# Introduction to microbial community profiling using amplicon sequencing
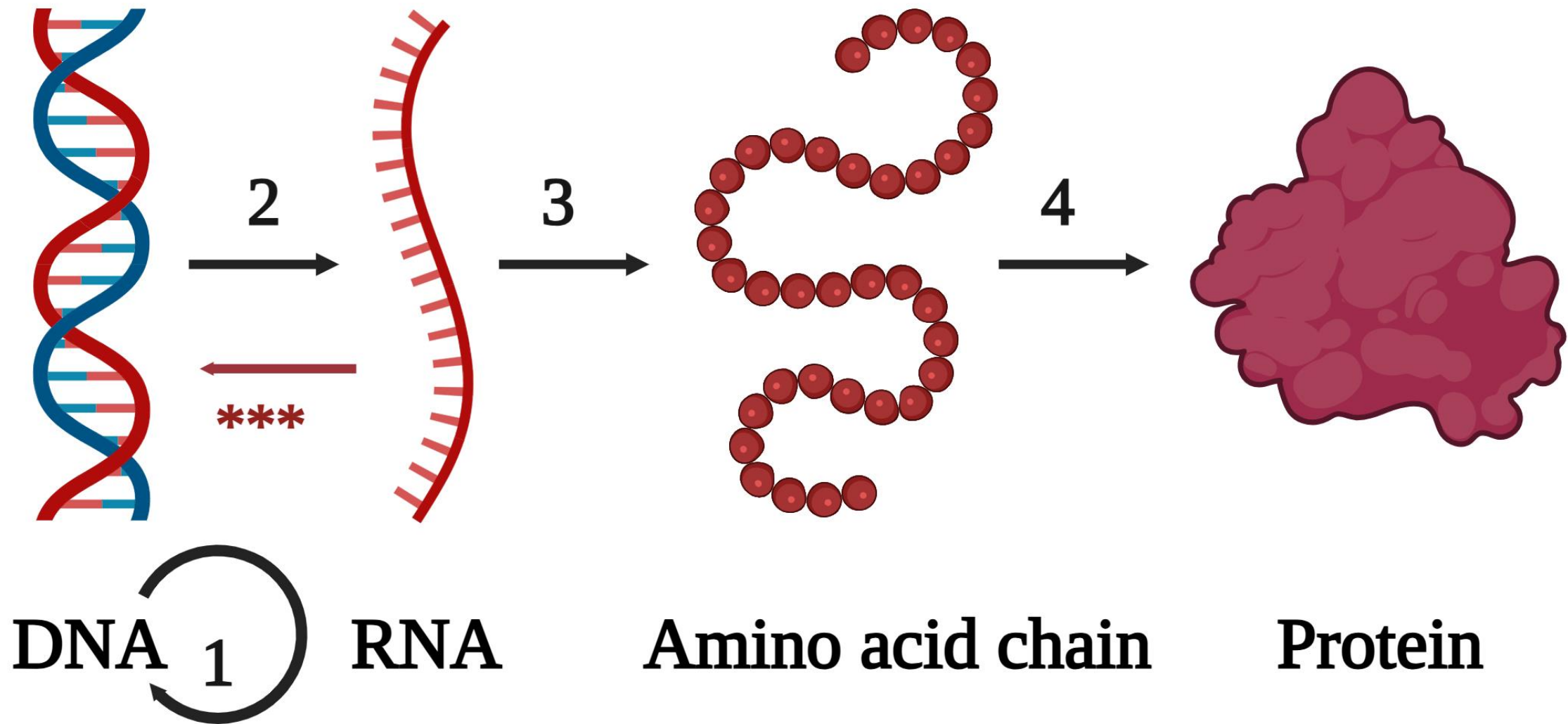
**Dr. Pankaj Singh**

