# YARN & MapReduce 2

Lucas BAKALIAN & Michaël HATOUM

August 15, 2021
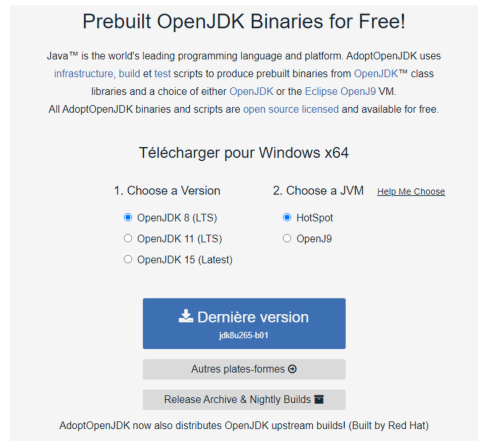
# 1   MapReduce JAVA

In this tutorial, we will use write MapReduce jobs using JAVA. Before that, we will install prerequirements, if you already have them, skip thoses parts.
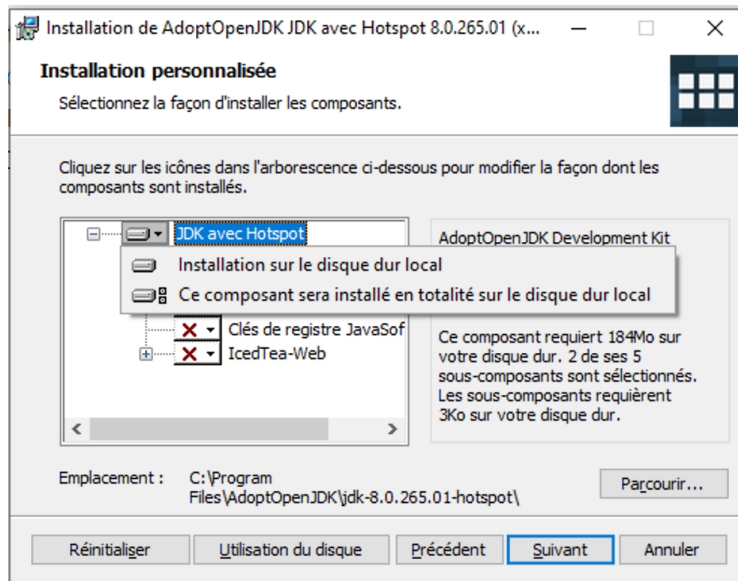
## 1.1   Install OpenJDK 8

### 1.1.1   Microsoft Windows

Download and install OpenJDK for Windows from here.



Start the installer and configure it to install everything:



### 1.1.2   Mac OS X

It is easier to install OpenJDK on OS X:

```
$ brew tap AdoptOpenJDK/openjdk
$ brew cask install adoptopenjdk8
```

## 1.2   Git

### 1.2.1   Microsoft Windows

Download and install 64-bit Git for Windows Setup from here.
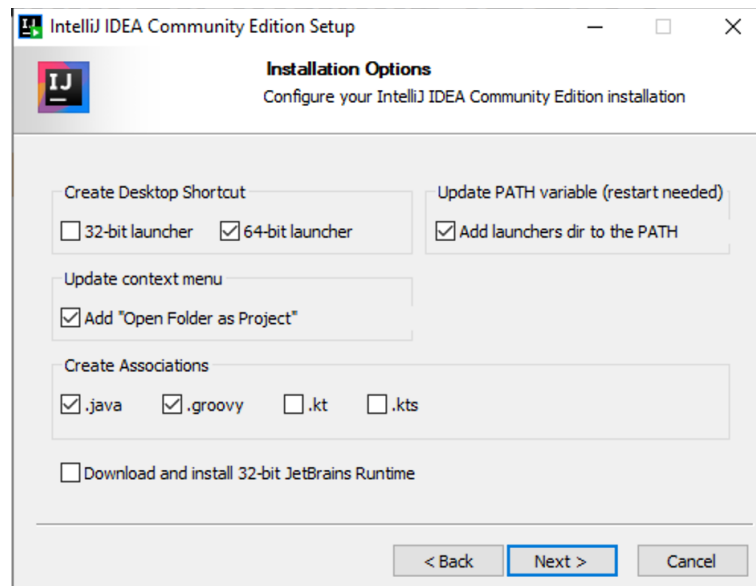
### 1.2.2   Mac OS X

