# Meta_Microbial Workshop

## Exploring Online Resources and Repositories **I**
- Practical example (Silva NGS)

**Catarina Magalhães**

# PLAN FOR THE NEXT 45'

- Silva NGS | Short Introduction
- Start a Case Study
  - Samples/Data Description
  - Silva NGS Registration
  - Execute a Project
  - Analyze the Outputs in Working Groups by Responding   to a Challenge
    Present Your Results to the Class

# SILVA NGS | SHORT INTRODUCTION

- SILVAngs is a data analysis service for ribosomal RNA gene (rDNA) amplicon reads generated by next-generation sequencing (NGS), based on an automatic software pipeline.

- It uses the SILVA databases to classify rDNA reads and provides several outputs, like taxonomy tables and multiple graphs for download.

- SILVAngs uses the de.NBI Cloud (German Network for Bioinformátics Infrastructure) to process all user projects. The de.NBI cloud provide free compute resources for academic users.

*Quast et al. 2013*
Developed by Frank Oliver Group
Max Planck Institute for Marine Microbiology

*Contact: ngs-contact@arb-silva.de*

# SILVA NGS | SHORT INTRODUCTION

The basic workflow of the pipeline can be divided into the following steps

- Alignment | Initial quality control

- Quality management | further quality filtering

- De-replication (identical sequences)

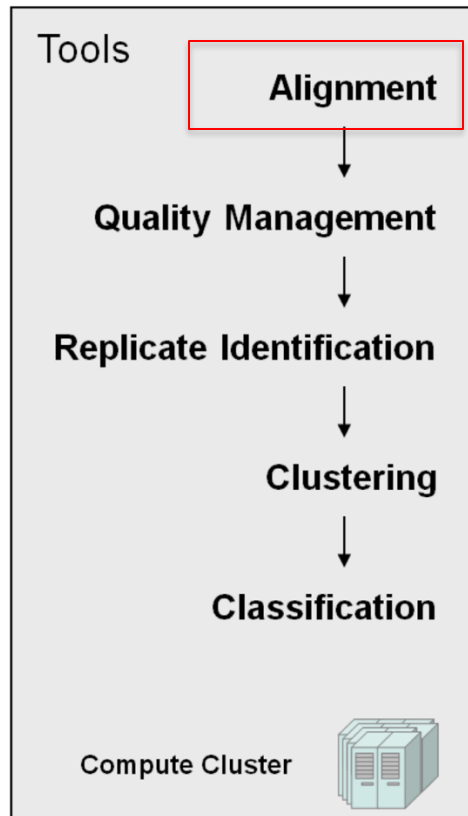- Clustering (OTU definition)

- Classification of the OTUs/reads

SILVAngs User Guide:

https://www.google.com/search?q=silva+ngs+user+guide&rlz=1C5CHFA_enPT857PT857&oq=silva+ngs+user+guide&gs_lcrp=EgZjaHJvbWUyBggAEEUYOTIGCAEQRRg80gEINDk2NWowajSoAgCwAgE&sourceid=chrome&ie=UTF-8

# 1º Alignment

## Tools

Alignment

↓

Quality Management

↓

Replicate Identification

↓

Clustering

↓

Classification

Compute Cluster

## Alignment | Initial quality control

All input reads are aligned by SILVA Incremental Aligner (SINA - http://www.arb-silva.de/aligner/sina-download/)

According to a number of parameters measured by SINA, problematic reads (such as PCR artefacts) or even contamination of the dataset with non-rRNA gene sequences are identified and not considered for further processing

*Pruesse et al. 2012*
*https://academic.oup.com/bioinformatics/article/28/14/1823/218226*

# 2º Quality management



Tools

**Alignment**
↓
**Quality Management**
↓
**Replicate Identification**
↓
**Clustering**
↓
**Classification**

Compute Cluster

**Further quality filtering**

All reads which have not been rejected by the previous alignment step undergo further quality filtering, including **length**, **ambiguity** (ambiguous bases) and **homopolymer** (series of consecutive identical bases) checks. The sequence length cut-off can be defined by the user, whereas for ambiguities and homopolymers thresholds of max. 2% are used.

# 2º Replication Identification



**De-replication | identification of identical sequences**

All remaining reads enter the de-replication stage of the pipeline. 100% identical reads, are identified and only one read is retained for further processing.

