



Microbiota Data Analysis Workshop

22 - 24 January, 2020

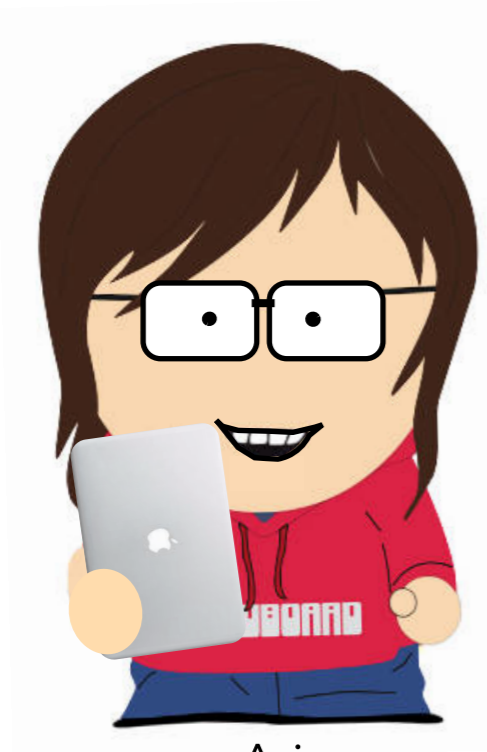
<https://www.gdc-docs.ethz.ch/MDA/site/>

Welcome



Genetic Diversity Centre
Zurich

<http://www.gdc.ethz.ch/>



Aria



Jean-Claude



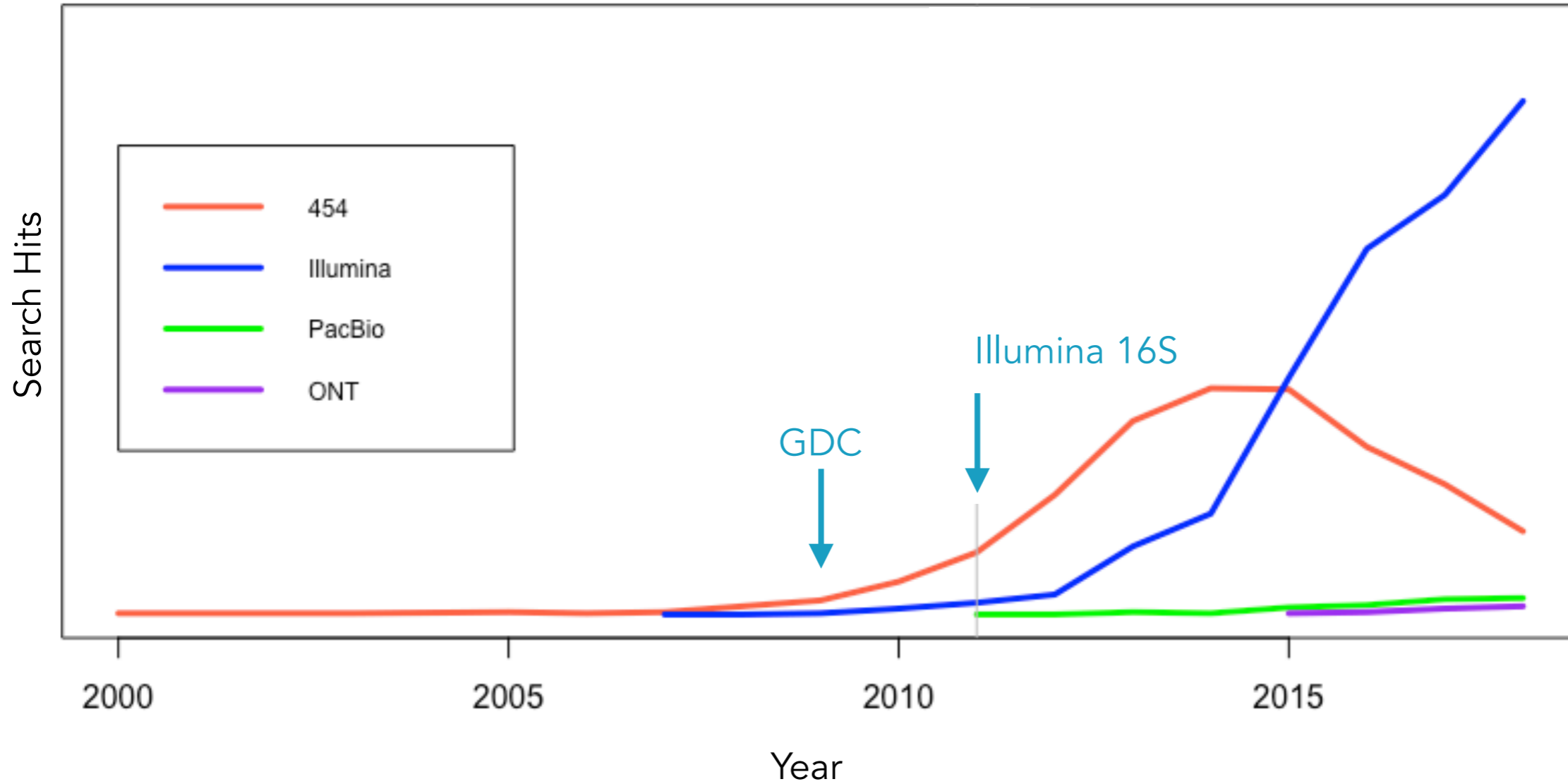
Silvia



Nik

The **Genetic Diversity Centre (GDC)** is a knowledge and technology platform of the D-USYS Department. The team GDC provides scientific and technical support for research related to genetic and genomic diversity in a wide range of non-model organisms.

