

# Metagenomics

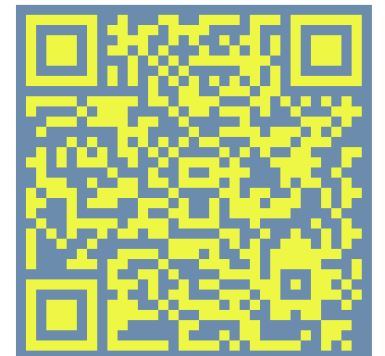
## Lecture 3

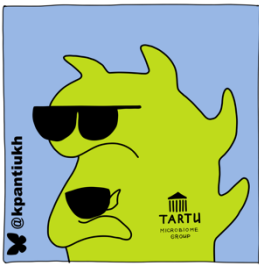
MWAS – METAGENOME WIDE  
ASSOCIATION STUDY

Kateryna Pantiukh

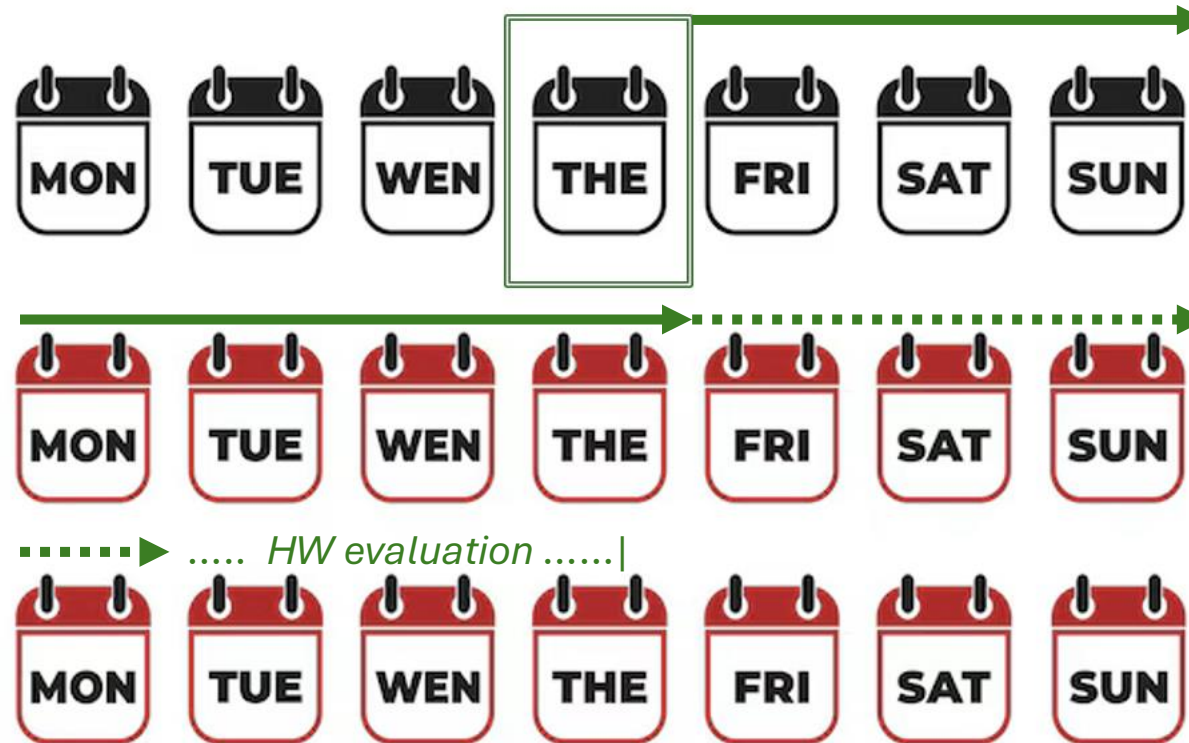
pantiukh@ut.ee

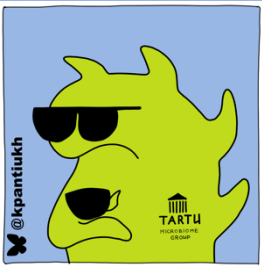
GitHub





# Revision of homework deadline





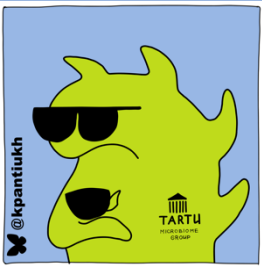
# Microbiome



A community of  
**microorganisms** that lives in  
a specific environment

- *Bacteria*
  - Primar degrades
- *Archaea*
  - Primary fermenters
- *Viruses*
  - Secondary fermenters
- *Microeucarites*
  - Sinks





# Microbiome



A community of  
**microorganisms** that lives in  
a specific environment

- *Bacteria*
  - Primar degrades
- *Archaea*
  - Primary fermenters
- *Viruses*
  - Secondary fermenters
- *Microeucarites*
  - Sinks

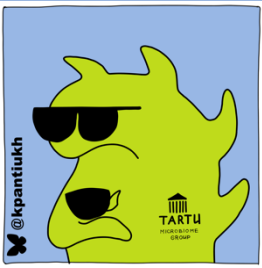
Community may have different level of complexity

... may be evaluated with  
Community **diversity indexes**

alpha-, beta-, gamma-







# Microbiome



A community of  
**microorganisms** that lives in  
a specific environment

- *Bacteria*
  - Primar degrades
- *Archaea*
  - Primary fermenters
- *Viruses*
  - Secondary fermenters
- *Microeucarites*
  - Sinks

Community may have different level of complexity

... may be evaluated with  
Community **diversity indexes**

alpha-, beta-, gamma-

... the first and simplest measures that can  
reveal differences between communities





# Microbiome

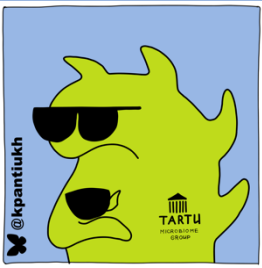


A community of **microorganisms** that lives in a specific environment

- *Bacteria*
- *Archaea*
- *Viruses*
- *Microeucarites*
- Primar degrades
- Primary fermenters
- Secondary fermenters
- Sinks

Associations between **community complexity** and the feature of interest





# Microbiome

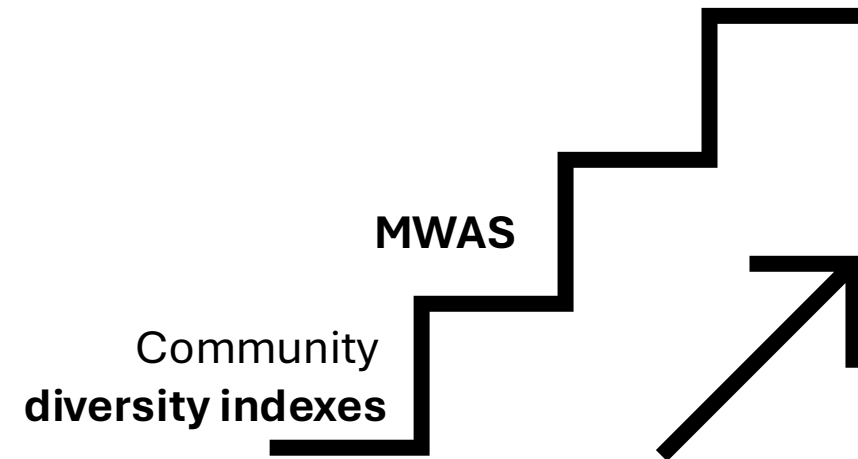


A community of  
**microorganisms** that lives in  
a specific environment

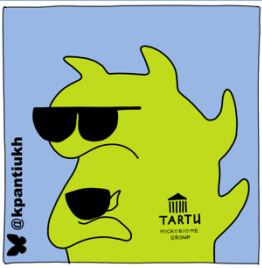
- *Bacteria*
- *Archaea*
- *Viruses*
- *Microeucarites*
- Primar degrades
- Primary fermenters
- Secondary fermenters
- Sinks