This is to reduce calculation time, since processing redundant reads is a waste of computing power.

# 4º Clustering



**Clustering (OTU definition)**

Clustering is done on a **97-99% identity level**. This is motivated by the fact that PCR and sequencing errors can easily introduce 3% artificial divergence in the sequences.

Technically, this is just another de-replication step to further reduce the number of reads that needs to be classified. Compared to previous de-replication,

97-99 % Identity
(can be adjusted)

# 4º Classification



**- Classification of the OTUs/reads**

In the classification step the representative reads (the longest read of each OTU) are compared to the SILVA reference datasets of the small- (16S/18S) and large (23S/28S) subunit rDNA with its corresponding SILVA taxonomy (Quast et al. 2013).

Only significant hits are considered; everything else is assigned to a class called 'No Relative'.

The classification result of each representative read is mapped back to all other reads of the OTU cluster and also to the corresponding identical reads from de-replication step.

Is recommended to avoid overinterpretation of the results especially for reads below 1200 bases. For amplicon illumina data SILVA's phylogeny-based taxonomy has a reliable resolution down to the genus level.

# 5º RESULTS

- Chart Gallery (e.g. Sequence lengh distribution; Rarefaction curves)

- Fingerprints (It offers the option to specify the taxonomic depth, filter by taxonomy,

- include/remove samples)

- Taxplotkrona (Interactive visualization)

- Archive Zip (with all the autputs)

- Report (HTML and PDF)

# SILVA NGS Pipeline | Let's Make It Happen!

Meta_Microbial
Workshop

1. Go to SILVA NGS Web Page at https://www.arb-silva.de/ngs/
2. Register at SILVA NGS platform and log in.
3. Start a project
    3.1 Go to "My Projects"
    3.2 Select "Create Project"
    3.3 Add a Name in "Project Name" – *No spaces or special characters*
    3.4 In Sequence Type chose 16S/18S, SSU
    3.5 In Sequence Technology select "Illumina (MiSeq/HiSeq)"
    3.6 In Expected Sequence Quantity write "300000"
    3.7 In Expected Read Length write "400"
    3.8.In project description write for example "Samples from Salt Pans"
4. In Upload Files "Upload your sequences". The FASTA files E1, E2, E3 and E4 available in:
https://drive.google.com/drive/folders/1dqUnQHbp1VJHUSSDCVdCoZIA3FLep9LS?usp=sharing
5. Select "Execute project".

*SILVAngs - rDNA-based microbial community analysis using next-generation sequencing (NGS) data -*
*User Guide_2020.* *https://ngs.arb-silva.de/silvangs/#*

Meta_Microbial Workshop

SAMPLES WE WILL ANALYZE

Olhão Salt Pans

Aveiro Salt Pans

E1 - Sal
E2 - Soil

E3 - Salt
E4 - Soil

# METABARCODING
# GENETIC MARKER | 16S rRNA GENE

**515F-Y/926R - primers**
**Hypervariable Regions V4 and V5 of the 16S rRNA Gene**

Number of Sequences
63,244

*Ficheiros **E1**; E2; E3; E4*

# CHALLENGE – Extract Insights from Your Data

## Form Groups of 4-5 Elements

In this challenge, we ask to conduct a visual and interpretative analysis of the SilvaNGS outputs to study the diversity of microbiomes in Olhão and Aveiro Salt Pans.

This platform will generate graphical visualizations and taxonomic tables are available at

https://drive.google.com/drive/folders/1d9I1ThGaiH-wB29PLyyCMNGCOVelPOG0?usp=sharing

**Analyze the results obtained in groups in order to answer one of the following questions**:

1. Which samples showed the greatest diversity of prokaryotes?

2. Which prokaryotic phyla were present in all samples?

3. Which prokaryotic phyla differed between the two studied salt pans (Olhão and Aveiro)?

4. Among the prokaryotes, which domain is better represented in the studied salt pan samples: Archaea or Bacteria? What differences were observed in terms of relative abundance of Archaea or Bacteria phyla identified in the different samples?

5. At a lower taxonomic level (genus), analyze a group of prokaryotes of your choice and compare this taxonomic group between the different analyzed samples.

6. Identify differences in the structure of microbial communities inhabiting the Salt and Sediment in the Aveiro Salt Pan.

7. Characterize the distribution of cyanobacteria in the different sampling stations.

8. Answer your own question...