Number of Publications on PubMed
(Search term: <SeqTech> and 16S)

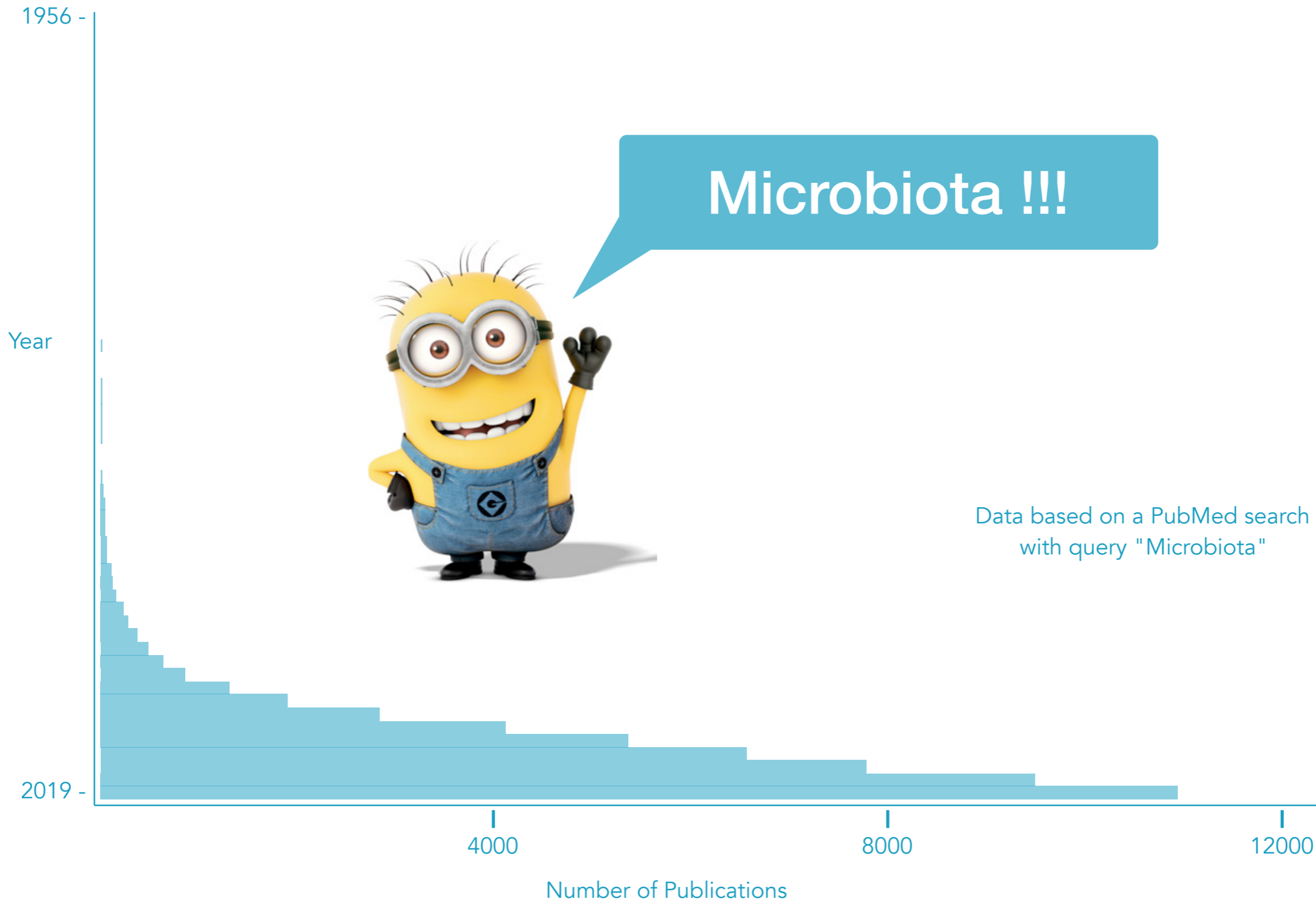


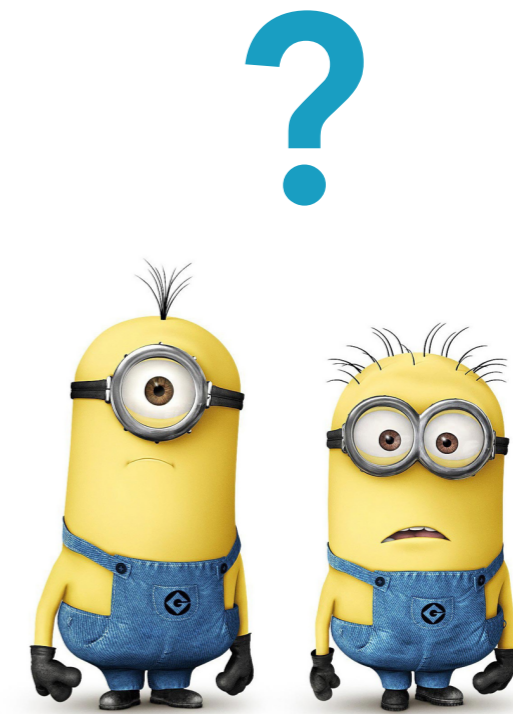
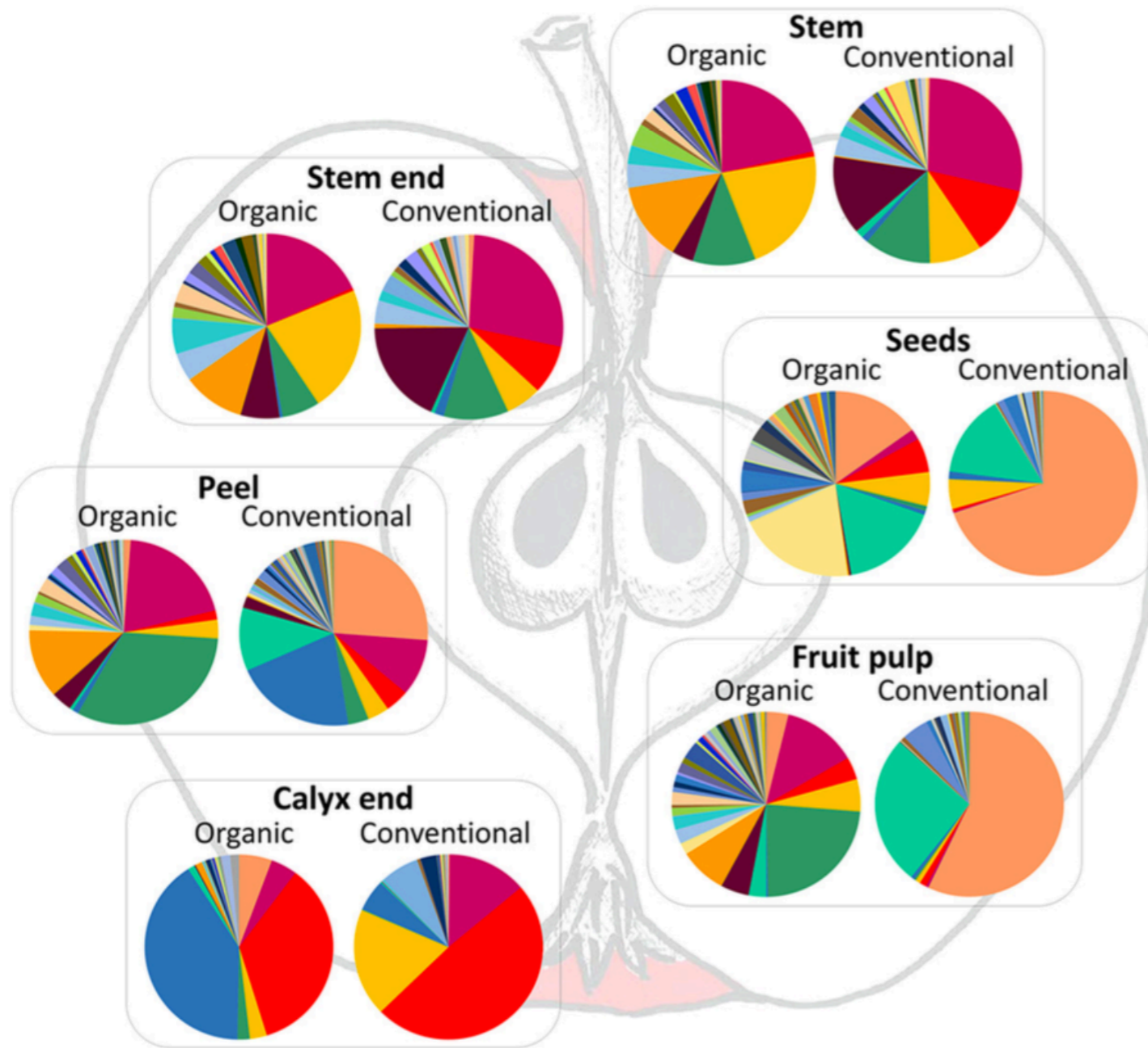
Announcements

