



# ggplot2 for microbiome analysis

Presenter: Huyha

# Table of contents

1. Why we need to visualize data? Anscombe quartet
2. Addition function
3. Microbiome composition + OTU and ASV term + microbiome data
4. Microbiome composition barplot practice

# Why we need to visualise data?

---

	mathscore	englishscore	sciencescore	state_of_living_in_usa	age	grade
1	64.4	97.8	94.5	California	18	D
2	89.4	72.7	84.6	New York	16	F
3	70.4	83.9	82.0	Florida	14	B
4	94.2	78.6	99.7	California	16	A
5	97.0	55.1	82.8	New York	14	A
6	52.3	95.0	85.4	Illinois	15	C
7	76.4	62.3	77.2	Florida	18	A
8	94.6	52.1	79.7	Texas	11	F
9	77.6	66.4	64.5	Illinois	14	A
10	72.8	97.7	57.4	California	17	B

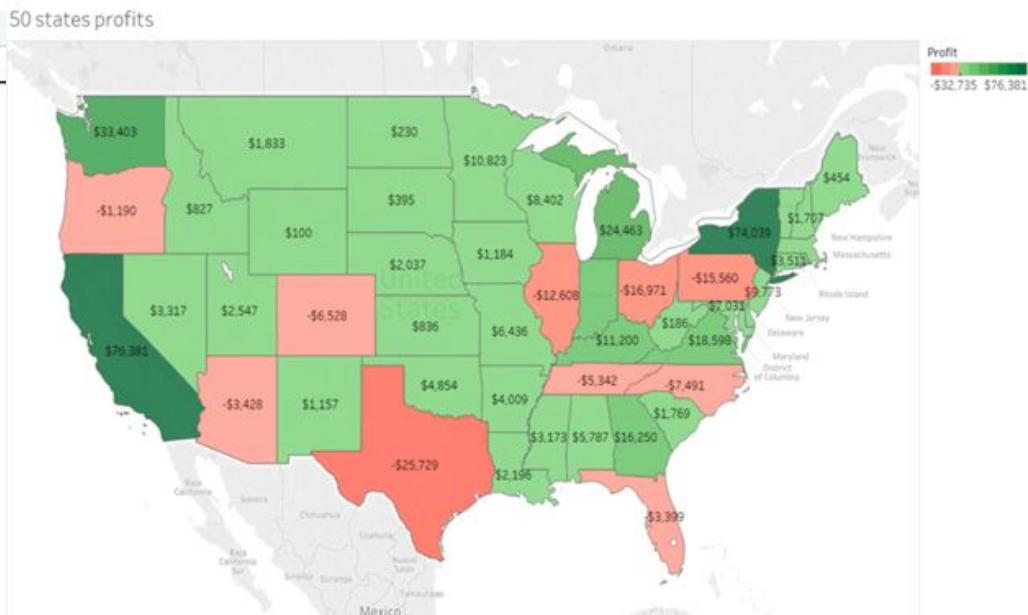
> |

# Why we need to visualise data?

R 4.3.3 · ~/

	mathscore	englishscore	sciencescore	state
1	64.4	97.8	94.5	
2	89.4	72.7	84.6	
3	70.4	83.9	82.0	
4	94.2	78.6	99.7	
5	97.0	55.1	82.8	
6	52.3	95.0	85.4	
7	76.4	62.3	77.2	
8	94.6	52.1	79.7	
9	77.6	66.4	64.5	
10	72.8	97.7	57.4	

> |



Map based on Longitude (generated) and Latitude (generated). Color shows sum of Profit. The marks are labeled by sum of Profit. Details are shown for Country and State. The view is filtered on Country, which keeps United States.

# Data [edit]

For all four datasets:

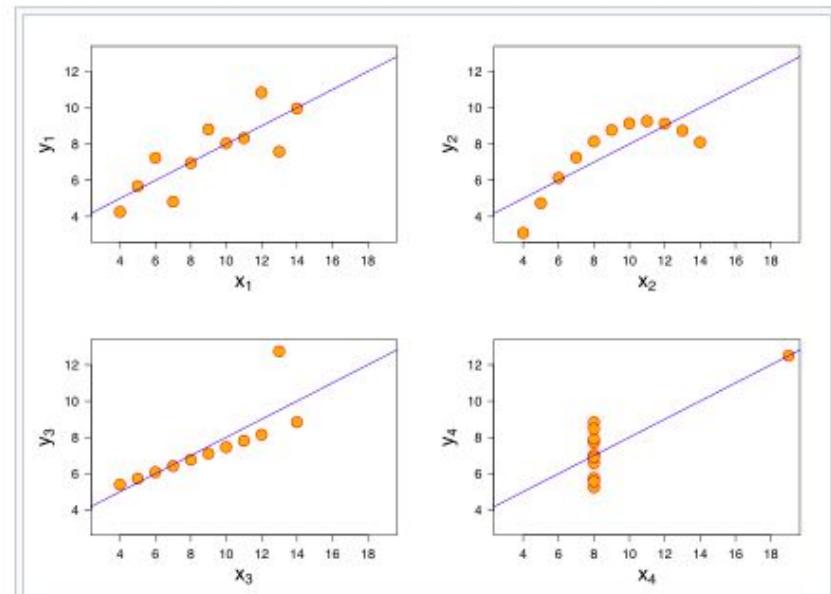
Property	Value				Accuracy			
Mean of x	9				exact			
Sample variance of x: $s_x^2$	x1	x2	x3	x4	y1	y2	y3	y4
Mean of y	1	10	10	10	8	8.04	9.14	7.46
Sample variance of y: $s_y^2$	2	8	8	8	8	6.95	8.14	6.77
Correlation between x and y	3	13	13	13	8	7.58	8.74	12.74
Linear regression line	4	9	9	9	8	8.81	8.77	7.11
Coefficient of determination of th	5	11	11	11	8	8.33	9.26	7.81
	6	14	14	14	8	9.96	8.10	8.84
	7	6	6	6	8	7.24	6.13	6.08
	8	4	4	4	19	4.26	3.10	5.39
	9	12	12	12	8	10.84	9.13	8.15
	10	7	7	7	8	4.82	7.26	6.42
	11	5	5	5	8	5.68	4.74	5.73
								6.89

<b>set</b> <code>&lt;chr&gt;</code>	<b>mean_x</b> <code>&lt;dbl&gt;</code>	<b>mean_y</b> <code>&lt;dbl&gt;</code>	<b>sd_x</b> <code>&lt;dbl&gt;</code>	<b>sd_y</b> <code>&lt;dbl&gt;</code>	<b>cor</b> <code>&lt;dbl&gt;</code>
1	9	7.500909	3.316625	2.031568	0.8164205
2	9	7.500909	3.316625	2.031657	0.8162365
3	9	7.500000	3.316625	2.030424	0.8162867
4	9	7.500909	3.316625	2.030579	0.8165214

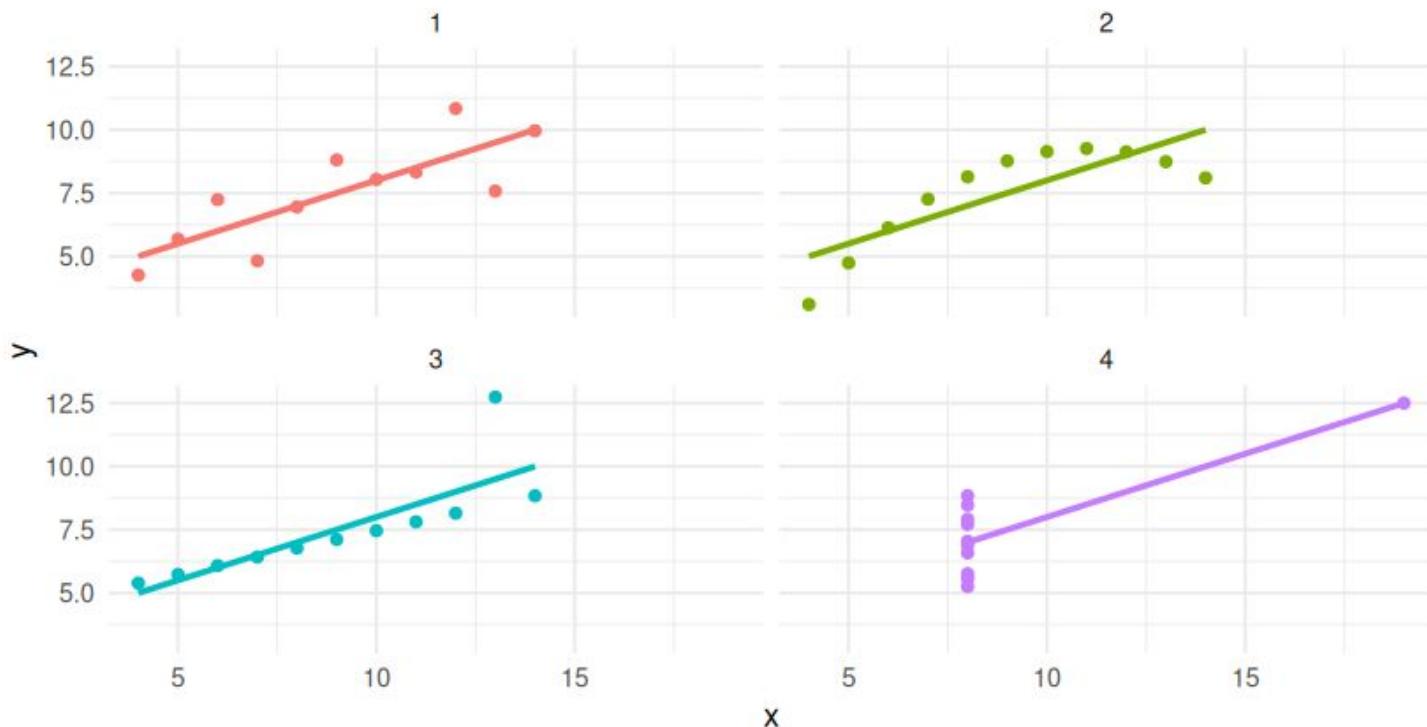
4 rows

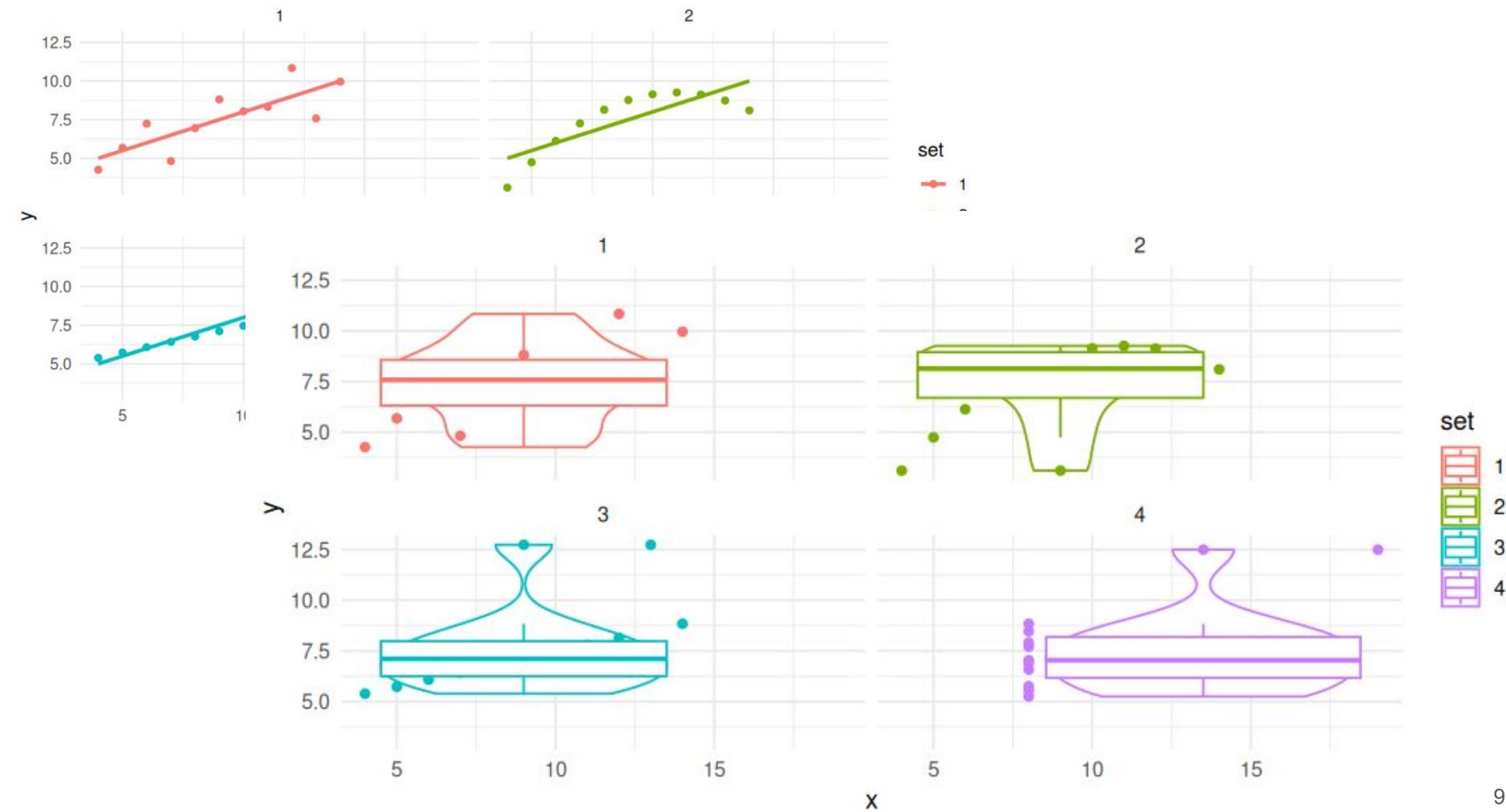
From Wikipedia, the free encyclopedia

**Anscombe's quartet** comprises four datasets that have nearly identical simple descriptive statistics, yet have very different distributions and appear very different when graphed. Each dataset consists of eleven  $(x, y)$  points. They were constructed in 1973 by the statistician Francis Anscombe to demonstrate both the importance of graphing data when analyzing it, and the effect of outliers and other influential observations on statistical properties. He described the article as being intended to counter the impression among statisticians that "numerical calculations are exact, but graphs are rough".<sup>[1]</sup>



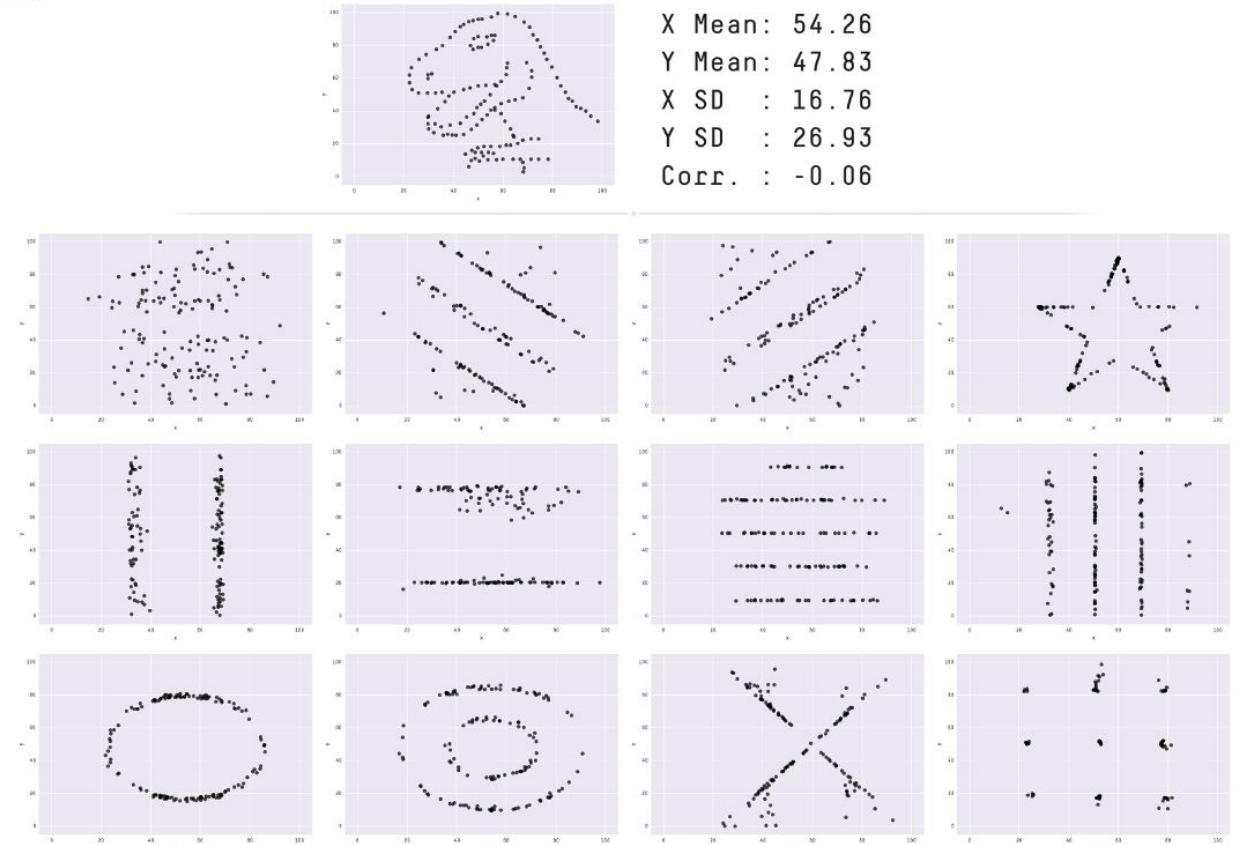
The four datasets composing Anscombe's quartet. All four sets have identical statistical parameters, but the graphs show them to be considerably different





# DATASAURUS

Recently, [Albert Cairo](#) created the [Datasaurus](#) dataset which urges people to “never trust summary statistics alone; always visualize your data.” While the data exhibits normal-seeming statistics, plotting the data reveals a picture of a dinosaur. Inspired by Anscombe’s Quartet and Datasaurus, we present, The Datasaurus Dozon (download.csv).



