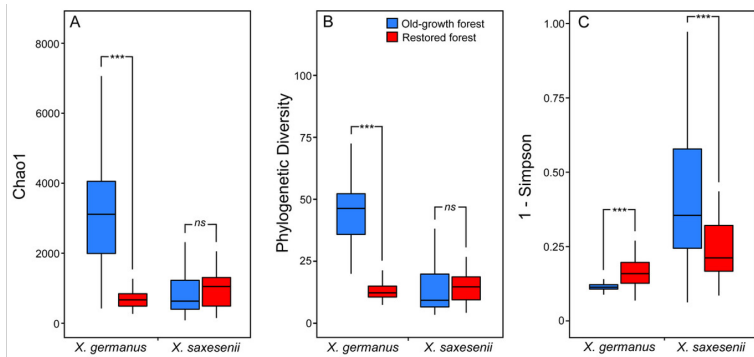
- ▶ **Richness** refers to how many different types of organisms are present in a sample
- ▶ **Evenness** tells us how even or uneven the distribution of species abundances are in a given environment
- ▶ **Diversity** is a measurement of species richness combined with evenness



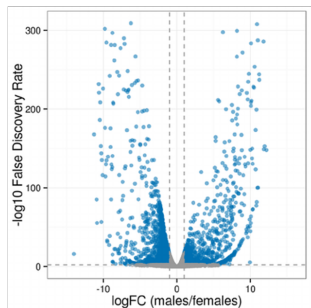
Lots of metrics!

- ▶ **Observed OTUs**: we simply count the OTUs that are observed in a given sample
- ▶ **Shannon index** assumes all species are represented in a sample and that they are randomly sampled
- ▶ **Simpson** index is a dominance index because it gives more weight to common or dominant species.
- ▶ **Chao1** index gives more weight to the low abundance species, only the singletons and doubletons are used to estimate the number of missing species
- ▶ **PD** is computed simply as the sum of the branch length in a phylogenetic tree that is "covered" or represented in a given sample.

Does the community diversity vary between treatments?



Who differs between treatments?



baseMean	log2FoldChange	lfcSE	stat	pvalue	padj	taxon
29.20535	1.91205	0.13432	14.23457	0.00000	0.00000	Clostridium difficile et rel.
51.65152	3.04116	0.28687	10.60107	0.00000	0.00000	Mitsuokella multiacida et rel.
12.39749	1.83825	0.18531	9.91994	0.00000	0.00000	Klebsiella pneumoniae et rel.
44.16494	1.78333	0.23072	7.72937	0.00000	0.00000	Megasphaera elsdenii et rel.
66.93783	1.68345	0.25330	6.64609	0.00000	0.00000	Escherichia coli et rel.
3.63459	1.53142	0.23140	6.61792	0.00000	0.00000	Weissella et rel.
5.74035	3.07334	0.47848	6.42308	0.00000	0.00000	Serratia
0.42171	1.70079	0.47147	3.60743	0.00031	0.00075	Moraxellaceae

More tools

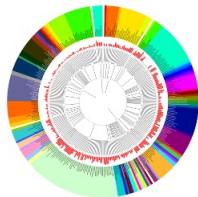
Metagenomics



Sample



magiSEQ™

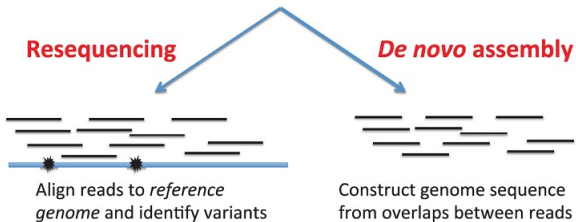


**Functional and taxonomic
assembly**

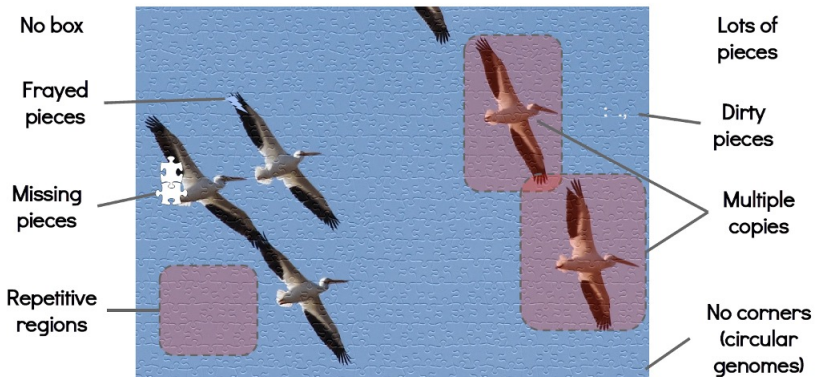
Metagenomics



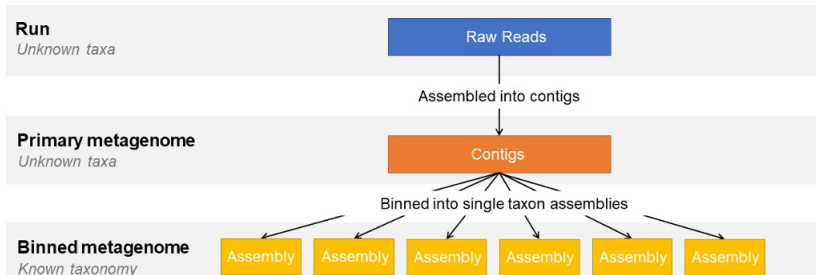
Let's take a step back: genomics










Let's take a step back: genomics






















Metagenomics



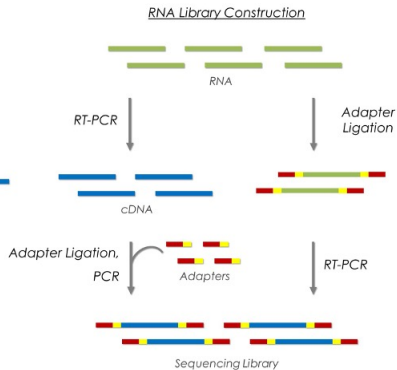
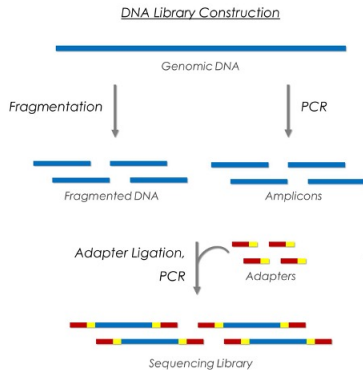
Metagenomics

			
	10	8	10
	4	2	14
	2	10	9
	1	5	88

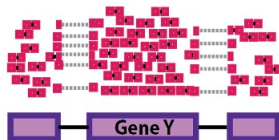
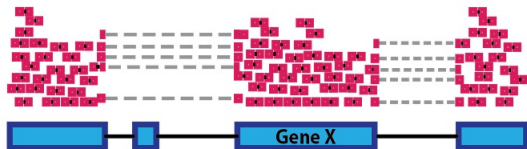
Metagenomics








Metatranscriptomics



Transcriptomics



Metatranscriptomics

			
	10	8	10
	4	2	14
	2	10	9
	1	5	88

Questions?