



Institut de Recherche
pour le Développement

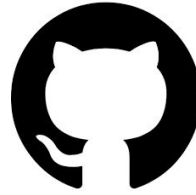
South Green
bioinformatics platform



plateau i-trop



www.southgreen.fr



<https://github.com/SouthGreenPlatform>



The South Green portal: a comprehensive resource for tropical and Mediterranean crop genomics, Current Plant Biology, 2016

Trainings 2018

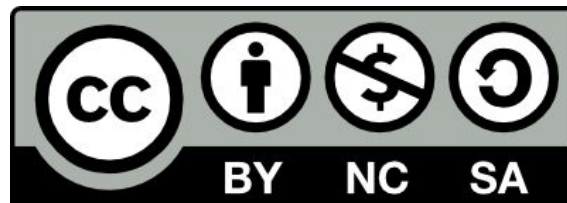


- TAll trainings here : <https://southgreenplatform.github.io/trainings/>
- Presentation & trainings : [HPC IRD](#)
- Work environment : [softwares to install](#)

HPC cluster Introduction

www.southgreen.fr

<https://southgreenplatform.github.io/trainings>



Objective

Knowing how to use the itrop HPC Cluster

Applications

- Knowing the architecture of the cluster
- Knowing the role of the different system partitions
- How to use SGE (qusb, qrsh, qhost, qacct, qstat, qqdel)
- Use the modules environment
- Do some basic scripting

ARCHITECTURE

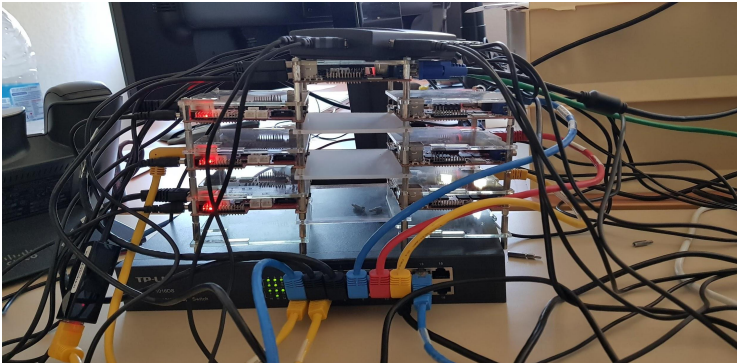
- A cluster is a logical unit composed of several servers
- Acts like a unique powerful server
- Allow to obtain a high computing performance
- A bigger storage
- A better reliability
- A high availability of the ressources

- A cluster is a logical unit composed of several servers
- Acts like a unique powerful server
- Allow to obtain a high computing performance
- A bigger storage
- A better reliability
- A high availability of the ressources



What is a cluster?

- A cluster is a logical unit composed of several servers
- Acts like a unique powerful server
- Allow to obtain a high computing performance
- A bigger storage
- A better reliability
- A high availability of the ressources



- Master node : Scheduler.
Handle the ressources and priorities of the jobs
- Computing nodes :
Ressources (CPU or RAM memory)
used by the master

COMPUTING



- Master node : Scheduler.
Handle the ressources and priorities of the jobs
- Computing nodes :
Ressources (CPU or RAM memory)
used by the master

COMPUTING



- Nas servers:
Store the users data and the analyses results

STORAGE





bioinfo-master
.ird.fr
91.203.34.148

Role : launch and schedule the jobs on the computing nodes
Accessible from the Internet
Connection : `ssh login@bioinfo-master.ird.fr`



**bioinfo-master
.ird.fr
91.203.34.148**

Role : launch and schedule the jobs on the computing nodes
Accessible from the Internet
Connection : `ssh login@bioinfo-master.ird.fr`



**25 noeuds
with
bioinfo-inter.ird.fr
91.203.34.150**

Role : Used by the master to execute jobs
Not accessible from the Internet
node0 à node25
Interactive node : `bioinfo-inter.ird.fr`
Connection : `ssh login@bioinfo-inter.ird.fr`



