# Genetic data and analysis on a cluster

FEB 2021

# Overview

I. Data
   1. Format
   2. IBD project

II. Cluster
   1. Uliege cluster
   2. Run an analysis
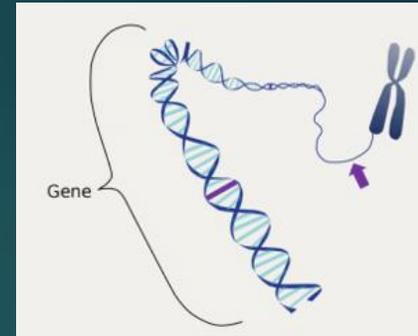
# I. Data
# 1. Format

PLINK: free open source command-line program for genomic analysis

Download: https://www.cog-genomics.org/plink/

Plink formats: https://www.cog-genomics.org/plink/1.9/formats#bed

- ▶ Bed: representation of genotype calls at biallelic variants, so the markers/SNPs. The file can't be open.

- ▶ Bim: variant information file (Chromosome code, variant identifier, alleles…)

- ▶ Fam: Sample information file (family ID, sex code, phenotype value…)

-> go together and represent the entire dataset.

# I. Data
# 1. Format

2. Obtain via –recode, loaded via -file

▶ Ped: The first six fields are the same as .fam. Then, variant information.

▶ Map: variant information file accompanying a .ped (chromosome code, variant identifier, …)

-> go together and represent the entire dataset.

# I. Data
# 1. Format

Basic plink functions for input filtering:

- removes all unlisted samples: --keep

- Remove all listed samples: --remove

- Extract a subset of SNP based on chromosomes: --chr

- removes all unlisted variants: --extract

- removes all listed variants: --exclude

- Linkage disequilibrium: --indep-pairwise

- Minimum allele frequency= --maf

- …

# I. Data
# 1. Format

Work with R from PLINK files:

❑ Change the format of the files using PLINK software so R can import them:

from bed, bim, fam to map and ped using option –recode

raw from –recodeA

-> can be read in R but will be huge

❑ Use specific R functions, for example read.bed()

# I. Data
# 1. IBD

Projects: Detect epistasis with multiple tools and same dataset

IBD: Inflammatory Bowel Disease.

Two main Datasets:

Same 66,280 individuals (~50% cases, +50% controls)

Same initial quality controls (LD, MAF, HWE…)

► Unfiltered

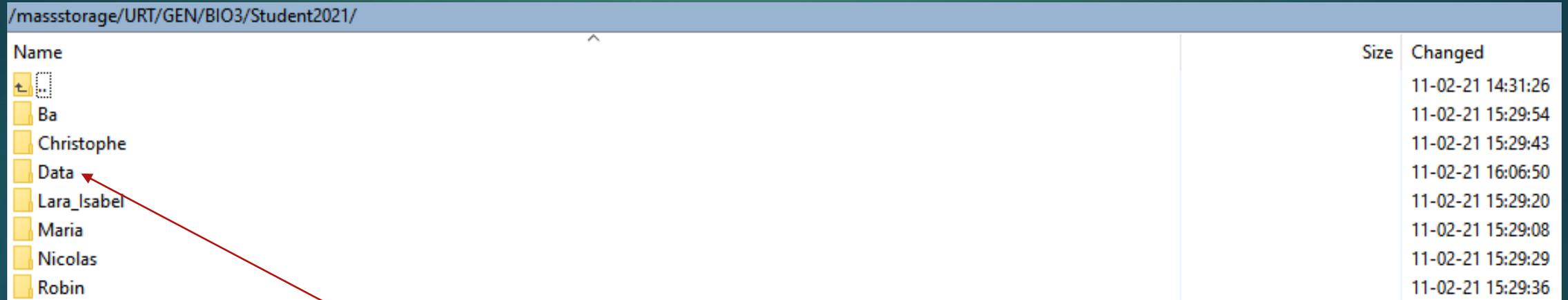► Functional: biological filters

# I. Data
# 1. IBD

For the 2 datasets, multiple variations: specific requirement of analysis

➢ More SNP filters (relief and epiblaster) for analysis that can't handle large amount of SNPs