- 👉 Lecture room: key, food and drinks
- 👉 Bathrooms
- 👉 Coffee / tea breaks
- 👉 Lunch
- 👉 Social dinner - Thursday evening > Doodle

- 👉 Power plugs
- 👉 Wi-Fi connection (eduroam)
- 👉 Access to Google drive
- 👉 R libraries
- 👉 Question(s)

Overview





Wassermann et al. (2019) An Apple a Day: Which Bacteria Do We Eat With Organic and Conventional Apples? *Frontiers in Microbiology*. Volume 10 | Article 1629.

Wednesday 22.01.20

Time	Subject
8:00 - 9:00	Registration & Installation
9:00-9:30	Introduction
9:30-10:30	Warm Up
10:30 - 11:00	Tea / Coffee Break
11:00-11:45	Data Preparation
11:45-12:30	eDNA
12:30 - 13:30	Lunch
13:30-15:00	Reproducible Science
15:00 - 15:30	Tea / Coffee Break
15:30-16:15	Reality Check
16:15-17:00	Q&A
17:00	End

Thursday 23.01.20

Time	Subject
8:00 - 9:00	Installation
9:00-10:30	Diversity Analysis
10:30 - 11:00	Tea / Coffee Break
11:00-12:30	Diversity Analysis
12:30 - 13:30	Lunch
13:30-15:00	Diversity Analysis
15:00 - 15:30	Tea / Coffee Break
15:30-17:00	Diversity Analysis
17:00	End

Friday 24.01.20

Time	Subject
8:00 - 9:00	Installation
9:00-10:30	Microbial Network
10:30 - 11:00	Tea / Coffee Break
11:00-12:30	Microbial Network
12:30 - 13:30	Lunch
13:30-15:00	Data Submission
15:00 - 15:30	Tea / Coffee Break
15:30-17:00	Random Forest
17:00	End

- Mahendra Mariadassou (INRA)
- Robert Edgar (drive5)
- Sam Holt (EMBL-EBI)
- Klaus Schläppi (University of Bern)
- Kristy Deiner (ETH Zürich)
- Jean-Claude Walser (GDC)
- Nik Zemp (GDC)

Dinner*

* Dinner with the speakers but everybody is welcome to join.

warmup



hops on the spot



side-to-side hops
single leg



hops on the spot



side-to-side hops
feet together



alt back expansions



chest expansions



arm circles (wide)



