Scientific computing resources

Maxime KERMARQUER Using the ICM cluster





CHERCHER, TROUVER, GUÉRIR, POUR VOUS & AVEC VOUS.

Introduction

> Goal

This training gets you started with the basics of connecting and running your code / applications on the cluster

> Prerequisites

- > This training assumes knowledge about the Linux command line
 - \succ ls , cd , pwd , tree , cat , tail , less , vim , sed , awk , |
 - Relative path / Absolute path
 - > data
 - /network/lustre/iss01/home/maxime.kermarquer/data
 - Bash environment variables
 - export VAR1=value
 - > echo \${VAR1}

For Topics

- What is a cluster
- Cluster storage
- Module environment
- Submitting jobs on the cluster with SLURM



2 Using the ICM cluster

Introduction

Agenda









3 Using the ICM cluster

02/06/2020

HPC What is it ? Do I need it ?





CHERCHER, TROUVER, GUÉRIR, POUR VOUS & AVEC VOUS.

High Performance Computing What is it?

Serial Computing

- A problem is broken down into instructions
- Instructions are executed sequentially
- >Executed on a single processor

Parallel Computing

- A problem is broken into parts that can be solved in parallel Each part is broken down into instructions
- \succ
- Instructions from each part execute simultaneously >
- Executed on multiple processors





High Performance Computing What is a cluster ?

> A **collection of computers** (nodes) connected by a network

- Lot of computers in the background for running tasks
- Most clusters run on Linux







High Performance Computing What is a cluster ?





High Performance Computing What can it do for me?





- Solve larger problems
- Use non-local resources



Who is using HPC ?

- Universities and government institutions
- Atmosphere, Earth, Environment
- Physics applied, nuclear, particle
- Mechanical Engineering
- Chemistry, Molecular Sciences
- Geology, Seismology
- Medical imaging and diagnosis
- Genetics
- ••

Industrial and Commercial

- Big Data, data mining
- Pharmaceutical design
- Financial and economic modeling

≻ ...



















02/06/2020

High Performance Computing What do I need for access ?

> A workstation on ICM network or a VPN connection open to computing resources

> Open an issue in category *Calculateur-07* to request an account

> SSH protocol used to connect to the cluster



2 IT Infrastructure What are the storage/computing resources?



CHERCHER, TROUVER, GUÉRIR, POUR VOUS & AVEC VOUS.

IT Infrastructure

ISSTITUT du Cerveau et de la Moelle épinière









Storage Structure



You can access Lustre through Windows (SMB), Linux (SMB/NFS) and Mac OS (NFS/SMB)



12 Using the ICM cluster

02/06/2020

Storage ISS01



Where your store all files relate to your teams and project

You can access Lustre through Windows (SMB), Linux (SMB/NFS) and Mac OS (NFS/SMB)



02/06/2020

3 How do I use the cluster environment?



CHERCHER, TROUVER, GUÉRIR, POUR VOUS & AVEC VOUS.

Cluster usage Process





Cluster usage Login node



Used for access to the cluster

- > Login
- Data transfers
- Job submission
- Editing scripts

Not used to compute





Cluster usage SSH connection

- SSH protocol used to connect to the cluster
- IP Address: 192.168.90.2
- Linux and MacOS
 - From a terminal
 - > ssh <username>@login02
 - \rightarrow -x option to activate Xforwarding
 - -Y option on Mac with XQuartz activated
- > Windows
 - > Putty
 - > Git Bash
 - Xming for Xforwarding







17 Using the ICM cluster



Cluster usage

Data transfer

File Transfer client

- > Linux : scp, rsync commands with your ICM credentials
 - > rsync <src> <dest>
 - > **Options:** -r -h --stats -progress -itemize-change

rsync -r data maxime.kermarquer@192.168.90.2:/network/lustre/iss01/maxime.kermarquer

