# Functional microbiome research – bioinformatics section
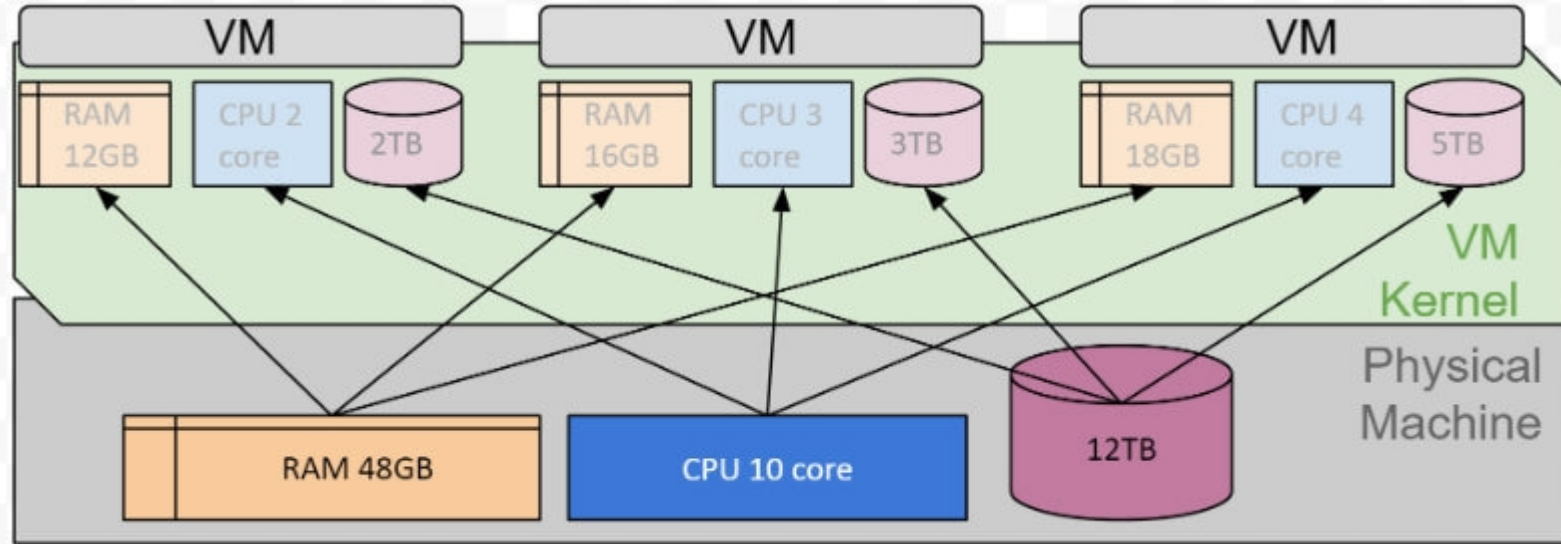
## Day 1 – intro to linux
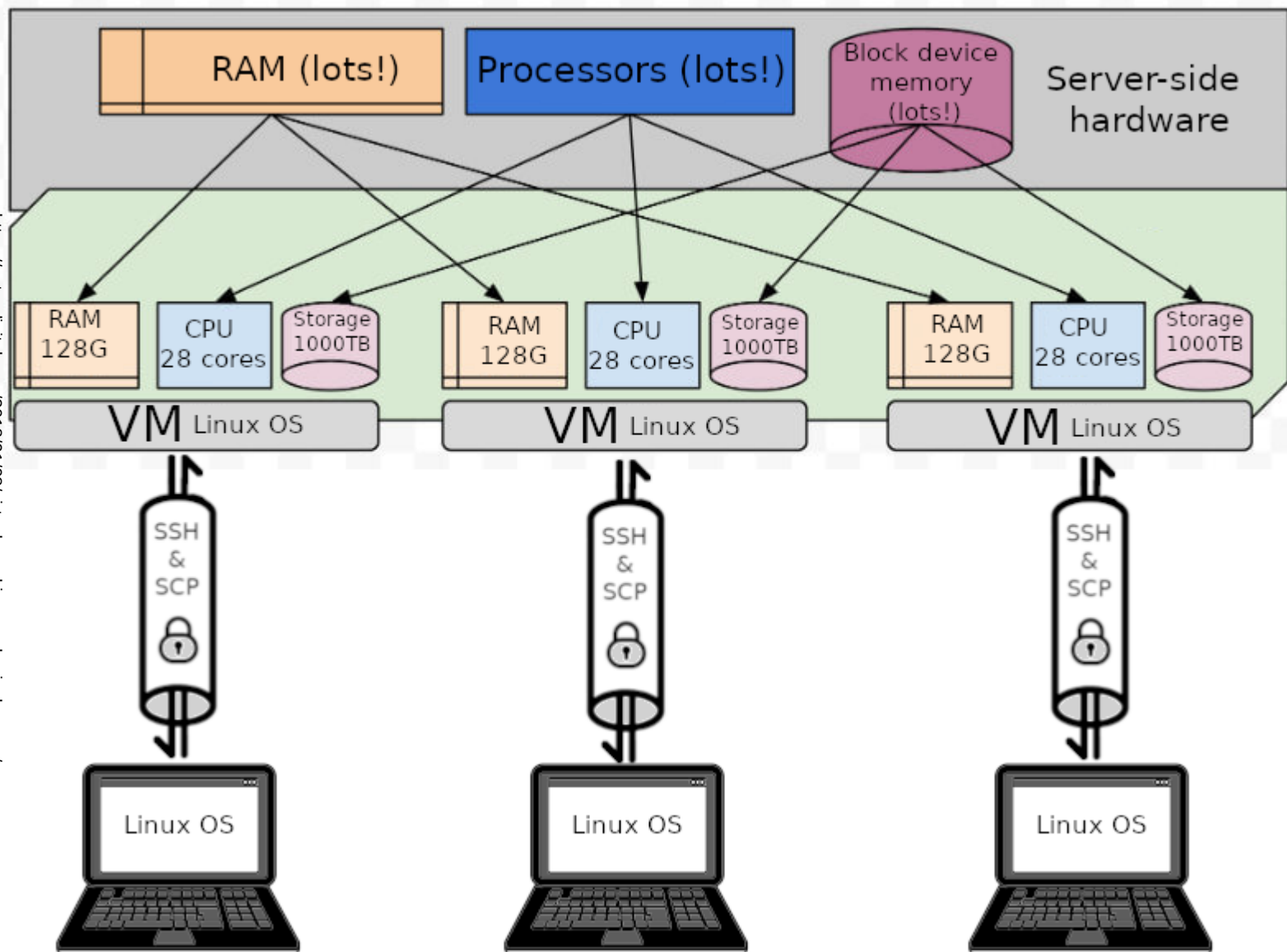
First, thanks to de.NBI. They are amazing.