arm circles



hops on the spot



side-to-side hops
single leg

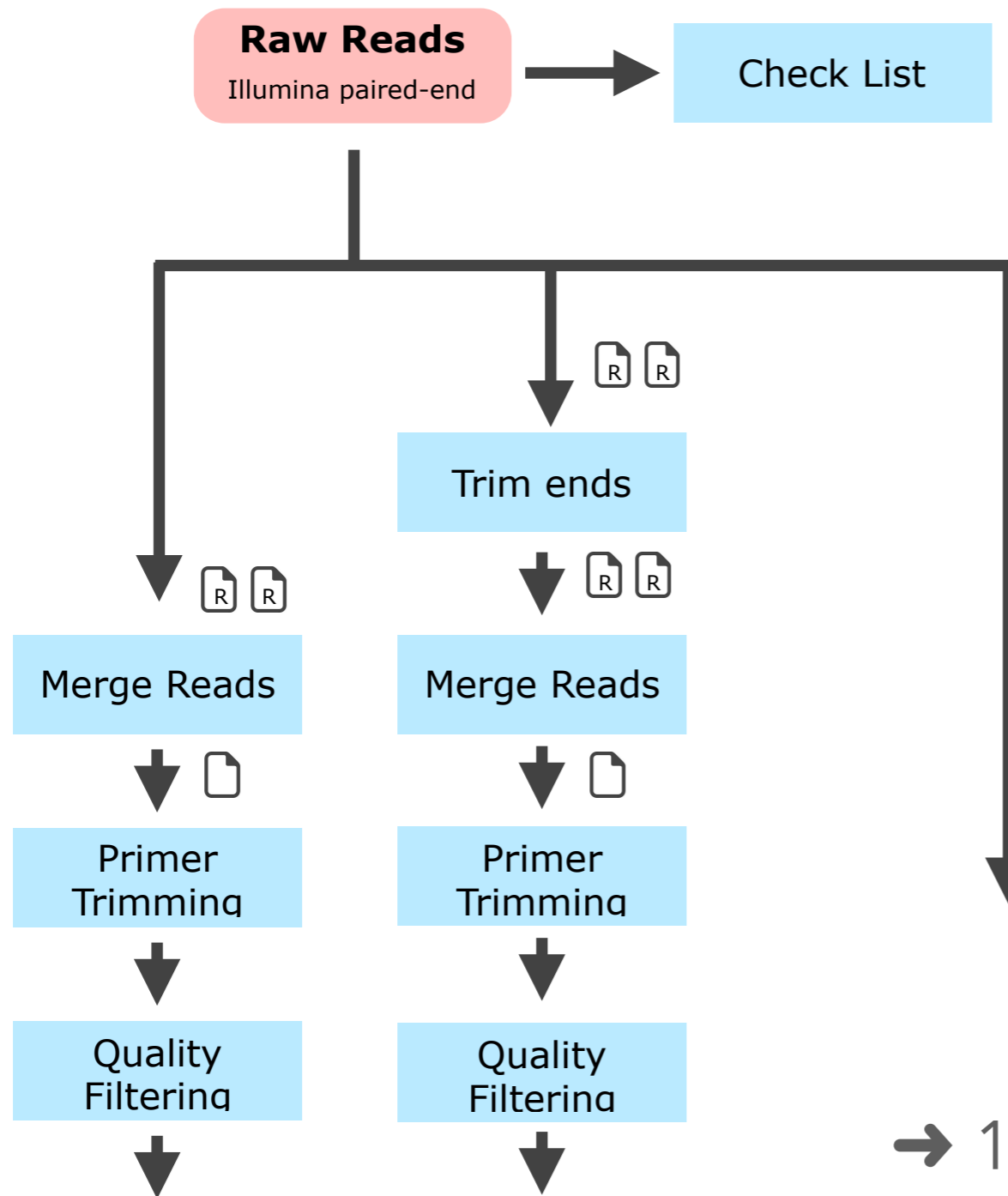


hip rotations



torso rotations

➔ 9:30-10:30 Jean-Claude Walser

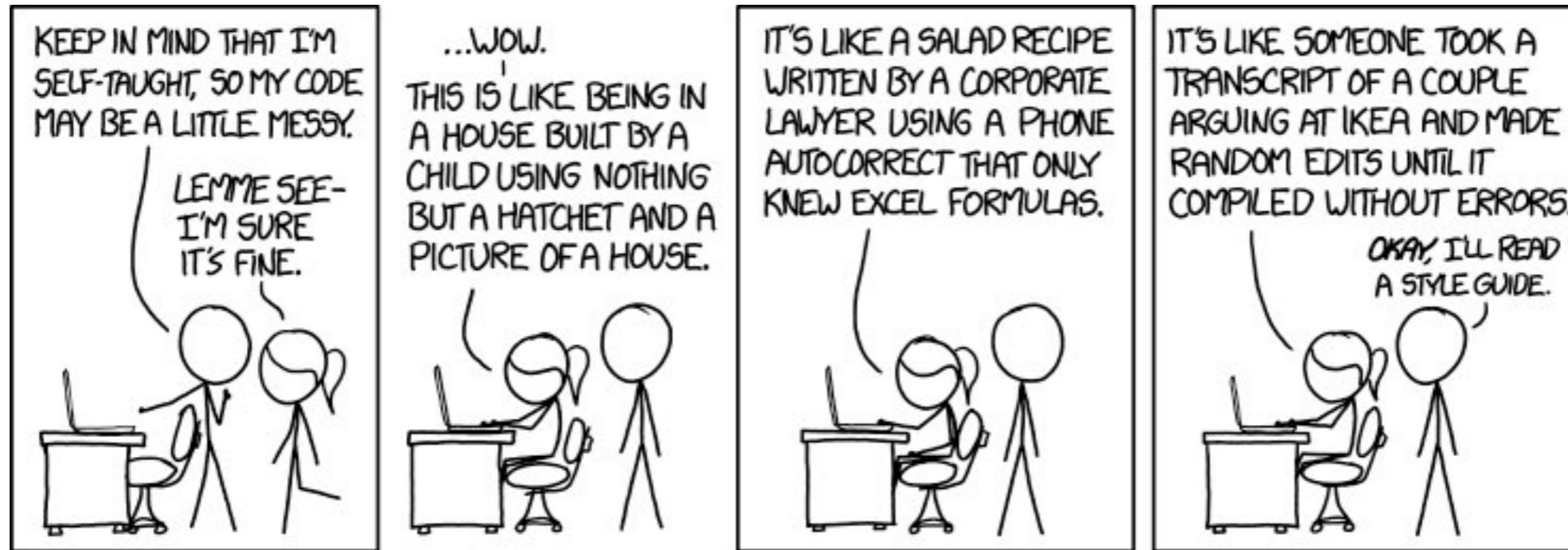


→ 11:00-11:45 Jean-Claude Walser

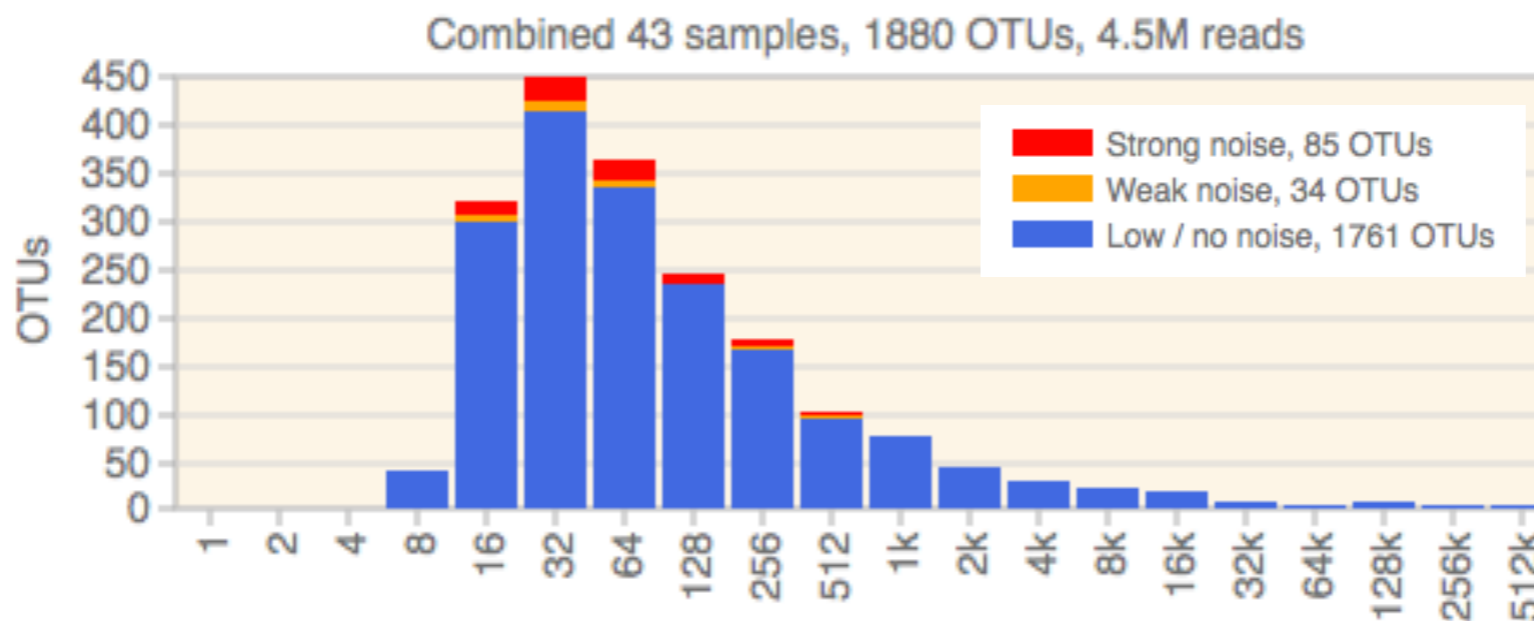


➔ 11:45 - 12:30 Kristy Deiner

Reproducible Science



➔ 13:30 - 15:00 Nik Zemp



➔ 15:30 - 16:15 Robert Edgar

➔ 16:15 - 17:00 Robert Edgar

Count / OTU Table

	Sample1	Sample2	Sample3	Sample4	Sample5	Sample6
OUT1	2519	2354	1074	2	452	1233
OUT2	2520	1158	1287	1	34	3184
OUT3	106	4	2	0	0	0
OUT4	490	82	148	0	0	22
OUT5	120	73	111	0	0	133
OUT6	13	1	6	0	0	0
OUT7	9	0	0	0	0	0
OUT8	813	415	142	0	0	808
OUT9	45	2	7	0	0	0

	Sample1	Sample2	Sample3	Sample4	Sample5	Sample6
OUT1	2519	2254	1074	2	452	1233
OUT2						4
OUT3						
OUT4						
OUT5						3