➢ Imputation (knn) for analysis that can't handle missing values

➢ Phenotypes:
   ▪ continuous
   ▪ binary

# I. Data
# 1. IBD

Folder structure

/massstorage/URT/GEN/BIO3/Student2021/

| Name | Size | Changed |
|---|---|---|
| .. | | 11-02-21 14:31:26 |
| Ba | | 11-02-21 15:29:54 |
| Christophe | | 11-02-21 15:29:43 |
| Data | | 11-02-21 16:06:50 |
| Lara_Isabel | | 11-02-21 15:29:20 |
| Maria | | 11-02-21 15:29:08 |
| Nicolas | | 11-02-21 15:29:29 |
| Robin | | 11-02-21 15:29:36 |

Input

# I. Data
# 1. IBD

Folder structure

/massstorage/URT/GEN/BIO3/Student2021/Data/

Name

📁 [..]
📁 GeneInformation
📁 Phenotypes
📁 SNP_to_gene_mapping
📁 SNPs
📄 README.txt

Two main SNP sets

/massstorage/URT/GEN/BIO3/Student2021/Data/SNPs/
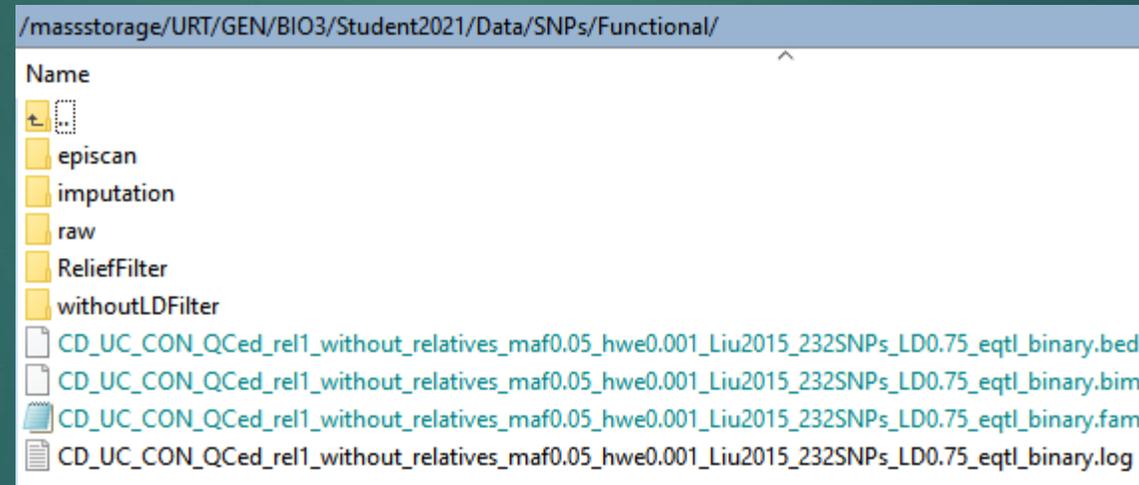
Name

📁 [..]
📁 Functional
📁 Unfiltered

# I. Data
# 1. IBD

Folder structure

Example: Functional dataset



Available options for specific requirements:
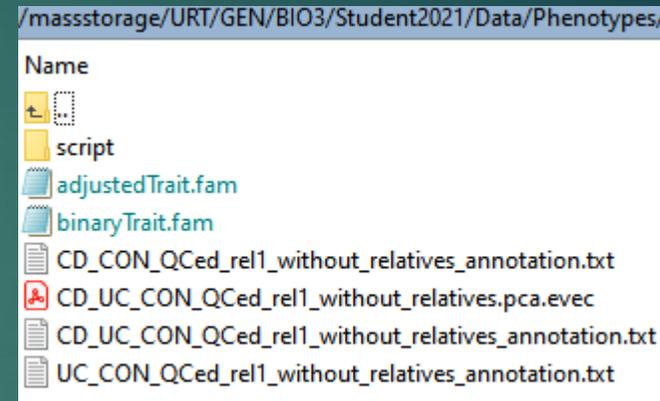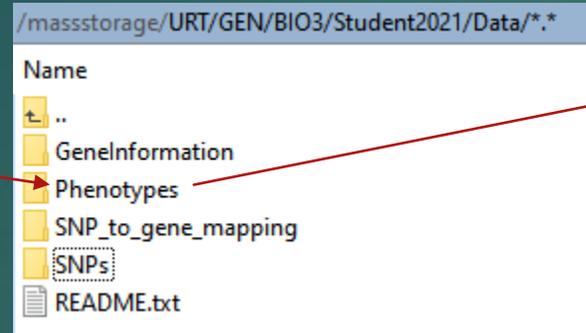    Imputation of 1. and 2.
    Reduction of the dataset via episcan and relief

# I. Data
## 1. IBD

Folder structure

Phenotypes



(Order is important)
1. If your tool can adjust for covariates: binary phenotypes and adjust for the first 7 PCs.
2. Else, if your tool can't include covariates but can handle continuous phenotypes: continuous phenotypes that are already adjusted for first 7 PCs
3. Else, if your tool can't include covariates and can't handle continuous phenotypes: binary phenotypes.