**MWAS** – microbiome wide association study

Associations between ***specific species***  
and the feature of interest







# MWAS

**MWAS** – microbiome wide association study

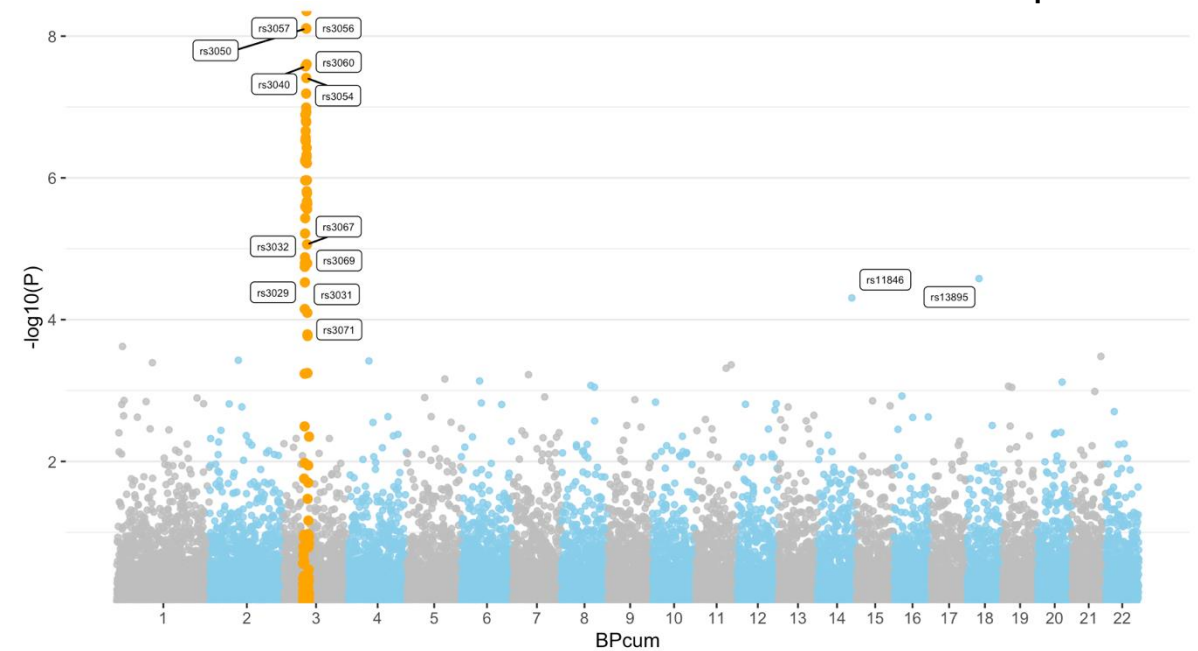
**GWAS** – genome wide association study

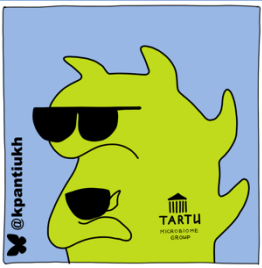
Associations between **specific genome variations** (SNP or indel) and the feature of interest

Genotyping



Manhattan plot

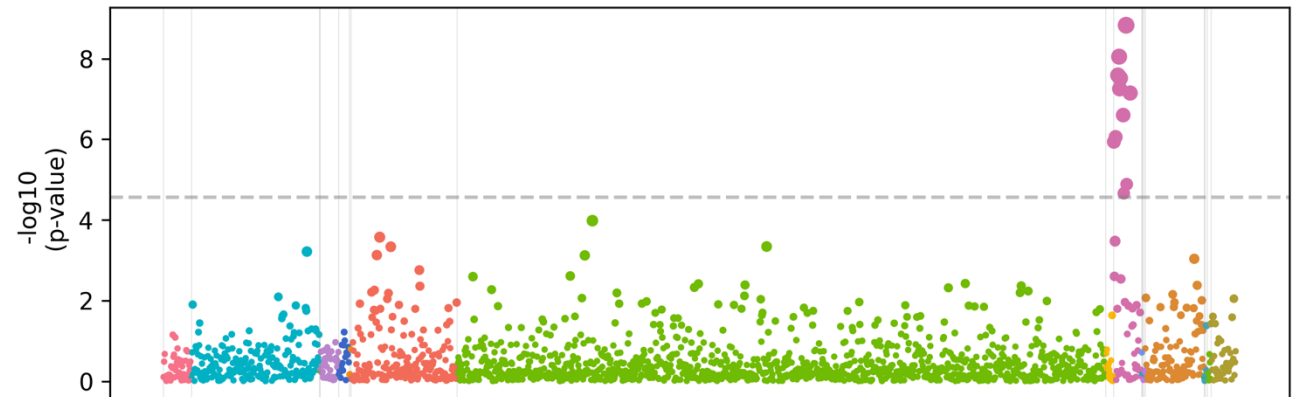




# MWAS

MWAS – microbiome wide association study

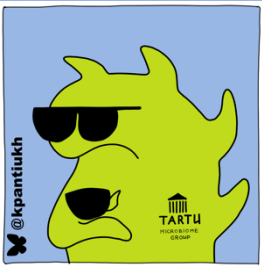
Bacteria sp. instead of SNP



- |                           |                           |                           |
|---------------------------|---------------------------|---------------------------|
| ● <i>Actinobacteriota</i> | ● <i>Bacteroidota</i>     | ● <i>Proteobacteria</i>   |
| ● <i>Bacillota</i>        | ● <i>Campylobacterota</i> | ● <i>Thermoplasmata</i>   |
| ● <i>Bacillota_A</i>      | ● <i>Cyanobacteriota</i>  | ● <i>Verucomicrobiota</i> |
| ● <i>Bacillota_B</i>      | ● <i>Desulfobacterota</i> | ● <i>Other phyla</i>      |
| ● <i>Bacillota_C</i>      | ● <i>Patescibacteria</i>  |                           |

Manhattan plot





# MWAS

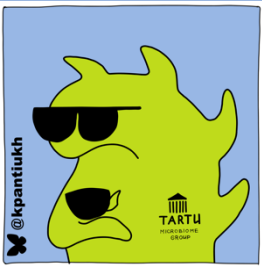
MWAS – microbiome wide association study

$$Y_i = \beta_0 + \beta_1 X_i + \sum_{j=2}^p \beta_j C_{ij} + \epsilon_i$$

Where:

- $Y_i$  is the phenotype (e.g., disease status, quantitative trait) for sample  $i$
- $X_i$  is the abundance (or presence/absence) of a microbial feature in sample  $i$
- $C_{ij}$  are covariates (age, sex, sequencing batch, etc.)
- $\beta_0$  is the intercept
- $\beta_1$  is the effect size of the microbial feature
- $\epsilon_i$  is the residual error





# MWAS

MWAS – microbiome wide association study



Heart  
Disease  
status




Intersept  
coefficient



$\Sigma$ (effect of species  $\times$   
abundance of species)