Windows: Filezilla with host as sftp://192.169.90.2 and your ICM credentials

E sftp://maxime.kermarquer@login01 - FileZilla		- 0 ×
Fichier Édition Affichage Transfert Serveur Favoris ? Nouvelle version disponible !		
Hôte : sftp://login01 Identifiant : me.kermarquer Mot de passe : ++++++++ Port : Connexion rapide	*	
Erreur : Impossible de récupérer le contenu du dossier		
Statut : Récupération du contenu du dossier "/network/lustre/iss01/home/maxime.kermarquer/Developpement"		
Statut : Listing directory /network/lustre/iss01/nome/maxime.kermarquer/Developpement Statut : Contenu du dossier "/network/lustre/iss01/home/maxime.kermarquer/Developpement" affiché avec succès		
Site local : C\Users\maxime.kermarquer\Desktop\	Site distant :	/network/lustre/iss01/home/maxime.kermarquer/Developpement
B Documents		- 2 bash
Downloads		- P Benchmarks
🔅 🛼 Favorites		
- IntelGraphicsProfiles		- Z Java
- I Links		A MATIAR
- Menu Démarrer	Nom de fich	Taille d_ Type de _ Dernière m_ Droits d'_ Propriét_
- Res documents	-	and an Observe serves and a serves
Nom de fichier Taille de Type de fic Demière modi	bash	Dossier 19/01/2018 drwxr-xr-x maxime
	Benchman	ks Dossier 05/03/2018 drwxr-xr-x maxime
Raccourcis I Dossier de f 17/05/2018 10	1 C	Dossier 12/01/2018 drwxr-xr-x maxime
a desktop.ini 282 Paramètres 17/06/2018 22:	🣜 Java	Dossier 30/11/2017 dnwxr-xr-x maxime
C Microsoft E 1 417 Raccourci 17/06/2018 22:	Licenses	Dossier 15/03/2018 drwxr-xr-x maxime
*** My ICM Intr 220 Raccourci I 09/07/2018 14:	MATLAB	Dossier 01/06/2018 drwxr-xr-xr maxime
	Python	Dossier 30/05/2018 drwxr-xr-x maxime
	R Seriete SI	UDSSIEF 22/03/2018
	StageM2	Dossier 14/03/2018 drawr wr w maxime
	bashrc	87 Fichier B., 23/01/2018., -rw-rr- maxime
	Ulesh2.0.	8.tgz 41.931 Fichier T 01/02/2018rw-rr- maxime
3 fichiers et 1 dossier. Taille totale : 1 919 octets	2 fichiers et 1	0 dossiers. Taille totale : 42 018 octets
Consum (Cables Land Direct Cables distant Taille Directed Parks		
Serveur / Fichier local Direc Fichier distant laine Phonte Statut		
Fichiers en file d'attente Transferts échoués Transferts réussis		
		🔒 🕜 File d'attente : vide

02/06/2020









19 Statistiques Cluster ICM - 2019

03/02/2020

Module Available software

- > Workstation
 - One or few users
 - A small selection of software
- > HPC cluster
 - Large number of users
 - > Each needing a different selection of software packages
 - Often with different versions and configurations
- Organized by the module tool
 - Manage multiple versions and configurations of software
 - Used by many HPC centers around the world
- Environment set for software when module is loaded
- Dependency automatically load



Module Commands

Display the available modules module avail

						Compilers					
cmake/3.10.2	gcc/5.5.0 gcc,	/6.4.0 gcc/7.3.0)(D) intel/2018	pgi/1710							
MATLAB/R2016	a MATLAB/R2017a	MATLAB/R2017b (L,	D) R/3.4.3 R/3.	5.0 (D) java/7	java/8 (D) la	gramming Language atex/3.14159265	s python/bioinfo	python/2.7 python/3.6 (D)			
						Neuroimaging					
AFNI/18.1.05 ANTs/2.2.0	FSL/5.0.1 FreeSurfer/5.3.0	FreeSurfer/6.0. ICA-AROMA/0.3	0 (D) ICA-AROMA/0. MNE/0.16	4 (D) MNE_py27/0 MRtrix3/no	0.16 o_gui_19_03_2018	PETPVC/1.2.1 SPM/12	TDT/3.991 clinica/05_04_2018	dcm2niix/v1.0.20171215 field-trip/20180625	fmriprep/1.0.11 phy/spyking_circus_dev	phy/spyking_circus_folder phy/1.0.0	s (D)
C3D/1.1.0	CUDA/8.0 CUDA/9.0	0 CUDA/9.1 (D)	ICU/58.2 ICU/60.	1 (D) ITK/4.13	fiji/20170530	libwebp/0.5.2	snakemake/4.8.0	svgo/1.0.5 zstd/1.3.3			
lmod/7.7.18	settarg/7.7.18				····· /usr/snare/	lmod/lmod/modulet	lles/core				
Où: L: Le modul D: Default	e est chargé Module										