There are several ways to install Git on a Mac. The easiest is probably to install the Xcode Command Line Tools. On Mavericks (10.9) or above you can do this simply by trying to run git from the Terminal the very first time.

```
$ git --version
```

## 1.3   IntelliJ IDEA

### 1.3.1   Microsoft Windows

Download and install IntelliJ IDEA for Windows (Community edition) from here. Here the installation options and restart your computer when asked.



### 1.3.2   Mac OS X

Download and install IntelliJ IDEA for Mac OSX (Community edition) from here. Open the .dmg folder and move the IDEA app into your applications folder.
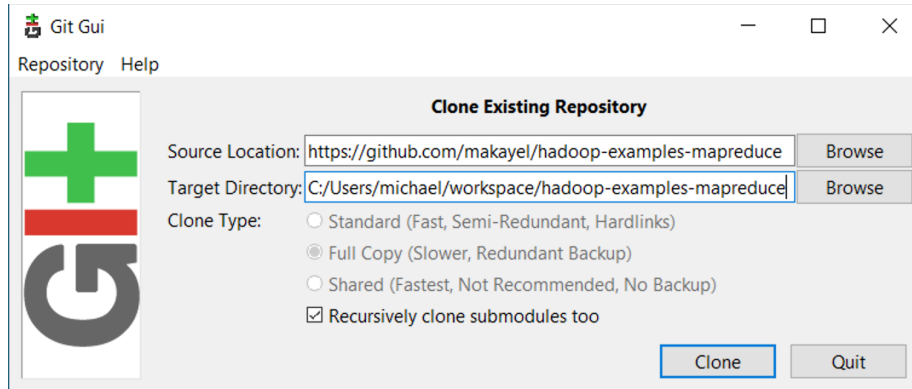
## 1.4   Clone the project

Here is a JAVA project hosted on github to start your lab.

### 1.4.1   Microsoft Windows

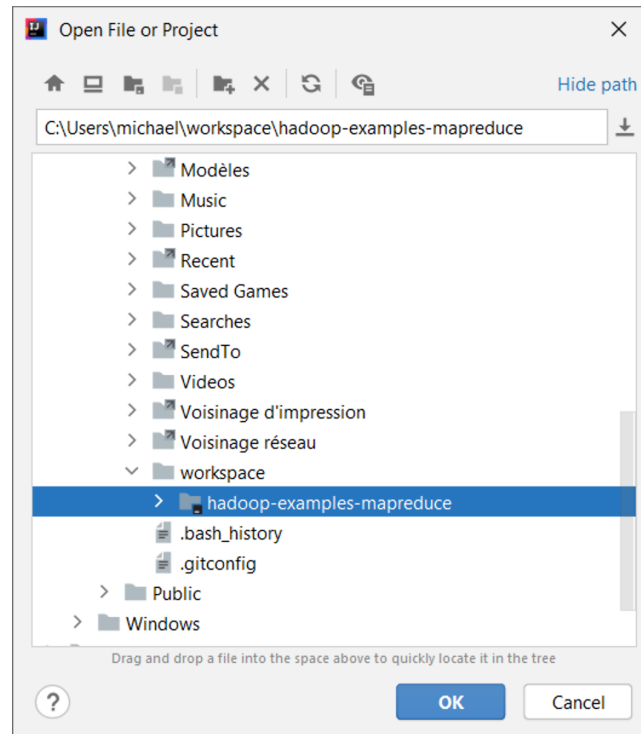You can use Git GUI for cloning the repository.