**bioinfo-nas3.ird.fr
91.203.34.180**



Role : Store data
Accessible from the Internet
Connection : filezilla or scp

**bioinfo-nas.ird.fr bioinfo-nas2.ird.fr
91.203.34.157 91.203.34.160**



**bioinfo-master.
ird.fr**
91.203.34.148

Rôle : Lancer et prioriser les jobs sur les nœuds de calcul

Accessible depuis Internet

Connexion : `ssh login@bioinfo-master.ird.fr`



**22 noeuds
avec
bioinfo-inter.ird.
fr**
91.203.34.150

Rôle : Utilisés par le maître pour exécuter des jobs

Pas accessibles depuis Internet

node0 à node22

Noeud interactif : `bioinfo-inter.ird.fr`

Connexion : `ssh login@bioinfo-inter.ird.fr`



bioinfo-nas3.ird.fr
91.203.34.180



Rôle : Stocker les données utilisateurs

Accessibles depuis Internet

Connexion : `filezilla` ou `scp`

bioinfo-nas.ird.fr
91.203.34.157

bioinfo-nas2.ird.fr
91.203.34.160

/home : Your personal folder
Quota 100Go
Hosted on : bioinfo-nas.ird.fr
Shared on all servers

/teams : project data shared between
users from the same team
Quota 200Go
Hosted on : bioinfo-nas.ird.fr
Shared on all servers

/data2 : project data shared between
several users
Quota 500Go à 1To
Hosted on: bioinfo-nas.ird.fr
Shared on all servers

/home : Your personal folder
Quota 100Go
Hosted on : bioinfo-nas.ird.fr
Shared on all servers

/data : project data shared between
several users
Quota 500Go à 1To
Hosted on : bioinfo-nas2.ird.fr
Shared on all servers

/teams : project data shared between
users from the same team
Quota 200Go
Hosted on : bioinfo-nas.ird.fr
Shared on all servers

/data3 :Données projet partagées
entre plusieurs utilisateurs
Quota 500Go à 1To
Hosted on : bioinfo-nas3.ird.fr
Shared on all servers

/data2 : project data shared between
several users
Quota 500Go à 1To
Hosted on: bioinfo-nas.ird.fr
Shared on all servers

/home : Your personal folder
Quota 100Go
Hosted on : bioinfo-nas.ird.fr
Shared on all servers

/data : project data shared between
several users
Quota 500Go à 1To
Hosted on : bioinfo-nas2.ird.fr
Shared on all servers

/teams : project data shared between
users from the same team
Quota 200Go
Hosted on : bioinfo-nas.ird.fr
Shared on all servers

/data3 :Données projet partagées
entre plusieurs utilisateurs
Quota 500Go à 1To
Hosted on : bioinfo-nas3.ird.fr
Shared on all servers

/data2 : project data shared between
several users
Quota 500Go à 1To
Hosted on: bioinfo-nas.ird.fr
Shared on all servers

/scratch : temporary working folder
1To à 5To
Hosted on : each node
Not shared local only
Data kept **3 weeks only**



/home : Votre répertoire personnel
Quota 100Go
Hébergée sur : bioinfo-nas.ird.fr
Partagée sur toutes les machines

/data : Données projet partagées
entre plusieurs utilisateurs
Quota 500Go à 1To
Hébergée sur : bioinfo-nas2.ird.fr
Partagée sur toutes les machines

/team : Données projet partagées
entre plusieurs utilisateurs
d'une même équipe
Quota 200Go
Hébergée sur : bioinfo-nas.ird.fr
Partagée sur toutes les machines

/data3 : Données projet partagées
entre plusieurs utilisateurs
Quota 500Go à 1To
Hébergée sur : bioinfo-nas3.ird.fr
Partagée sur toutes les machines