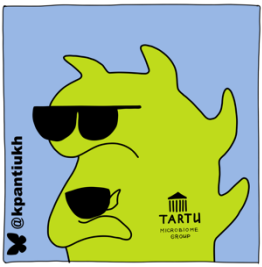
{ Age, sex, BMI,  
Stool type  }

$\Sigma$ (effect of covariates  
 $\times$  covariates)



error





# MWAS

MWAS – microbiome wide association study



Heart  
Disease  
status

=

Intersept  
coefficient

+

$\Sigma(\text{effect of species} \times$   
abundance of species)

+

$\Sigma(\text{effect of covariates}$   
 $\times$  covariates)

+

error



Heart  
Desease  
status

=

Intersept  
coefficient

+

$\Sigma(\text{effect of species} \times$   
abundance of species)

+

$\Sigma(\text{effect of covariates}$   
 $\times$  covariates)

+

error



Heart  
Desease  
status

=

Intersept  
coefficient

+

$\Sigma(\text{effect of species} \times$   
abundance of species)

+

$\Sigma(\text{effect of covariates}$   
 $\times$  covariates)

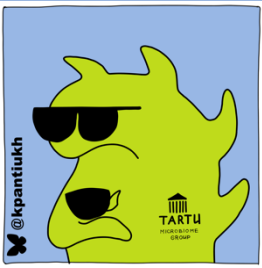
+

error



!! Correction for  
multiple testing





# MWAS

MWAS – microbiome wide association study



!! Correction for  
multiple testing

## Bonferoni correction

Significance level / number of tests

$$0.05 / 3 = 0.0166$$

Corrected significance level = 0,0166





# MWAS

MWAS – microbiome wide association study



Heart  
Disease  
status



Intersept  
coefficient



$\Sigma(\text{effect of species} \times \text{abundance of species})$



$\left\{ \begin{array}{l} \text{Age, sex, BMI,} \\ \text{Stool type} \end{array} \right\}$

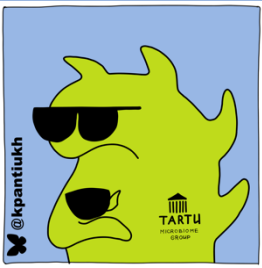
$\Sigma(\text{effect of covariates} \times \text{covariates})$



error

Main outcome: p-value and beta coefficient



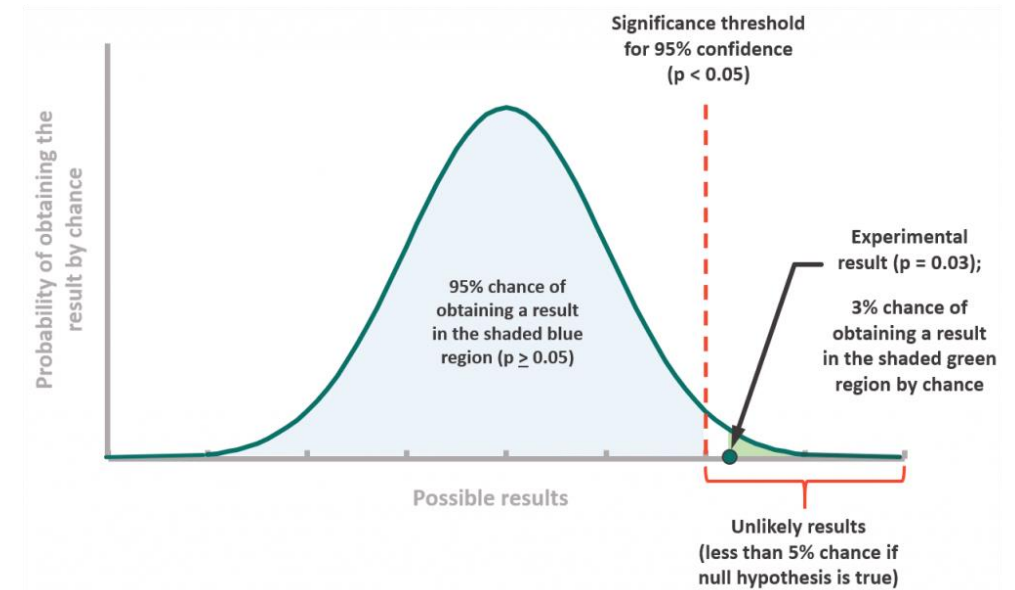


# p-value

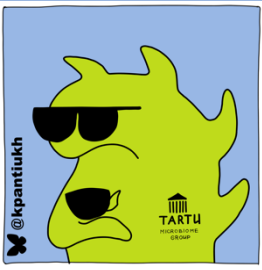
MWAS – microbiome wide association study

**null hypothesis:** the true effect is zero

The p-value is the probability of observing a result **at least as extreme as the one you got**, assuming that null hypothesis is true.







# p-value

MWAS – microbiome wide association study

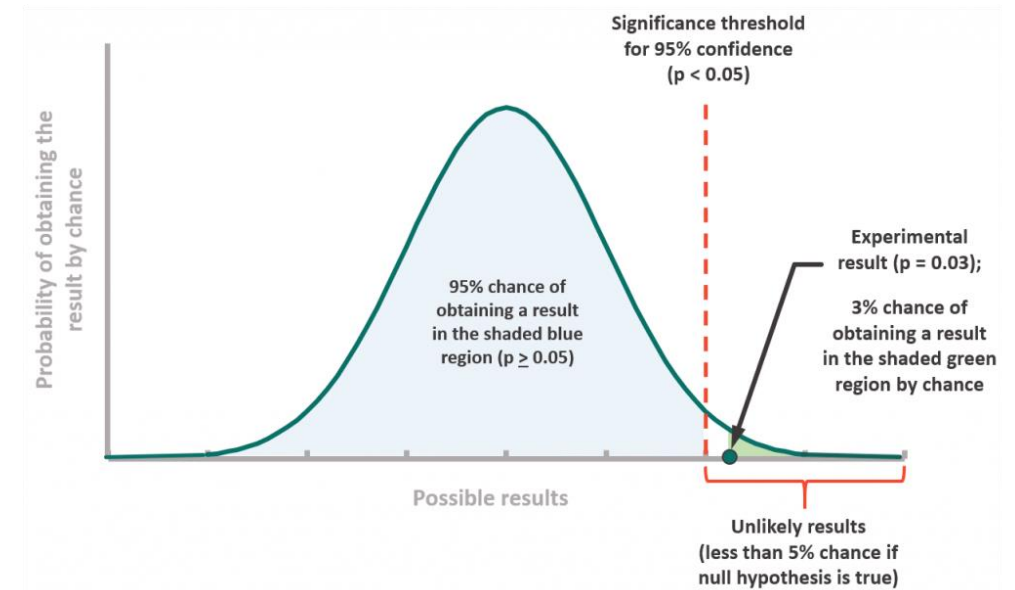
What it tells you:

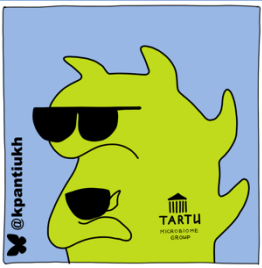
- A **small p-value** means your result would be unlikely if the effect were truly zero.
- A **large p-value** means your data are quite compatible with no effect.

What it does *not* tell you:

- It is **not** a measure of effect size or importance!!!