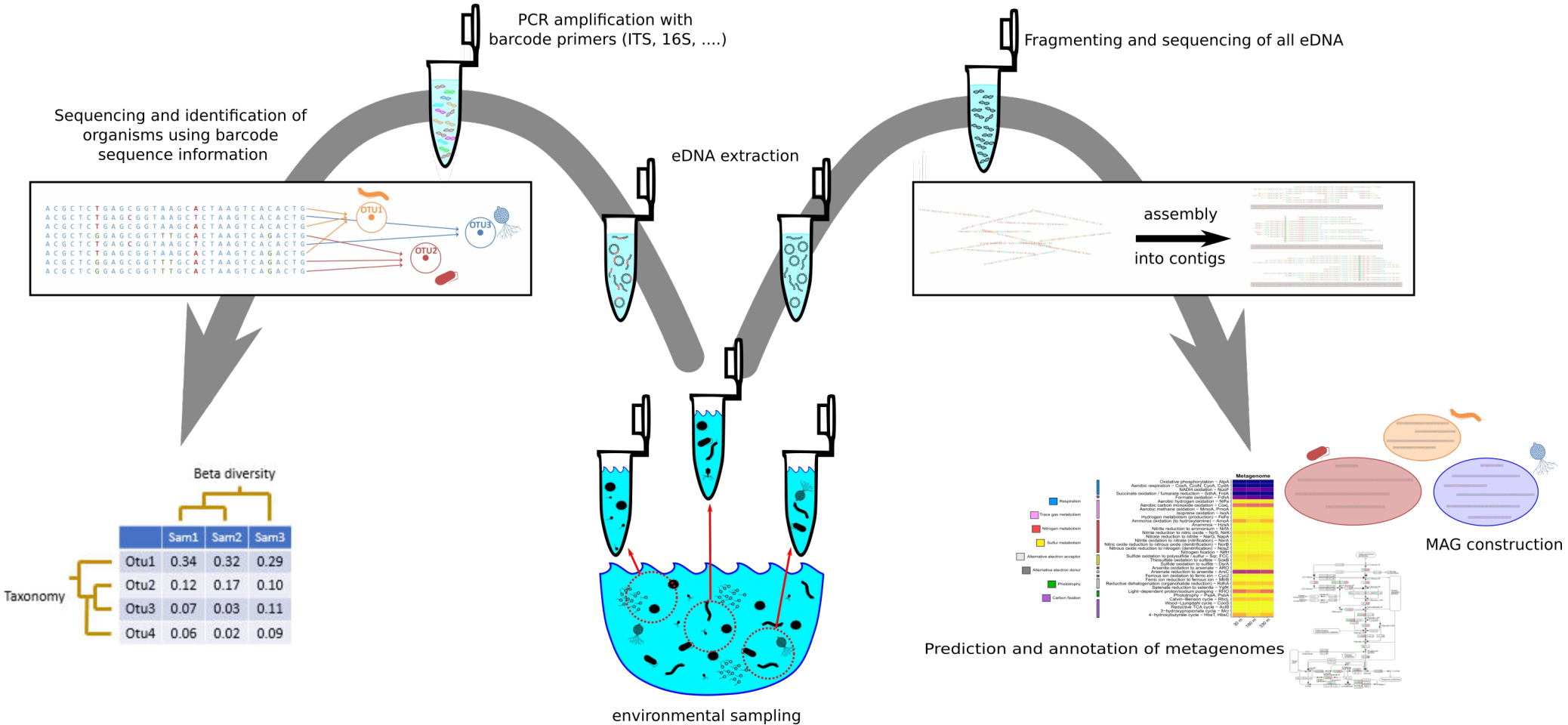Distributed computing and compute clusters are essential tools for scientists.