<https://www.research.autodesk.com/publications/same-stats-different-graphs/>

## 2. How to draw in ggplot

1. Import-data
2. Data manipulation
3. Basic plot
4. Extend more detail
5. Save file

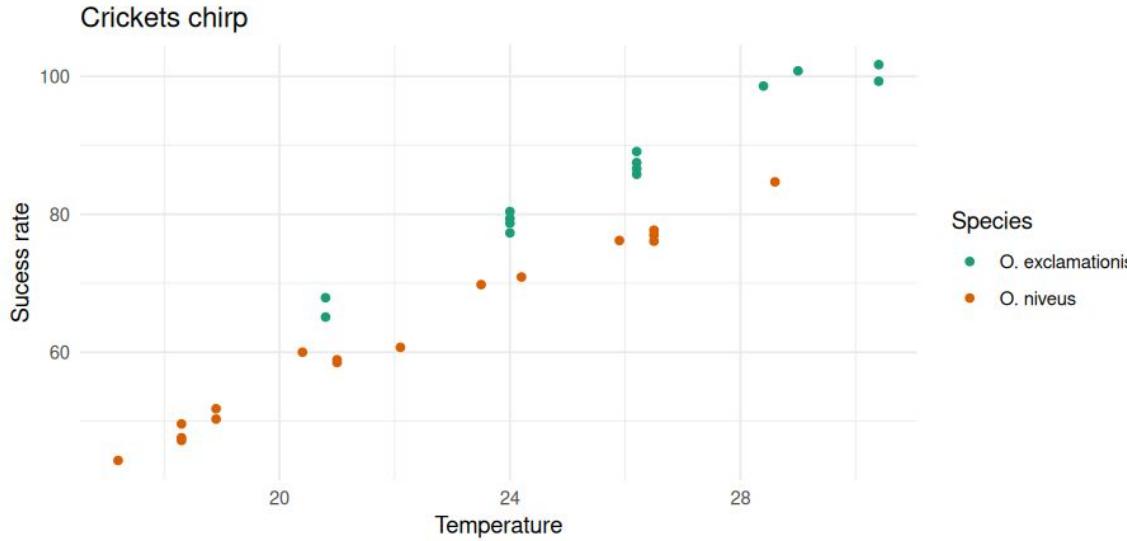
# Crickets chirp relate to temp

	species	temp	rate
	<fct>	<dbl>	<dbl>
1	O. exclamationis	20.8	67.9
2	O. exclamationis	20.8	65.1
3	O. exclamationis	24	77.3
4	O. exclamationis	24	78.7
5	O. exclamationis	24	79.4
6	O. exclamationis	24	80.4
7	O. exclamationis	26.2	85.8
8	O. exclamationis	26.2	86.6
9	O. exclamationis	26.2	87.5
10	O. exclamationis	26.2	89.1

# i 21 more rows



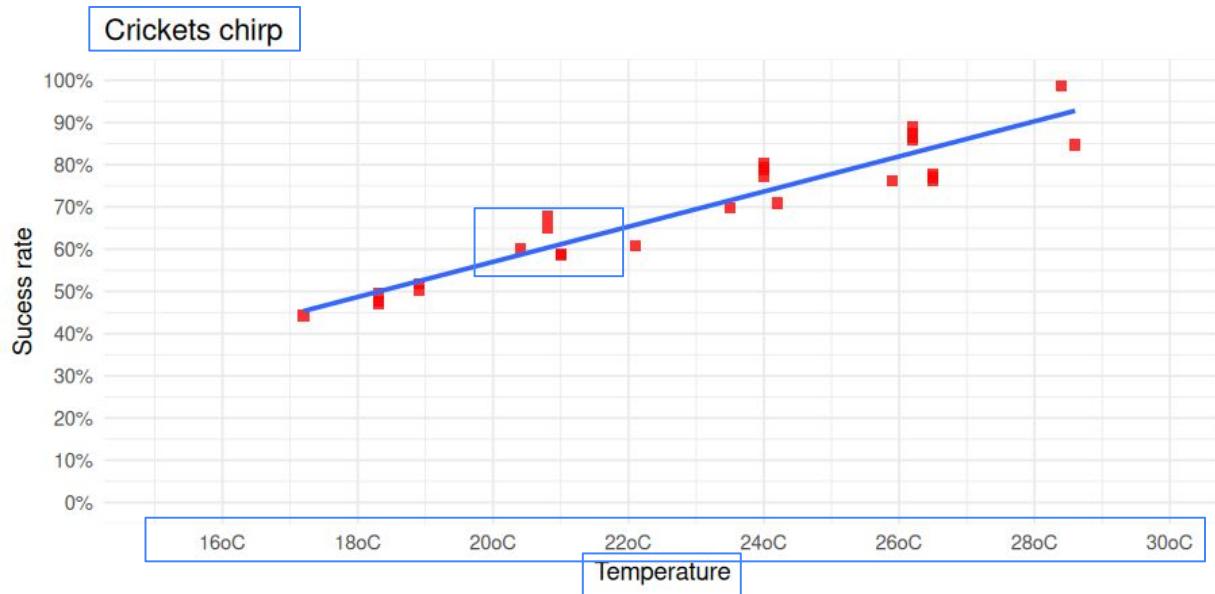
## 2. Basic plot



```
> ggplot(crickets, aes(x=temp,y=rate,
+ color=species)) +
+   geom_point() +
+   labs(x="Temperature", y="Success rate",
+        color="Species",
+        title="Crickets chirp") +
+   scale_color_brewer(palette = "Dark2")
```

## 2. Extend more detail

```
ggplot(crickets, aes(x=temp,y=rate,))+  
  geom_point(color="red",  
            size=2,  
            shape="square",  
            alpha=0.8) +  
  labs(x="Temperature", y="Success rate",  
       color="Species",  
       title="Crickets chirp") +  
  scale_x_continuous(  
    limits = c(15,30),  
    breaks = seq(10,30, by =2),  
    labels = function(x) paste0(x,"oC")  
  )+  
  scale_y_continuous(  
    limits = c(0,100),  
    breaks = seq(0,100, by=10),  
    labels = function(x) paste0(x,"%")  
  )+  
  geom_smooth(method="lm",se=FALSE)
```



# How to learn ggplot2 by yourself?

1. Cheatsheet
2. Help function in Rstudio
3. Use AI to help
4. Book - (recommended)

[Introduction to ggplot\(recommend\)](#)



# 2. Helps in Rstudio

?<command>

Example: ?labs

The screenshot shows the RStudio interface with the 'Help' menu open, displaying the documentation for the 'labs' function from the 'ggplot2' package. The top navigation bar includes tabs for Environment, History, Connections, Tutorial, Files, Plots, Packages, Help, Viewer, and Presentation. Below the tabs are icons for search, refresh, and file operations. The main content area shows the title 'R: Modify axis, legend, and plot labels' and a search bar labeled 'Find in Topic'. The page title is 'labs (ggplot2)' and the right sidebar indicates 'R Documentation'. The 'Description' section explains the importance of good labels for accessibility and provides guidance on using title, subtitle, caption, and tag arguments. The 'Usage' section shows the syntax for using labs(), including xlab(), ylab(), and ggtitle(). The 'Arguments' section details the '...' argument as a list of name-value pairs and the 'title' argument as text for the title.

R: Modify axis, legend, and plot labels

Find in Topic

labs (ggplot2)

R Documentation

## Description

Good labels are critical for making your plots accessible to a wider audience. Always ensure the axis and legend labels display the full variable name. Use the `plot.title` and `subtitle` to explain the main findings. It's common to use the `caption` to provide information about the data source. `tag` can be used for adding identification tags to differentiate between multiple plots.

## Usage

```
labs(  
  ...,  
  title = waiver(),  
  subtitle = waiver(),  
  caption = waiver(),  
  tag = waiver(),  
  alt = waiver(),  
  alt_insight = waiver()  
)  
  
xlab(label)  
  
ylab(label)  
  
ggtitle(label, subtitle = waiver())
```

## Arguments

...	A list of new name-value pairs. The name should be an aesthetic.
title	The text for the title.

# 3. USE AI to learn

ChatGPT ▾

results.

```
 Ø {r}
#Plot_2
dat %>%
  ggplot(aes(x = sample_id, y = count)) +
  facet_grid(~ fraction + soil, scales = "free_x", s
  geom_bar(aes(fill = Phylum), stat = "identity", po
  scale_fill_brewer(palette = "Paired")
i want to draw black line to distingquis bar area to
```

To draw a black line around each bar to distinguish bar areas from each other in a `ggplot2` bar plot, you can use the `color` argument inside the `geom\_bar` function. This will add an outline to each bar.

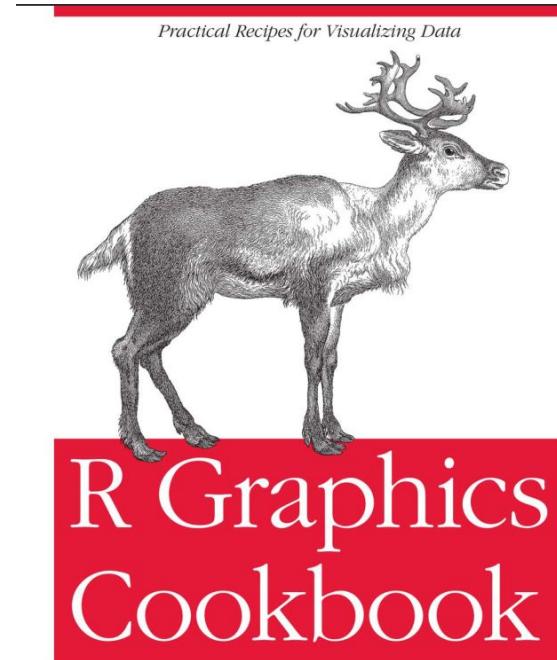
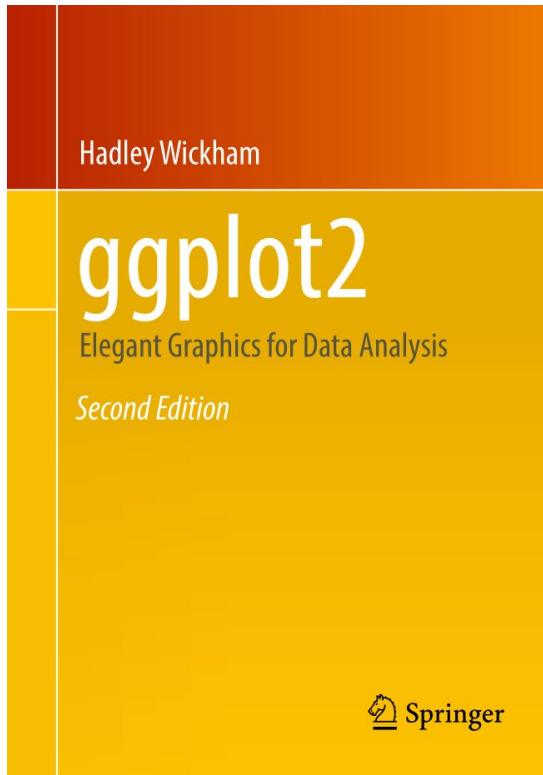
Here's your modified code with the black line added around each bar:

```
 r
Copy code
# Assuming you have the necessary data frame `dat` and the ggplot2 package loaded
library(ggplot2)
library(dplyr)

# Create the plot
dat %>%
  ggplot(aes(x = sample_id, y = count)) +
  facet_grid(~ fraction + soil, scales = "free_x", space = "free_x") +
  geom_bar(aes(fill = Phylum), stat = "identity", position = "fill", width = 1, col
  scale_fill_brewer(palette = "Paired",
```

Message ChatGPT

## 4. Book to learn



O'REILLY®

Winston Chang



Practice is the key to  
success

## II. Some addition function

Reshape from Long  
to Wide and back

country	year	metric
x	1960	10
x	1970	13
x	2010	15
y	1960	20
y	1970	23
y	2010	25
z	1960	30
z	1970	33
z	2010	35

```
pivot_wider(names_from = "year",  
            names_prefix = "yr",  
            values_from = "metric")
```

country	yr1960	yr1970	yr2010
x	10	13	15
y	20	23	25
z	30	33	35

```
pivot_longer(cols = yr1960:yr2010,  
            names_to = "year",  
            names_prefix = "yr"  
            values_to = "metric")
```

	subject	measurement1	measurement2	measurement3
1	A	10	15	20
2	B	20	25	30
3	C	30	35	40
>				

	subject	measurement	value
1	A	measurement1	10
2	B	measurement1	20
3	C	measurement1	30
4	A	measurement2	15
5	B	measurement2	25
6	C	measurement2	35
7	A	measurement3	20
8	B	measurement3	30
9	C	measurement3	40

```
# Reshape from wide to long format form tidyR
df_long1 <- df %>%
  pivot_longer(
    cols = starts_with("measurement"),
    names_to = "measurement",
    values_to = "value"
  )
```

```
# Using melt from package reshape2
df_long2 <- melt(df, id.vars = "subject",
  variable.name =
  "measurement",
  value.name = "value")
print(df_long2)
```

### III. OTU vs AVS

Source: <https://www.youtube.com/watch?v=azl9taCIDhQ>

The screenshot shows a video player interface with a presentation slide. The slide has a yellow and green curved background. At the top left, it says "Microbiome Informatics: OTUs vs. ASVs | Zymo Research". At the top right, there is a "To exit full screen, press Esc" button and a Zymo Research logo with the tagline "The Beauty of Science is to Make Things Simple". The main title of the slide is "OTUs vs. ASVs" in large bold black font. Below the title, the subtitle "Understanding the key differences" is displayed. The bottom of the slide shows a progress bar indicating the video is at 0:13 of 10:57, and the section is "Intro >". The bottom right corner of the slide contains the name "Michael Weinstein" and his email "mweinstein@zymoresearch.com".