\* - Effect size - beta





# MWAS

MWAS – microbiome wide association study

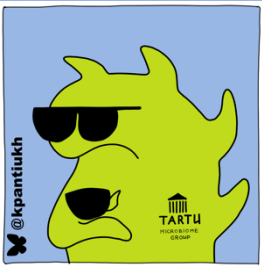


If  $p\text{-value} < \text{level of significance}$   
We consider the association  
Statistically significant

- Positive beta – positive correlation
- Negative beta – negative correlation

pheno	name	bacteria	p-value	beta
N97	Female infertility	A0002_Methanobrevibacter_A_smithii_A	8,53E-06	0,00073
N97	Female infertility	H0023_Alistipes_communis	2,84E-07	0,003378
K21	Gastro-esophageal reflux disease	H0092_CAG-41_sp900066215	2,09E-06	0,000291
M13	Other arthritis	H0117_Scotosoma_sp900555925	9,53E-07	0,001968
H40	Glaucoma	H0220_Ruminiclostridium_E_sp900539195	7,09E-06	0,001756
G43	Migraine	H0224_Merdimorpha_sp002314265	1,43E-06	0,000376
I48	Atrial fibrillation and flutter	H0237_Dysosmobacter_welbionis	1,58E-06	0,000956
M13	Other arthritis	H0262_UMGS692_sp900544545	8,77E-06	0,006968
F41	Other anxiety disorders	H0280_Enterocloster_sp000431375	2,32E-06	0,000527





# MWAS

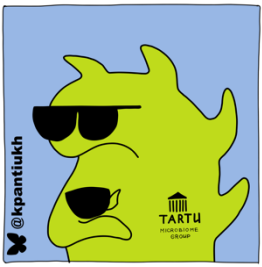
MWAS – microbiome wide association study

level of significance =  $5 * 10^{-6}$

pheno	name	bacteria	p-value	beta
N97	Female infertility	A0002_Methanobrevibacter_A_smithii_A	8,53E-06	0,00073
N97	Female infertility	H0023_Alistipes_communis	2,84E-07	0,003378
K21	Gastro-esophageal reflux disease	H0092_CAG-41_sp900066215	2,09E-06	0,000291
M13	Other arthritis	H0117_Scariosoma_sp900555925	9,53E-07	0,001968
H40	Glaucoma	H0220_Ruminiclostridium_E_sp900539195	7,09E-06	0,001756
G43	Migraine	H0224_Merdimorpha_sp002314265	1,43E-06	0,000376
I48	Atrial fibrillation and flutter	H0237_Dysosmobacter_welbionis	1,58E-06	0,000956
M13	Other arthritis	H0262_UMGS692_sp900544545	8,77E-06	0,006968
F41	Other anxiety disorders	H0280_Enterocloster_sp000431375	2,32E-06	0,000527

*Alistipes putredinis* is positively associated with female infertility status, indicating higher relative abundance in diagnosed individuals compared with controls





# MWAS

MWAS – microbiome wide association study

level of significance =  $5 * 10^{-6}$

pheno	name	bacteria	p-value	beta
N97	Female infertility	A0002_Methanobrevibacter_A_smithii_A	8,53E-06	0,00073
N97	Female infertility	H0023_Alistipes_communis	2,84E-07	0,003378
K21	Gastro-esophageal reflux disease	H0092_CAG-41_sp900066215	2,09E-06	0,000291
M13	Other arthritis	H0117_Scariosoma_sp900555925	9,53E-07	0,001968
H40	Glaucoma	H0220_Ruminiclostridium_E_sp900539195	7,09E-06	0,001756
G43	Migraine	H0224_Merdimorpha_sp002314265	1,43E-06	0,000376
I48	Atrial fibrillation and flutter	H0237_Dysosmobacter_welbionis	1,58E-06	0,000956
M13	Other arthritis	H0262_UMGS692_sp900544545	8,77E-06	0,006968
F41	Other anxiety disorders	H0280_Enterocloster_sp000431375	2,32E-06	0,000527

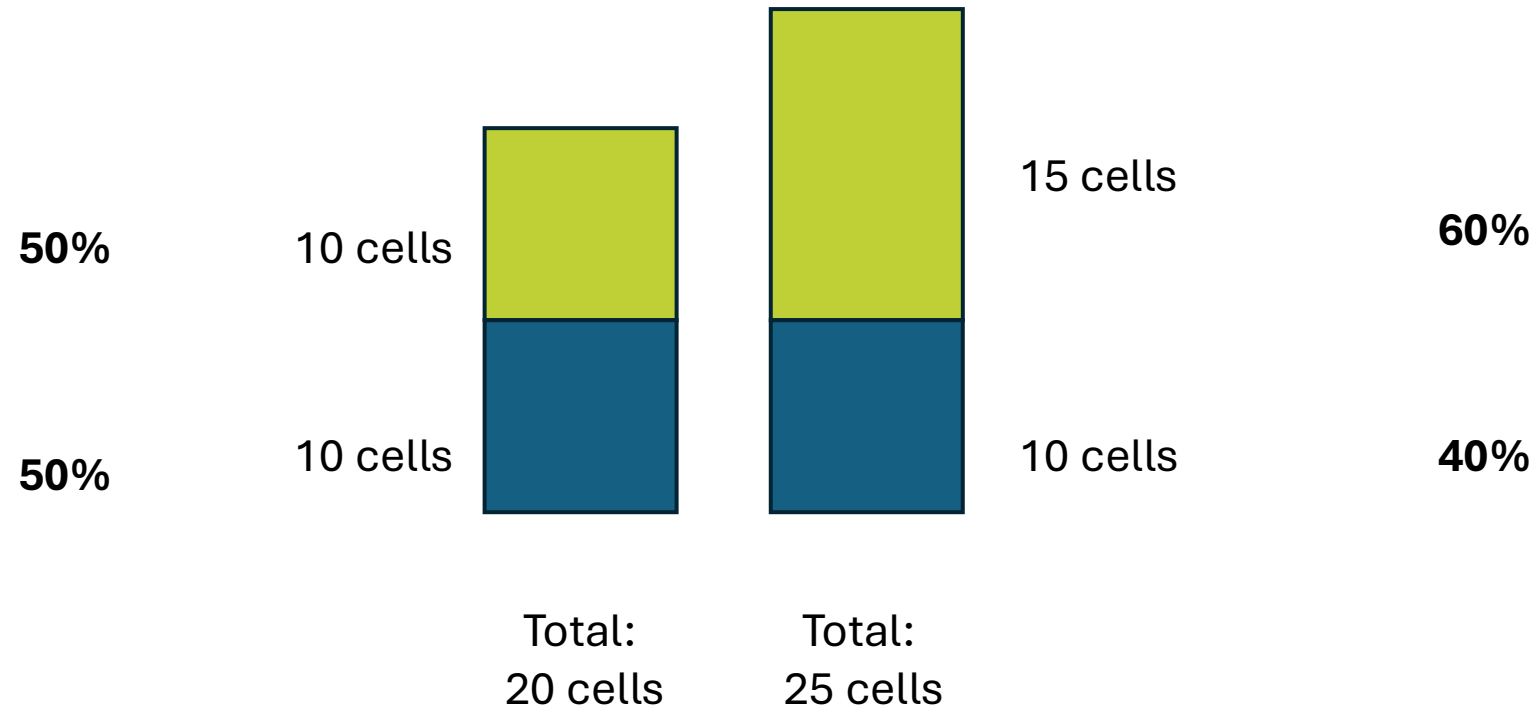
*Alistipes putredinis* is positively associated with female infertility status, indicating higher relative abundance in diagnosed individuals compared with controls, however the estimated effect size is small