**School of Agriculture and Environment**

**University of Western Australia**

THE UNIVERSITY OF
WESTERN
AUSTRALIA

# Central Dogma

**2**

**3**

**4**

DNA **1** RNA

Amino acid chain

Protein

1 DNA replication
2 Transcription
3 Translation
4 Protein folding

*** Reverse transcription

**Population** –  group of individuals <u>of the same species</u> living in the same area, potentially interacting

**Community** – group of **populations** of <u>different species</u> living in the same area, <u>potentially interacting</u>

**What are some ecological interactions?**

# Why are ecological interactions important?

Ecological interaction can shape up distribution/ diversity/abundance of any organism within a population

# Diversity, richness and evenness

❑ Diversity indicates like how many different type of species are in present within a community

Alpha diversity--diversity on a local scale, describing the species diversity (richness) within a functional community

Beta diversity--describes the rate at which species composition changes across a region

❑ Richness quantifies how many species does a population contains.

❑ Evenness refers to how closer the total number of individual species are present within a population. Lower the dominance of individual species better will the evenness of any populations
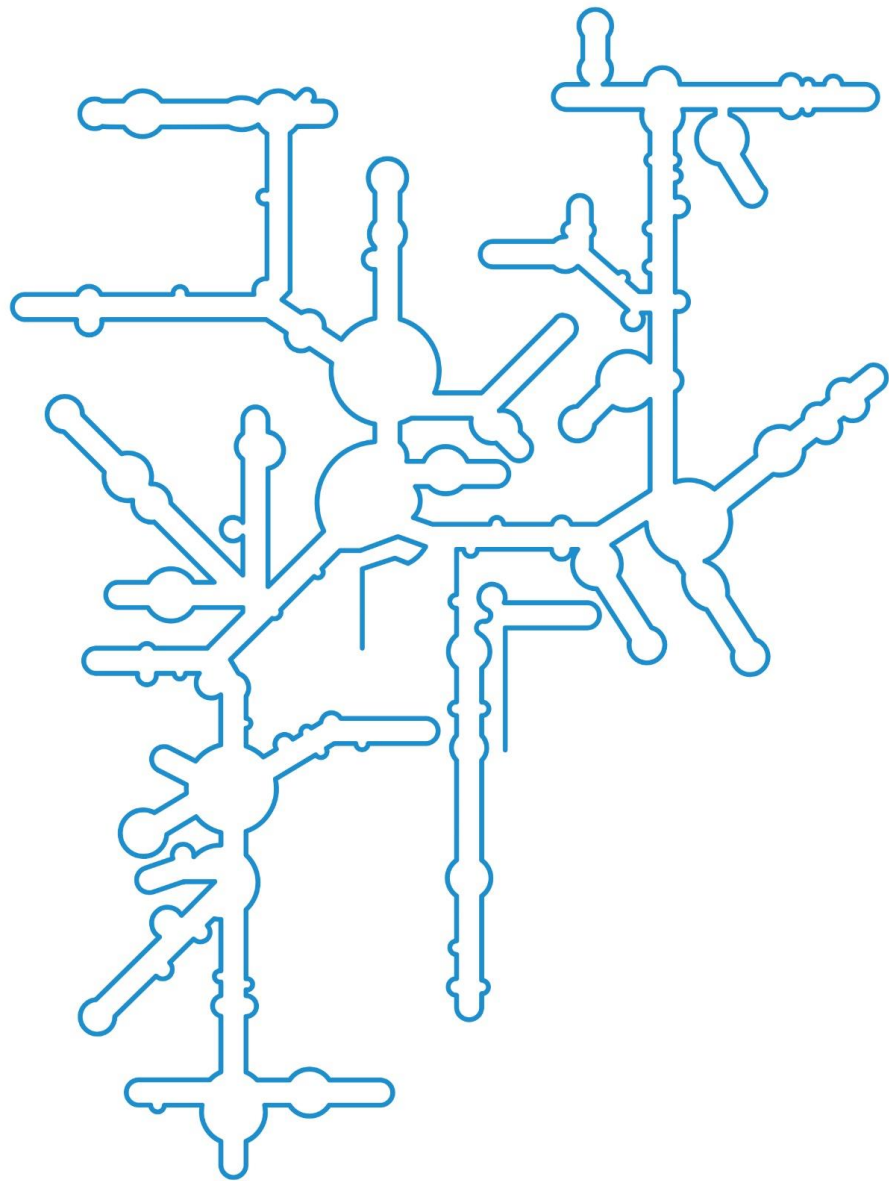
# Symbiosis and co-evolution

❑ Researchers have challenged Darwinism on the basis of theory of symbiosis and co-evolution, **Lynn Margulis** was one of them.