Microbiome Informatics: OTUs vs. ASVs | Zymo Research

To exit full screen, press Esc

ZYMO RESEARCH  
The Beauty of Science is to Make Things Simple

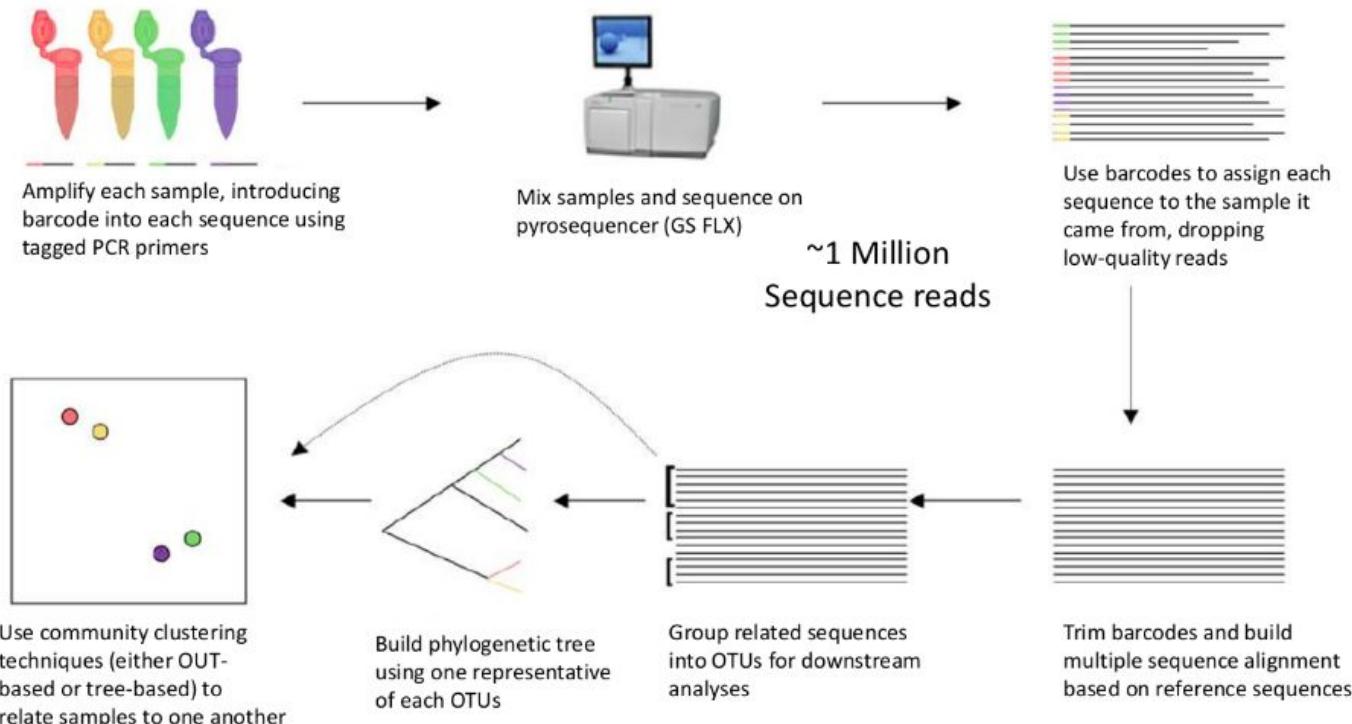
# OTUs vs. ASVs

Understanding the key differences

0:13 / 10:57 • Intro >

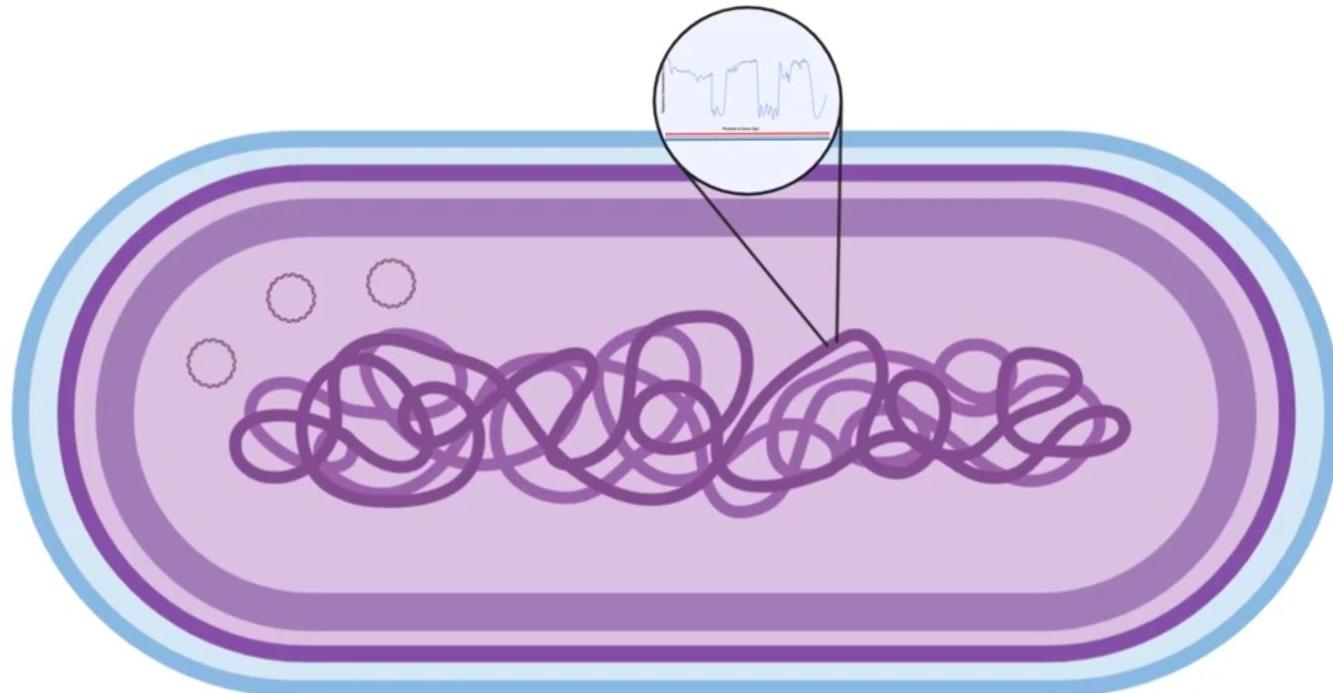
Michael Weinstein  
mweinstein@zymoresearch.com

# Overview of workflow with amplicon 16S metagenomics



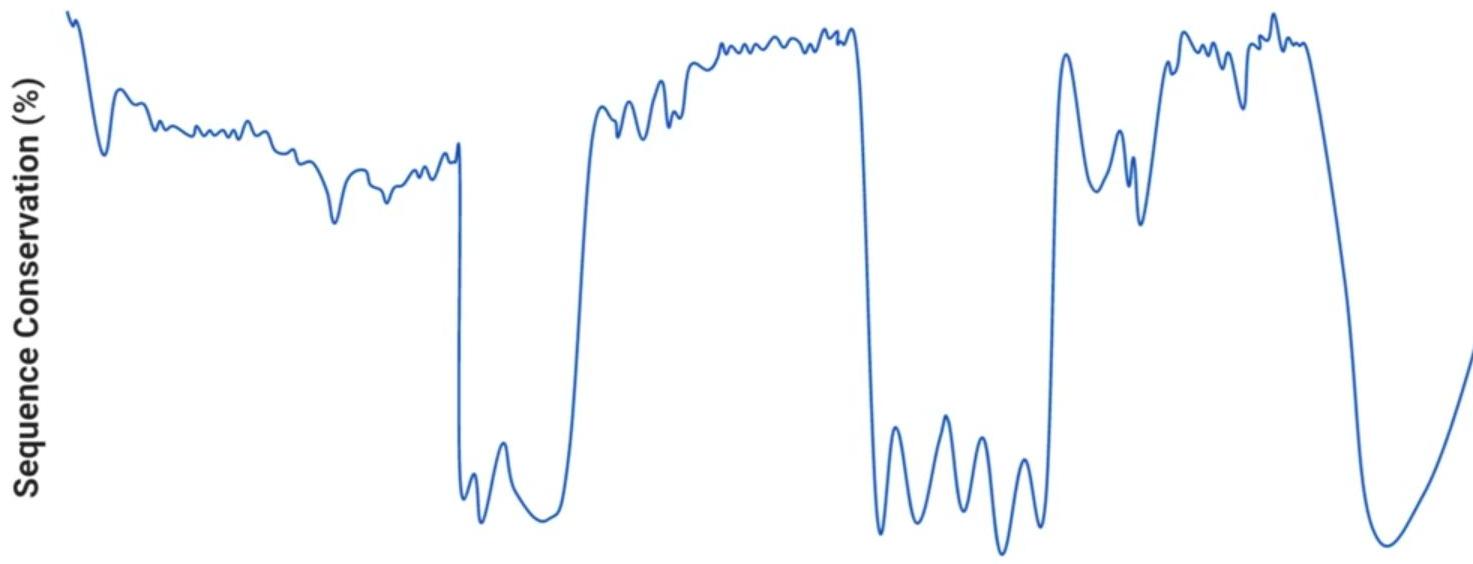
Hamady & Knight 2009 Genome Res. 19: 1141

# Targeted Sequencing



ZYMO RESEARCH

# A Potential Target for Sequencing

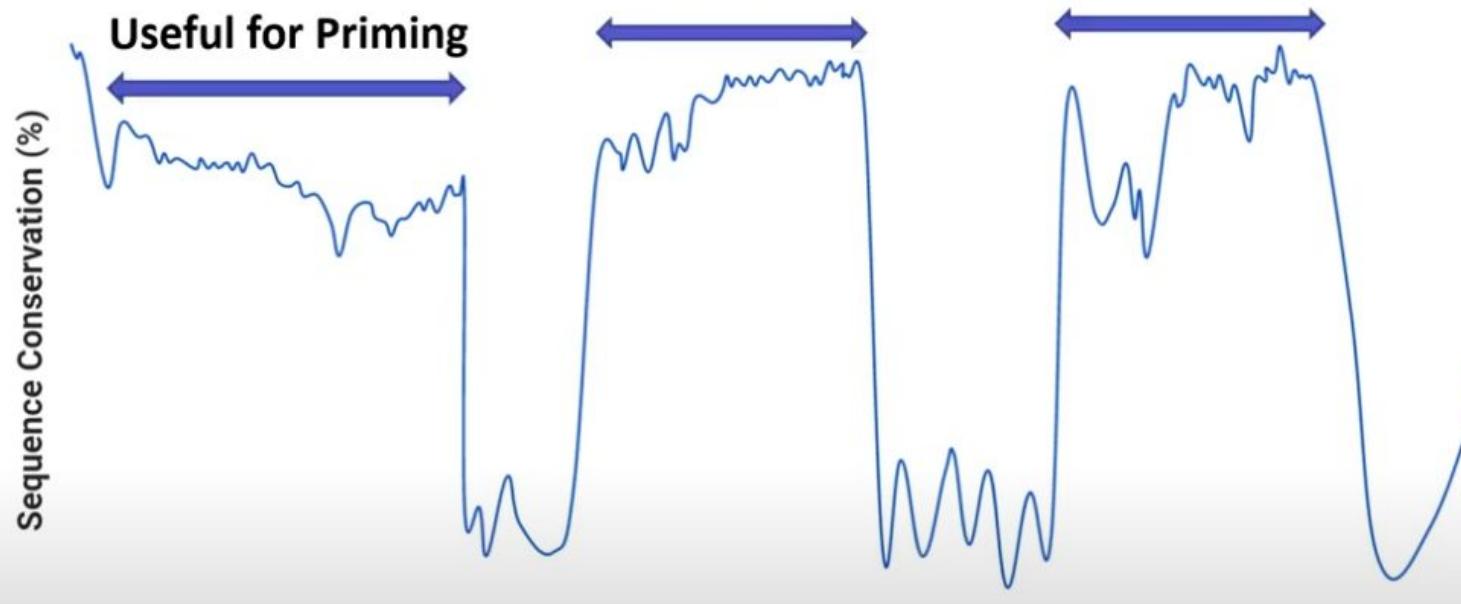


Position in Gene (bp)



ZYMO RESEARCH

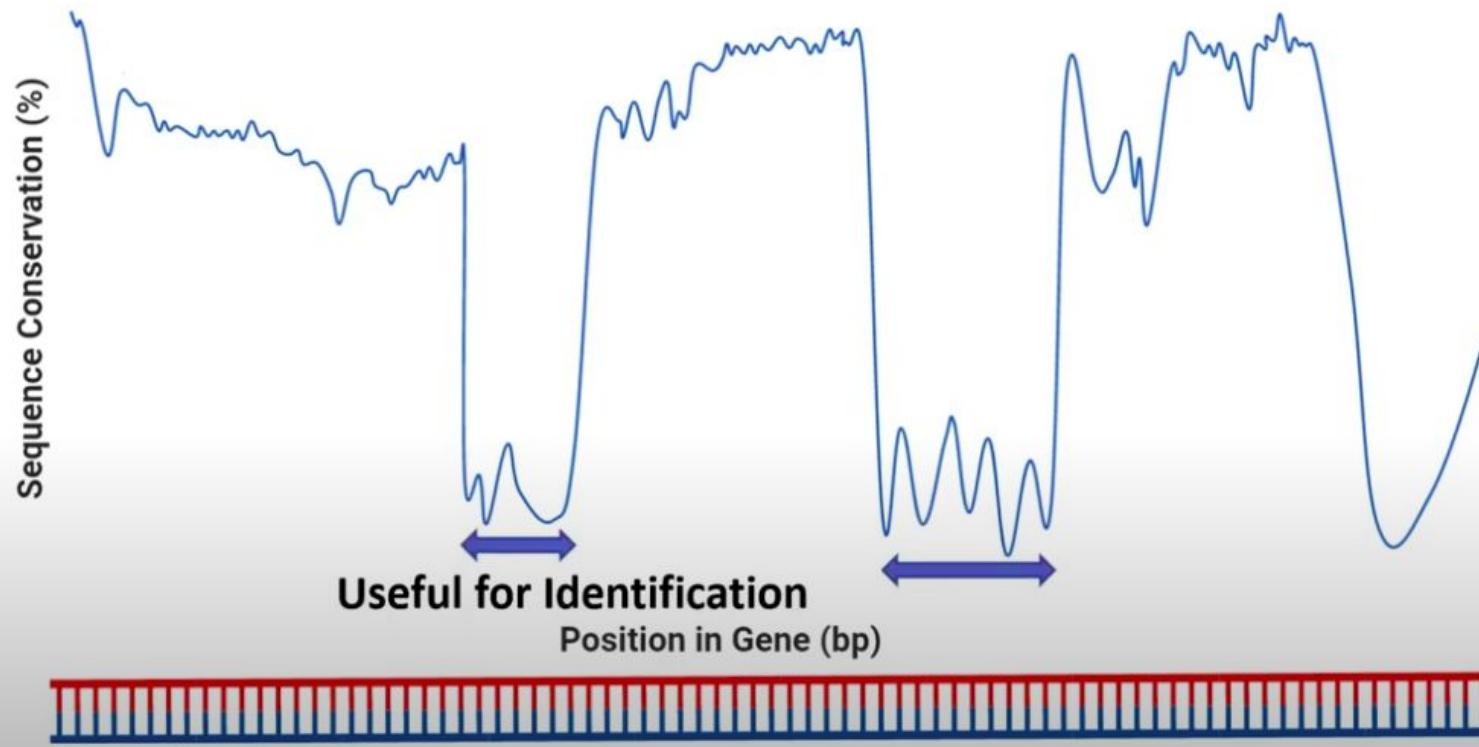
# A Potential Target for Sequencing



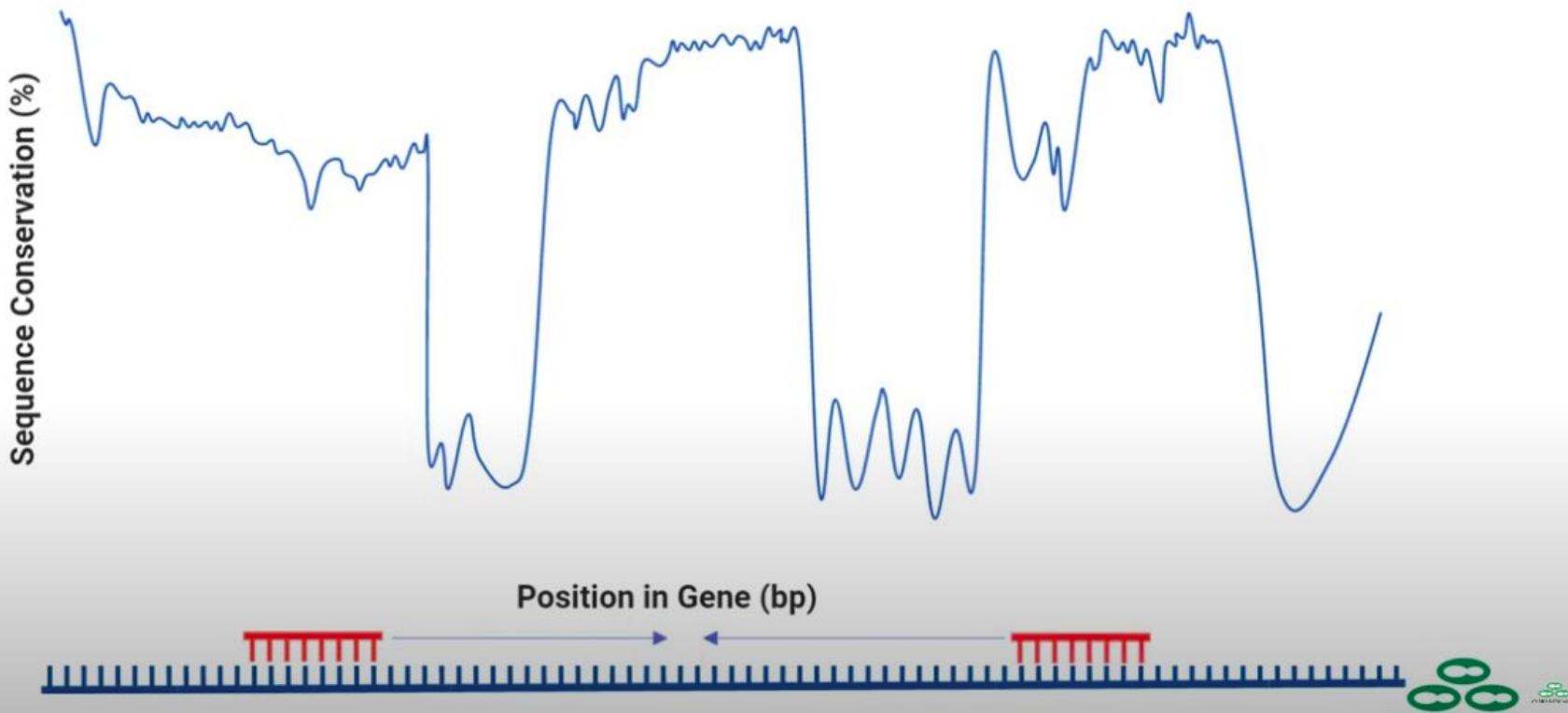
Position in Gene (bp)



# A Potential Target for Sequencing



# A Potential Target for Sequencing



# Finding an Ideal Target Sequence

## What to look for in a targeted sequencing gene

- Sequencable highly conserved regions surrounding variable regions
- Present in potential target species
- Reasonably well-characterized and understood
- Large existing database of reference sequences already available

# Finding an Ideal Target Sequence

## What to look for in a targeted sequence gene

- Sequencable highly conserved regions surrounding variable regions
- Present in potential target species
- Reasonably well-characterized and understood
- Large existing database for analysis available