We use our local computers to talk to shared remote servers that have lots of high-powered processors, and storage and RAM.
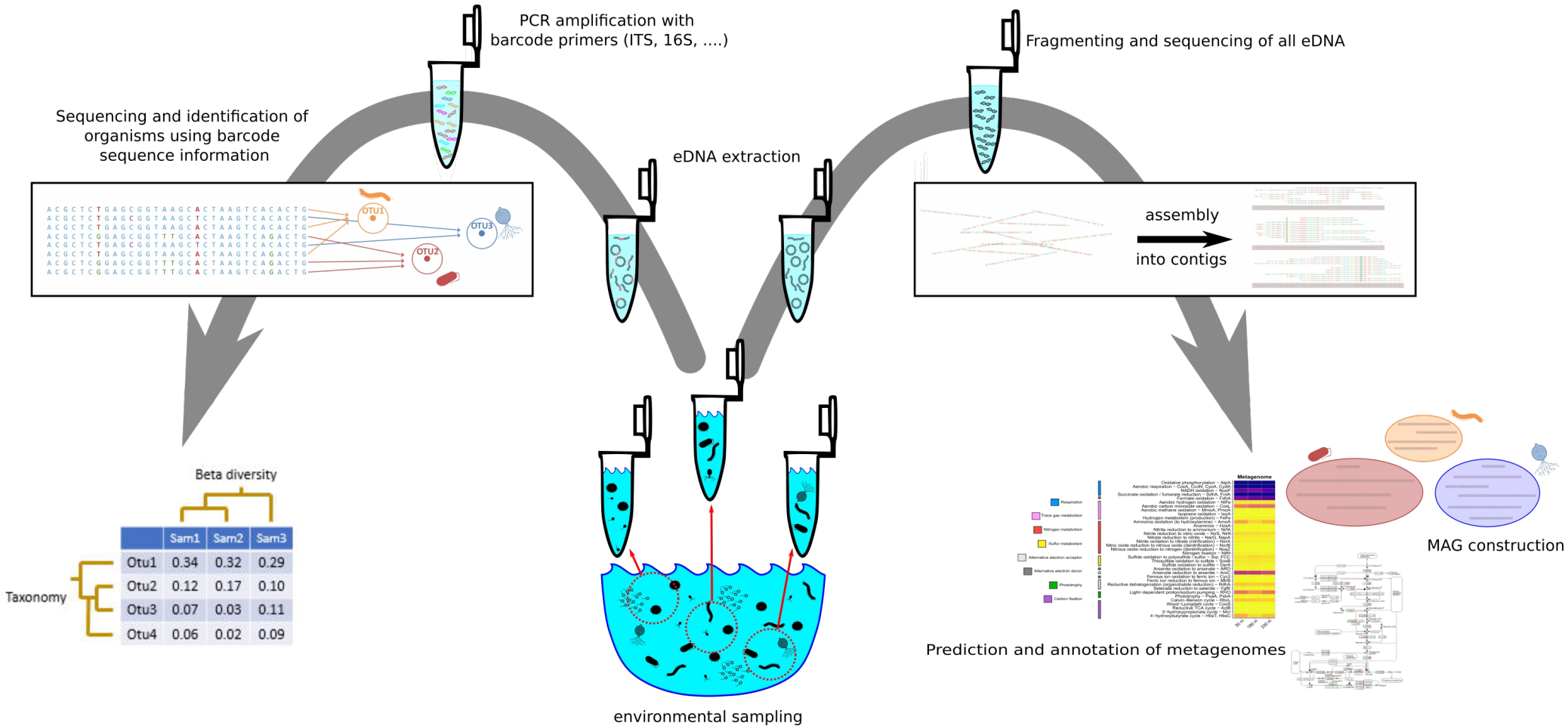
RAM (lots!)

