# Hands-on with S. *cerevisiae* (~12 Mb, 16 chromosomes)

You will split in two teams (or more).

Your missions is to perform, compare and give information about different assemblies :
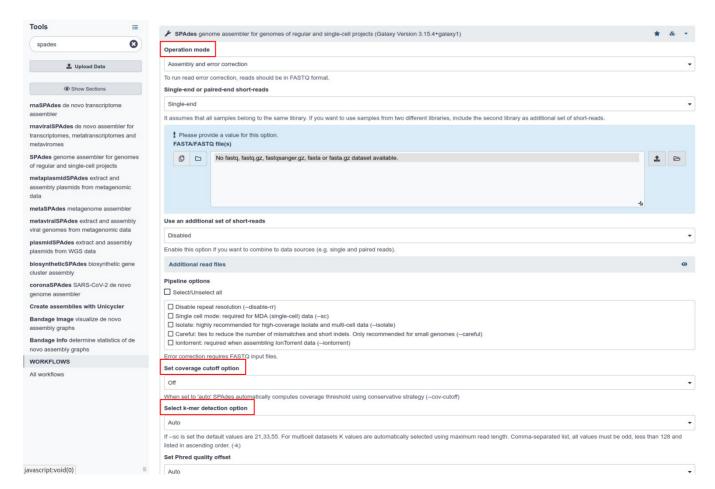
- using different data types (short reads (MiSeq PE), ONT reads (R9.4), HiFi reads),
- different coverages (from 5X to 60X),
- different software packages (SPAdes, flye, wtdbg2, Hifiasm).

You will compare the results using quast and dgenies (https://dgenies.toulouse.inra.fr/).
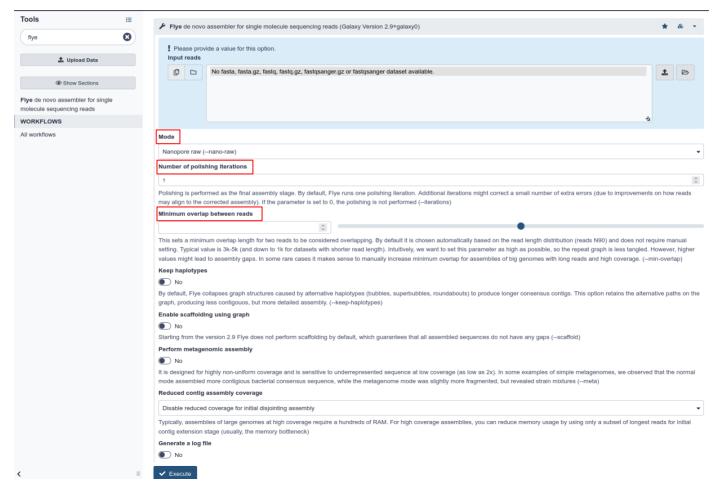
You will fill the table : https://lite.framacalc.org/f952lz03mb-9vey

You will present you findings (impact of data types, coverage, software package,...) : 5 minutes per team
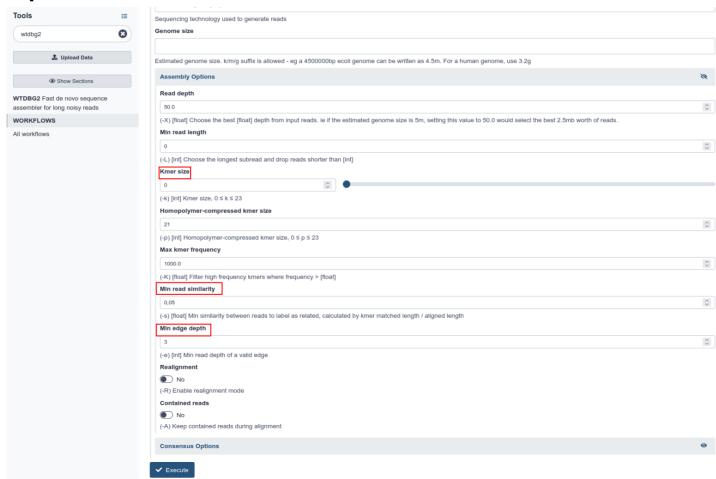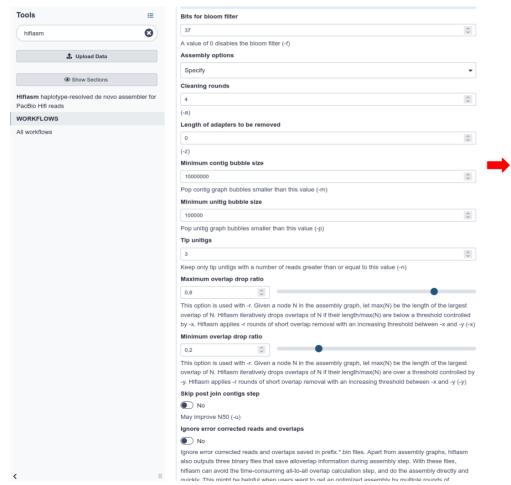
# SPAdes parameters
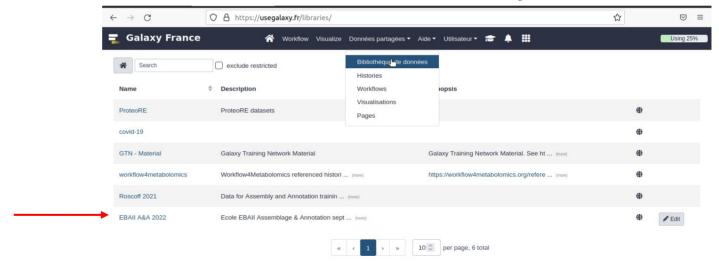
# Flye parameters

# wtdbg2 parameters

# Hifiasm parameters



5

# And

Processing will be performed using : https://usegalaxy.fr/

The data files are located at : shared libraries/EBAII A&A 2022/Assembly



You can ask question using : https://semestriel.framapad.org/p/abgi5b9vdm-9vey?lang=fr

# Let's go!