C A G A U C

GUIGCACAAUGAUCAUCGUCUAGAU

THIS PART

# 16S Sequencing Challenges

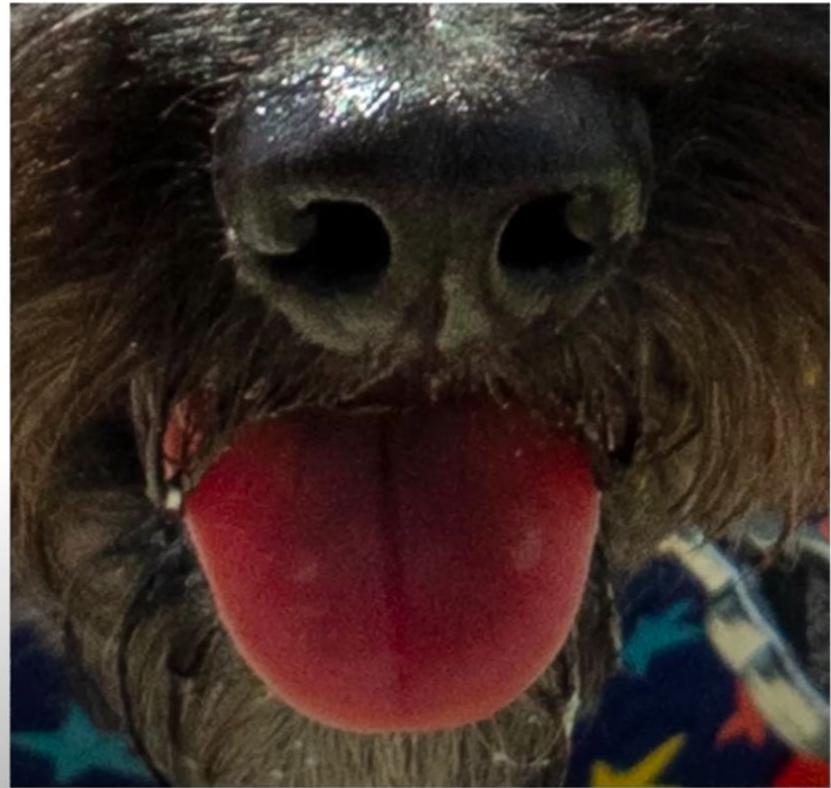
## **Targeted sequence with a few bases differentiating species**

- Sequencing is imperfect
  - Illumina usually makes some base call errors
  - Nanopore makes more
  - Errors are not necessarily evenly-distributed
- We do not want errors to be confused with real diversity/new species

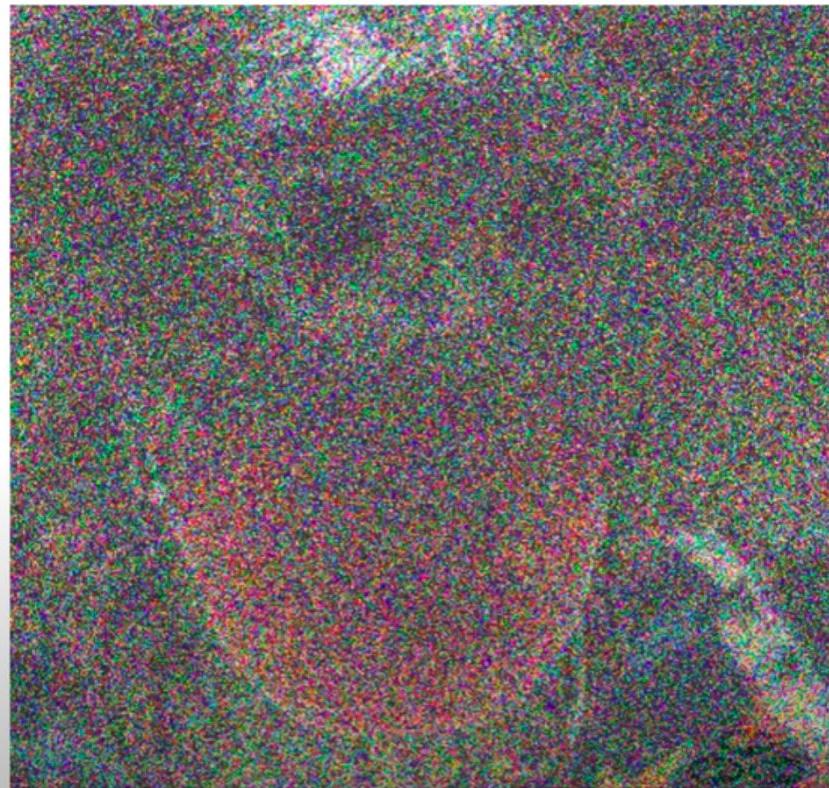
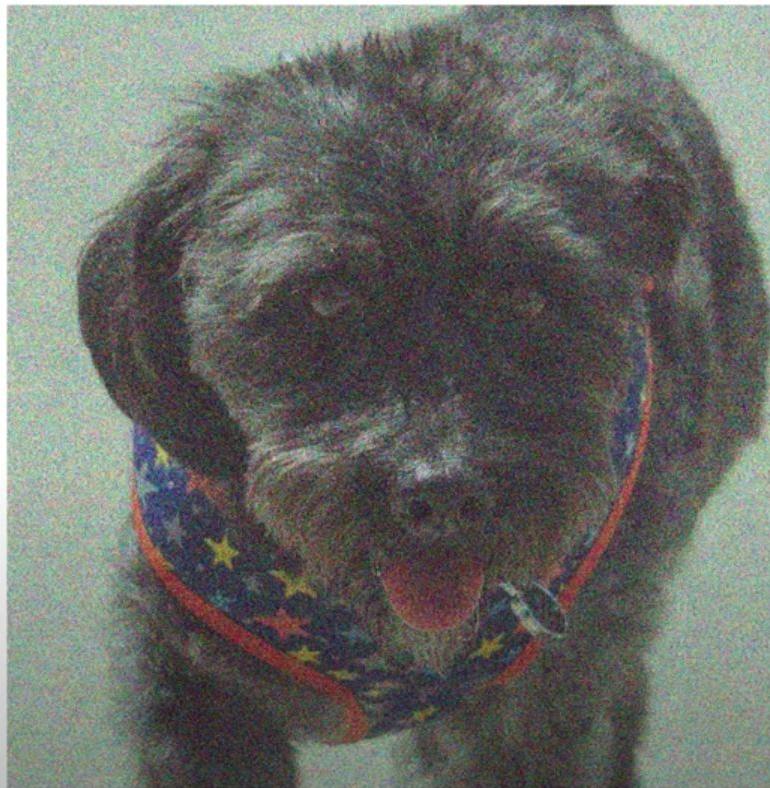
# The True Image



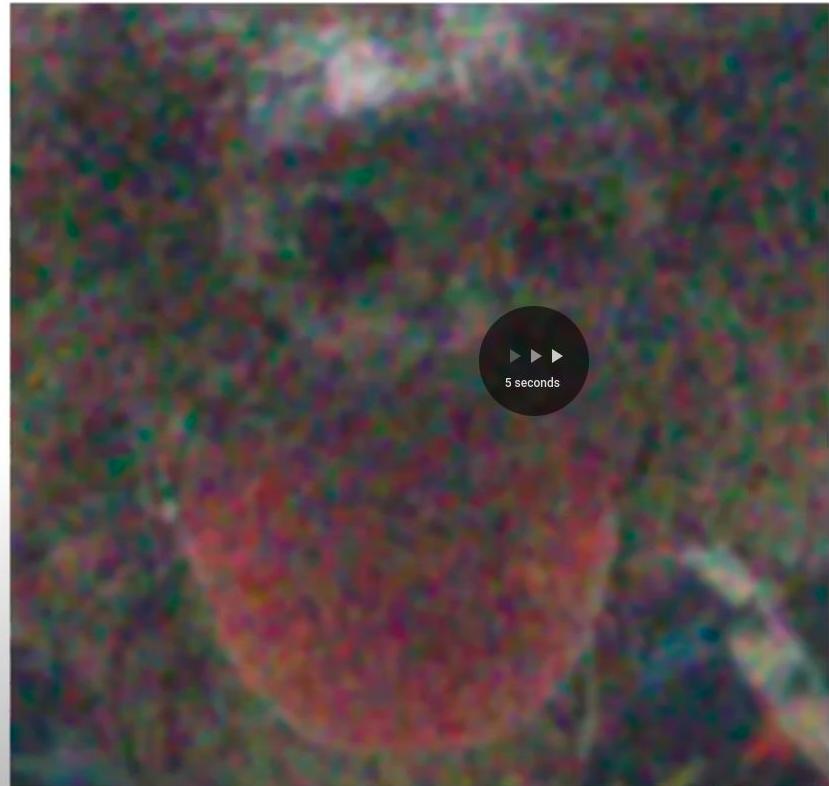
# The True Image



# The Addition of Noise



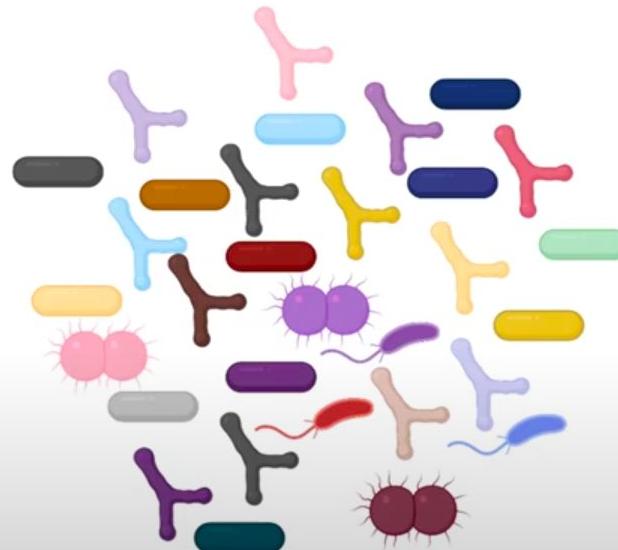
# Blur: Loss of Detail and Noise



# The True Image



# Microbes Need to Be Organized/grouped



# A Quick Note

**FOR OUR PURPOSES:**



= **1 sequence and one species**



# A Quick Note

**FOR OUR PURPOSES:**



= **1 sequence and one species**

**REAL BACTERIA**



**One species can have multiple, different copies of a gene**

# A Quick Note

**FOR OUR PURPOSES:**



= **1 sequence and one species**

**REAL BACTERIA**



**One species can have multiple, different copies of a gene**



**Two similar species may share an identical sequence**

## A Quick Note

## **FOR OUR PURPOSES:**



= 1 sequence\* and one species  
\*May contain errors



## REAL BACTERIA

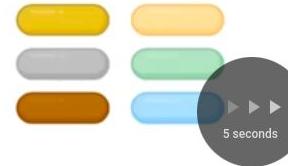
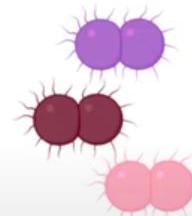


**One species can have multiple, different copies of a gene**

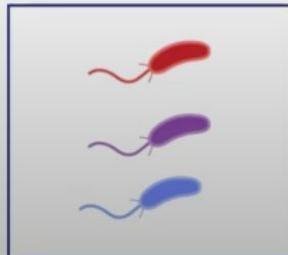
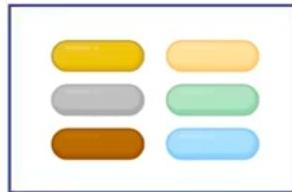
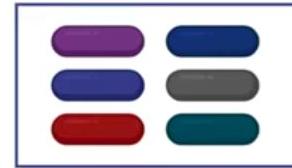
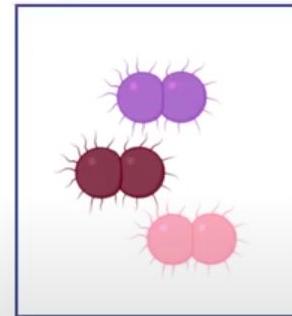


**Two similar species may share an identical sequence**

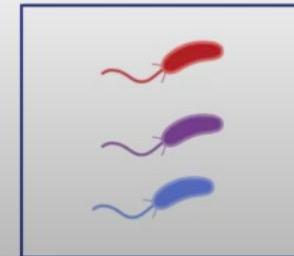
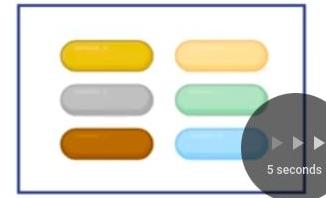
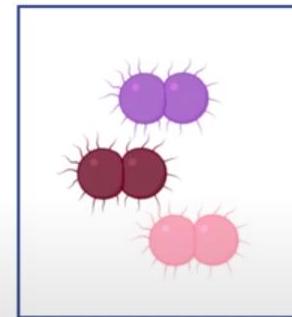
# Clustering (Blurring)



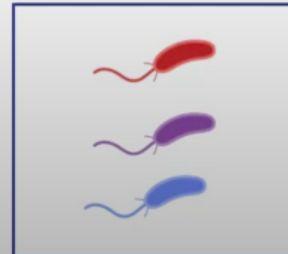
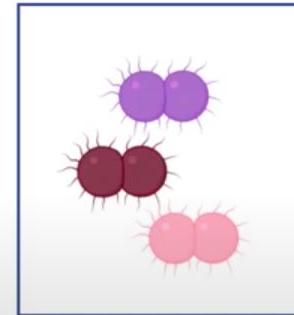
# Clustering



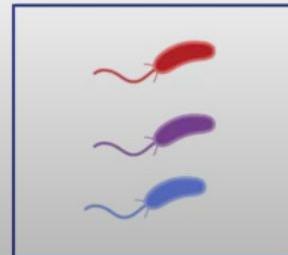
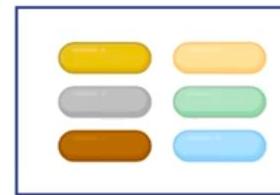
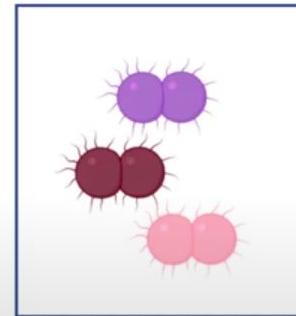
## Clustering (*De Novo*)



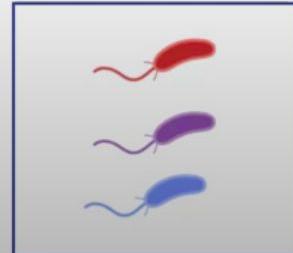
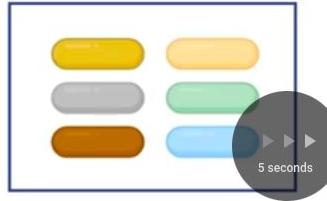
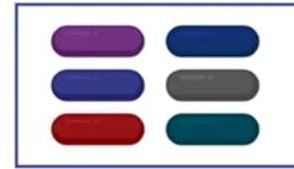
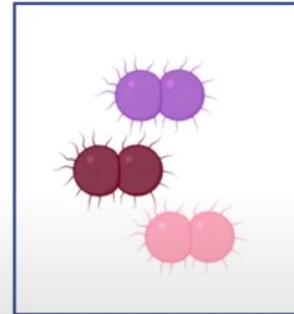
# Clustering (*De Novo*)



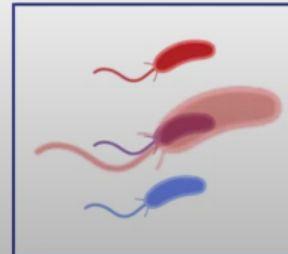
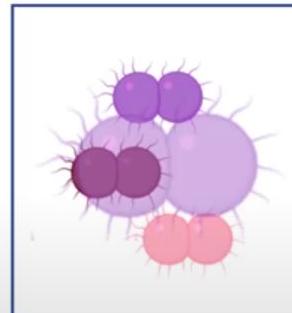
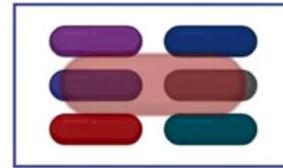
# Clustering (*De Novo*)



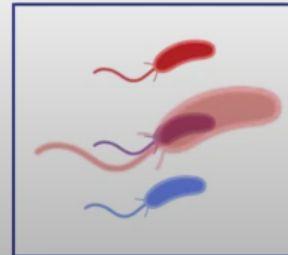
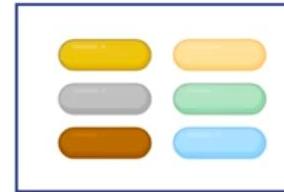
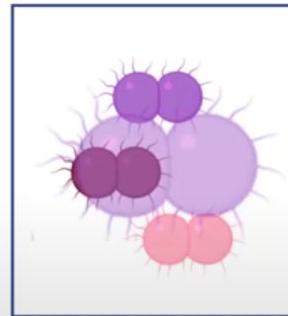
# Clustering (Open Reference)



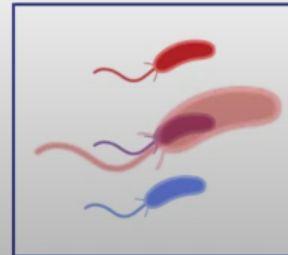
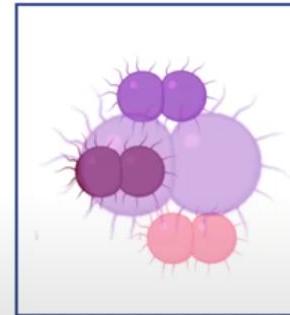
# Clustering (Open Reference)



# Clustering (Closed Reference)



# Clustering (Closed Reference)



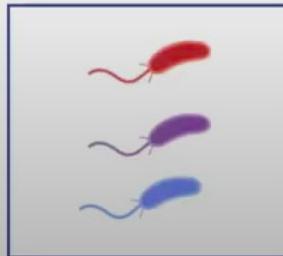
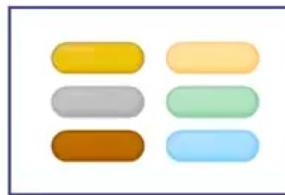
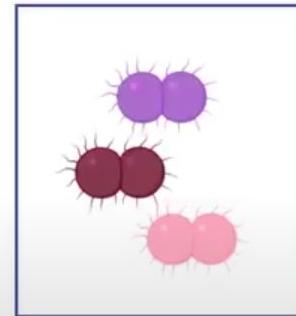
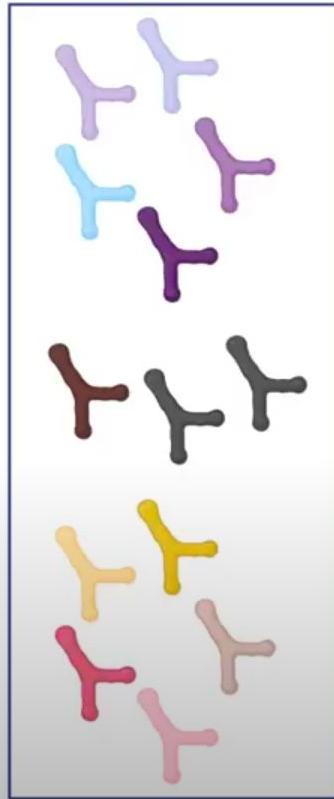
# Operational Taxonomic Unit (OTU) Approach

**We know some of these sequences arose from error/artifact.**

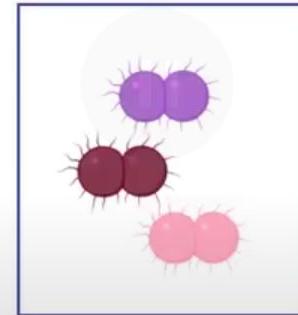
Combine extremely similar sequences (usually 97% or more identity) to minimize the effects of observed errors. Then treat each OTU as a representative sequence.