Processors (lots!)

Block device memory (lots!)

Server-side hardware

de.NBI / Universität Bielefeld

RAM 128G    CPU 28 cores    Storage 1000TB    RAM 128G    CPU 28 cores    Storage 1000TB    RAM 128G    CPU 28 cores    Storage 1000TB

VM Linux OS    VM Linux OS    VM Linux OS

SSH & SCP    SSH & SCP    SSH & SCP

Linux OS    Linux OS    Linux OS

Us.

https://actusdigital.com/2018/01/08/virtual-machine-or-physical-server/

Review – what do we do with DNA sequence data?

# Why do we need distributed computing to do this?



PCR amplification with barcode primers (ITS, 16S, ....)

Fragmenting and sequencing of all eDNA

Sequencing and identification of organisms using barcode sequence information

eDNA extraction

assembly

into contigs

Beta diversity

|      | Sam1 | Sam2 | Sam3 |
|------|------|------|------|
| Otu1 | 0.34 | 0.32 | 0.29 |
| Otu2 | 0.12 | 0.17 | 0.10 |
| Otu3 | 0.07 | 0.03 | 0.11 |
| Otu4 | 0.06 | 0.02 | 0.09 |

Taxonomy

MAG construction

Prediction and annotation of metagenomes

environmental sampling

Because some steps take ~100 gig of RAM (or more!) and lots of processors.

PCR amplification with barcode primers (ITS, 16S, ....)

Fragmenting and sequencing of all eDNA

Sequencing and identification of organisms using barcode sequence information

eDNA extraction

assembly

into contigs

Beta diversity

Metagenome

MAG construction

Taxonomy

| | Sam1 | Sam2 | Sam3 |
|---|---|---|---|
| Otu1 | 0.34 | 0.32 | 0.29 |
| Otu2 | 0.12 | 0.17 | 0.10 |
| Otu3 | 0.07 | 0.03 | 0.11 |
| Otu4 | 0.06 | 0.02 | 0.09 |

Prediction and annotation of metagenomes

environmental sampling

## Dataset 1: ZymoBIOMICS HMW DNA Standard.



| Species | Theoretical Composition (%) | | | | |
|---|---|---|---|---|---|
| | Genomic DNA | 16S Only[1] | 16S & 18S[1] | Genome Copy[2] | Cell Number[3] |
| *Pseudomonas aeruginosa* | 14 | 5.1 | 4.6 | 7.8 | 7.9 |
| *Escherichia coli* | 14 | 12.4 | 11.2 | 10.9 | 11.0 |
| *Salmonella enterica* | 14 | 12.7 | 11.4 | 11.2 | 11.2 |
| *Enterococcus faecalis* | 14 | 12.1 | 10.9 | 18.8 | 18.8 |
| *Staphylococcus aureus* | 14 | 19 | 17.1 | 19.6 | 19.6 |
| *Listeria monocytogenes* | 14 | 17.3 | 15.6 | 17.8 | 17.9 |
| *Bacillus subtilis* | 14 | 21.4 | 19.2 | 13.2 | 13.2 |
| *Saccharomyces cerevisiae* | 2 | NA | 10 | 0.63 | 0.32 |