❑ Its based upon interaction of two species (components) and their evolution for countering each other or for existence in a symbiotic way.

❑ This contradicts the Darwinian view that evolution occurs mainly as a result of competition between species.

❑ The organisms form a symbiotic partnership, typically by one engulfing the other– a process known as endosymbiosis. Dramatic evolutionary changes result.
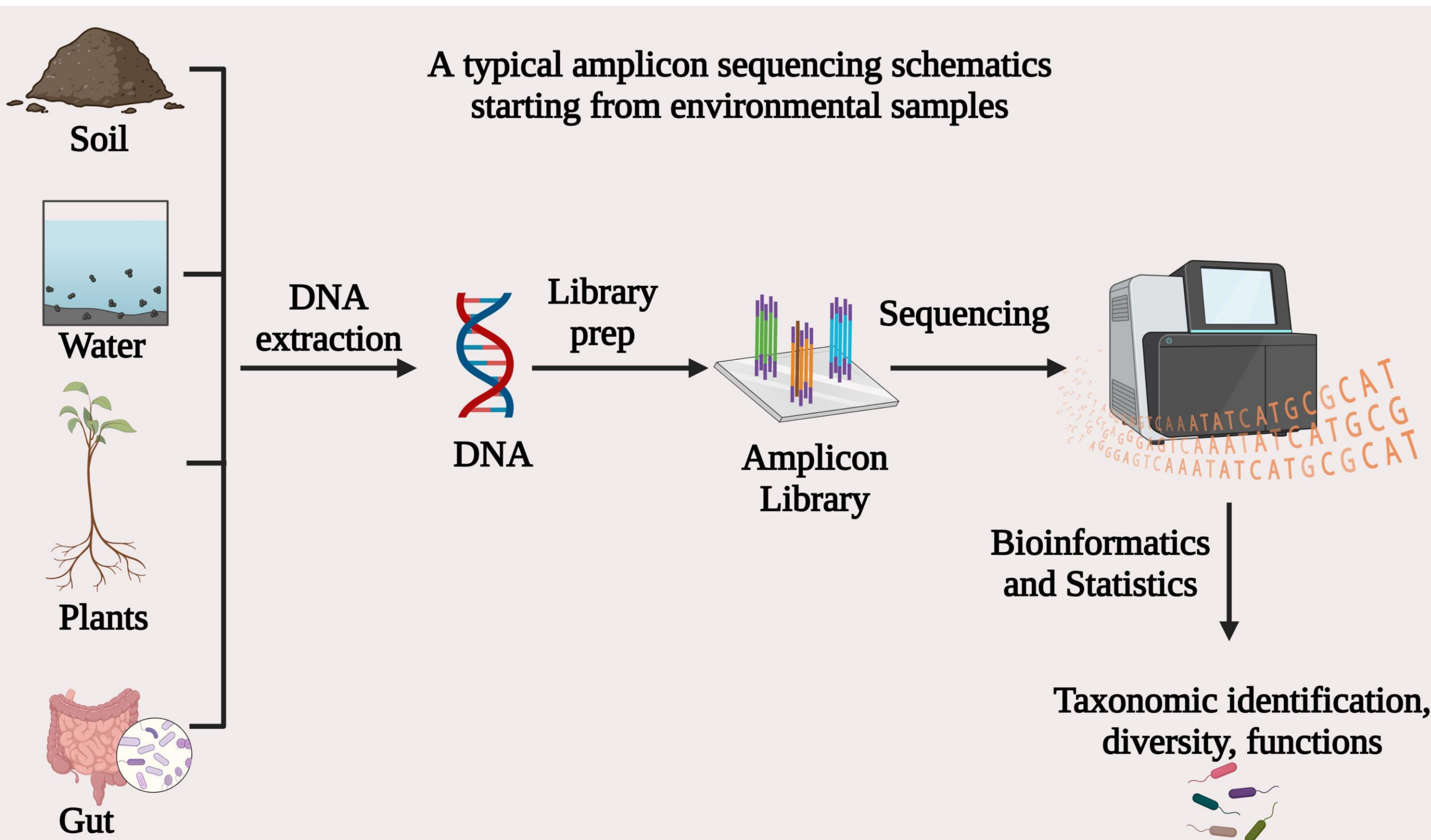
# Why is amplicon sequencing exciting ?