# I. Cluster
## 1. Uliege cluster

What is a cluster?

Set of connected computers that work together.

Why are we using a cluster?

- Big dataset, big analysis -> improve performance and availability
- Legal agreement

Advice:

Create and try your code on a small dataset* on your own computer. Then, run the real analysis on the cluster once you made sure your code is ok.

Why: Easier and faster to find errors.

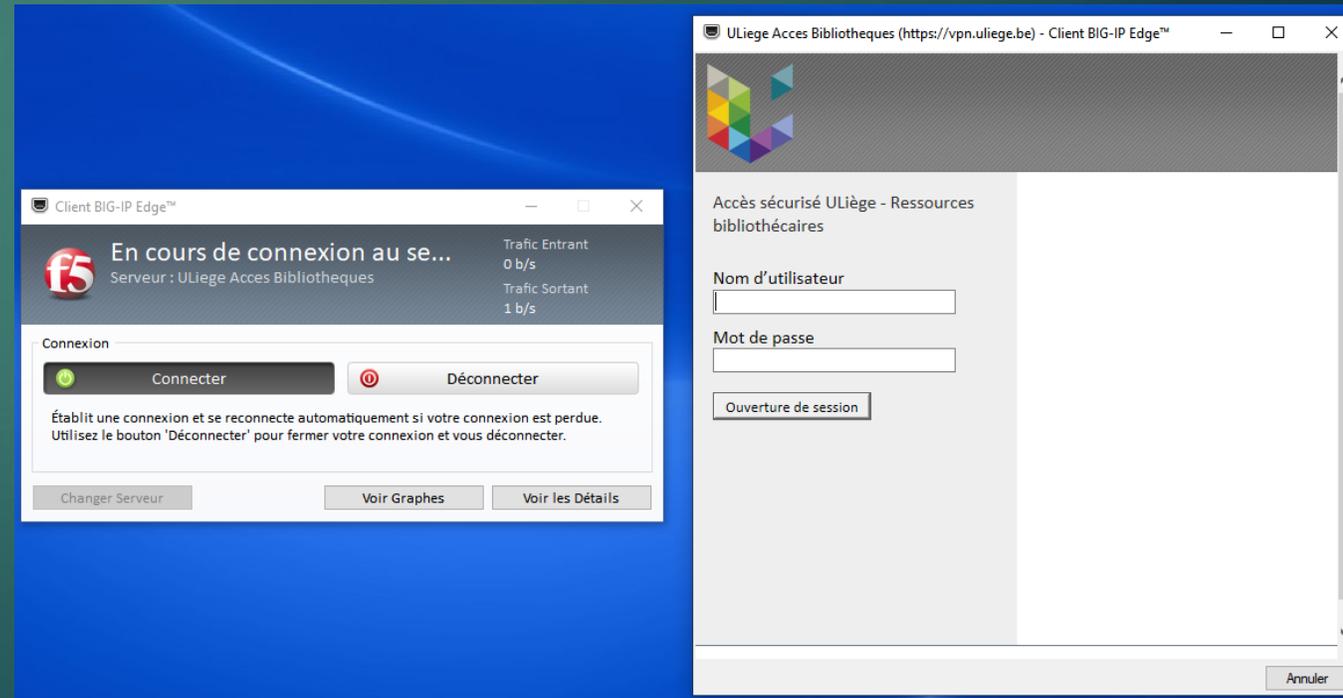*a public dataset, not the IBD one which can't be downloaded.