OPEN

# Oxford Nanopore R10.4 long-read sequencing enables the generation of near-finished bacterial genomes from pure cultures and metagenomes without short-read or reference polishing
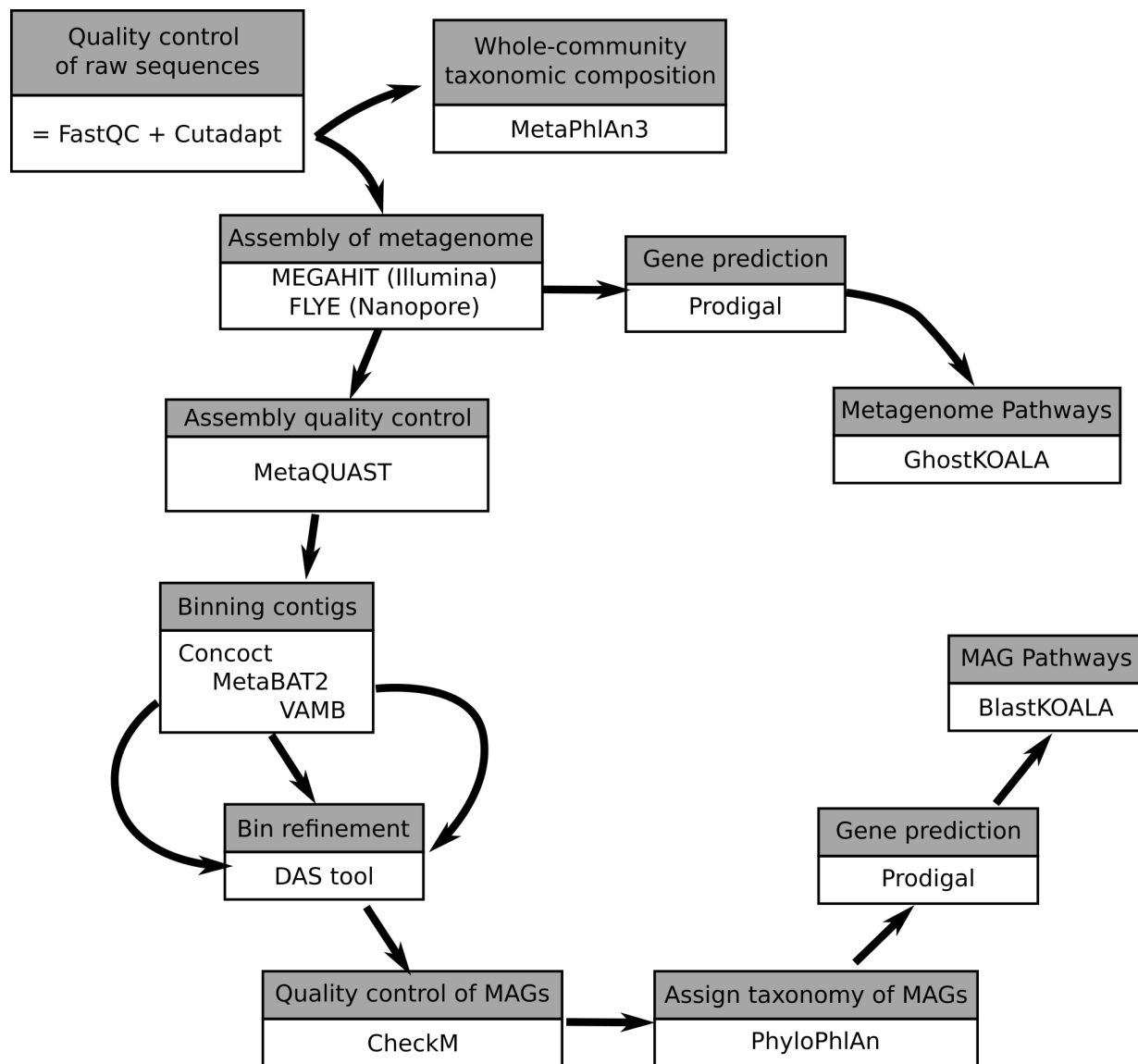
Mantas Sereika [1,4], Rasmus Hansen Kirkegaard [1,2,4], Søren Michael Karst [1], Thomas Yssing Michaelsen[1], Emil Aarre Sørensen [1], Rasmus Dam Wollenberg[3] and Mads Albertsen [1]✉