/data2 : Données projet partagées
entre plusieurs utilisateurs
Quota 500Go to 1To
Hébergée sur : bioinfo-nas.ird.fr
Partagée sur toutes les machines

/scratch : Répertoire temporaire de travail
1To à 5To
Hébergée sur : **chaque noeud**
Pas partagée mais uniquement en local
Données conservées **3 semaines**



SUN GRID ENGINE (SGE)

- SGE (SUN Grid Engine) est un gestionnaire de ressources de calcul sous linux, capable de gérer de deux à des milliers de serveurs et des centaines de clusters de plusieurs nœuds à la fois.
- An opensource tool
- 3 main functions :
 - Allocates ressources (CPU,RAM) to users to allow them to launch their analyses
 - Provides a frame to launch,execute et monitore the jobs on the whole allocated nodes
 - Deals with jobs in queue wait

Bioinfo.q : default queue

Nodes: node2, node8, node9, node10,
node11,node12,node13,
,node14,node15,node16,node17,
node19,node20
RAM: from 48Go to 64Go
Cores: from 12 to 20 cores

~~dynadiv.q : priority for the dynadiv
team~~

Nodes: node2, node10
RAM: 48Go
Cores: 12 cores

/scratch of 5To for node10

dynadiv2.q : priority for thomas
Couvreur

Nodes: node20
RAM: 64Go
Cores: 20 cores

smrtportal.q : priority for the smrtportal
software

Nodes: node17, node18
RAM: 64Go
Cores: 12 cores

alizon.q : priority for the samuel Alizon
team

Nodes: node8, node9, node12
RAM: 48Go
Cores: 12 cores

Bioinfo.q : default queue

Nodes: node2, node8, node9, node10,
node11,node12,node13,
,node14,node15,node16,node17,
node19,node20
RAM: from 48Go to 64Go
Cores: from 12 to 20 cores

dynadiv.q : priority for the dynadiv
team

Nodes: node2, node10
RAM: 48Go
Cores: 12 cores
/scratch of 5To for node10

dynadiv2.q : priority for thomas
Couvreur

Nodes: node20
RAM: 64Go
Cores: 20 cores

smrtportal.q : priority for the smrtportal
software

Nodes: node17, node18
RAM: 64Go
Cores: 12 cores

alizon.q : priority for the samuel Alizon
team

Nodes: node8, node9, node12
RAM: 48Go
Cores: 12 cores

r900.q : queue with **DELL nodes**

Nodes: node5, node21
RAM: 32Go
Cores: 16 cores

longjob.q : long jobs or > to10 jobs

Nodes: node0, node1, node11
RAM: 48Go
Cores: 12 cores

bigmem.q : memory need

Nodes: node3
RAM: 96Go
Cores: 12 cores

highmem.q :gig need of memory

Nodes: node4 et node7
RAM: 144Go
Cores: 12 cores

Bioinfo.q : queue par défaut
Noeuds: node2, node8, node9, node10,
node11,node12,node13,
,node14,node15,node16,node17,
node19,node20
RAM: de 48Go à 64Go
Cœurs: de 12 à 20 cœurs

~~dynadiv.q : priorité pour l'équipe~~
dynadiv
Noeuds: node2, node10
RAM: 48Go
Cœurs: 12 cœurs
/scratch de 5To pour node10

dynadiv2.q : priorité pour thomas
Couvreur
Noeuds: node20
RAM: 64Go
Cœurs: 20 cœurs

smrtportal.q : priorité pour le logiciel
smrtportal
Noeuds: node17, node18
RAM: 64Go
Cœurs: 12 cœurs

alizon.q : priorité pour l'équipe de
samuel Alizon
Noeuds: node8, node9, node12
RAM: 48Go
Cœurs: 12 cœurs

r900.q : queue avec noeud **DELL**
Noeuds: node5, node21
RAM: 32Go
Cœurs: 16 cœurs

longjob.q : jobs longs ou > à 10 jobs
Noeuds: node0, node1, node11
RAM: 48Go
Cœurs: 12 cœurs

bigmem.q : besoin de mémoire
Noeuds: node3
RAM: 96Go
Cœurs: 12 cœurs

highmem.q : besoin de mémoire
Noeuds: node4 et node7
RAM: 144Go
Cœurs: 12 cœurs

Actions

- Reserve a core on a node in an interactive way
- Reserve a core on a particular node
- Reserve X core on a node