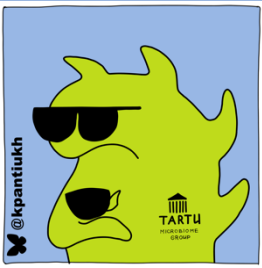




# Abundance metrics

## relative abundance **PROBLEM**





# MWAS

MWAS – microbiome wide association study

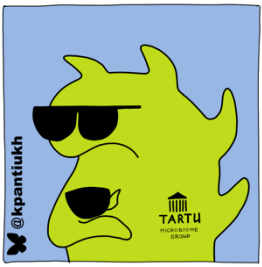
INPUT:  
relative abundance

	Tom	Mary
Species 1	0.1	3.7
Species 2	0.0	0.2
Species 3	2.3	0.0

INPUT:  
presence-absence

	Tom	Mary
Species 1	1	1
Species 2	0	1
Species 3	1	0



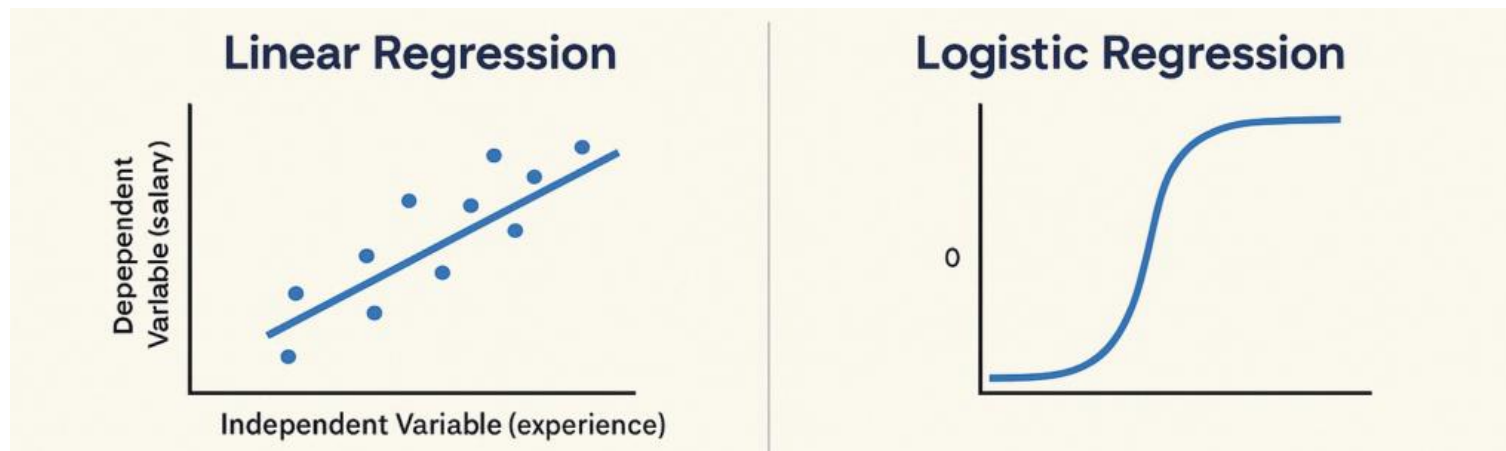


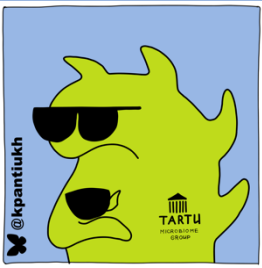
# MWAS

MWAS – microbiome wide association study

INPUT:  
relative abundance

INPUT:  
presence-absence





# Effect size - odds ratios

MWAS – microbiome wide association study

## Example:

- Predictor: Alistipes putredinis present vs absent
- Outcome: Female infertility
- Logistic regression gives OR = 2.0



## Interpretation:

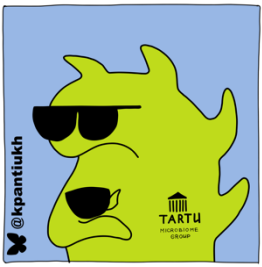
Individuals with Alistipes putredinis present have twice the odds of being diagnosed with female infertility compared to individuals without this bacterium, holding other variables constant.

If OR = 0.5, the interpretation flips:

Individuals with Alistipes have half the odds of infertility compared to those without it.





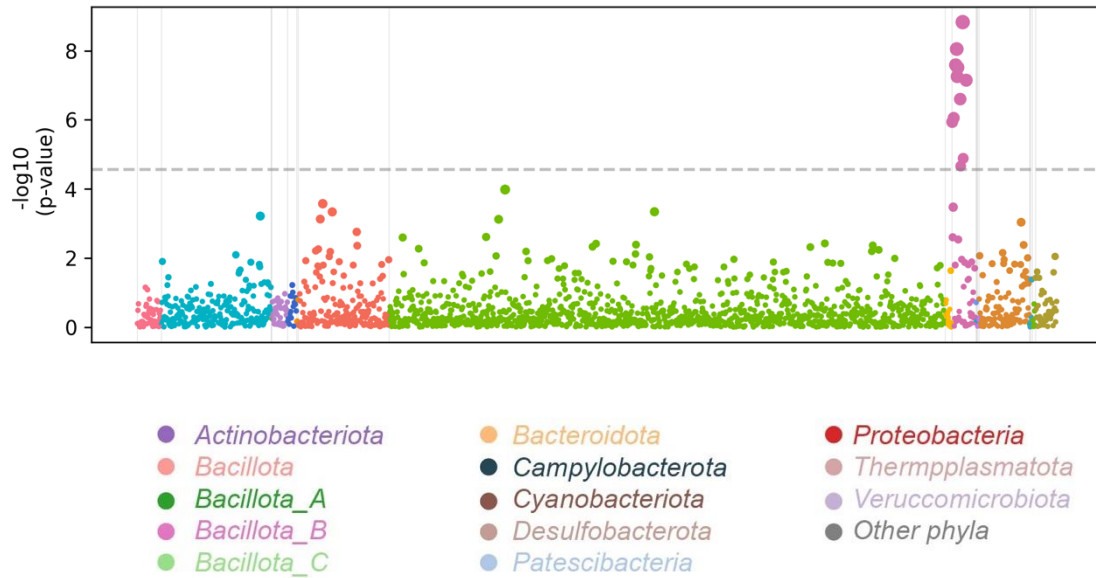


# MWAS

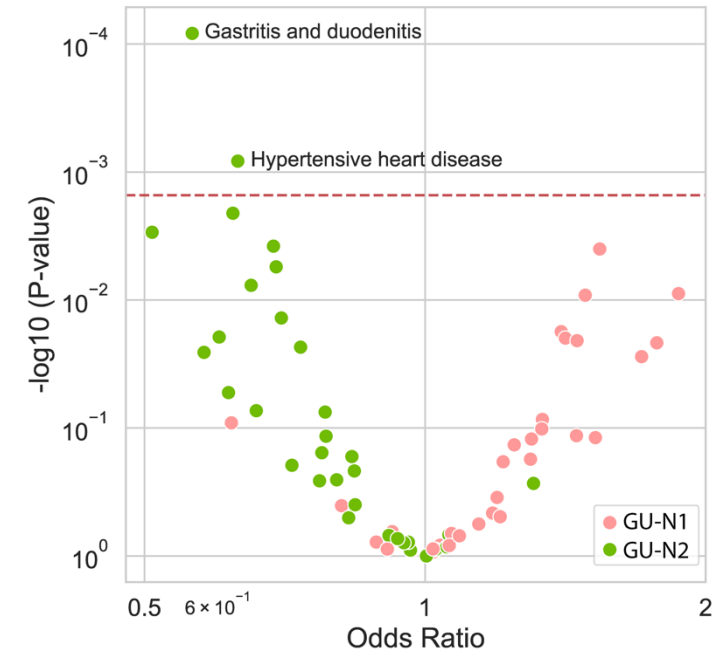
Visualisation example

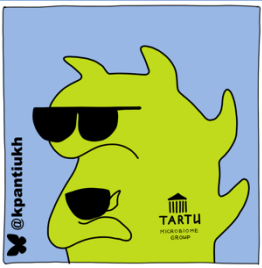
Linear regression  
(relative abundance input)