- OTUs can be represented by a representative consensus sequence
- Closed reference OTUs are fast to create, but are subject to reference bias
- *De novo* OTUs are free of reference bias, but computationally expensive and can change with changed samples
- Open reference OTUs are in between, with sequence similar to reference behaving like closed reference and more novel sequences behaving like *de novo*

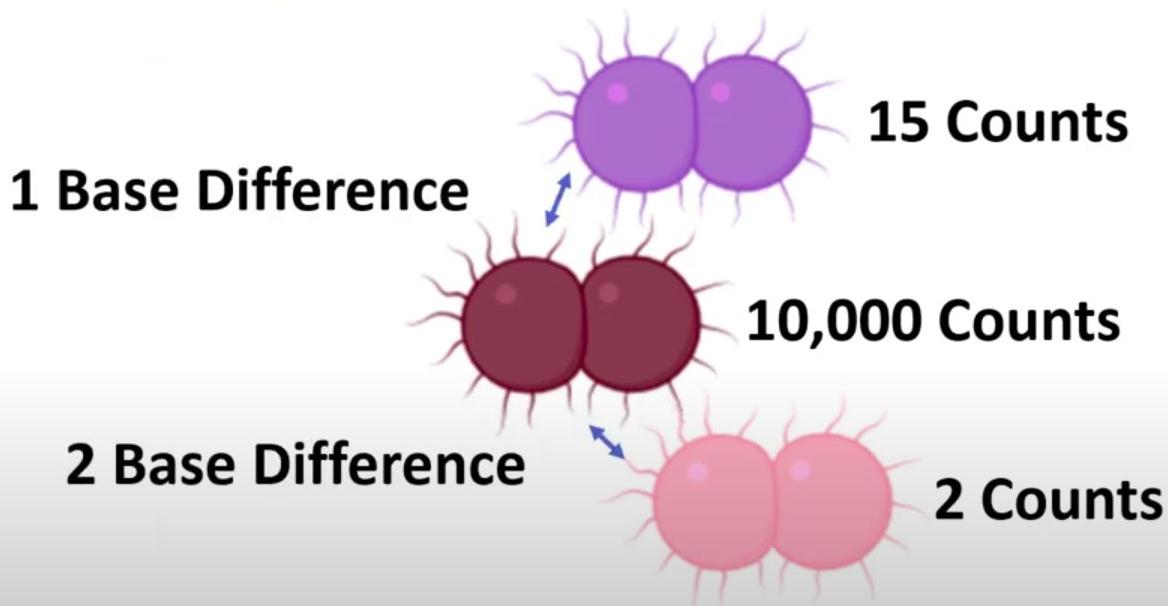
# But...



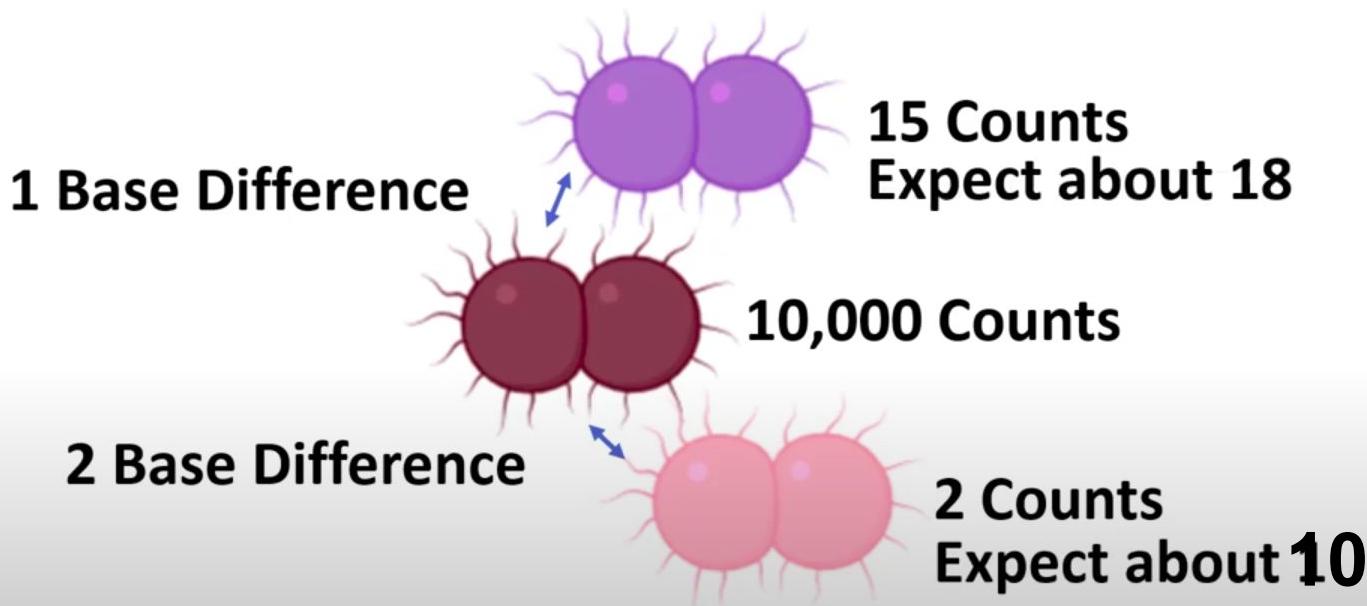
# What if...



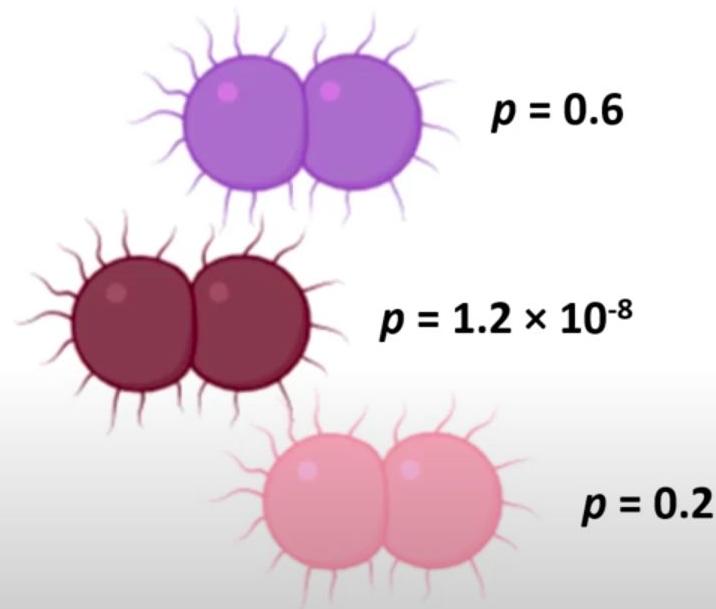
# And Differences...



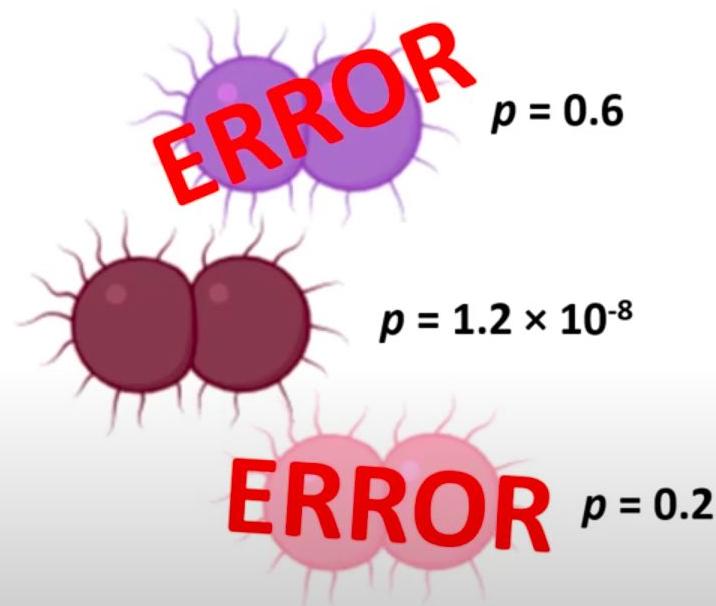
# And an Error Model for the Run...



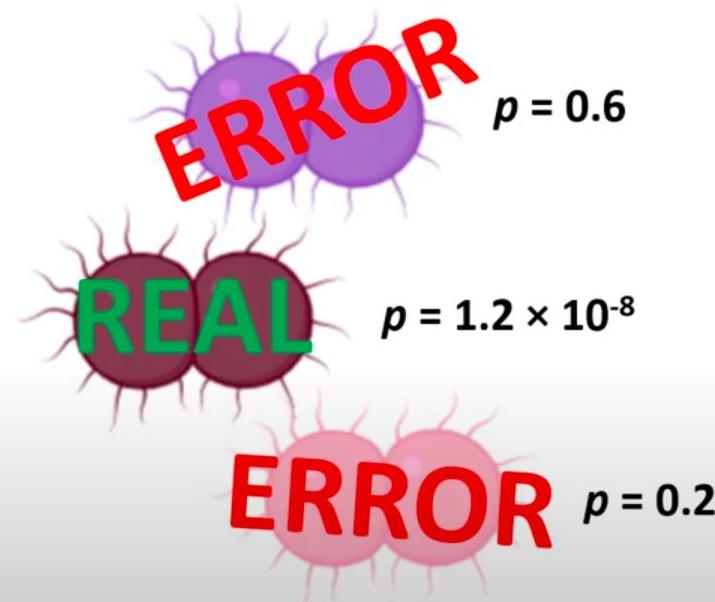
# We Could Determine Our Confidence In a Sequence



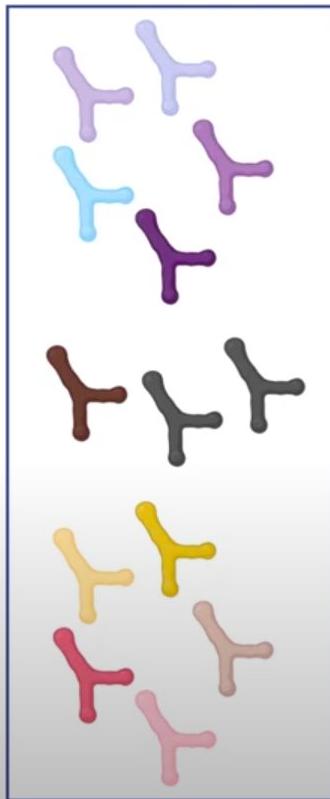
# And Drop Sequences that Are Explained by Error



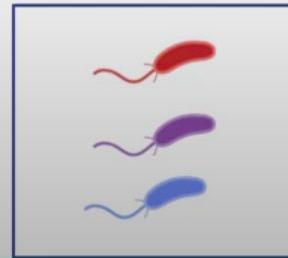
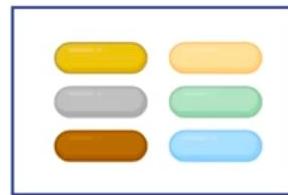
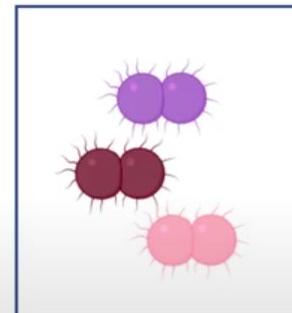
# While Retaining Those In Which We Are Confident



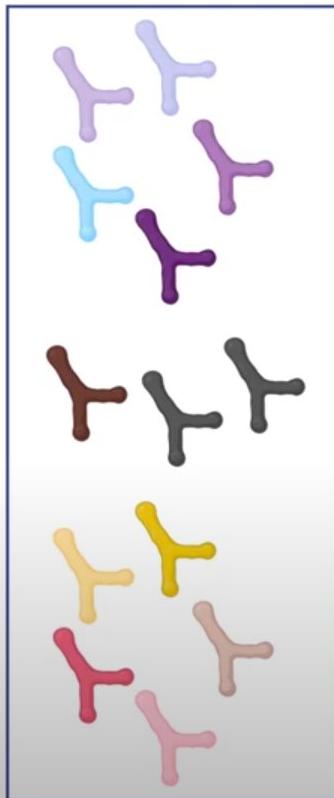
# The Outcome Looks Like Our Clustering...



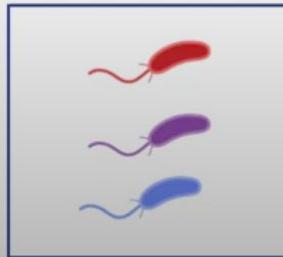
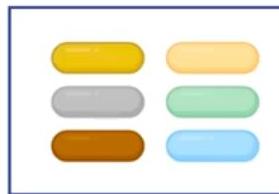
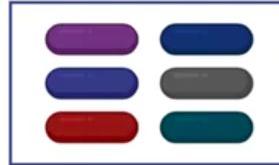
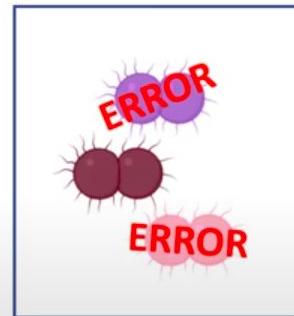
O  
p  
e  
r  
a  
t  
i  
n  
g  
  
T  
a  
x  
o  
n  
o  
m  
i  
c  
c  
o  
d  
e  
s  
  
U  
n  
i  
t  
s



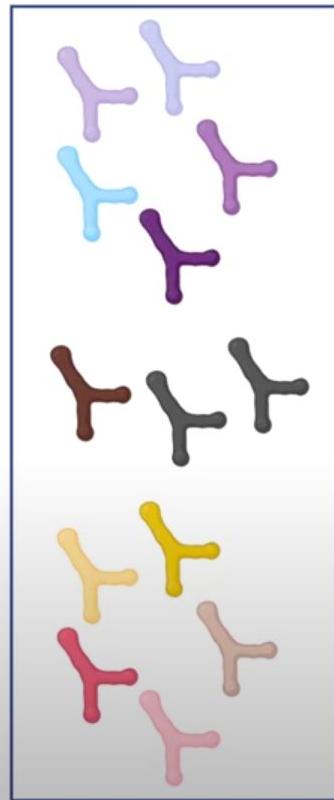
# ...What Changed?