Commands

\$ qrsh

\$ qrsh -l hostname=nodeX

With X the node number

\$~ qrsh -pe ompi X

With X : number of processors from 0 to 12

Actions

- Launch a script in batch mode
- Propagate the load environment to the node
- Name your job
- Use several processors
- Have a certain amount of memory
- Have a particular node
- Directly launch a command with qsub

Commands

`$qsub + script.sh`

`$qsub -V script.sh`

`$~ qsub -N job_name script.sh`

`$~ qsub -pe ompi X script.sh`

Avec X le nombre de coeurs à utiliser

`$~ qsub -l mem_free=XG script.sh`

Avec X le montant de mémoire à réserver

`$~ qsub -l hostname=nodeX script.sh`

`$~ qsub -b y command`

Actions

- Informations on nodes
- Watch your jobs state
- Informations on running jobs
- Informations on completed jobs
- Global informations on queues

Commands

\$ qhost

\$~ qstat

\$~ qstat -j <JOB_ID>

With JOB_ID :the job number

\$~ qacct -j <JOB_ID>

With JOB_ID :the job number

\$~ qstat -g c

Actions

- Deletion of a job

Commandes

`$~ qdel <JOB_ID>`

With JOB_ID : the job id

1



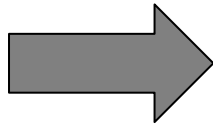
Training: Launch a blast analysis in a interactive way