➢ Only about 1-2 % of microbes can be cultured using conventional laboratory practices

➢ That means majority of microbial flora remains unidentified. Hence their role and functions remain unresolved .

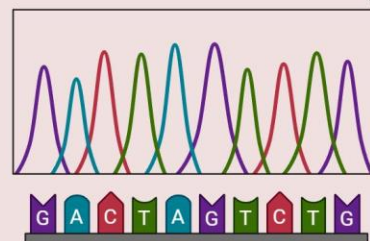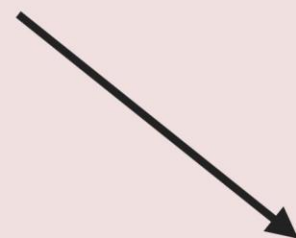➢ NGS techniques allow the reading of DNA from uncloned samples.

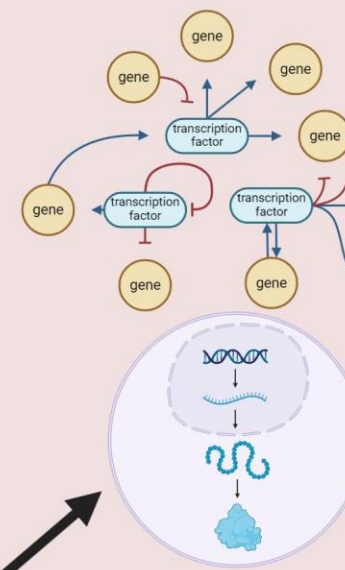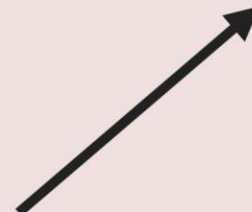# What to use for bacterial identification

16S rRNA a molecular barcode

- Universal
- Undergone less mutation
- Horizontal gene transfer is not an issue
- Conserved and has multiple variable regions for targeted amplification (V1 to V9)
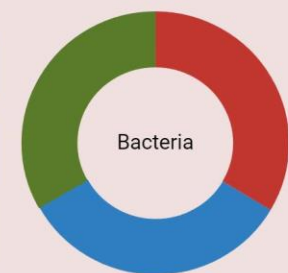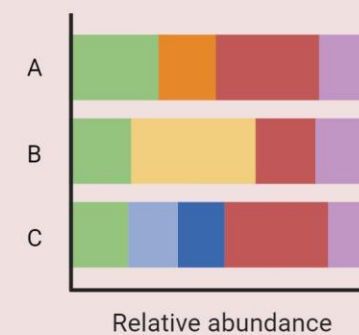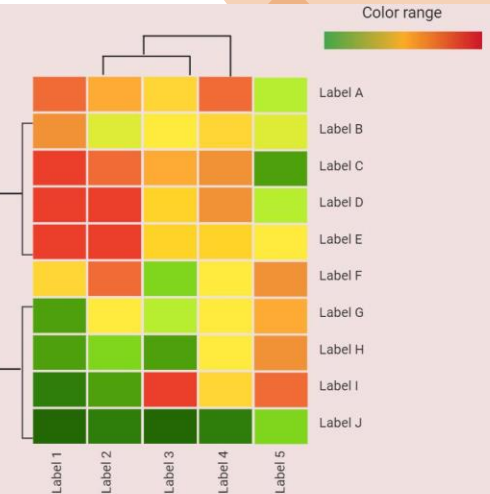
A typical amplicon sequencing schematics starting from environmental samples