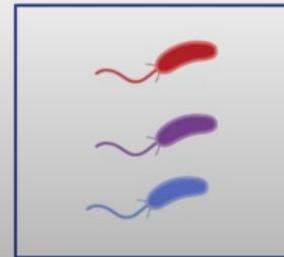
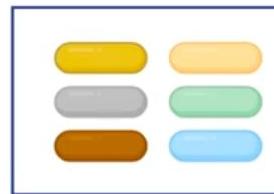
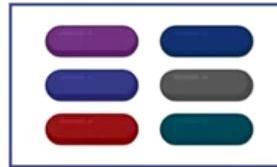
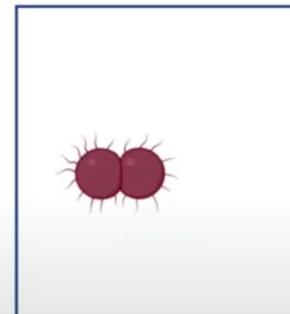
O  
perating  
T  
axonomic  
U  
nits



# ...What Changed?

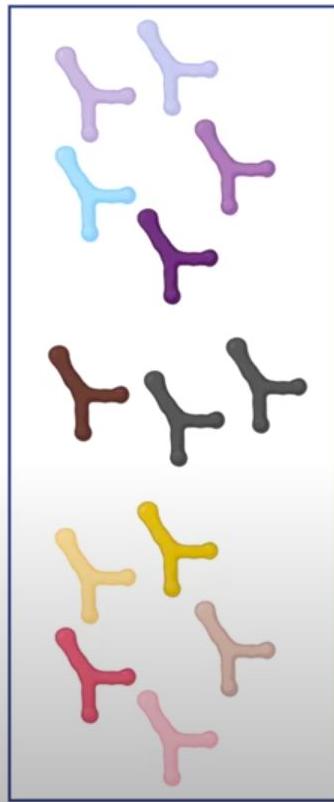


Operating  
Taxonomic  
Units

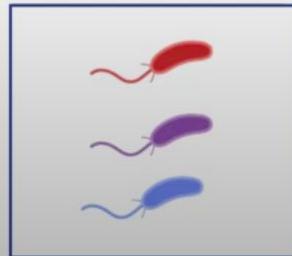
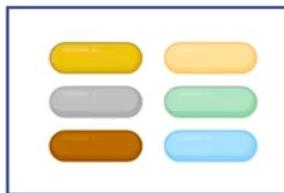


ZYMO RESEARCH

# ...What Changed?

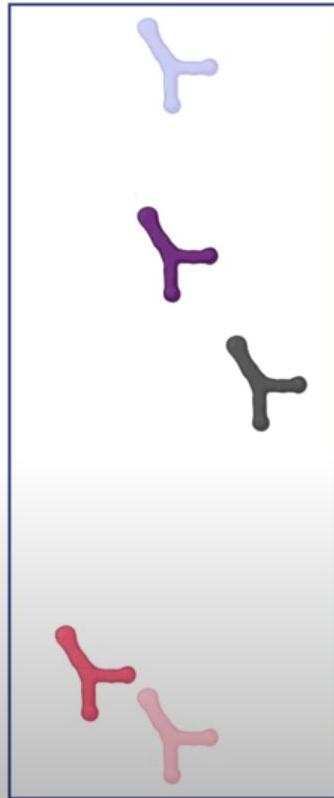


Amplicon  
Sequence  
Variant

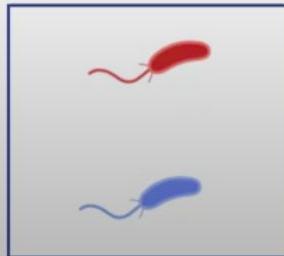
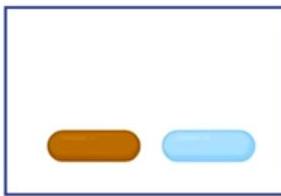
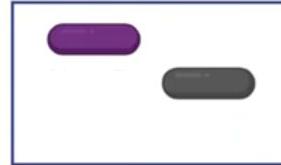
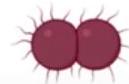


ZYMO RESEARCH

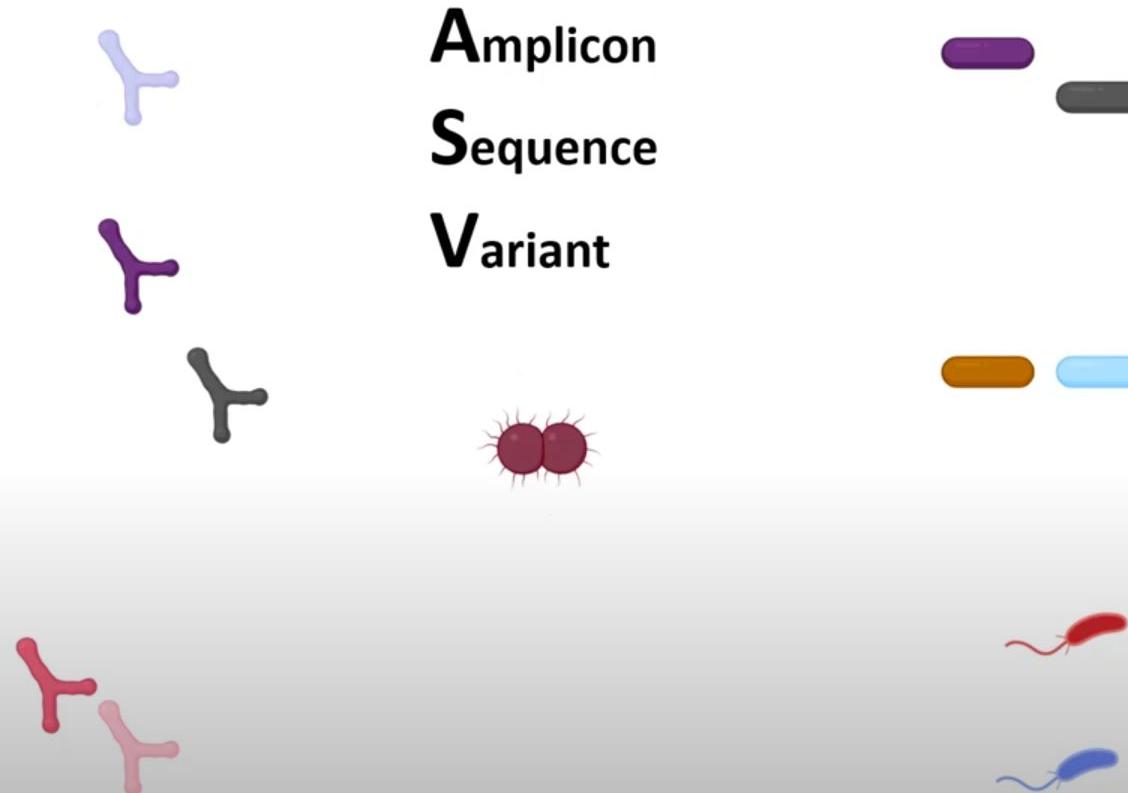
# ...What Changed?



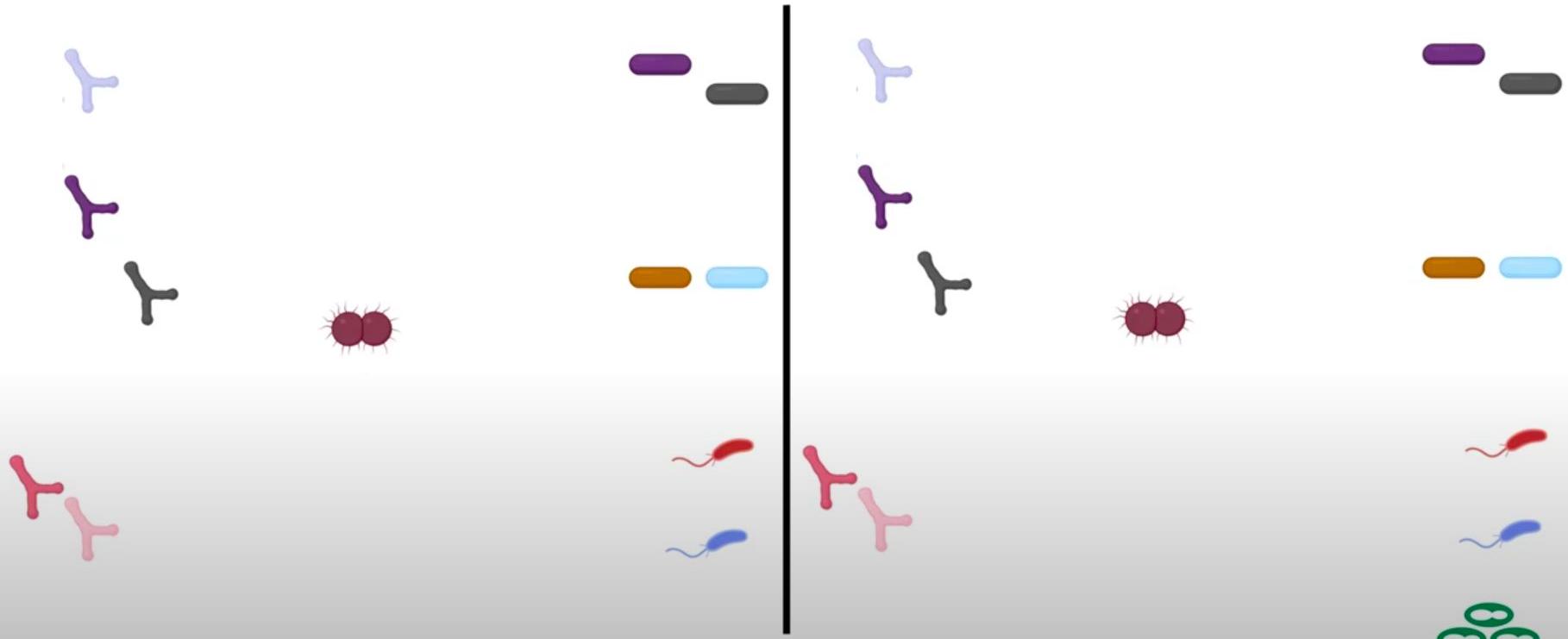
Amplicon  
Sequence  
Variant



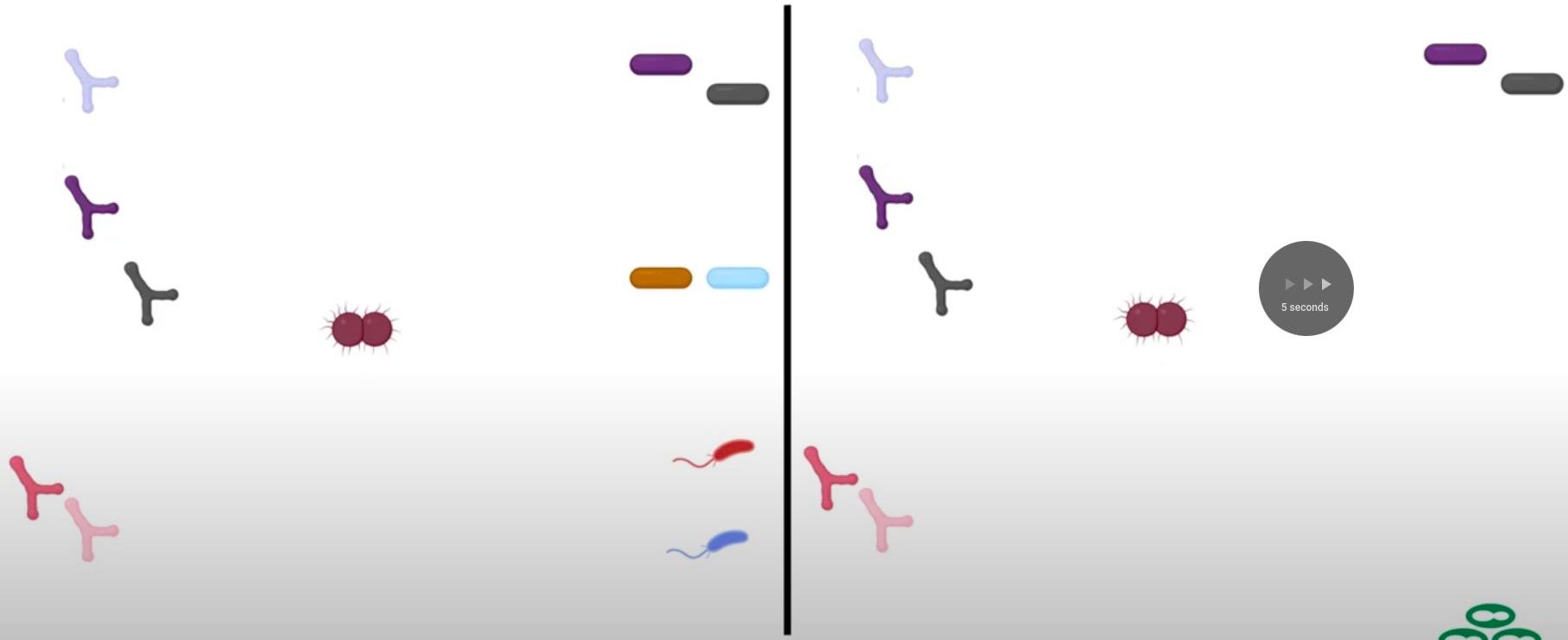
# ...What Changed?



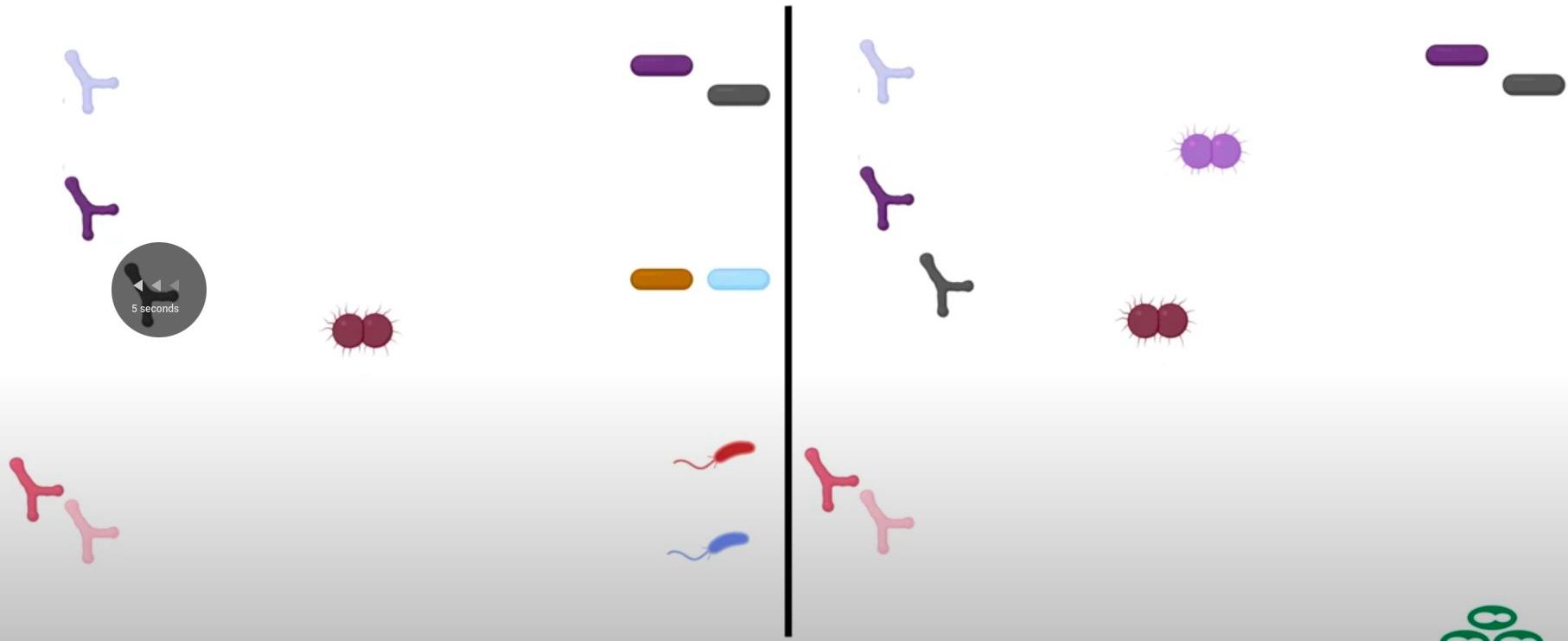
# Adding More Samples



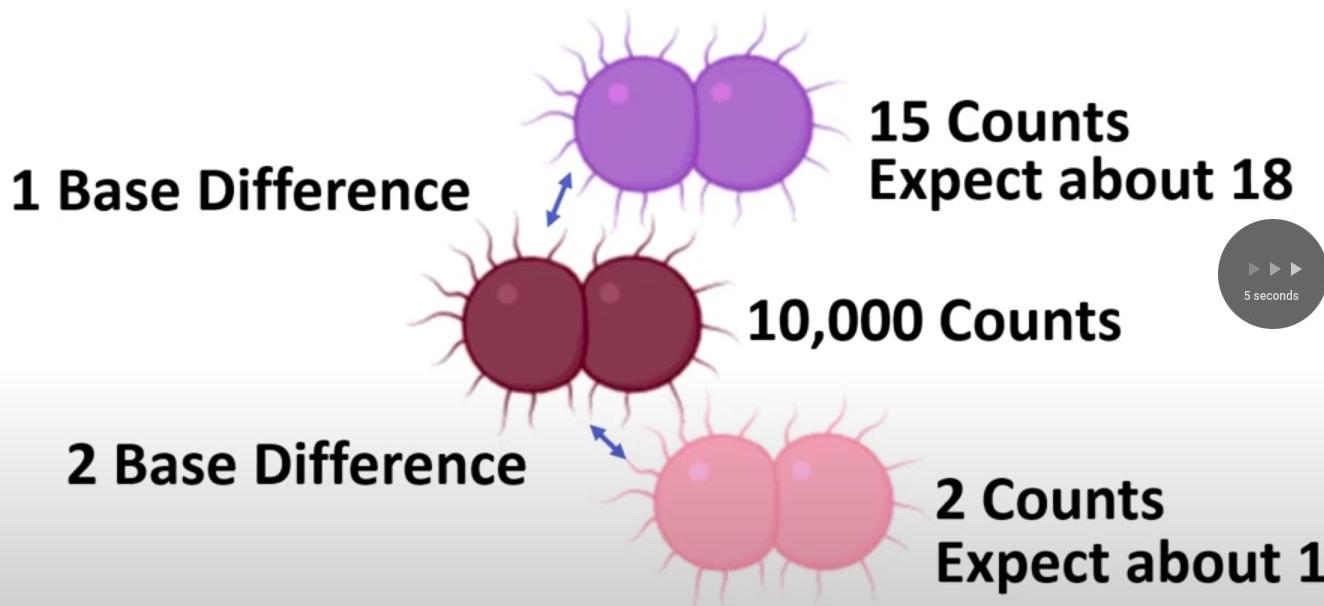
# Adding More Samples



# Adding More Samples

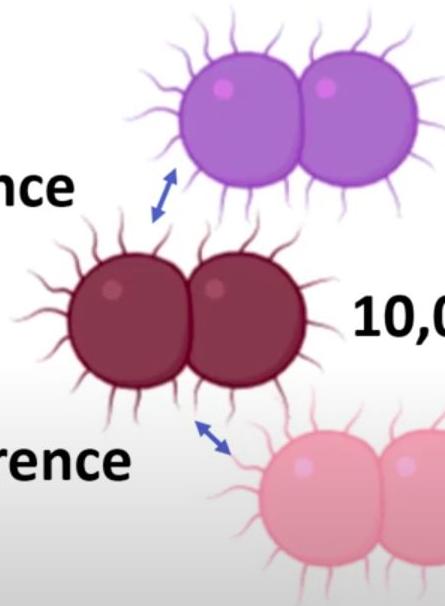


# What Else Can Change?



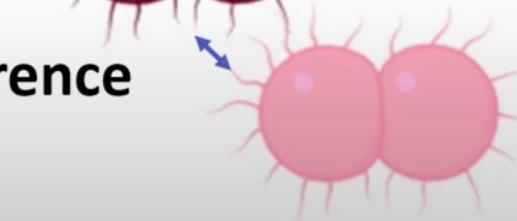
# What Else Can Change?