Manhattan plot



Logistic regression  
(presence-absence input)

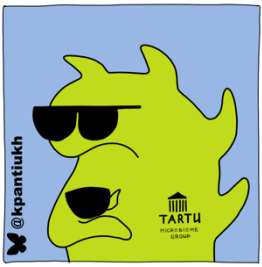




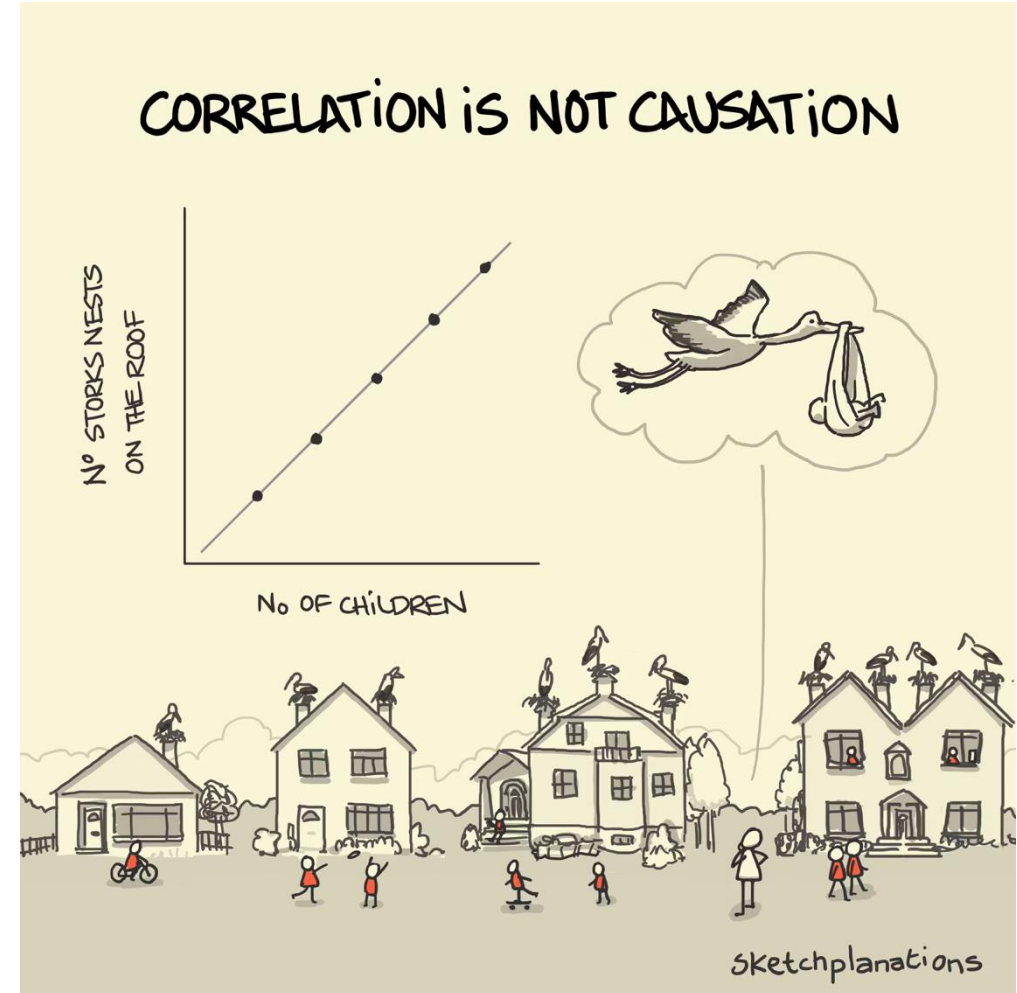
# Causation issue

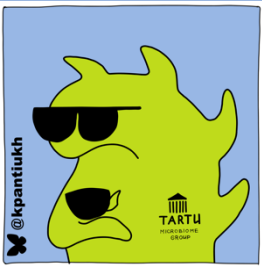
Correlation != Causation



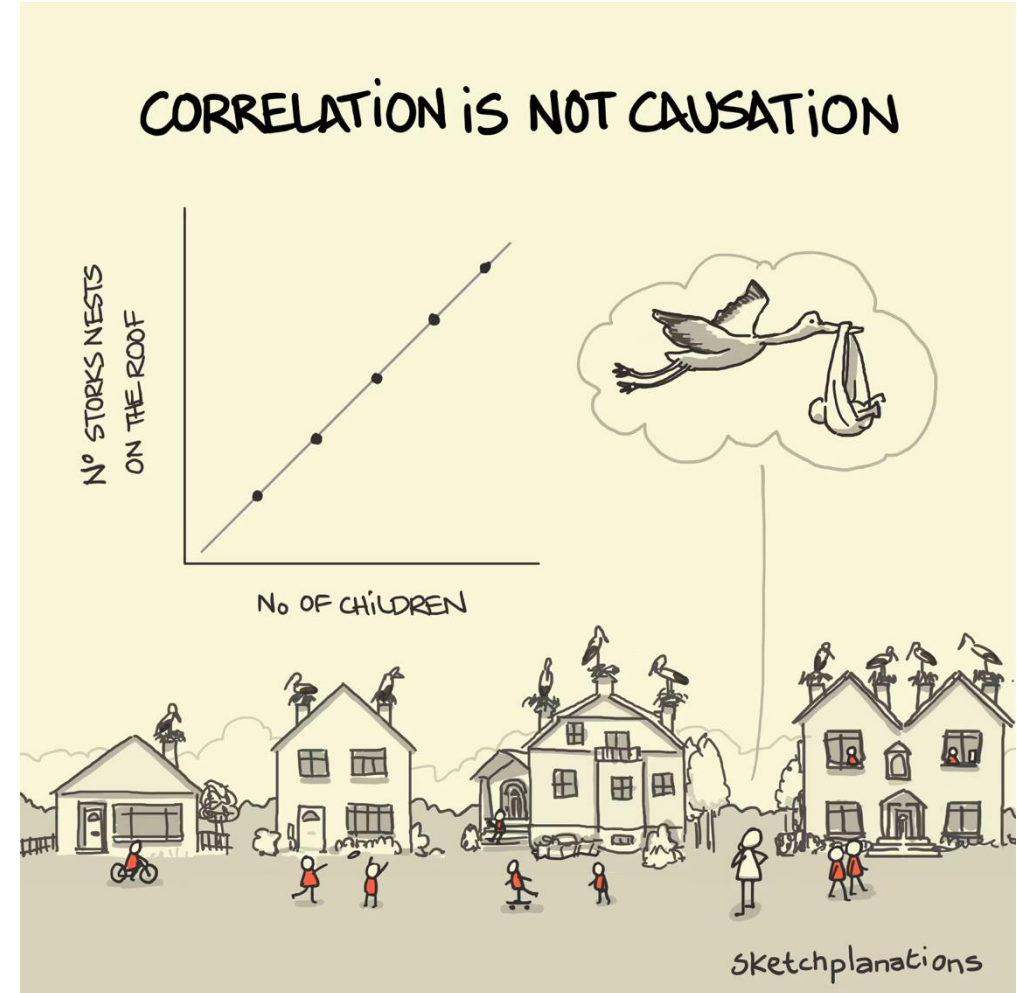
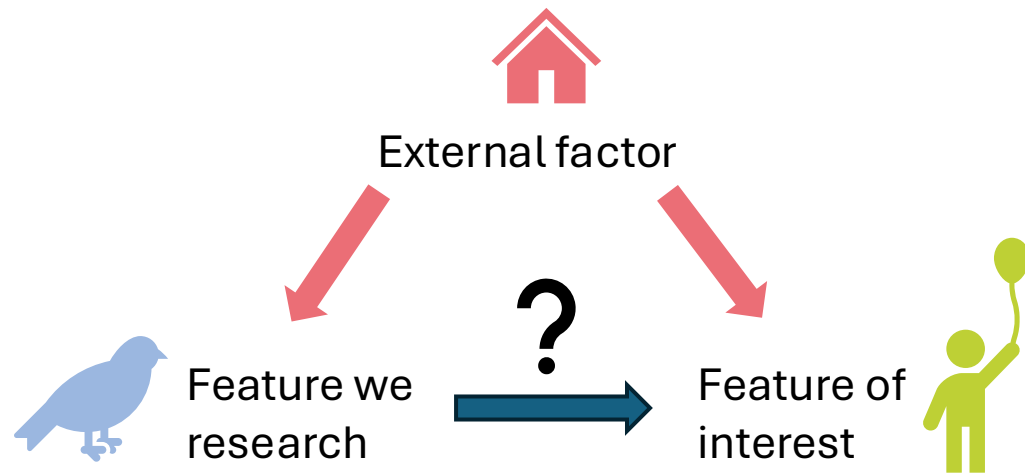