> Load / unload a specific module

- > module load <modulefile>
- > module unload <modulefile>

module load MATLAB/R2017b



Module Commands

- Display the available modules
 - > module list



- Display a brief description
 - module what-is <modulefile>

Version: R2017b Keywords: MATLAB

URL: https://fr.mathworks.com/products/matlab.html

Description: MATLAB combines a desktop environment tuned for iterative analysis and design processes with a programming language that expresses matrix and array mathematics directly.



Module Installation of new software

- Into your own space storage
- Keep control over software yourself
- No special privileges required
- Cannot be used by others users by default
- > You can request to create module to share with all users
- > Into a new module
 - > Can be used by other users
 - > Installation requires special privileges
 - > Open an issue in *Calculateur-07* category





SLURM How do I compute my data?



CHERCHER, TROUVER, GUÉRIR, POUR VOUS & AVEC VOUS.

SLURM Presentation

- > SLURM stands for Simple Linux Utility for Resource Management
- > Manages and shares hardware resources between users
- The most used scheduler for Linux cluster
- > Resources: logical or hardware unit necessary to run software on a cluster, e.g.:
 - o Walltime
 - Compute node
 - Memory
 - Accelerator cards
 - o ...





SLURM Presentation

- 1. Users write **job description** and **needed resources**
- 2. SLURM finds matching resources
- 3. Scheduler tries to make optimal use of the resources
- 4. No resource for a job wait in a queue
- 5. Priority determined by usage of the system in recent past and job size



SLURM Most used commands

Check the cluster state > sinfo

> cluster_state

> Running jobs

- > sbatch to submit a job script
- > srun to run a job (wait the ressources)
- > salloc to allocate resources (wait the ressources)

> Monitoring jobs

- > squeue to check the jobs in queue
- > sacct to check the job history





> 29 compute nodes form the computing cluster

- Based on Intel Xeon E5-2680 v4 processors
- > 22 nodes with 128GB memory and 28 cores
- 4 nodes with 256GB memory and 28 cores
- I GPU node with 4 NVIDIA Tesla P100 graphics cards with nvlink





- > A partitions/queue is a logical view of all compute node
- > The partitions divide the nodes using characteristics:
 - Compute node memory
 - Number of processor
 - Accelerator cards
 - ...
- > Partitions are:

	Nodes	Core/node	Mem/node	Processor
normal	22	28	128G	Broadwell
bigmem	4	28	256G	Broadwell
gpu	1	24	128G	Broadwell



Slurm sinfo

cluster_state or sinfo command display nodes state \succ

maxime.k	ermarqu	er@login01	~ > 9	sinfo –	
PARTITION	AVAIL	TIMELIMIT	NODES		IODELIST
debug	up	infinite	4	idle	mw[001-004]
bigmem	up	infinite	4	mix	imb[001-004]
normal*	up	infinite	1	drain*	.mb020
normal*	up	infinite	3	drain	.mb[021,025,030]
normal*	up	infinite	13	mix	mb[022-024,026,028,035-036,039,041-045]
normal*	up	infinite	8	alloc	mb[027,033-034,037-038,040,046-047]
gpu	up	infinite	1	idle	.mgpu01

- \succ
 - State are the following: down : unavailable node for use
 - mix : node has some of its CPUs ALLOCATED while others are IDLE •
 - alloc : node has been allocated to one or more jobs •
 - idle : node free •
 - drain : the node is in maintenance •



SLURM Most used commands

Check the cluster state > sinfo

> cluster_state

> Running jobs

- > sbatch to submit a job script
- > srun to run a job (wait the ressources)
- > salloc to allocate resources (wait the ressources)