**1 Base Difference**



**15 Counts  
Expect about 18**

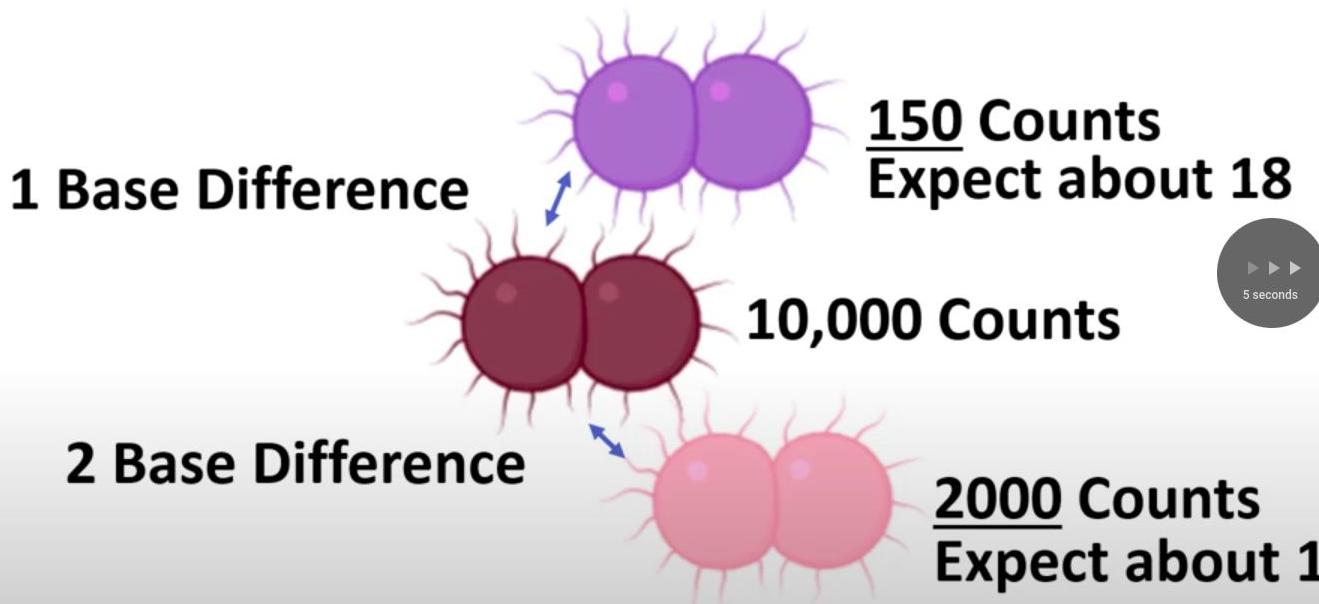
**2 Base Difference**



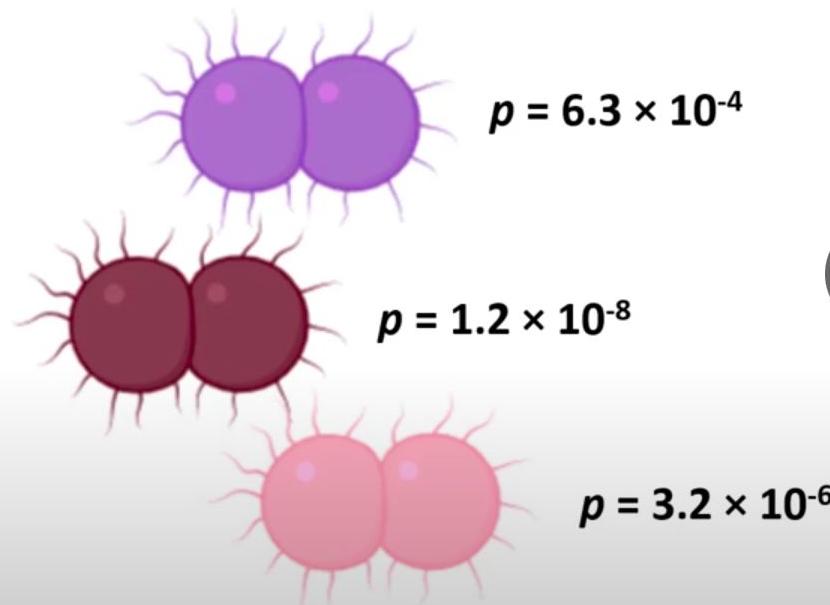
**10,000 Counts**

**2 Counts  
Expect about 1**

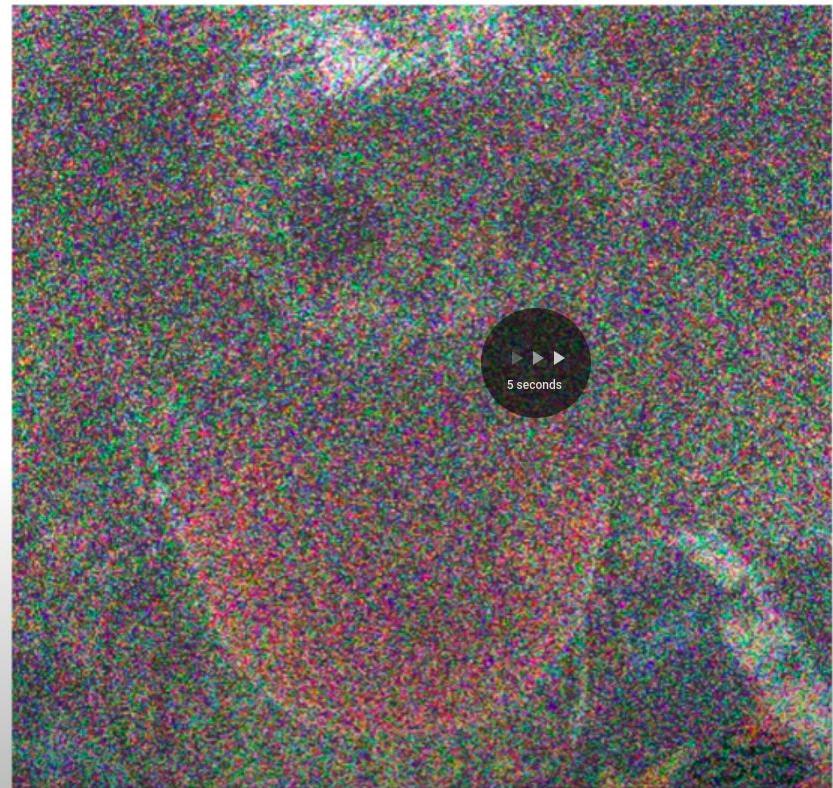
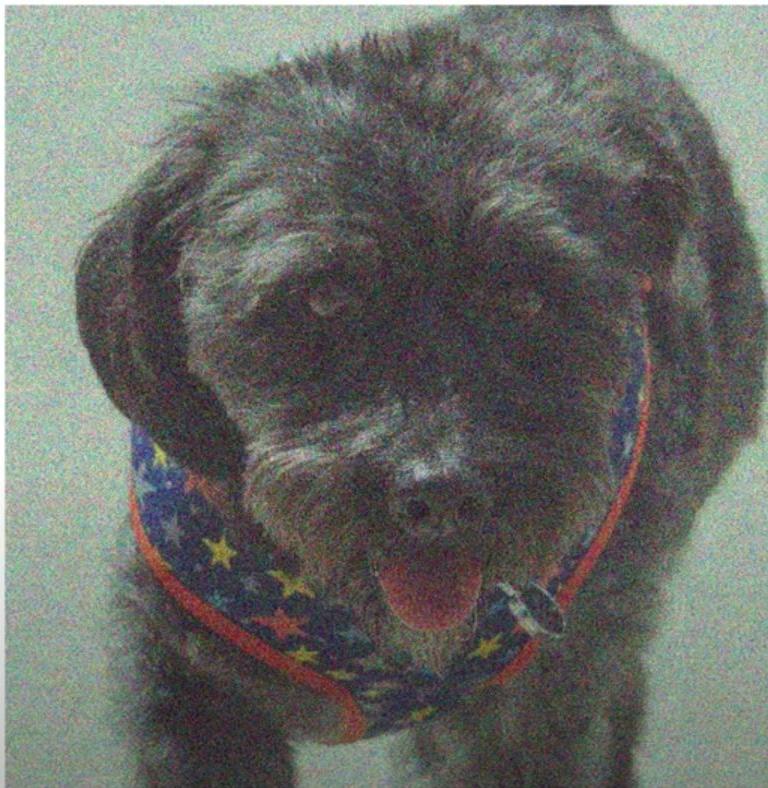
# What If The Counts Were Far Higher?



# What If The Counts Were Far Higher?



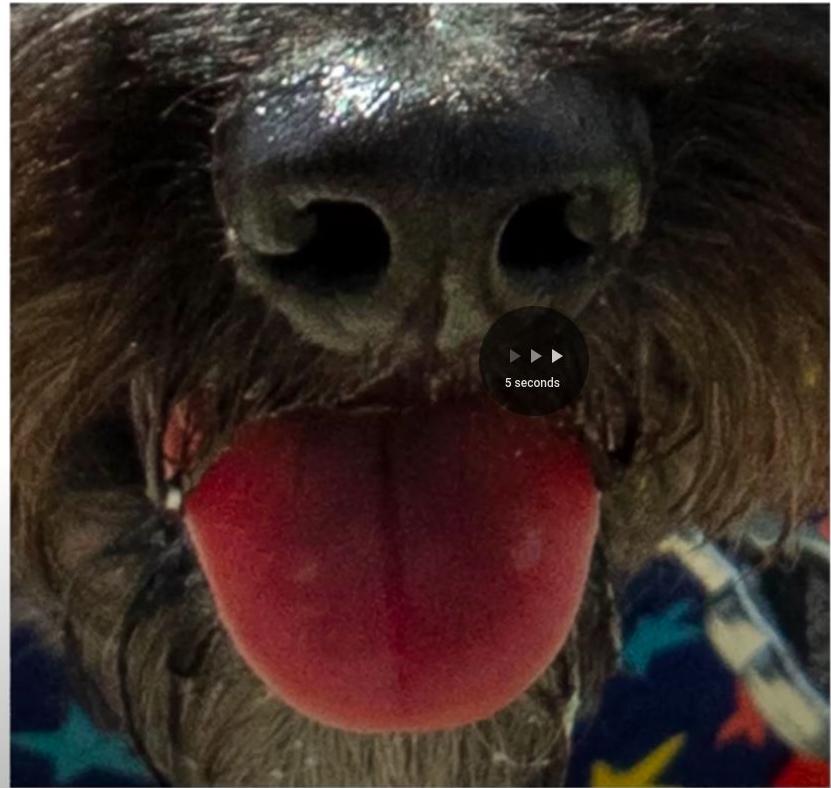
# Can This Noise Be Predicted?



ZYMO RESEARCH

73

# The True Image



ZYMO RESEARCH

# Amplicon Sequence Variant (ASV) Approach

## What is the statistical support for each sequence's existence?

Throw out amplicon sequences that lack strong statistical support for not being artifacts of sequencing. Cost: potential loss of real sequence that was present at very low levels.

- No representative consensus sequence
  - Each sequence is supported as being present in the sample
- Potentially higher resolution, no potential to combine multiple “real” sequences into an abstract
- Generated without the use of a reference, no risk of reference bias.
- May also be called an ESV (exact sequence variant) or zOTU (zero-radius OTU)



# OTU vs. ASV

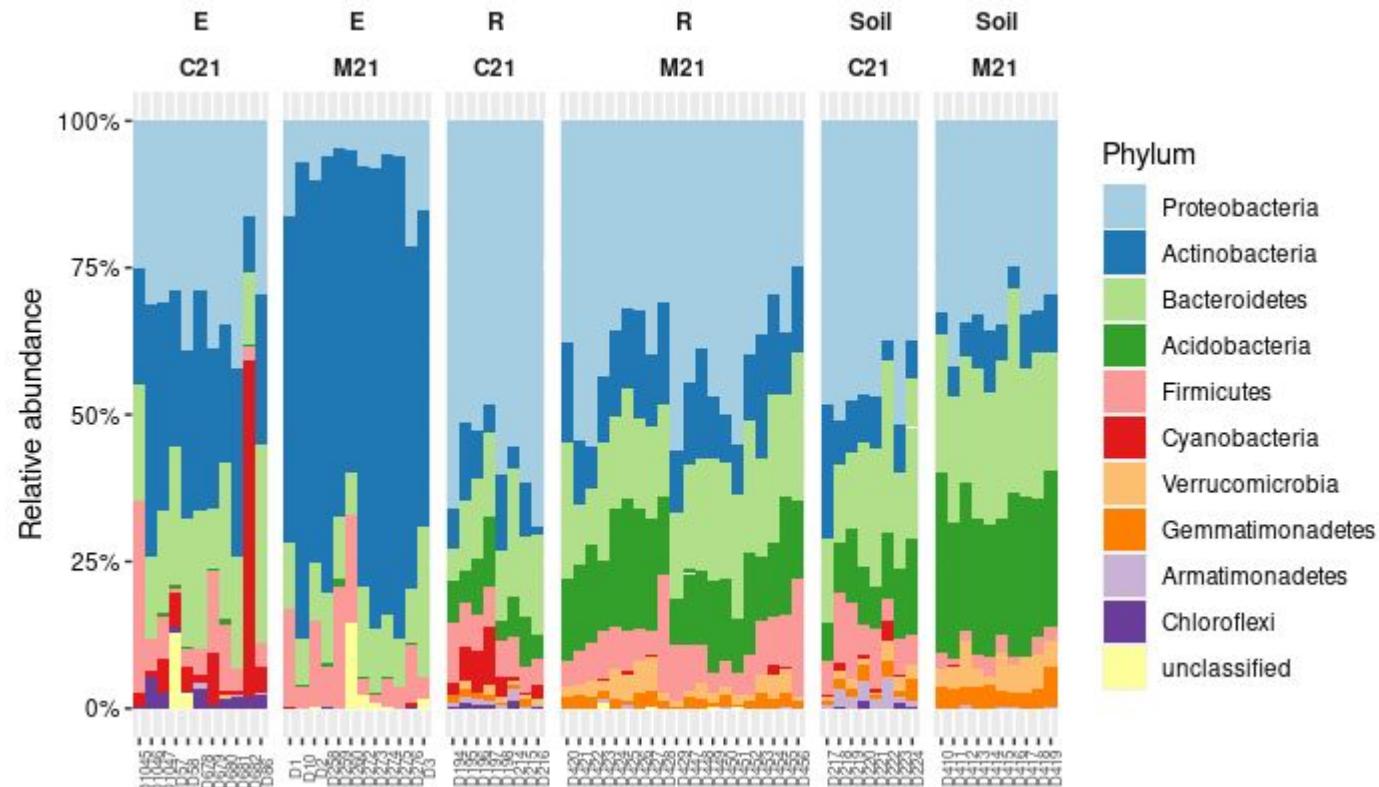
OTU	ASV
Can be subject to reference bias	Reference is not used until taxonomy assignment
OTU tables cannot be combined between studies	ASV tables can be compared across studies
Represented by a consensus sequence	Represented by an exact sequence
Can represent multiple species with different sequences	If it represents multiple species, it is because they share the sequence
Subject to chimeric sequences	Subject to chimeric sequences
Chimera detection can be complex and may require reference bias	Chimera detection is simple and reference-free



# OTU vs. ASV

OTU	ASV
Can be subject to reference bias	Reference is not used until taxonomy assignment
OTU tables cannot be combined between studies	ASV tables can be compared across studies
Represented by a consensus sequence	Represented by an exact sequence
Can represent multiple species with different sequences	If it represents multiple species, it is because they share the sequence
Subject to chimeric sequences	Subject to chimeric sequences
Chimera detection can be complex and may require reference bias	Chimera detection is simple and reference-free
Constantly changed to “OUT” by autocorrect	Generally left alone