OUT6						
OUT7						
OUT8	815	415	142	0	0	808
OUT9	45	2	7	0	0	0

Microbiome data are ...

(a) **compositional** (multiple parts of nonnegative numbers).

(b) **high dimensional** (few data points and many features)

and **underdetermined** (the number of OTUs much greater than the number of samples).

(c) **overdispersed** (variance of the counts of read is larger than expected).

(d) often spares with **many zeros** (zero-inflated).

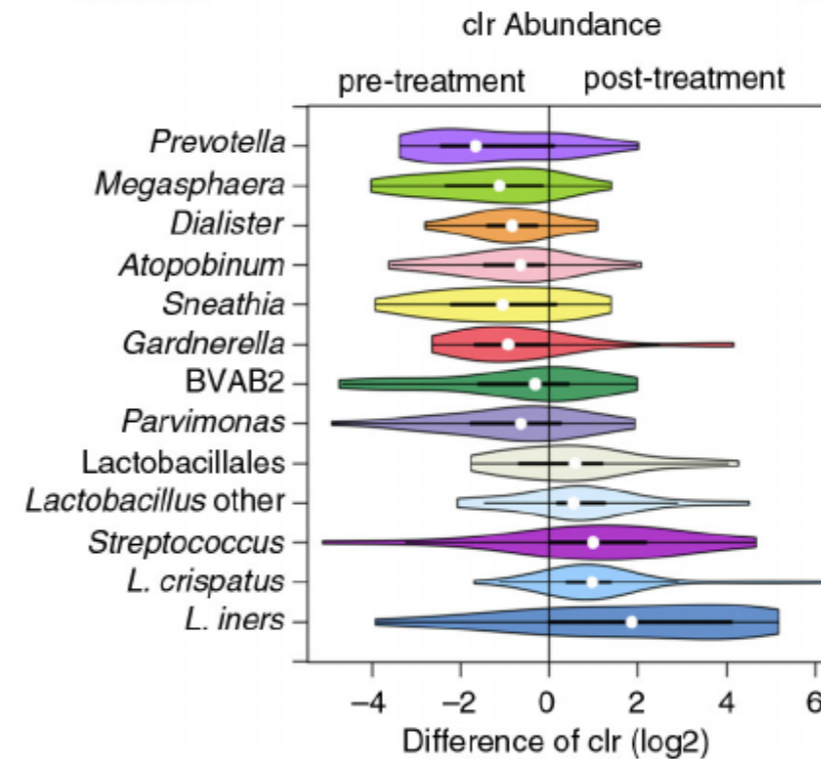
Metadata / Mapfile

	pH	TP	NO2	NO3	NH4	Shannon
Sample1	8.3	13.3	17.7	1.6	67.4	3.5
Sample2	8.2	13.3	17.8	1.6	67.9	3.3
Sample3	8.2	13.2	17.7	1.6	68.6	3.3
Sample4	8.2	14.3	17.3	1.5	75.2	3.4
Sample5	9.1	14.2	16.6	1.5	75.1	3.5
Sample6	9.1	14.1	17.5	1.6	94.2	3.4
Sample7	9.1	15.3	22.3	1.5	11.8	3.6
Sample8	7.9	14.3	22.3	1.4	11.3	3.6
Sample9	7.9	16.0	22.4	1.4	9.3	3.7