# I. Cluster
# 1. Uliege cluster

**Connect** to the cluster

If not onsite (wifi of university of Liège), download the VPN:

https://my.segi.uliege.be/cms/c_116507 35/fr/mysegi-new-vpn
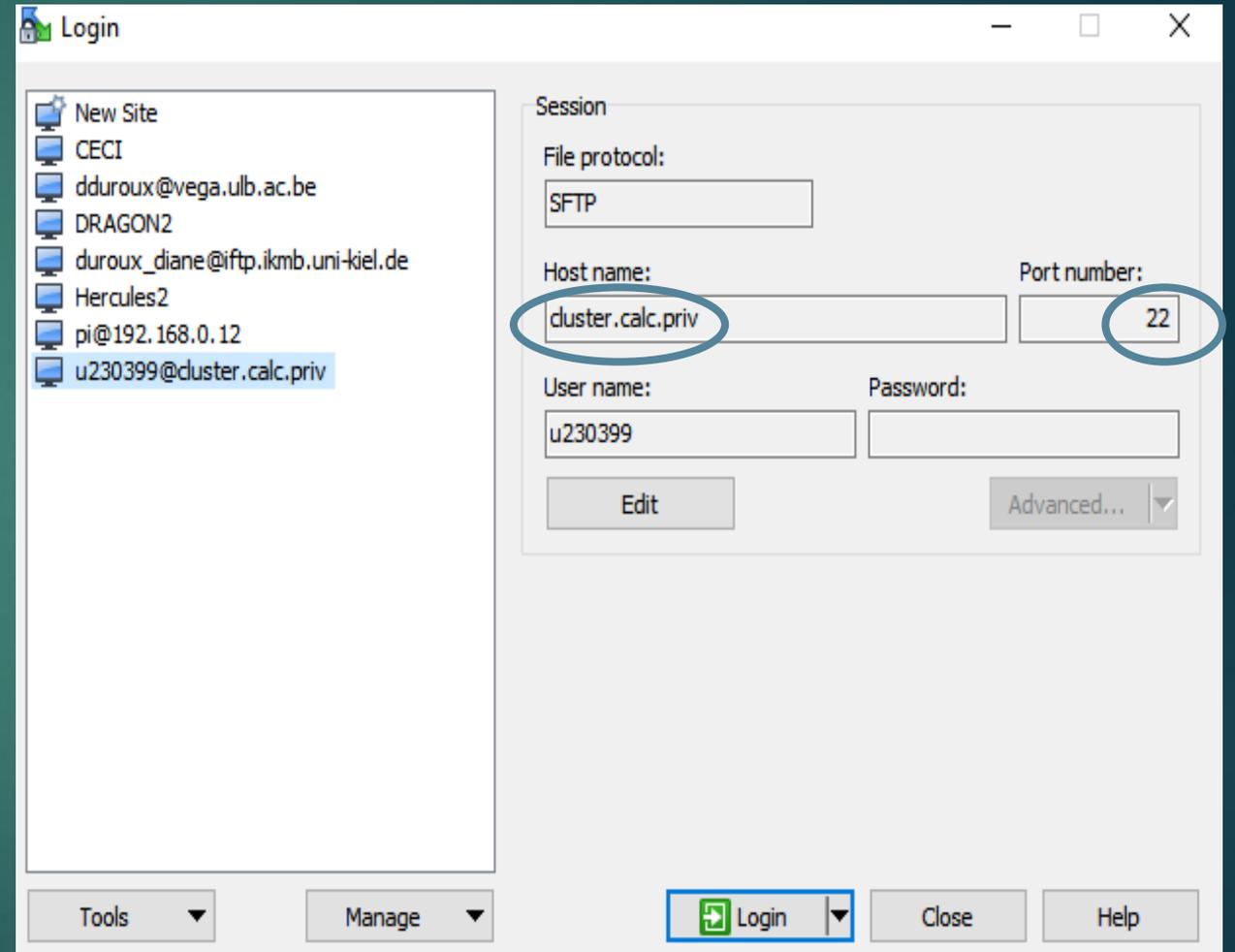
Enter your id and passwork to connect

# I. Cluster
## 1. Uliege cluster

Access and **visualize** your individual folder and the data

Windows:
https://winscp.net/eng/index.php

# I. Cluster
# 1. Uliege cluster

Access and **visualize** your individual folder and the data

Windows:
https://winscp.net/eng/index.php

# I. Cluster
## 1. Uliege cluster

Access and **visualize** your individual folder and the data

Windows:
https://winscp.net/eng/index.php

Linux: ssh command
https://docs.oracle.com/en/cloud/paas/big-data-cloud/csbdi/connect-cluster-node-secure-shell-ssh.html#GUID-E6F4421D-3D7F-415B-ABD6-D3CC0C870947