### III. Microbiome composition abundance visualize tutorial



# 1 - Otu table

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
	otu	X10n	X10p	X11n	X11p	X120n	X120p	X121n	X121p	X122n	X122p	X125n	X125p	X126n	X126p	X127n	X13n	X13p	X140n	X140p	X141n
1																					
2	Otu001	13679	6292	42	2500	18850	5	43	7138	9432	10541	9	9772	1388	7	31538	38	2338	23	9	1358
3	Otu002	18	7134	38	9830	45	61420	182	23751	36	11	4535	3502	11018	5473	26	14411	38	19018	12	3080
4	Otu003	9939	8983	31	13	24620	19	19	16	12502	3831	4621	2240	9924	4052	9292	18	0	37	7	3680
5	Otu004	3675	4234	24	22	11	16	32967	35	6	18	6908	5	16	8702	24	11	37717	0	25	4196
6	Otu005	0	5	0	7	0	8	0	16	20166	0	0	2	5	8	2	16	0	13	0	0
7	Otu006	0	8	0	0	0	8	0	0	5	3	3	0	0	9	0	5	4	0	0	3
8	Otu007	4587	518	4	386	8775	5	6	1102	14336	0	0	3626	51	0	6	12	0	10	0	395
9	Otu008	1	8	2	4408	3	29	6	12355	0	0	0	0	0	9	3	1588	0	6	3	3
10	Otu009	115	914	3	325	0	629	1	834	5	0	1354	2108	1117	67	0	2010	1897	11227	1	3
11	Otu010	780	8	23810	12	3279	0	12	7	3027	0	2	4156	0	0	18	0	0	0	0	0
12	Otu011	0	3	2	2	0	13	5	5	4	7	3081	11	4	6804	0	3	11	0	5	0
13	Otu012	0	0	0	6	0	0	0	16	3	0	0	0	0	0	0	17	0	6	0	0
14	Otu013	6321	2471	2	0	12	3	0	0	4	20272	0	15	9	0	5	0	11	0	14	0
15	Otu014	0	82	4	3304	1	1667	4	9233	13	3	0	2707	0	0	3	4806	9	3	5	0
16	Otu015	0	12	0	3	7	25	1	6	10	0	4	2772	1	3	0	2	0	10	13	8052
17	Otu016	1	0	0	9	5	0	0	14	0	0	0	0	2654	0	0	6	1	1	0	0
18	Otu017	0	0	0	0	0	0	0	17	8	0	0	0	0	0	0	17	24	48	35210	4
19	Otu018	1	0	9	911	0	0	15	2702	6	4	342	2217	606	0	13	3846	4	6	8513	1
20	Otu019	0	0	13	0	0	0	29	0	0	0	0	0	0	0	11	0	0	5	4	0
21	Otu020	425	0	1	0	1706	0	8447	1	0	0	0	0	0	26	0	0	3490	0	2620	0
22	Otu021	0	4	0	0	0	10	0	0	0	0	2	0	0	4	0	0	0	0	0	4
23	Otu022	0	0	0	4987	0	0	0	6	90	1	1	524	0	467	0	4	8	6198	0	1
24	Otu023	4	0	1	0	0	3	0	0	0	0	0	0	3351	3	0	3910	1	2	3	0
25	Otu024	0	0	0	0	0	0	0	0	0	0	0	0	17	0	0	0	0	2	1	0
26	Otu025	69	0	0	0	290	0	0	0	21	0	118	2	9	513	2	0	0	2	0	0
27	Otu026	0	2	0	0	0	3	0	0	0	1	0	0	0	0	0	0	0	0	0	0
28	Otu027	6	2304	0	0	5	0	0	0	57	4	0	14529	9597	2	6	0	0	0	0	0

# 2 - Taxonomy Table

	A	B	C	D	E	F	G	H	I
1	otu	Domain	Supergroup	Division	Class	Order	Family	Genus	
2	Otu001	Eukaryota	Archaeplastida	Chlorophyta	Mamiellophyceae	Mamiellales	Bathycoccaceae	Ostreococcus	
3	Otu002	Eukaryota	Hacrobia	Haptophyta	Prymnesiophyceae	Prymnesiophyceae_X	Braarudosphaeraceae	UCYN_A1_host	
4	Otu003	Eukaryota	Archaeplastida	Chlorophyta	Mamiellophyceae	Mamiellales	Bathycoccaceae	Bathycoccus	
5	Otu004	Eukaryota	Alveolata	Dinophyta	Dinophyceae	Dinophyceae_X	Dinophyceae_X	Prorocentrum	
6	Otu005	Eukaryota	Stramenopiles	Ochrophyta	Bacillariophyta	Mediophyceae	Mediophyceae_X	Thalassiosira	
7	Otu006	Eukaryota	Stramenopiles	Ochrophyta	Bacillariophyta	Bacillariophyceae	Bacillariophyceae_X	Pseudo_nitzschia	
8	Otu007	Eukaryota	Stramenopiles	Ochrophyta	Pelagophyceae	Pelagophyceae_X	Pelagophyceae_X	Pelagomonas	
9	Otu008	Eukaryota	Alveolata	Dinophyta	Dinophyceae	Dinophyceae_X	Dinophyceae_X	Dinophyceae_X	
10	Otu009	Eukaryota	Hacrobia	Haptophyta	Prymnesiophyceae	Prymnesiales	Chrysochromulinaceae	Chrysochromulina	
11	Otu010	Eukaryota	Opisthokonta	Metazoa	Craniata	Craniata_X	Craniata_XX	Craniata_XX_unclassified	
12	Otu011	Eukaryota	Stramenopiles	Ochrophyta	Chrysophyceae	Chrysophyceae_X	Chrysophyceae_Clade_C	Chrysophyceae_Clade_C_X	
13	Otu012	Eukaryota	Alveolata	Dinophyta	Dinophyceae	Dinophyceae_X	Dinophyceae_X	Gonyaulax	
14	Otu013	Eukaryota	Alveolata	Dinophyta	Syndiniales	Syndiniales_Group_III	Syndiniales_Group_III_X	Syndiniales_Group_III_X	
15	Otu014	Eukaryota	Stramenopiles	Ochrophyta	Chrysophyceae	Chrysophyceae_X	Chrysophyceae_Clade_G	Chrysophyceae_Clade_G_X	
16	Otu015	Eukaryota	Alveolata	Dinophyta	Dinophyceae	Dinophyceae_X	Dinophyceae_X	Dinophyceae_X	
17	Otu016	Eukaryota	Hacrobia	Centroheliozoa	Centroheliozoa_X	Pterocystida	Pterocystida_X	Pterocystida_X	
18	Otu017	Eukaryota	Opisthokonta	Fungi	Basidiomycota	Agaricomycotina	Agaricomycetes	Hypodontia	
19	Otu018	Eukaryota	Stramenopiles	Ochrophyta	Dictyochophyceae	Dictyochophyceae_X	Pedinellales	Pedinellales_X	
20	Otu019	Eukaryota	Opisthokonta	Fungi	Basidiomycota	Agaricomycotina	Agaricomycetes	Iteronilla	
21	Otu020	Eukaryota	Hacrobia	Haptophyta	Prymnesiophyceae	Prymnesiophyceae_X	Braarudosphaeraceae	Braarudosphaera	
22	Otu021	Eukaryota	Alveolata	Dinophyta	Dinophyceae	Dinophyceae_X	Dinophyceae_X	Dinophyceae_X	
23	Otu022	Eukaryota	Hacrobia	Haptophyta	Prymnesiophyceae	Prymnesiophyceae_X	Prymnesiophyceae_X	Syracosphaera	
24	Otu023	Eukaryota	Stramenopiles	Ochrophyta	Bacillariophyta	Bacillariophyceae	Bacillariophyceae_X	Bacillariophyceae_X	
25	Otu024	Eukaryota	Archaeplastida	Streptophyta	Klebsormidiophyceae	Klebsormidiophyceae_X	Klebsormidiophyceae_XX	Klebsormidium	
26	Otu025	Eukaryota	Archaeplastida	Chlorophyta	Mamiellophyceae	Mamiellales	Mamiellaceae	Micromonas	
27	Otu026	Eukaryota	Stramenopiles	Ochrophyta	Bacillariophyta	Bacillariophyceae	Bacillariophyceae_X	Cylindrotheca	
28	Otu027	Eukaryota	Alveolata	Dinophyta	Dinophyceae	Suessiales	Suessiales_X	Karlodinium	
29	Otu028	Eukaryota	Hacrobia	Haptophyta	Prymnesiophyceae	Isochrysidales	Noelaerhabdaceae	Emiliania	
30	Otu029	Eukaryota	Opisthokonta	Fungi	Ascomycota	Saccharomycotina	Saccharomycetales	Debaryomyces	
31	Otu030	Eukaryota	Hacrobia	Cryptophyta	Cryptophyceae	Cryptophyceae_X	Cryptomonadales	Teleaulax	
32	Otu031	Eukaryota	Alveolata	Dinophyta	Syndiniales	Syndiniales_Group_I	Syndiniales_Group_I_Clade_1	Syndiniales_Group_I_Clade_1_X	
33	Otu032	Eukaryota	Archaeplastida	Chlorophyta	Prasinophyceae	Prasinophyceae_VII	Prasinophyceae_VII_X	Prasinophyceae_VII_A	Prasinophyceae_VII_A_X

# 3 - Sample table (metadata)

A	B	C	D	E	F	G	H	I	J	K	L
sample	fraction	Select_18S_nifH	total_18S	total_16S	total_nifH	sample_number	transect	station	depth	latitude	longitude
X10n	Nano	Yes	53230	8772	36	10	1	81	140	-27.42	-44.72
X10p	Pico	Yes	47390	4448	6241	10	1	81	140	-27.42	-44.72
X11n	Nano	No	24007	6193	3772	11	1	85	110	-26.8	-45.3
X11p	Pico	Yes	31899	14	10201	11	1	85	110	-26.8	-45.3
X120n	Nano	Yes	70455	5292	93	120	2	96	5	-27.39	-47.82
X120p	Pico	Yes	76162	53272	23147	120	2	96	5	-27.39	-47.82
X121n	Nano	Yes	52401	5958	26638	121	2	96	30	-27.39	-47.82
X121p	Pico	Yes	71785	10993	23706	121	2	96	30	-27.39	-47.82
X122n	Nano	Yes	78740	11730	15543	122	2	96	50	-27.39	-47.82
X122p	Pico	Yes	37364	11817	11045	122	2	96	50	-27.39	-47.82
X125n	Nano	Yes	27381	9	14331	125	2	98	5	-27.59	-47.39
X125p	Pico	Yes	55179	10419	21461	125	2	98	5	-27.59	-47.39
X126n	Nano	Yes	65714	15	16929	126	2	98	50	-27.59	-47.39
X126p	Pico	Yes	30406	3	10140	126	2	98	50	-27.59	-47.39
X127n	Nano	Yes	60610	9	11493	127	2	98	85	-27.59	-47.39
X13n	Nano	Yes	48001	33	21316	13	1	86	105	-26.33	-45.41
X13p	Pico	Yes	59626	7217	11954	13	1	86	105	-26.33	-45.41
X140n	Nano	Yes	48126	10428	25286	140	2	101	5	-27.79	-46.96
X140p	Pico	Yes	48569	10448	12301	140	2	101	5	-27.79	-46.96
X141n	Nano	Yes	30081	6394	21302	141	2	101	60	-27.79	-46.96
X141p	Pico	Yes	64221	11318	10428	141	2	101	60	-27.79	-46.96
X142n	Nano	Yes	85219	23243	11753	142	2	101	110	-27.79	-46.96
X142p	Pico	Yes	89797	9553	17156	142	2	101	110	-27.79	-46.96
X155n	Nano	Yes	54162	8237	20674	155	2	106	5	-28.12	-46.17
X155p	Pico	Yes	50782	7384	66172	155	2	106	5	-28.12	-46.17
X156n	Nano	Yes	55065	11371	14447	156	2	106	60	-28.12	-46.17
X156p	Pico	Yes	43917	9665	16093	156	2	106	60	-28.12	-46.17
X157n	Nano	Yes	29078	4978	15532	157	2	106	100	-28.12	-46.17
X157p	Pico	Yes	51848	9139	15204	157	2	106	100	-28.12	-46.17
X15n	Nano	Yes	22468	2887	2678	15	1	87	105	-26.22	-45.48
X15p	Pico	Yes	78390	13813	1033	15	1	87	105	-26.22	-45.48
X165n	Nano	Yes	50732	15337	14706	165	2	114	5	-28.65	-44.99
X165p	Pico	Yes	48514	10902	39918	165	2	114	5	-28.65	-44.99
X166n	Nano	Yes	51112	51112	51112	166	2	114	50	-28.65	-44.99

# Work in R\_notebook

1.-Import data

2.-Analyse-manipulate data

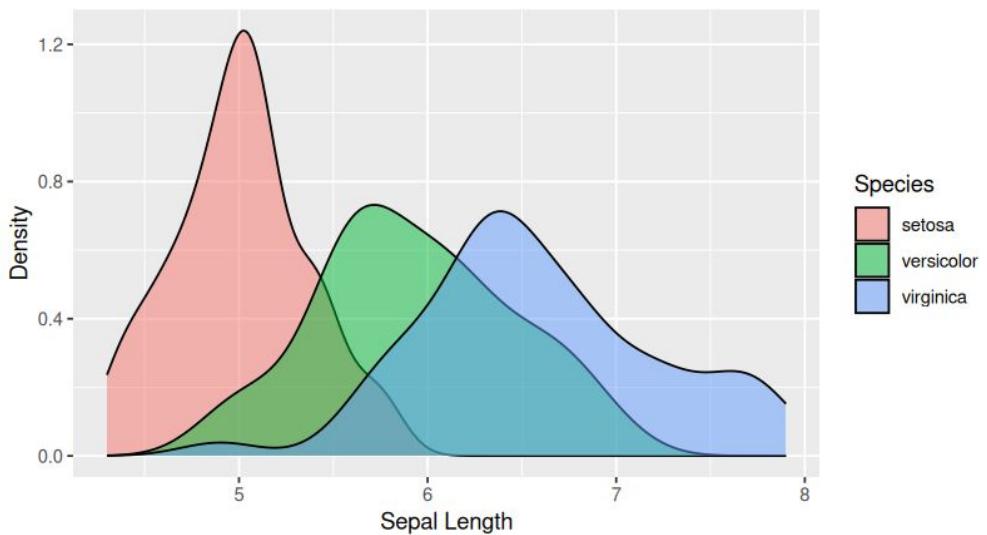
3-Draw basic\_plot

4-Extend more detail

# Summary

1. Visualise data is IMPORTANT!-quick note about install R packages in Linux
2. Reshape2
3. Learn by yourself
4. OTU and ASV
5. Microbiome barplot tutorial
6. Homework

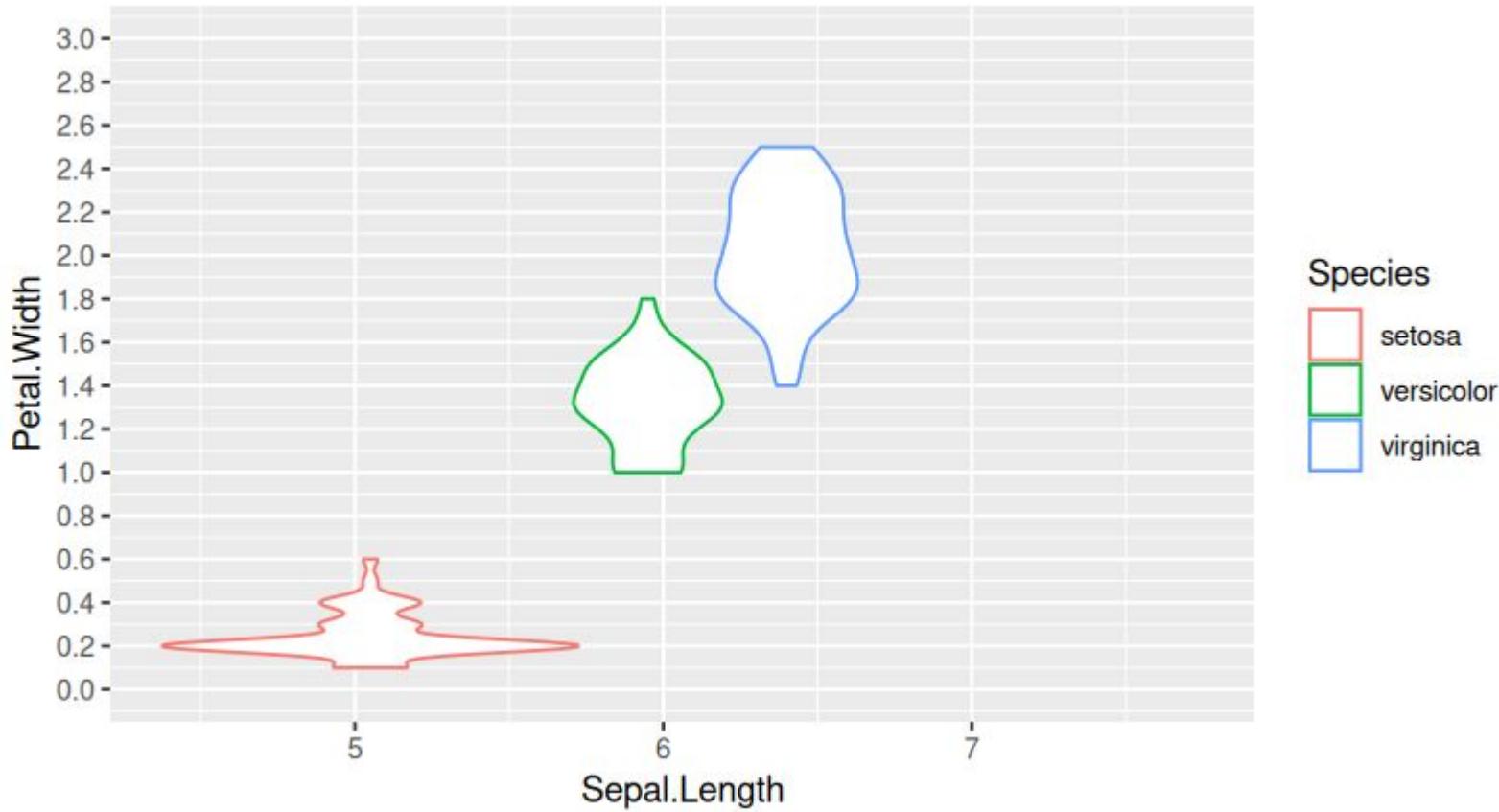
# Homework



```
# Generate the sample  
data frame
```

```
data("iris")
```

```
print(data)
```



Correlation between Sepal width and length  
iris data