Data
transfer
from PC to
nas servers

Step 1

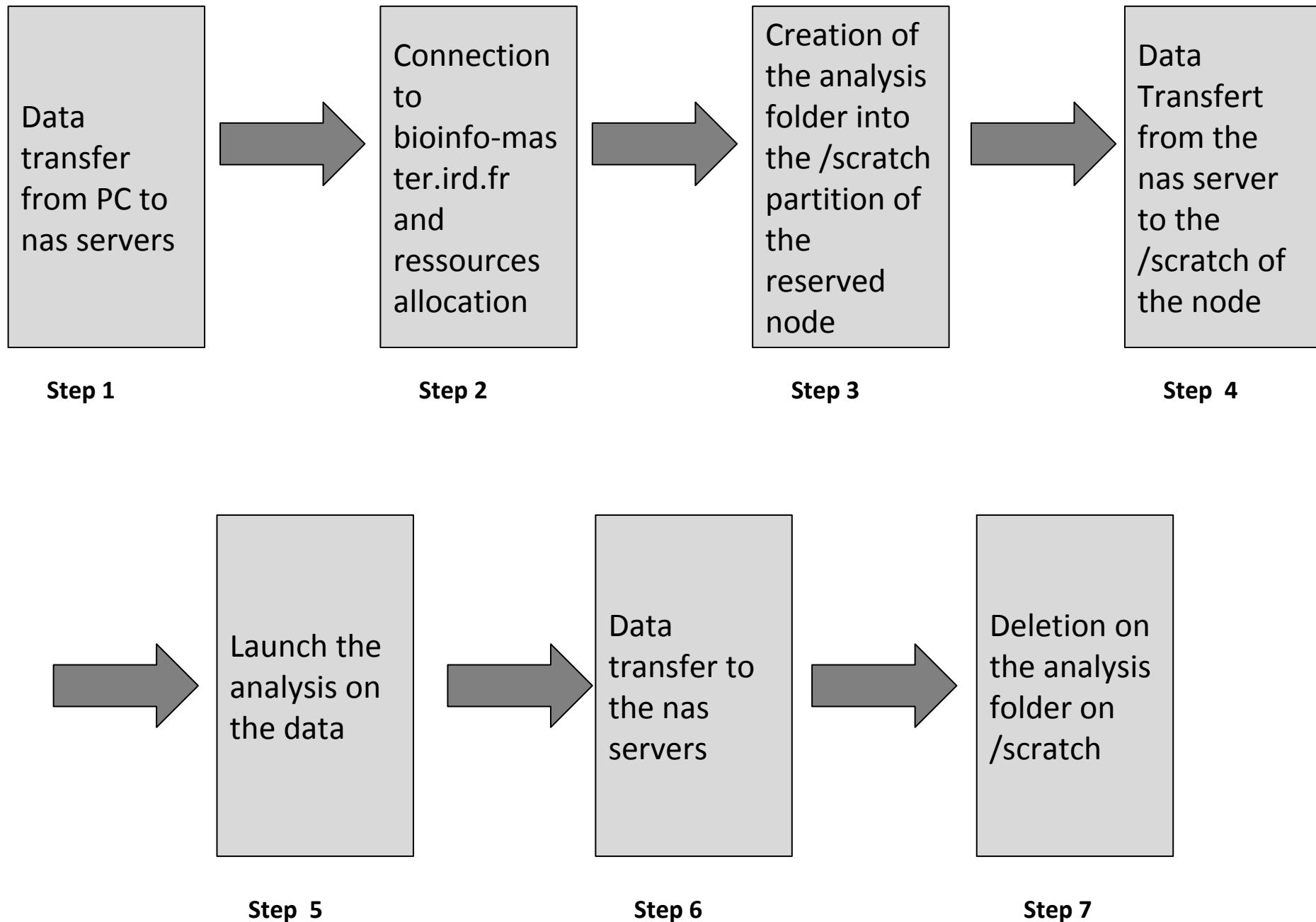
Data
transfer
from PC to
nas servers

Step 1



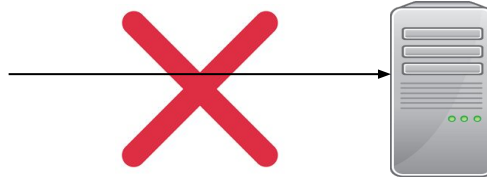
Connection
to
bioinfo-mas
ter.ird.fr
and
ressources
allocation