> Monitoring jobs

- > squeue to check the jobs in queue
- > sacct to check the job history



Running jobs Commands to submit jobs

> 3 ways allow to submit jobs with SLURM



- > salloc + srun / ssh
- > Each job is identified by a **job id** useful to follow its progress
- > You need to specify a set of resources in which your job will run
- Exceeding the limit resources will, stop your application/job



Running jobs GPU server example





33 Using the ICM cluster

02/06/2020



Type of compute node used normal, bigmem, gpu --partition=<partition name> Or -p <partition name>

Time limit of a job with this format for time HH:MM:SS or days-hours:minutes:seconds --time=<time> Or -t <time>

The memory size usable by the job with this format 1024 , 1024MB or 1G --mem=<count>

> The number of CPUs per process --cpus-per-task=<count> Or -c <number>

The number of GPUs (only for gpu partition) --gres=gpu:<count>





The job name (the default command name) --job-name=<jobname> or -J <jobname>

> To choose the working directory
--chdir=<path_workdir>

> To redirect the error output
--error=<filename pattern> or -e <filename pattern>

> To redirect the standard output --output=<filename pattern> or -o <filename pattern>



Running jobs sbatch

> A job script is shell file, starts with #!/bin/bash followed by #SBATCH and SLURM resources options

> Sample script:







Running jobs MATLAB example

```
#!/bin/bash
#SBATCH --job-name=run_MATLAB
#SBATCH --partition=normal
#SBATCH --time=03:00:00
#SBATCH --mem=10G
#SBATCH --cpus-per-task=1
#SBATCH --chdir=/network/lustre/iss01/home/maxime.kermarquer
#SBATCH --output=file_output_%j.log
#SBATCH --error=file_output_%j.log
```

module load MATLAB
matlab -nojvm -nodesktop -nodisplay < mycode.m</pre>

You can find other examples from gitlab at https://gitlab.com/maxime.kermarquer/examples-slurm-scripts



Running jobs GPU allocation example

#!/bin/bash #SBATCH --partition=gpu #SBATCH --time=02:00:00 #SBATCH --mem=20G #SBATCH --gres=gpu:2 #SBATCH --cpus-per-task=8 #SBATCH --cpus-per-task=8 #SBATCH --workdir=. #SBATCH --output=tensorflow_gpu_job_%j.log #SBATCH --error=tensorflow_gpu_job_%j.log #SBATCH --job-name=test_tensorflow

module load python/2.7 module load CUDA/9.0 python tensorflow_example.py



Running jobs GPU allocation example





Running jobs Slurm environment variables

> SLURM defines some environment variable you can use to tune your script:

Variable	Function	Example
SLURM_JOBID	Job number	4242
SLURM_SUBMIT_DIR	Workdir	/network/iss01/home/maxime.kermarquer
SLURM_JOB_NAME	Job name	run_MATLAB
SLURM_NODELIST	Node list	lmb021
SLURM_CPUS_PER_TASK	CPUs by tasks	2
SLURM_MEM_PER_NODE	Memory	4096M



Running jobs

Slurm environment variables – Usage example

```
#!/bin/bash
#SBATCH --partition=normal
#SBATCH --cpus-per-task=8
#SBATCH --mem=32G
#SBATCH --time=20:00:00
#SBATCH --chdir=.
#SBATCH --output=outputs-slurm/fmriprep-20.0.0-test-%j.txt
#SBATCH --error=outputs-slurm/fmriprep-20.0.0-test-%j.txt
#SBATCH -- job-name=fmriprep-20.0.0-test
module load fmriprep
export BIDS DIR=../bids-example/ds001 BIDS
export OUTPUT DIR=outputs
export PARTICIPANT=participant
export WORK DIRECTORY=workdir
export OPTIONS="--verbose --work-dir ${WORK DIRECTORY} --nthreads ${SLURM CPUS PER TASK}
--mem mb ${SLURM MEM PER NODE}"
```

fmriprep \${BIDS_DIR} \${OUTPUT_DIR} \${PARTICIPANT} --work-dir \${WORK_DIRECTORY} \${OPTIONS}



Running jobs srun

> You can directly run a command on the cluster without write a script with srun command.