### 1.4.2   Mac OS X
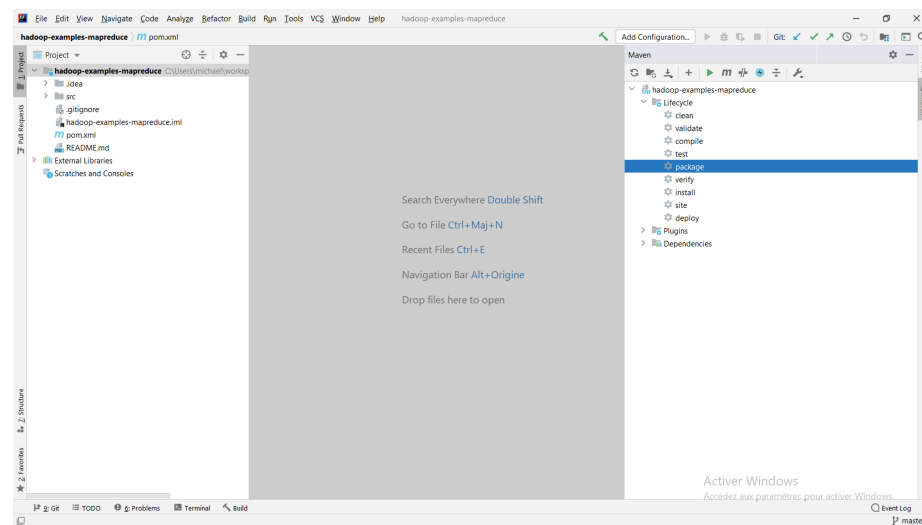
Use CLI to clone the repository.

```
$ cd /Users/michael/workspace/IdeaProjects
$ git clone https://github.com/makayel/hadoop-examples-mapreduce
```
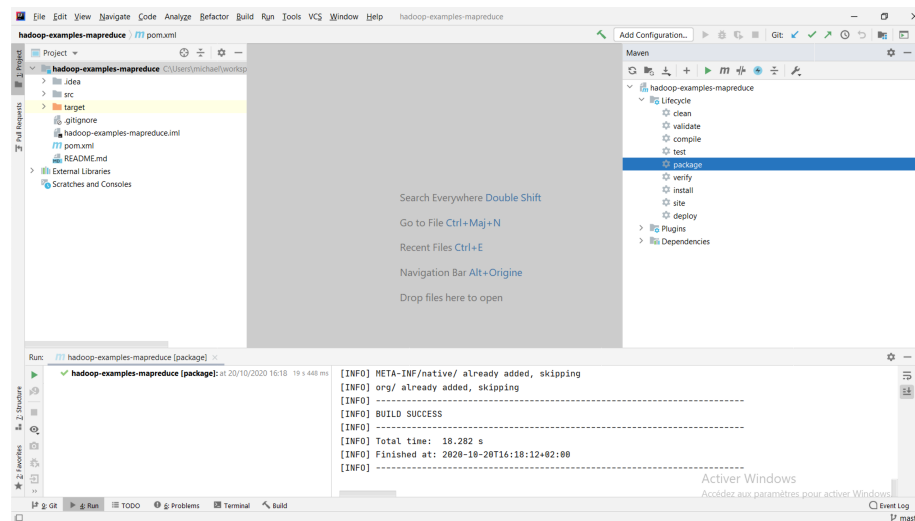
## 1.5   Import the project