Thursday, 23.01.2020

	Sample1	Sample2	Sample3	Sample4	Sample5	Sample6
OUT1	2519	2354	1074	2	452	1233
OUT2	2520	1158	1287	1	34	3184
OUT3	106	4	2	0	0	0
OUT4	490	82	148	0	0	22
OUT5	120	73	111	0	0	133
OUT6	13	1	6	0	0	0
OUT7	9	0	0	0	0	0
OUT8	813	415	142	0	0	808
OUT9	45	2	7	0	0	0

	pH	TP	NO2	NO3	NH4	Shannon
Sample1	8.3	13.3	17.7	1.6	67.4	3.5
Sample2	8.2	13.3	17.8	1.6	67.9	3.3
Sample3	8.2	13.2	17.7	1.6	68.6	3.3
Sample4	8.2	14.3	17.3	1.5	75.2	3.4
Sample5	9.1	14.2	16.6	1.5	75.1	3.5
Sample6	9.1	14.1	17.5	1.6	94.2	3.4
Sample7	9.1	15.3	22.3	1.5	11.8	3.6
Sample8	7.9	14.3	22.3	1.4	11.3	3.6
Sample9	7.9	16.0	22.4	1.4	9.3	3.7

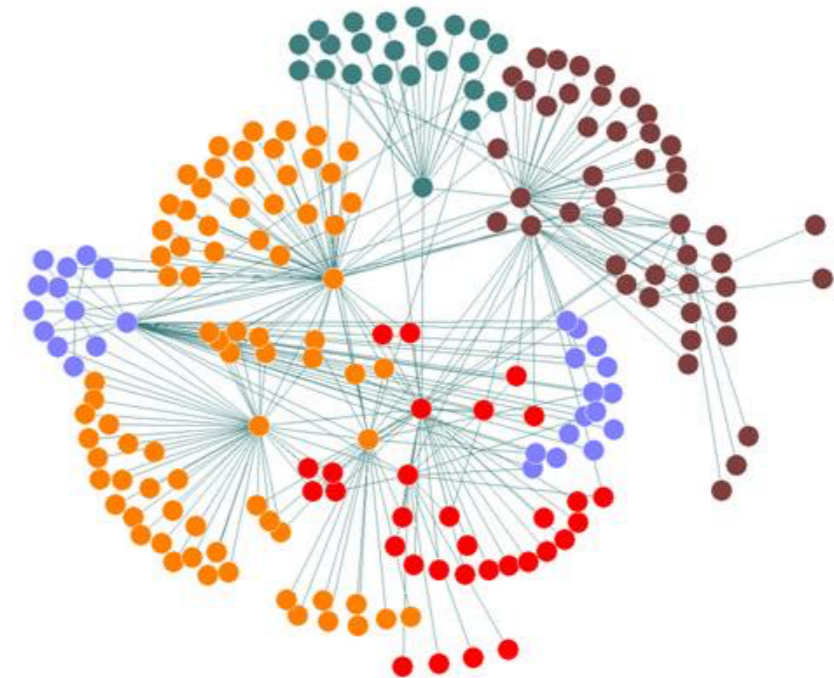


→ Tomorrow: 9.00 - 17.00 Mahendra Mariadassou

Friday, 24.01.2020

	Sample1	Sample2	Sample3	Sample4	Sample5	Sample6
OUT1	2519	2354	1074	2	452	1233
OUT2	2520	1158	1287	1	34	3184
OUT3	106	4	2	0	0	0
OUT4	490	82	148	0	0	22
OUT5	120	73	111	0	0	133
OUT6	13	1	6	0	0	0
OUT7	9	0	0	0	0	0
OUT8	813	415	142	0	0	808
OUT9	45	2	7	0	0	0

	pH	TP	NO2	NO3	NH4	Shannon
Sample1	8.3	13.3	17.7	1.6	67.4	3.5
Sample2	8.2	13.3	17.8	1.6	67.9	3.3
Sample3	8.2	13.2	17.7	1.6	68.6	3.3
Sample4	8.2	14.3	17.3	1.5	75.2	3.4
Sample5	9.1	14.2	16.6	1.5	75.1	3.5
Sample6	9.1	14.1	17.5	1.6	94.2	3.4
Sample7	9.1	15.3	22.3	1.5	11.8	3.6
Sample8	7.9	14.3	22.3	1.4	11.3	3.6
Sample9	7.9	16.0	22.4	1.4	9.3	3.7



➔ Friday 9:00 - 10:30

Jean-Claude Walser

➔ Friday 11:00 - 12:30

Klaus Schlaeppli

Friday, 24.01.2020

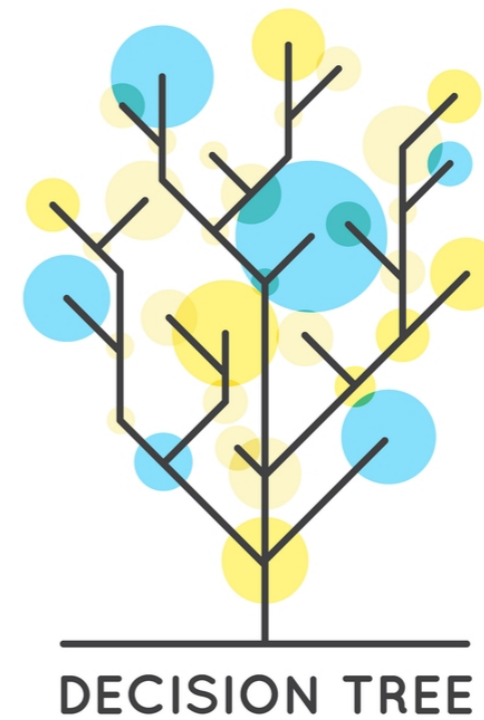


→ Friday 13:30 - 15:00 Sam Holt

Friday, 24.01.2020

	Sample1	Sample2	Sample3	Sample4	Sample5	Sample6
OUT1	2519	2354	1074	2	452	1233
OUT2	2520	1158	1287	1	34	3184
OUT3	106	4	2	0	0	0
OUT4	490	82	148	0	0	22
OUT5	120	73	111	0	0	133
OUT6	13	1	6	0	0	0
OUT7	9	0	0	0	0	0
OUT8	813	415	142	0	0	808
OUT9	45	2	7	0	0	0

	pH	TP	NO2	NO3	NH4	Shannon
Sample1	8.3	13.3	17.7	1.6	67.4	3.5
Sample2	8.2	13.3	17.8	1.6	67.9	3.3
Sample3	8.2	13.2	17.7	1.6	68.6	3.3
Sample4	8.2	14.3	17.3	1.5	75.2	3.4
Sample5	9.1	14.2	16.6	1.5	75.1	3.5
Sample6	9.1	14.1	17.5	1.6	94.2	3.4
Sample7	9.1	15.3	22.3	1.5	11.8	3.6
Sample8	7.9	14.3	22.3	1.4	11.3	3.6
Sample9	7.9	16.0	22.4	1.4	9.3	3.7