Step 2





Personal
Computer

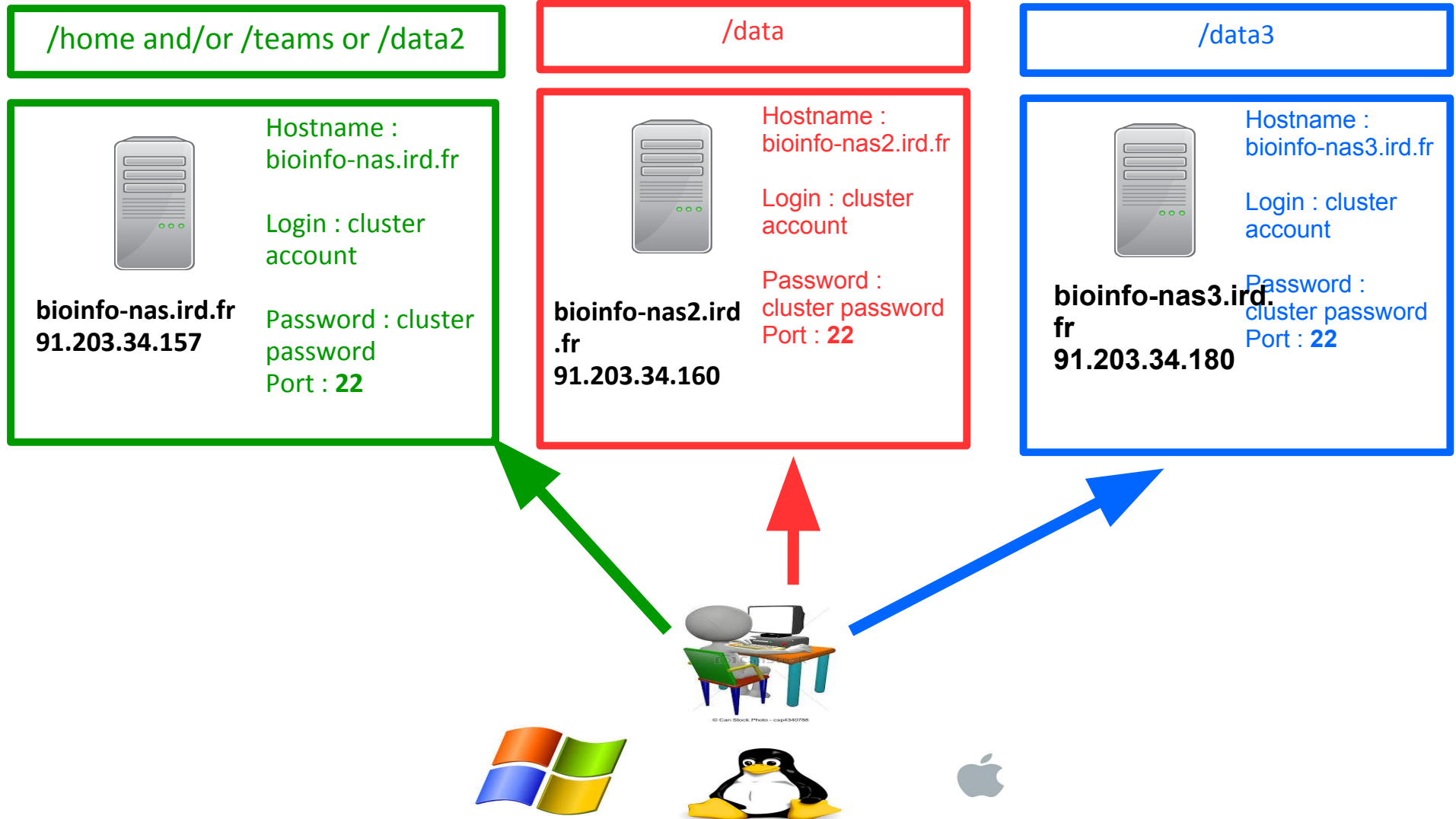


**direct transfert
through filezilla
is forbidden**



**bioinfo-master.ird.fr
91.203.34.148**

Data transfer on the cluster itrop



Open filezilla and retrieve the file « HPC_en.pdf »
Hosted in /data/projects/tp-cluster/training_2018

Open filezilla and retrieve the file « HPC_en.pdf »
Hosted in /data/projects/tp-cluster/training_2018

Enter the following parametres:

Hostname : **bioinfo-nas2.ird.fr**

Login : votre login

Password : votre login

Port :**22**

Navigate into the right window through /data/projects/tp-cluster/training_2018

Retrieve the file HPC_en.pdf with a drag-and-drop

Launch scripts to several nodes



bioinfo-master.ird.fr
91.203.34.148

Use the
qsub
command

Hostname :
bioinfo-master.ird.f
r

Login : cluster
account

Password : cluster
password
Port : 22

Test your script(s)



bioinfo-inter.ird.fr
91.203.34.150

Hostname :
bioinfo-inter.ird.fr

Login : cluster
account

Password : cluster
password

Port : 22

Or use the qrun command on
bioinfo-master.ird.fr



With Putty
Use parameters above



© Can Stock Photo - csp4340788



With terminal
Use the ssh command

Connect to bioinfo-master.ird.fr via ssh

Type :

\$~ssh [login@bioinfo-master.ird.fr](https://login.bioinfo-master.ird.fr) onApple or Linux

Under windows : Download Mobaxterm to :