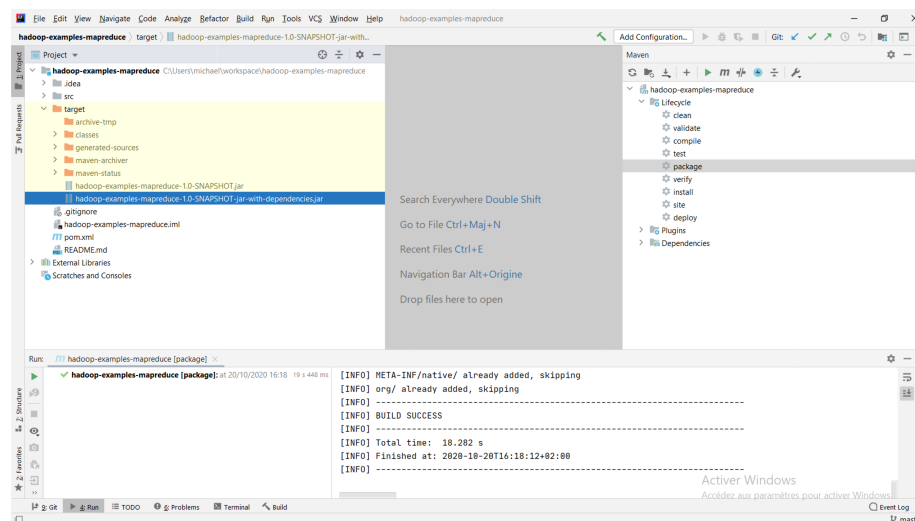Open a new project (import) from the source cloned.



Generate the JAR using maven lifecycle `package`.

On the bottom of the window, you will see the building process, wait for an `[INFO] BUILD SUCCESS` message.



You will see a new folder `package`, inside this folder you will find the JAR. We will use the JAR `*-with-dependencies.jar`.

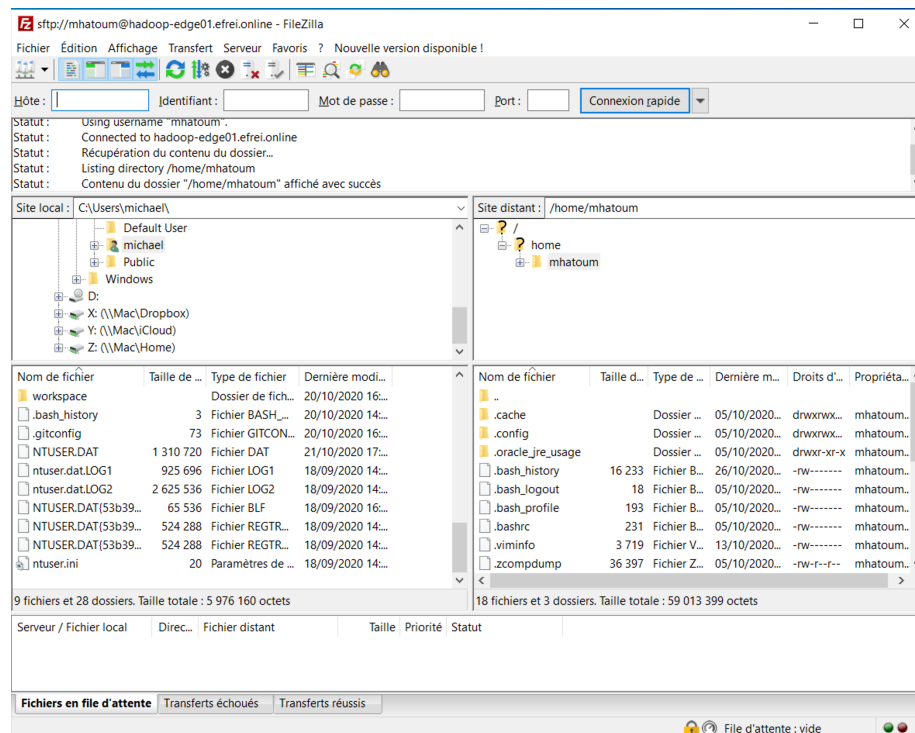## 1.6   Send the JAR to the edge node

### 1.6.1   Microsoft Windows

Download and install FileZilla client for Windows from here. Connect to the edge node using thoses parameters:

- *Host*: `sftp://hadoop-edge01.efrei.online`

- *Username*: `ssh username`

- *Password*: `ssh password`

- *Port*: `22`
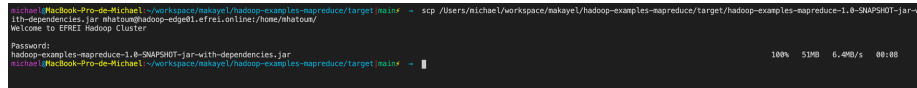
You can slide files from left (your directories) to the right (edge node).