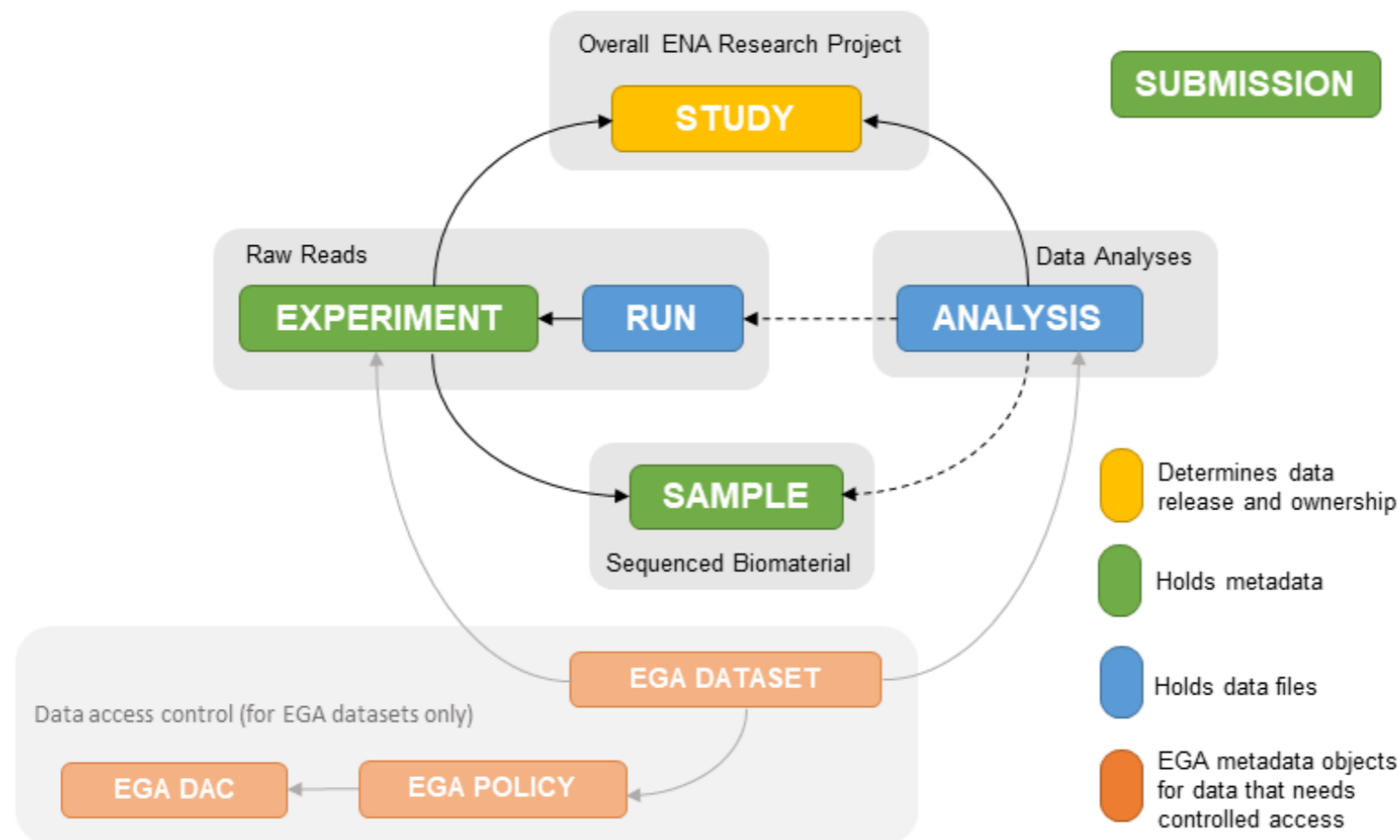
➔ Friday 15:30 - 17:00 Jean-Claude Walser

Introduction

Cluster-Game



European Nucleotide Archive (ENA) accepts sequence reads and associated analyses. Once public, data submitted to ENA is exchanged between International Nucleotide Sequence Database Collaboration (INSDC) partners: NCBI and DDBJ.



MIxS – Minimum information about any (x) sequence



MIxP – Minimum information about any (x) participants

*The **MIxS** is a unified standard developed by the Genomic Standards Consortium (GSC) for reporting of minimum information about any (x) nucleotide sequence. It consists of MIGS (Minimum Information about a Genome Sequence), MIMS (Minimum Information about a Metagenome Sequence) and **MIMARKS** (Minimum Information about a MARKer gene Sequence) standards and describes fourteen environments.

Questions:

1. What **data type** are you working with?
2. From which biome(s) do your samples come from?
3. What species-diversity are you interested in?
4. What sequencing platform/technology are you using?
5. What is the read / sequence length?
6. Do you prepare the library yourself or do you outsource it?