Soil

Water

Plants

Gut

DNA extraction

DNA

Library prep

Amplicon Library

Sequencing

Bioinformatics and Statistics

Taxonomic identification, diversity, functions

Color range

Label A
Label B
Label C
Label D
Label E
Label F
Label G
Label H
Label I
Label J

Label 1 Label 2 Label 3 Label 4 Label 5

gene

transcription factor

Functional analysis
(Functional diversity,
metabolic pathways, MAGs)

GACTAGTCTG

High quality reads

A
B
C

Relative abundance

Bacteria

Taxonomic analysis
(Composition, diversity,
differential abundance
analysis)

# File formats in NGS

- ❑ SRF
- ❑ HDF5

- ❑ FASTQ
- ❑ FASTA

- ➢ Helicos
- ➢ PacBio, Applied Biosystems, Oxford Nanophore

- ➢ Most common(Illumina and others)
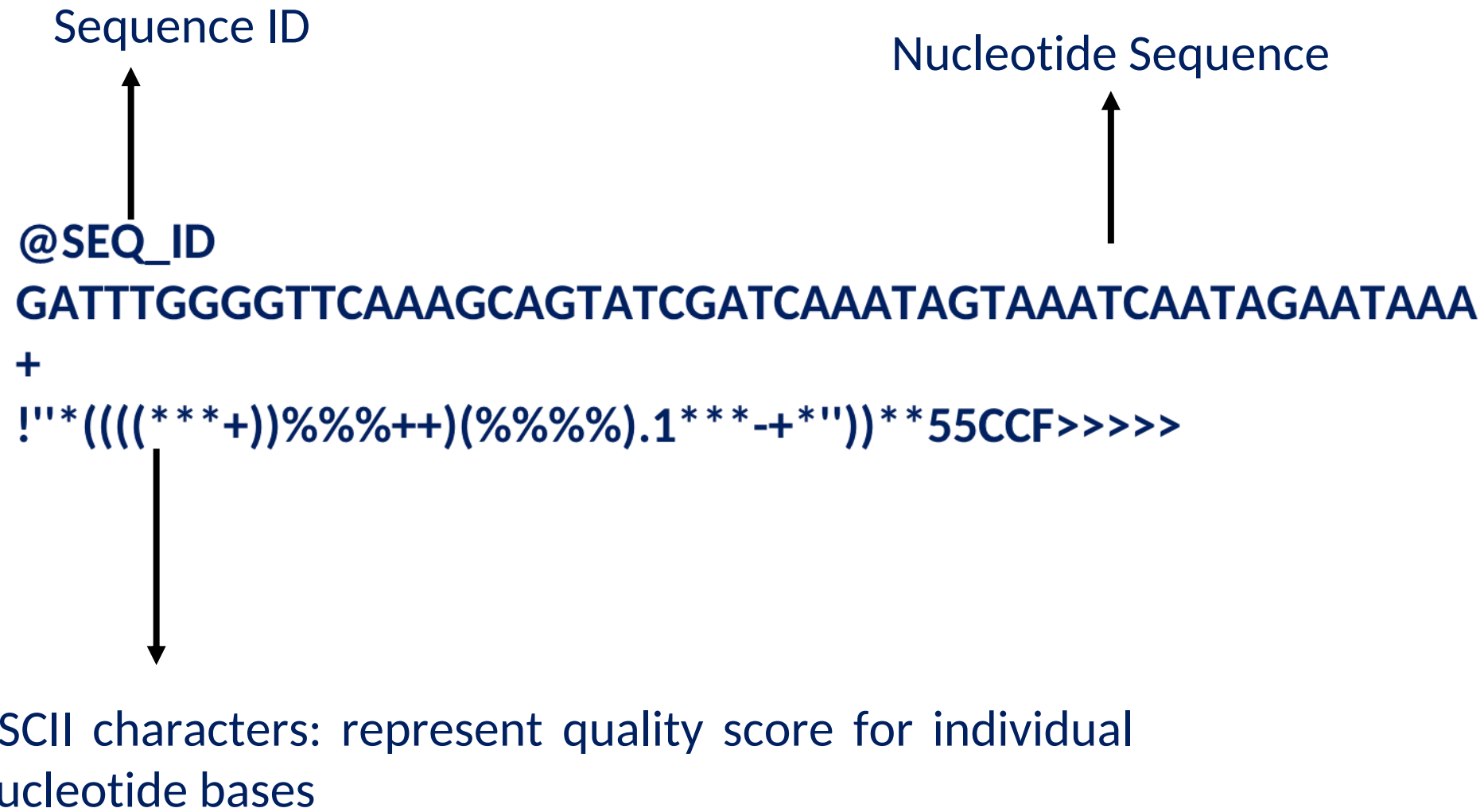- ➢ No quality information