### 1.6.2   Mac OS X

On Unix based operating system, you can use the `scp` command. Here an example:

```
$ scp myjar.jar username@hadoop-edge01.efrei.online:/home/username/
```



### 1.6.3   Run the job

Start the `wordcount` job using the JAR that we send to the edge node. Referer to previous labs for the command.

## 1.7    How will the lab work?

### 1.7.1    Important 1

This lab will last at least 2 weeks and there are all kinds of exercises, easy and hard. So start with all the easy exercises, you will do the hard ones next time.

If you start with the difficult ones and you have no experience, you are going to have problems. *You do as you want, it's up to you.*

For each of these projects, if you can't do everything that is asked for, try to get closer. The scoring takes into account your attempts.

### 1.7.2    Important 2

Don't forget to write unit test if you want bonus points. Here some documentation about it.

### 1.7.3    Important 3

If you don't have a GitHub account, you must create one. You must clone the base repository following this tutorial (keep your repository public).

At the end of each lab session, you must add on Moodle a file with your commit ID:

```
$ <name1>-<name2>-<commit-id>.md
```

## 1.8   Remarkable trees of Paris

You are going to write some *MapReduce jobs* on the remarkable trees of Paris using this dataset.

Download the file and `put` it in your HDFS home directory.

*Remember to ignore the first row in every mapper on this dataset.*

### 1.8.1   Districts containing trees (very easy)

Write a MapReduce job that displays the list of distinct containing trees in this file. Obviously, it's twenty or less different arrondissements, but exactly how many?

You just need to put rounding as a key and put any value or NullWritable as a value, output from the mapper.

The reducer will just have to output the keys and values it receives, it doesn't even have to loop through the values since it ignores them.

### 1.8.2   Show all existing species (very easy)

Write a MapReduce job that displays the list of different species trees in this file.

### 1.8.3   Number of trees by kinds (easy)

Write a MapReduce job that calculates the number of trees of each kinds.

For example, there are 3 Acer, 19 Platanus, etc. How will you define the keys and values passed between the mapper and the reducer?

The mapper must extract the kind of tree. The reducer gets the pairs (`keyI`, `valueI`) with the same key, so this key must be the kind of tree; the value being the number of such trees.

### 1.8.4   Maximum height per kind of tree (average)

Write an MapReduce job that calculates the height of the tallest tree of each kind.

For example, the tallest Acer is 16m, the tallest Platanus is 45m, etc.

### 1.8.5   Sort the trees height from smallest to largest (average)

Write an MapReduce job that sort the trees height from smallest to largest.

### 1.8.6   District containing the oldest tree (difficult)

Write a MapReduce job that displays the district where the oldest tree is.

The mapper must extract the age and district of each tree.

The problem is, this information can't be used as keys and values (why?). You will need to define a subclass of Writable to contain both information.

The reducer should consolidate all this data and only output district.

### 1.8.7   District containing the most trees (very difficult)

Write a MapReduce job that displays the district that contains the most trees.

The problem is that the program will almost certainly display a list of pairs (district number, number of trees) not ordered by number.

If we apply this program to real big data not limited to the arrondissements of Paris, we would recover a huge list, unusable and the classification could take hours. How do you keep only the best answer?

The solution probably goes through a second MapReduce phase, in which the Mapper retrieves the pairs (district, number) from the first phase, and makes them pairs whose key is unimportant (NullWritable) and the values themselves are the same input pairs.

The Reducer receives them all and must keep the best pair.