Answers:

1. Amplicon ⇨ 16S

2. Sediment ⇨ Deep sea sediment

3. Bacteria ⇨ Cyanobacteria

4. Illumina ⇨ MiSeq

5. paired-end 2 x 300

6. NO (of course not)

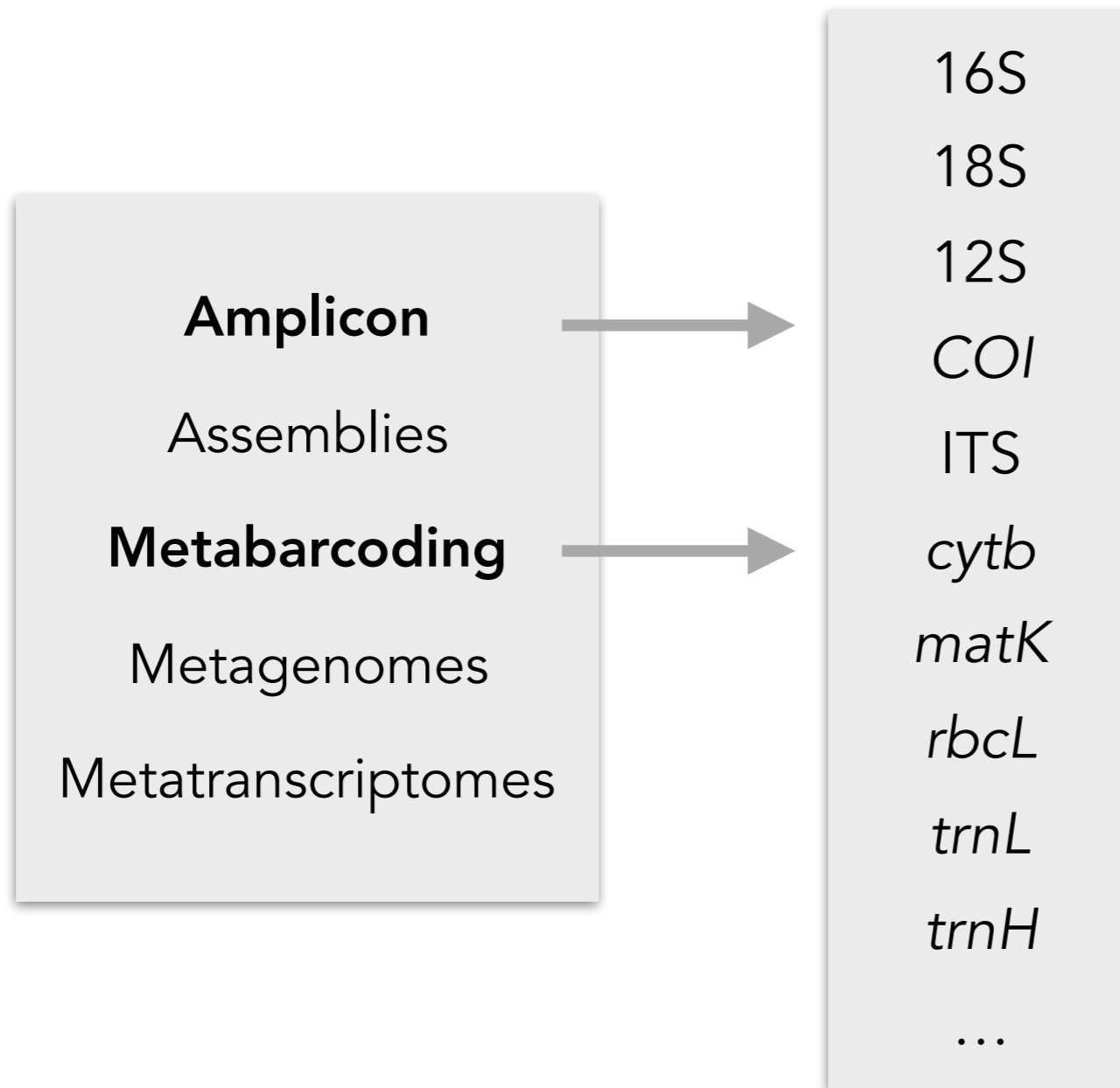
1

Data Type

- Amplicon
- Assemblies
- Metabarcoding
- Metagenomes
- Metatranscriptomes

1_b

Data Sub-Types



2

Environment-specific Descriptors

(Genomics Standards Consortium (GSC) MixS Checklists)

- Air
 - Host-associated
 - Human-associated
 - Microbial mat/biofilm
 - Plant-associated
- Sediment
- Soil
- Wastewater / Sludge
- Water

2_bEnvironment-specific **Sub-Descriptors**

(Genomics Standards Consortium (GSC) MixS Checklists)

Air
Host-associated
Human-associated
Microbial mat/biofilm
Plant-associated
Sediment
Soil
Wastewater / Sludge
Water



Air ⇔ AC filter system
Host ⇔ daphnia
Human ⇔ baby gut
Biofilm ⇔ rubber ducks
Plant ⇔ leave
Sediment ⇔ deep sea sediment
Soil ⇔ farm land
Sludge ⇔ WWTPs
Water ⇔ river systems

3 Target Species

Bacteria
Archaea
Fungi
Plants
Eukaryote

...

4 NGS Data Formats

ILLUMINA
Roche LS454
PacBio
IonTorrent
Oxford Nanopore

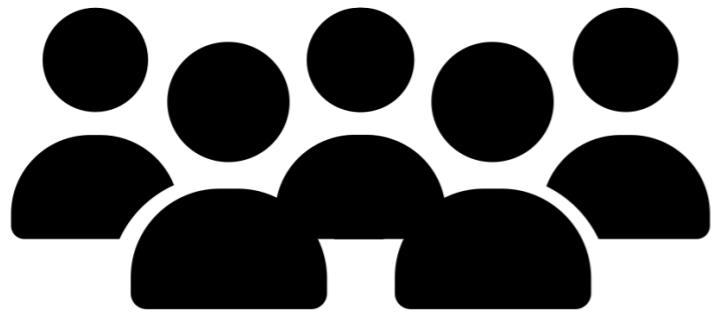
4b

NGS Data Formats - Read Length



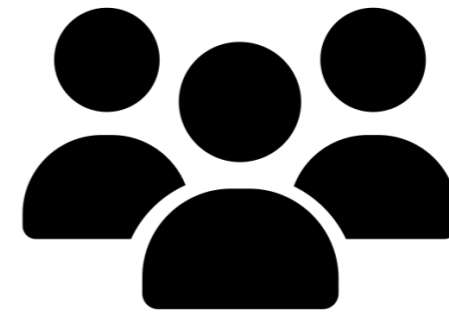
5 Library Prep

Yes / No



Answers:

1. **Amplicon** ⇨ 16S



Answers:

1. **Metagenomics**

Lets Cluster!