# fastq sequence format

Sequence ID

Nucleotide Sequence

**@SEQ_ID**
**GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCAATAGAATAAA**
**+**
**!''*((((***+))%%%++)(%%%%).1***-+*''))**55CCF>>>>>**

ASCII characters: represent quality score for individual
nucleotide bases

# Phred quality (Q) Scores

❑ Quality (Q) scores represent probability of erroneous of a base call

For e.g. $Q = -10log_{10}P \longrightarrow P = 10^{-Q/10}$

Q= 20 means error probability of $P = 10^{-2}$ = 1 in 100

Q= 30 means error probability of $P = 10^{-3}$ = 1 in 1000

❑ Better the Q Score lesser the chances of error better will be the data quality

❑ Q Score ranges from 0-93

# FASTA FILE FORMAT

Just like @ indicates start of
new sequence

Sequence nameC

>ERR010482.1 FT9FZH301B6YPS/3
ATCAACACATTAGGACTTACACGAATCAGGCATTCGTTACCAT
CAGTATGTCGAT

>ERR010482.2 FT9FZH301ARSRC/3
ATGCTTGCTCGGCCGACGTGAGCGTTATTCGAGCAGGGCTCG
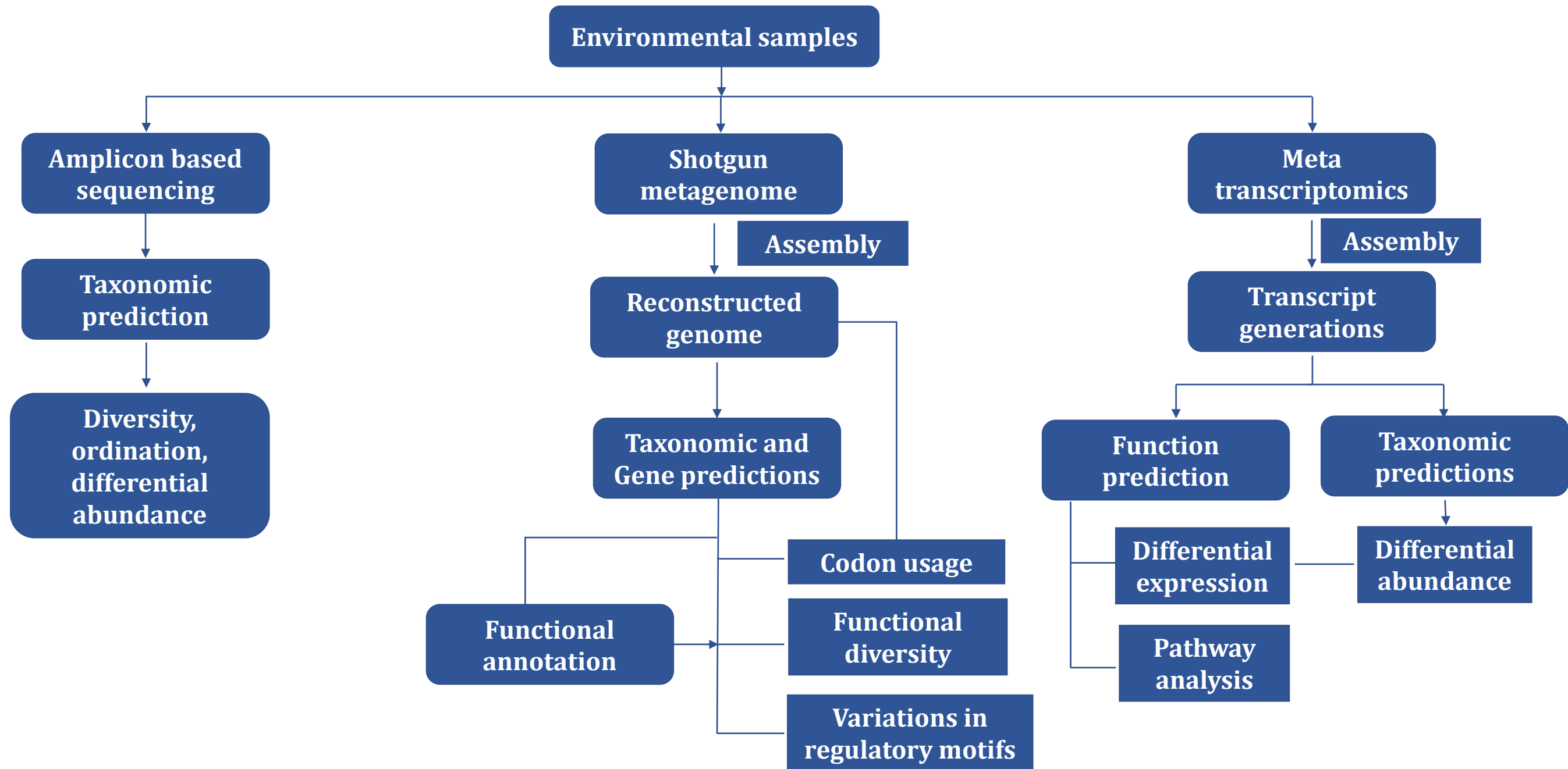GATGGTAGTTAGCGATCCAAAGGGGAGTC

*There is no quality information*

# Quality Control

❑ Quality control is an important step in data processing

❑ Low quality reads below a defined threshold are removed from the dataset

❑ Only high-quality reads are processed further for bioinformatics

❑ General threshold is defined as Q=30, but it varies according to the source of samples and sequencing results