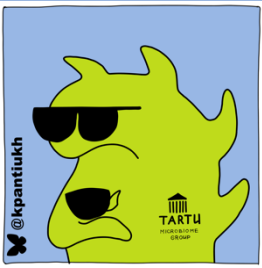
# Causation issue



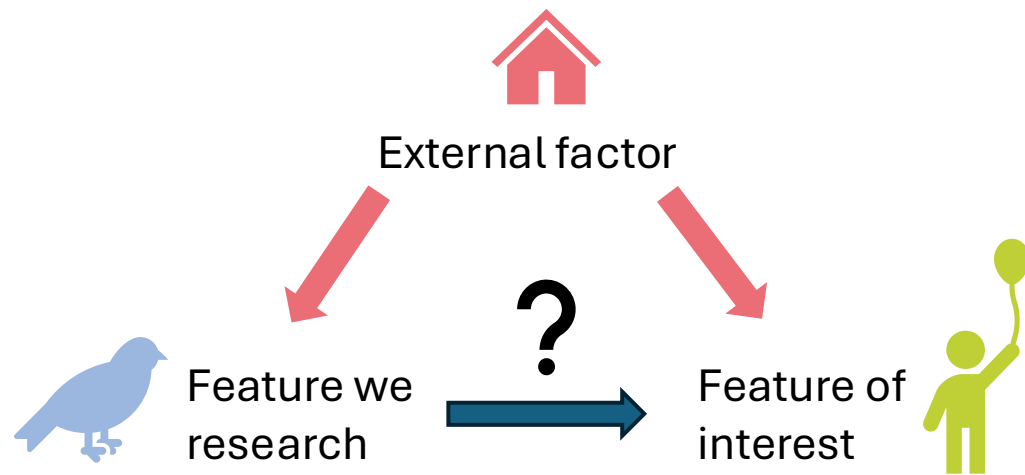


# Causation issue





# Causation issue



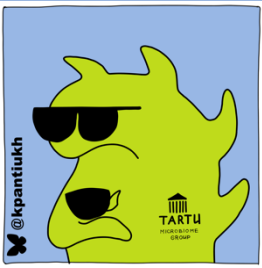
## WHAT WE CAN DO ABOUT IT?

1. Even when causation is questionable, correlation may be used as a predictor

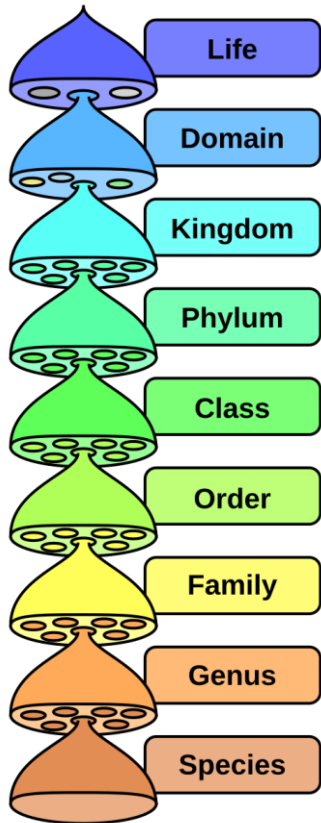
Useful for early diagnosis

2. We can design additional experiments to check causation





# Different level abundance tables



	Tom	Mary
Genus 1	0.1	3.7
Genus 2	0.0	0.2
Genus 3	2.3	0.0

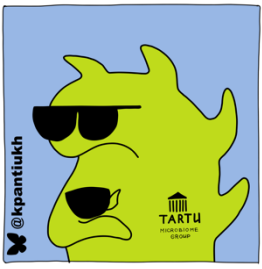
Low dimensionality → fewer multiple-testing problems.

	Tom	Mary
Species 1	0.1	3.7
Species 2	0.0	0.2
Species 3	2.3	0.0

Species & strain level:  
can link associations to  
particular functions or  
pathogenic potential.

BUT: group size matters!





# How to decide what taxonomic level to use?

## 1. Cohort size and statistical power

- species/strain have many rare or zero-count taxa.
- Smaller cohorts may not provide enough observations per taxon to detect associations reliably.
- Broader levels (phylum, class, family) aggregate taxa, increasing counts and power but may hide specific effects.

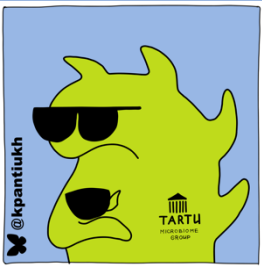
## 2. Expected prevalence of taxa

- Rare taxa are often absent in many samples.
  - Presence/absence models can work for very rare taxa, but effect size estimates become unstable.
  - Focus on taxa that occur in a meaningful fraction of samples (e.g., >10–20% prevalence).
- Tip:** Check prevalence at different levels before deciding; sometimes grouping into higher levels improves coverage.

## 3. Biological interpretability

- Broad levels show general trends (e.g., Bacteroidetes increase) but often lack actionable insights.
- Genus or species level allows linking findings to metabolic pathways, pathogenicity, or prior literature.
- Strain-level associations are most informative for functional or mechanistic hypotheses but require high-resolution sequencing.





# How to decide what taxonomic level to use?

## 4. Multiple testing burden

- Species/strain levels increase the number of tests, requiring stricter p-value correction and reducing statistical power.
- Consider pre-filtering low-abundance taxa or focusing on taxa with prior evidence to reduce false negatives.

## 5. Sequencing depth and MAGs availability

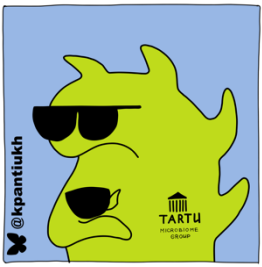
- Low-depth sequencing may not resolve species or strains reliably.
- High-resolution analysis is only meaningful if the data can support it; otherwise, stick to genus/family.

## 6. Hierarchical approach

- Start broad to detect global shifts, then zoom in to finer levels for taxa showing signals.
- This balances power, interpretability, and control of multiple testing.