Fredericia – Centralrenseanlægget is the second largest wastewater treatment plant in Denmark

**Quality control of raw sequences**

= FastQC + Cutadapt

**Whole-community taxonomic composition**

MetaPhlAn3

**Assembly of metagenome**

MEGAHIT (Illumina)
FLYE (Nanopore)

**Gene prediction**

Prodigal

**Metagenome Pathways**

GhostKOALA

**Assembly quality control**

MetaQUAST

**Binning contigs**

Concoct
MetaBAT2
VAMB

**Bin refinement**

DAS tool

**Quality control of MAGs**

CheckM

**Assign taxonomy of MAGs**

PhyloPhlAn

**Gene prediction**

Prodigal

**MAG Pathways**

BlastKOALA

We will use various scientific software packages to process our sequence data into a metagenome and Metagenome-Assembled-Genomes (MAGs).

**Quality control of raw sequences**
= FastQC + Cutadapt

**Whole-community taxonomic composition**
MetaPhlAn3

**Assembly of metagenome**
MEGAHIT (Illumina)
FLYE (Nanopore)

**Gene prediction**
Prodigal

**Metagenome Pathways**
GhostKOALA

**Assembly quality control**
MetaQUAST

**Binning contigs**
Concoct
MetaBAT2
VAMB

**Bin refinement**
DAS tool

**MAG Pathways**
BlastKOALA

**Gene prediction**
Prodigal

**Quality control of MAGs**
CheckM

**Assign taxonomy of MAGs**
PhyloPhlAn

We will do these steps with both illumina and nanopore data, for our learning process with the Zymo Mock Community data.

## Intro to command line. Open your terminal and try these:

**pwd** – print the current working directory

**ls** – list the contents of the current directory

**cd** – change directory

**cp** – copy a file

**mv** – move or rename a file

**mkdir** – make a new directory

**echo** – print something out to the display

**cat** – concatenate two or more files

**rm** [-r] – remove a file. Essentially deletes a file forever. Be careful!!

**sudo** – execute another command that requires root privileges

## Intro to command line. Open your terminal and try these:

**top** – opens up a real time display that shows you the "busiest" processes ongoing in your machine. After you start the program, press `i` to see only active processes, press `e` to make memory units more readable, and press `1` then `t` to keep a close eye on your cores.

**man** – print the reference manual for a utility, command, or program, if available. If not, try the `--help` flag for any given program without a manual.

**head** – show top several lines from a text file

**tail** – show last several lines from a text file

**less** – interactively open a text file for reading

## Intro to command line. Open your terminal and try these:

Special symbols are very important to the linux command line.

~ home directory (try with cd)

. current directory (try with cd). Also filenames that start with "." are "hidden".

.. parent directory (try with cd)

> direct the output of a process to a file

$ variable expansion

| pipe, connects the output from one command to another

; command separator

/ directory separator (compare to Windows!)

\ escape a character, so that it is no longer a "special character"

& send a process to run in the background

= variable assignment

* wildcard for multiple characters (try with ls)

? wildcard for single character (try with ls)

# comment (place it before a command and see what happens)

## Intro to command line. Open your terminal and try these:

**top** – opens up a real time display that shows you the "busiest" processes ongoing in your machine. After you start the program, press `i` to see only active processes, press `e` to make memory units more readable, and press `1` then `t` to keep a close eye on your cores.

**man** – print the reference manual for a utility, command, or program, if available. If not, try the `--help` flag for any given program without a manual.

**head** – show top several lines from a text file

**tail** – show last several lines from a text file

**less** – interactively open a text file for reading

Special command: **wget**

We'll keep track of what we do in scripts that will be kept in an online github repository. Let's practice grabbing the latest version:

**wget** https://raw.githubusercontent.com/danchurch/FunctionalMicrobiomePractical2022/main/funmic2023/funBASHterminalScript.txt

(Maybe better to follow the link in the skript PDF on page 20, or download this presentation, then copy-paste)

| First task: connect to VMs | |
| --- | --- |
| Milan | ssh -p 30284  -i /path/to/your/privateKey ubuntu@129.70.51.6 |
| Julia | ssh -p 30381  -i /path/to/your/privateKey ubuntu@129.70.51.6 |
| Maria | ssh -p 30355  -i /path/to/your/privateKey ubuntu@129.70.51.6 |
| Anja | ssh -p 30452  -i /path/to/your/privateKey ubuntu@129.70.51.6 |
| Susanne | ssh -p 30461  -i /path/to/your/privateKey ubuntu@129.70.51.6 |
| Rebekka | ssh -p 30372  -i /path/to/your/privateKey ubuntu@129.70.51.6 |
| Or try with a shell variable: | key=/home/daniel/.ssh/funmic2023<br>ssh -p 30500 -i $key ubuntu@129.70.51.6 |