Garbage in

Garbage out

Quality of the data matters the most

# Amplicon sequencing, metagenomics and meta-transcriptomics

**Environmental samples**

**Amplicon based sequencing**

**Shotgun metagenome**

Assembly

**Meta transcriptomics**

Assembly

**Taxonomic prediction**

**Reconstructed genome**

**Transcript generations**

**Diversity, ordination, differential abundance**

**Taxonomic and Gene predictions**

**Function prediction**

**Taxonomic predictions**

**Codon usage**

**Differential expression**

**Differential abundance**

**Functional annotation**

**Functional diversity**

**Pathway analysis**

**Variations in regulatory motifs**

# Taxonomic identification of processed reads

❑ The general rational behind taxonomic classification of sequences is based on the sequence similarity /homology

❑ Prokaryotes and eukaryotes are classified on different molecular barcodes such as 5S rRNA, 16S rRNA (Prokaryotes) 18SrRNA ,23S rRNA ITS (Internal Transcribed Spacer)

❑ Taxonomic annotation is performed against reliable databases. For eg
Bacteria : SILVA, RDP
Fungi : UNITE
Specific databases based on samples: Anaerobic digestors : MIDAS database
https://www.midasfieldguide.org/guide

Molecular barcode and proper primer selection for sequencing is very critical.