<https://mobaxterm.mobatek.net/download-home-edition.html>

Then connect to bioinfo-master.ird.fr

We can reserve a core of a node to launch an analysis
Through a limited time using the `qrsh` command
Type the `qstat` command and analyse the result

We can reserve a core of a node to launch an analysis
Through a limited time using the qrsh command
Type the qstat command and analyse the result

Type :
\$~qrsh
Check on wich node you are with the command
\$~ uname -a
\$~ qstat

Go into /scratch
Create a folder to host your data

Go into /scratch
Create a folder to host your data

Type the commands :
\$~cd /scratch
\$~ mkdir login (with login the name folder of your choice)

Data transfer with scp

Being connected to A

Remote folder to transfer : /data/projects/tp-cluster/training_2018

Login : login

Destination folder on the node : /scratch/tando

Copy remote folder from server B to local server A



**Destination
ServerA**



**Source
ServerB**

Data transfer with scp

Being connected to A

Remote folder to transfer : /data/projects/tp-cluster/training_2018

Login : login

Destination folder on the node : /scratch/tando

Copy remote folder from server B to local server A

```
scp -r login@source_server:/remote_path
```



**Destination
ServerA**



**Source
ServerB**

/data/projects/tp-cluster/training_2018

Data transfer with scp

Being connected to A

Remote folder to transfer : /data/projects/tp-cluster/training_2018

Login : login

Destination folder on the node : /scratch/tando

Copy remote folder from server B to local server A

```
scp -r login@source_server:/remote_path local_folder
```

/scratch/tando



**Destination
ServerA**



**Source
ServerB**

/data/projects/tp-cluster/training_2018

Data transfer with scp

Being connected to A

Remote folder to transfer : /data/projects/tp-cluster/training_2018

Login : login

Destination folder on the node : /scratch/tando

Copy remote folder from server B to local server A

```
scp -r login@source_server:/remote_path local_folder
```

/scratch/tando



**Destination
ServerA**



**Source
ServerB**

/data/projects/tp-cluster/training_2018

```
scp -r login@serverB:/data/projects/tp-cluster/training_2018
```

Data transfer with scp

Being connected to A

Remote folder to transfer : /data/projects/tp-cluster/training_2018

Login : login

Destination folder on the node : /scratch/tando

Copy remote folder from server B to local server A

```
scp -r login@source_server:/remote_path local_folder
```

/scratch/tando



**Destination
ServerA**



**Source
ServerB**

/data/projects/tp-cluster/training_2018

```
scp -r login@serverB:/data/projects/tp-cluster/training_2018 /scratch/tando
```

Copy the folder `/data/projects/tp-cluster/training_2018/Blast` to `/scratch/login`

Copy the folder /data/projects/tp-cluster/training_2018/Blast to /scratch/login

Type the commands :

```
$~cd /scratch/login
```

```
$~ scp -r login@bioinfo-nas2.ird.fr :/data/projects/tp-cluster/training_2018/Blast  
/scratch/login
```

Go into the folder `/scratch/login/Blast`
List the files of the folder

Go into the folder /scratch/login/Blast
List the files of the folder

Type :
\$~cd /scratch/login/Blast
\$~ls -ali

- Allow to choose the version of software you want to use
- 2 types of softwares :
 - bioinfo : includes all the bioinformatics softwares (example BEAST)
 - system : includes all the system softwares(example JAVA)
- Overcome the environnement variables
- 5 types of commands :

See the available modules : `module avail`

Obtain infos on a particular module: `module whatis + module name`

For example `module whatis bioinfo/blast/2.4.0+`

Load a module : `module load + modulename`

For example `module load bioinfo/blast/2.4.0+`

List the loaded module : `module list`

Unload a module : `module unload + modulename`

For example `module unload bioinfo/blast/2.4.0+`

Unload all the modules :