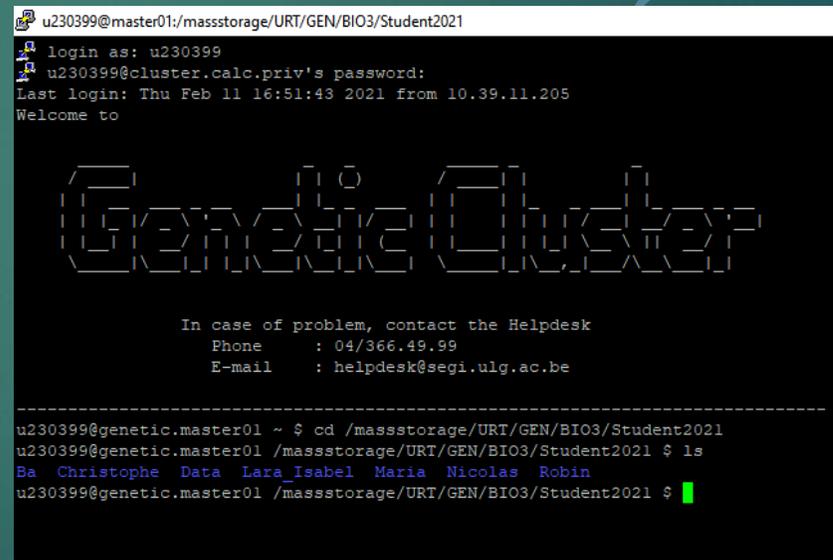# 1. Cluster
# 2. Run an analysis

How to **communicate** with the cluster (sofware)

▶ Windows: puTTY software, open source SSH client

[https://www.putty.org/](https://www.putty.org/)



▶ Linux: directly in terminal

# 1. Cluster
# 2. Run an analysis

How to **communicate** with the cluster (language)

Slurm: ressource manager / job scheduler

Goal: organize ressource sharing on a supercomputer

How: Users submit jobs, which are scheduled and allocated resources (CPU time, memory, etc.)

# 1. Cluster
# 2. Run an analysis

Basic commands to navigate in your folders and check your files (bash, shell):

▶ Cd folderName: change directory (go into another directory)

▶ Ls: display what's in a directory

▶ Head fileName: See the top of the file

▶ Tail filename: see the end of the file

▶ wc –l: count the number of rows in a file

▶ du –sh folderName: check the size of a folder

▶ rm fileName: delete a file

▶ mkdir folderName: create a new folder

▶ …

More info: https://www.educative.io/blog/bash-shell-command-cheat-sheet

# 1. Cluster
# 2. Run an analysis

**Run** an analysis / a script / a job on the cluster:

Create a .sh file (for example: run.sh).

This file has a specific structure so the cluster understands whit it needs to do

Header: must to start with #
Specify the ressource required

Load softwares needed

Analysis: here call an external R script

```
#!/bin/bash
#SBATCH --ntasks=1 #each job has one task
#SBATCH --cpus-per-task=1 # each task uses 1 cpu
#SBATCH --partition=kosmos
#SBATCH --mem-per-cpu=8000 #8GB


module load R/3.2.4


R CMD BATCH pathToFile/FileName.R
```

# 1. Cluster
# 2. Run an analysis

Header is very important:

If too much ressources asked: will never start

If not enough: job will stop before the end


Need to investigate the ressources needed: time, nb of CPUs…


Example:

```
#!/bin/bash
#SBATCH --ntasks=1
#SBATCH --cpus-per-task=6
#SBATCH --partition=kosmos
#SBATCH --mem-per-cpu=8000
--time=01:00:00
```

Always required

Number of core per task

Each task uses 6 cpus

Select a partition

8GB required

Time limit for the job.

More info: https://ubccr.freshdesk.com/support/solutions/articles/5000688140-submitting-a-slurm-job-script

# I. Cluster
# 2. Run an analysis

Some basic slurm commands:

- Submit/start a job: sbatch pathToFile/FileName.sh.

  What does it do: You ask permission to run a job on the cluster.

  If resources are available, it will start.

  If not, it will wait in the queue until enough resources are available.

- Ask if a program is running or pending: squeue -u yourUsername

- Get more info about the cluster: sinfo

- Stop a job: scancel jobNumber

  More info: https://support.ceci-hpc.be/doc/_contents/QuickStart/SubmittingJobs/SlurmTutorial.html

# 1. Cluster
# 2. Run an analysis

Tips:

▶ Tests or debugging:

  Slurm jobs are normally batch jobs: they are run unattended.

  If you want to have a direct view on your job, run: srun –pty bash

▶ NEVER work on the master node of the cluser (ie without srun or sbatch)

▶ Always google your problem (stackOverflow, mathOverflow, …)

# Thank you