# Databases for taxonomic annotation

https://www.arb-silva.de/



https://rnacentral.org/expert-database/rdp



p.s. Greengenes has not been updated for longtime avoid using it

# Common bioinformatics software for amplicon sequencing data analysis

❑ QIIME2 : Quantitative Insights Into Microbial Ecology 2
   One of the popular tools
   Unix/Linux depend, cannot use straightway on Windows
   Need Docker or Virtual Machine on Windows
   Good for visualizations
   Requires good knowledge of Bash scripting and python
❑ USEARCH
   Unix/Linux depend, cannot use straightway on Windows
   Need Docker or Virtual Machine on Windows
   Need good skill of Bash scripting
❑ Mothur
   Platform independent
   Light and works on Windows
   Does not require installation

# What do we get out of bioinformatics analysis

❑ OTU  file : Classifying sequenced reads into Amplicon Sequencing Variants(ASVs) or Operational Taxonomic Units

❑ Information regarding how many time every sequence in the form of OTU/ASV has appeared-----this relates to abundance

ASV number

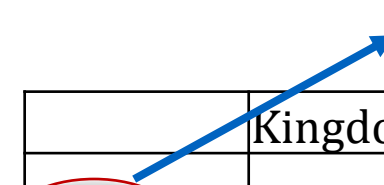| | ASV1 | ASV2 | ASV3 | ASV4 | ASV5 | ASV6 | ASV7 | ASV8 | ASV9 |
|-----|------|------|------|------|------|------|------|------|------|
| S01 | 1020 | 1544 | 325 | 845 | 2100 | 3215 | 2154 | 120 | 0 |
| S02 | 2590 | 454 | 1214 | 21 | 2121 | 785 | 445 | 549 | 423 |
| S03 | 3101 | 021 | 4785 | 196 | 352 | 268 | 124 | 412 | 563 |
| S04 | 3580 | 954 | 12 | 687 | 51 | 0 | 14 | 75 | 945 |
| S05 | 1257 | 758 | 352 | 635 | 487 | 753 | 951 | 852 | 159 |

Abundance

Sample ID

# What do we get out of bioinformatics analysis cont….

Taxonomy file

❑ Taxonomic information for every OTU/ASV corresponding to the OTU file

❑ Taxonomic information starting from Domain to Genus (sometimes species)

ASV number

Taxonomic information of individual ASVs

| | Kingdom | Phylum | Class | Order | Family | Genus |
|------|---------|--------|-------|-------|--------|-------|
| ASV1 | Bacteria | Firmicutes | Bacilli | Bacillales | Bacillaceae | NA |
| ASV2 | Bacteria | Firmicutes | Bacilli | Bacillales | Bacillaceae | NA |
| ASV3 | Bacteria | Firmicutes | Bacilli | Bacillales | Planococcaceae | Sporosarcina |
| ASV4 | Bacteria | Firmicutes | Bacilli | Bacillales | Bacillaceae | NA |
| ASV5 | Bacteria | Firmicutes | Bacilli | Bacillales | Bacillaceae | NA |
| ASV6 | Bacteria | Firmicutes | Bacilli | Bacillales | Bacillaceae | Bacillus |
| ASV7 | Bacteria | Proteobacteria | Alphaproteobacteria | Sphingomonadales | Sphingomonadaceae | Sphingomonas |

# Most important file : Metadata file

❏ Most important part, designed by user
❏ Has all variable and factors from the experiment
❏ Can include plant, soil parameters, etc.

| SampleID | SampleType | Plant | Stress | Timepoint |
|----------|------------|-------|--------|-----------|
| S01 | Bulk Soil | Ryegrass | Control | T1 |
| S02 | Rhizosphere | Ryegrass | Heat | T1 |
| S03 | Root | Ryegrass | Control | T2 |
| S04 | Root | Lucerne | Control | T2 |

# What basic stats to do post analysis

❑ Taxonomic composition of dataset

❑ Alpha diversity, richness indices based on metadata(experimental factors and variables)

❑ Beta diversity : Ordination plots

❑ Responsive ASVs/bacteria to factors : Indicator species analysis
Differential abundance analysis

❑ Statistical modelling

❑ Application of machine learning : Random Forest

❑ Many more

*All depends on your dataset and hypothesis*