`Module purge`

Load the blast module version 2.4.0+
Use the `blastn` command to launch a blast analysis
Hat will produce the result file called `blastn.out`

Load the blast module version 2.4.0+
Use the blastn command to launch a blast analysis
That will produce the result file called blastn.out

Type :
\$~ module load bioinfo/blast/2.4.0+
\$~ blastn -db All-EST-cofea.fasta -query sequence-NMT.fasta -out blastn.out

Edit the blastn.out file with the nano tool

Edit the blastn.out file with the nano tool

Type :
\$~ nano blastn.out

Copy the file blastn.out to your home folder
Check that the file has been copied

Copy the file blastn.out to your home folder
Check that the file has been copied

Type :
\$~scp blastn.out login@bioinfo-nas.ird.fr:/home/login
\$~ls -ali /home/login

Go into the /scratch folder
Delete your working directory

Type:
\$~cd /scratch
\$~ rm -rf *login*



**TP: Launch a bwa in a
interactive way**

- Follow the steps from the last training and adapt them to this one
 - The folder to copy: /data/projects/training_2018/bwa
 - Bwa version to use: 0.7.12
 - Commands to launch:
bwa index referenceIrgin.fasta
bwa mem referenceIrgin.fasta irigin1_1.fastq irigin1_2.fastq >mapping.sam
 - Retrieve the file mapping.sam and place it in your /home/

Cf solution: [practice2](#)



**Training: Launch an analyse
via a bash script**

- Execute a bash script via sge
- We use the command:

```
$~ qsub script.sh
```

With `script.sh` : the script name

First part of the script (in green): sge execution options with the key word #

```
#!/bin/sh

##### SGE CONFIGURATION #####
# Ecrit les erreurs dans le fichier de sortie standard
#$ -j y

# Shell que l'on veut utiliser
#$ -S /bin/bash

# Email pour suivre l'exécution
#$ -M prenom.nom@ird.fr ##### Mettre son adresse mail

# Type de message que l'on reçoit par mail
# - (b) un message au démarrage
# - (e) à la fin
# - (a) en cas d'abandon
#$ -m bea

# Queue que l'on veut utiliser
#$ -q bioinfo.q

# Nom du job
#$ -N Nom_a_choisir
#####
```


In the 2nd part of the script: the command to execute

```
path_to_dir="/data/projects/rep_a_choisir";  
path_to_tmp="/scratch/nom_rep_a_choisir-$JOB_ID"  
  
##### Create the temporary folder on the node and load the blast module  
module load bioinfo/blastn/2.4.0+  
mkdir $path_to_tmp  
scp -rp nas2:$path_to_dir/* $path_to_tmp # choisir nas pour/home, /data2 et /teams ou nas2 pour /data ou nas3 pour /data3  
echo "transfert donnees master -> noeud";  
cd $path_to_tmp  
  
##### Program execution  
cmd="blastn -db All-EST-cofea.fasta -query sequence-NMT.fasta -num_threads $NSLOTS -out blastn1-$JOB_ID.out";  
echo "Commande executee : $cmd";  
$cmd;  
  
##### Data transfer from node to nas  
scp -rp $path_to_tmp/ nas:$path_to_dir/  
echo "Transfert donnees node -> master";  
  
#### Deletion of the tmp folder  
rm -rf $path_to_tmp  
echo "Suppression des donnees sur le noeud";
```

- Using the Training 1 create a script to launch a blastn analysis
- Make the script launchable with

\$~ chmod 755 script.sh

- Launch the script with qsub:

\$~ qsub script.sh

- Check the running script with the command: watch qstat

[solution script blastn](#)

Use the dos2unix command when the script has been written under
Windows



- Using the Training 1 create a script to launch a bwa analysis
- Make the script launchable with

\$~ chmod 755 script.sh

- Launch the script with qsub:
\$~ qsub script.sh
- Check the running script with the command: watch qstat

[solution script bwa](#)



Use the dos2unix command when the script has been written under Windows