> srun [options to define needed resource] [your command]

> By example

srun -p gpu -t 02:00:00 --mem=20G --gres=gpu:1 -c 1
-J "run_tensorflow" python tensorflow_example.py

This command blocks your terminal until your command ends or is cancelled.

02/06/2020

> The workdir is the current directory.

Running jobs salloc + srun

> If you use srun command, the resources allocated for the job will be released at his end.

> Using salloc command allows to keep the allocation at the end of the command and to start an other command without waiting for the allocation of resources.

exit command release the allocation

> By example:

```
salloc -p gpu -t 02:00:00 --mem=20G --gres=gpu:1 -N 1 -n 1 -c 1
srun python script1.py
srun python script2.py
exit
```



Running jobs

salloc + ssh





SLURM Most used commands

Check the cluster state > sinfo

> cluster_state

> Running jobs

- > sbatch to submit a job script
- > srun to run a job (wait the ressources)
- > salloc to allocate resources (wait the ressources)

> Monitoring jobs

- > squeue to check the jobs in queue
- > sacct to check the job history



Monitoring jobs squeue

> The squeue command displays the status of your jobs. It gives information on :

- The job ID
- The job partition \succ
- The name of the job \succ
- The user >
- The state
- The running time
- . . .

Can be filtered on:

- squeue -u \$USER # To display his own jobs, for new users you can use the alias "squeueme" # To display the running jobs
- squeue -t R
- squeue -t PD
 - # To display the pending jobs
- squeue -u \$USER -t R # To display his own running jobs



Monitoring jobs squeue

JOBID	PARTITION	NAME	USER	ST	TIME	CPU	MIN_MEMO	SUBMIT_TIME	NODES	NODELIST(REASON)
1278770	gpu	preprocessing	elina.thibeausutre	PD	0:00	1	10G	2018-11-27T16:59:24	1	(Resources)
1282652	normal	snakejob.bam2bw.2.sh	justine.guegan	R	4:29:54	6	80000M	2018-11-28T09:53:44	1	lmb027
1282915	normal	snakejob.mirdeep2.5.sh	thomas.gareau	R	46:47	1	8G	2018-11-28T13:36:52	1	lmb026
1282913	normal	snakejob.mirdeep2.18.sh	thomas.gareau	R	47:23	1	8G	2018-11-28T13:36:16	1	lmb026
1282923	normal	snakejob.mirdeep2.17.sh	thomas.gareau	R	38:34	1	8G	2018-11-28T13:45:04	1	lmb027
1282921	normal	snakejob.mirdeep2.14.sh	thomas.gareau	R	39:10	1	8G	2018-11-28T13:44:28	1	lmb027
1282941	normal	snakejob.mirdeep2.25.sh	thomas.gareau	R	30:08	1	8G	2018-11-28T13:53:30	1	lmb047
1282933	normal	snakejob.mirdeep2.6.sh	thomas.gareau	R	32:45	1	8G	2018-11-28T13:50:54	1	lmb030
1282946	normal	<pre>snakejob.mirdeep2_inputs.32.sh</pre>	thomas.gareau	R	26:55	1	20G	2018-11-28T13:56:44	1	lmb047
1282951	normal	snakejob.mirdeep2.11.sh	thomas.gareau	R	20:31	1	8G	2018-11-28T14:03:08	1	lmb026
1282950	normal	snakejob.mirdeep2.3.sh	thomas.gareau	R	20:44	1	8G	2018-11-28T14:02:55	1	lmb026
1282959	normal	<pre>snakejob.mirdeep2_inputs.33.sh</pre>	thomas.gareau	R	15:20	1	20G	2018-11-28T14:08:19	1	lmb030
1282957	normal	snakejob.mirdeep2.23.sh	thomas.gareau	R	16:43	1	8G	2018-11-28T14:06:56	1	lmb026
1282955	normal	snakejob.mirdeep2.24.sh	thomas.gareau	R	17:43	1	8G	2018-11-28T14:05:56	1	lmb026
1282954	normal	snakejob.mirdeep2.13.sh	thomas.gareau	R	17:56	1	8G	2018-11-28T14:05:43	1	lmb026
1282967	normal	<pre>snakejob.mirdeep2_inputs.41.sh</pre>	thomas.gareau	R	10:31	1	20G	2018-11-28T14:13:07	1	lmb047
1282966	normal	snakejob.mirdeep2.22.sh	thomas.gareau	R	11:06	1	8G	2018-11-28T14:12:32	1	lmb030
1282964	normal	snakejob.mirdeep2.15.sh	thomas.gareau	R	12:06	1	8G	2018-11-28T14:11:32	1	lmb047
1282963	normal	<pre>snakejob.mirdeep2_inputs.34.sh</pre>	thomas.gareau	R	13:31	1	20G	2018-11-28T14:10:07	1	lmb028
1282961	normal	<pre>snakejob.mirdeep2_inputs.37.sh</pre>	thomas.gareau	R	14:20	1	20G	2018-11-28T14:09:19	1	lmb030
1282971	normal	snakejob.mirdeep2.1.sh	thomas.gareau	R	5:17	1	8G	2018-11-28T14:18:22	1	lmb047
1282970	normal	snakejob.mirdeep2.2.sh	thomas.gareau	R	6:06	1	8G	2018-11-28T14:17:33	1	lmb047
1282969	normal	snakejob.mirdeep2.4.sh	thomas.gareau	R	9:30	1	8G	2018-11-28T14:14:08	1	lmb047
1282975	normal	snakejob.mirdeep2.10.sh	thomas.gareau	R	1:17	1	8G	2018-11-28T14:22:21	1	lmb047
1282973	normal	snakejob.mirdeep2.20.sh	thomas.gareau	R	1:28	1	8G	2018-11-28T14:22:10	1	lmb047
1282974	normal	<pre>snakejob.picard_rnaseq_metrics</pre>	thomas.gareau	R	1:28	1	4G	2018-11-28T14:22:10	1	lmb047
1282972	normal	snakejob.mirdeep2.21.sh	thomas.gareau	R	4:53	1	8G	2018-11-28T14:18:45	1	lmb047
1277638	gpu	deformetrica benchmark	benoit.martin	R	23:03:01	24	120G	2018-11-27T15:20:38	1	lmgpu01
1282649	normal	InteractiveSession	marieconstance.corsi	R	5:36:49	28	120G	2018-11-28T08:46:50	1	lmb046