# First task: connect to VMs



**cd /vol/danBot**

## Once you are in, try to find your storage volume.

# Second task: get a file from your VM to your computer



You will often need to get files from your de.NBM VM, to inspect them locally either in your linux or windows evironment.

de.NBI virtual machine → Local linux environment → Windows → Browsers, Excel, other GUIs…

# Second task: get a file from your VM to your computer



We will use the **scp** program for this.

Second task: get a file from your VM to your computer

```
fileOnVM="/vol/danBot/helloWorld.txt"
mySSHkey="/home/daniel/.ssh/funmic2023"

scp -i $mySSHkey -P 30500 ubuntu@129.70.51.6:$fileOnVM .
```

Try **scp** , to get "helloWorld.txt" from your de.NBI VM to your local machine.

de.NBI virtual machine ⟹ Local linux environment ⟹ Windows ⟹ Browsers, Excel, other GUIs…

Second task: get a file from your VM to your computer

```
fileOnVM="/vol/danBot/helloWorld.txt"
mySSHkey="/home/daniel/.ssh/funmic2023"

scp -i $mySSHkey -P 30500 ubuntu@129.70.51.6:$fileOnVM .
```

(Can copy and paste from scripts, but make sure you understand!)

de.NBI virtual machine → Local linux environment → Windows → Browsers, Excel, other GUIs…

Second task: get a file from your VM to your computer

Once you have your file in your local linux environment, you will probably want to look at it in your windows environment, too.

| de.NBI virtual machine | ⇒ | Local linux environment | ⇒ | Windows | ⇒ | Browsers, Excel, other GUIs… |

Second task: get a file from your VM to your computer

If you are using WSL in windows, the try placing the file here:

cd /mnt/c/

See if you can find it in windows by looking at your C:\ directory. Considering making a special folder for this kind of file transfer.

de.NBI virtual machine → Local linux environment → Windows → Browsers, Excel, other GUIs…

Last task of the day: start your nanopore assembly.



Assembly of metagenomes is one of the most computationally intensive tasks we will ask our de.NBI machines to do for us. We will review metagenome assembly a bit romorrow, but let's let's get our first assembly started now, so it can as long as it needs to.

This will be the first of two assemblies we do, since we have both an illumina and a nanopore dataset from the same Zymo mock community. We will start with the nanopore data, since it takes much more time and computing resources.

Last task of the day: start your nanopore assembly.



This also gives a chance to see a complex command in action, that requires us to use much of what we just learned about the linux command line.

Last task of the day: start your nanopore assembly.

From my de.NBI machine, the command looks like this:

```
## get into the right conda environment
conda activate flye

## make a directory for our new assembly:
mkdir -p /vol/danBot/assemblies/zymoMC/nanopore/

## set my variables
reads=/vol/danBot/datasets/zymoMC/nanopore/ERR7287988_shortened.fastq
outdir=/vol/danBot/assemblies/zymoMC/nanopore/


## run it, using nohup to make sure it doesn't die when we hang up the connection
nohup flye --meta \
    --threads 25 \
    --out-dir $outdir \
    --nano-hq $reads &

## this is the same as:
nohup flye --meta --threads 25 --out-dir $outdir --nano-hq $reads &
```

(You can see the full help file by typing flye --help)

Last task of the day: start your nanopore assembly.

Did that work?



Let's check with the top program.

And final note for the day: keep a script of what you do!!!



With scientific computing, it can sometimes take hours to build, debug, and successfully execute a single program, sometimes just one line of code.

And final note for the day: keep a script of what you do!!!



You will often want to return to your commands days or months later. Other scientists will want to know precisely what commands and settings you used.

And final note for the day: keep a script of what you do!!!



notepad-plus-plus.org

www.vim.org

www.gnu.org/software/emacs/tour/

Keep a running text-file script of all that you do. Comment it with your own mental notes so you can remember what you did and why in human language.
Pick a good **text editor** (not word processor!!) and learn to use it.

**And final note for the day: keep a script of what you do!!!**



https://xkcd.com/378