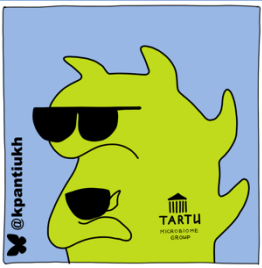


# How to decide what taxonomic level to use?

*It's about finding the right balance and testing many different options before discovering what works*

- Small cohort, rare taxa, low sequencing depth → use higher taxonomic levels (family/genus).
- Large cohort, common taxa, high-resolution data → species or strain level can be explored.
- Always consider prevalence, expected group size, and interpretability.





# Preparing an abundance table

## 1. Filter low-abundance and rare taxa

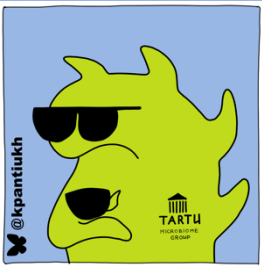
- Remove taxa that are present in very few samples ( $<1\%$  prevalence)
- Remove taxa with extremely low relative abundance (not popular)

*This reduces sparsity, improves statistical power, and decreases the multiple-testing burden.*

## 2. Handle zeros & Compositional transformation

- Zero counts are common in microbiome data.
- Perform CLR transformations





# Bacterial species

## main characteristics

**Ecological niche / lifestyle**

**Functional traits / metabolism**

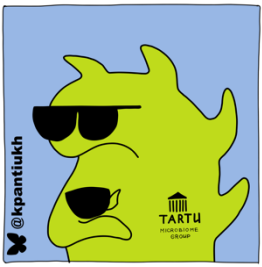
**Genomic**

**Morphology and structural features**

**Interaction with other microbes**

**Clinical or industrial relevance**





# Bacterial species

## main characteristics

### Ecological niche / lifestyle

### Functional traits / metabolism

### Genomic

### Morphology and structural features

### Interaction with other microbes

### Clinical or industrial relevance

- **Habitat:** gut, oral cavity, soil, water, skin, etc.
- **Host association:** commensal, symbiont, opportunistic pathogen, obligate pathogen.
- **Temperature preference:** psychrophile, mesophile, thermophile.
- **pH tolerance** and other environmental tolerances.

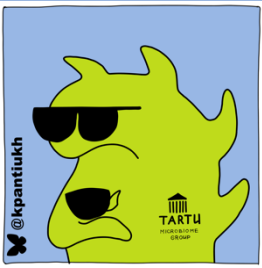
<https://metatraits.embl.de>

metaTraits

Databases ▾

Try the family "M





# Bacterial species

## main characteristics

Ecological niche / lifestyle

Functional traits / metabolism

Genomic

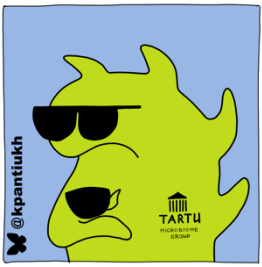
Morphology and structural features

Interaction with other microbes

Clinical or industrial relevance

- **Carbon source utilization:** sugars, proteins, lipids.
- **Energy generation:** respiration, fermentation, photosynthesis, chemolithotrophy.
- **Nitrogen/sulfur cycling capabilities:** nitrate reduction, sulfate reduction, ammonia oxidation.
- **Secondary metabolite production:** antibiotics, bacteriocins, signaling molecules.





# Bacterial species

## main characteristics

Ecological niche / lifestyle

Functional traits / metabolism

Genomic

Morphology and structural features

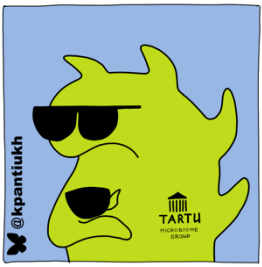
Interaction with other microbes

Clinical or industrial relevance

- **Genome size** and GC content.
- **Plasmid presence** or mobile genetic elements.
- **Virulence genes** or toxin production.
- **Antibiotic resistance genes.**

<https://gtdb.ecogenomic.org>





# Bacterial species

## main characteristics

**Ecological niche / lifestyle**

**Functional traits / metabolism**

**Genomic**

**Morphology and structural features**

**Interaction with other microbes**

**Clinical or industrial relevance**

- **Cell shape:** cocci, rods, spirals.
- **Motility structures:** flagella, pili.
- **Surface structures:** capsule, S-layer, biofilm-forming ability.
- **Sporulation ability.**

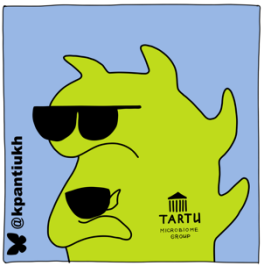
<https://metatraits.embl.de>

metaTraits

Databases ▾

Try the family "M"





# Bacterial species

## main characteristics

**Ecological niche / lifestyle**

**Functional traits / metabolism**

**Genomic**

**Morphology and structural features**

**Interaction with other microbes**

**Clinical or industrial relevance**

- **Symbiosis or antagonism:** production of inhibitory compounds, mutualistic relationships.
- **Biofilm formation:** ability to form communities on surfaces.
- **Quorum sensing / communication:** signaling mechanisms.

<https://pubmed.ncbi.nlm.nih.gov>

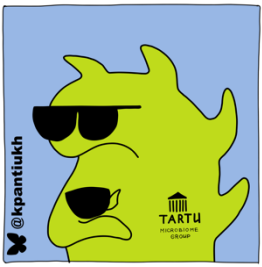
<https://www.biorxiv.org>



**bioRxiv**  
THE PREPRINT SERVER FOR BIOLOGY







# Bacterial species

## main characteristics

**Ecological niche / lifestyle**

**Functional traits / metabolism**

**Genomic**

**Morphology and structural features**

**Interaction with other microbes**

**Clinical or industrial relevance**

- Pathogenicity to humans, animals, or plants.
- Probiotic potential.
- Industrial applications: fermentation, bioremediation, enzyme production.

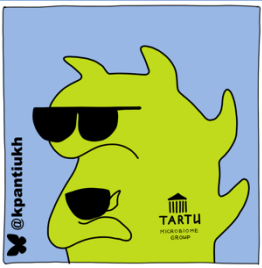
<https://pubmed.ncbi.nlm.nih.gov>

<https://www.biorxiv.org>



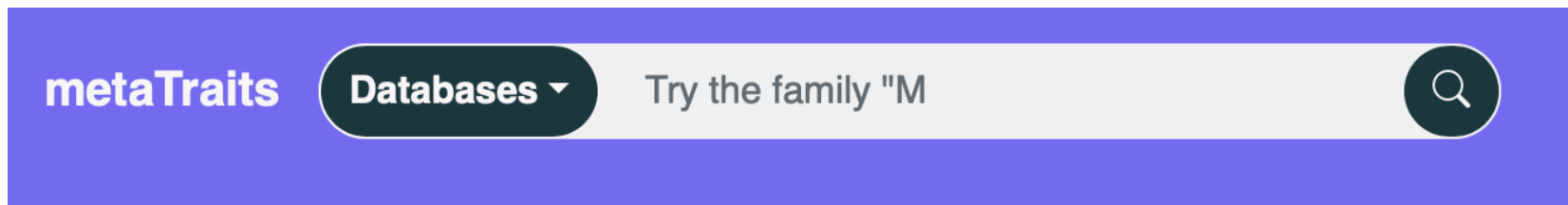
**bioRxiv**  
THE PREPRINT SERVER FOR BIOLOGY

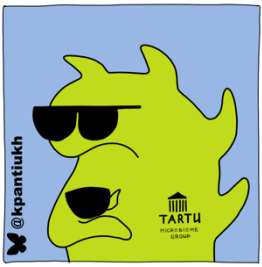




# MetaTraits demo

<https://metatraits.embl.de>





# GTDB demo

<https://gtdb.ecogenomic.org>