Monitoring jobs

scontrol

> To display all job information

> scontrol show job <job id>

JobId=1284207 JobName=hello_world_matlab

UserId=maxime.kermarquer(20108) GroupId=utilisa. du domaine(10513) MCS label=N/A Priority=7994 Nice=0 Account=dsi QOS=qos1 JobState=RUNNING Reason=None Dependency=(null) Requeue=1 Restarts=0 BatchFlag=1 Reboot=0 ExitCode=0:0 RunTime=00:00:04 TimeLimit=01:00:00 TimeMin=N/A SubmitTime=2018-11-28T18:48:57 EligibleTime=2018-11-28T18:48:57 StartTime=2018-11-28T18:48:58 EndTime=2018-11-28T19:48:58 Deadline=N/A PreemptTime=None SuspendTime=None SecsPreSuspend=0 LastSchedEval=2018-11-28T18:48:58 Partition=normal AllocNode:Sid=login01:18177 ReqNodeList=(null) ExcNodeList=(null) NodeList=lmb030 BatchHost=lmb030 NumNodes=1 NumCPUs=1 NumTasks=1 CPUs/Task=1 ReqB:S:C:T=0:0:*:* TRES=cpu=1,mem=10G,node=1,billing=1 Socks/Node=* NtasksPerN:B:S:C=0:0:*:* CoreSpec=* MinCPUsNode=1 MinMemoryNode=10G MinTmpDiskNode=0 Features=(null) DelayBoot=00:00:00 Gres=(null) Reservation=(null) OverSubscribe=OK Contiguous=O Licenses=(null) Network=(null) Command=/network/lustre/iss01/home/maxime.kermarguer/hello world/job.sh WorkDir=/network/lustre/iss01/home/maxime.kermarguer/hello world/. StdErr=/network/lustre/iss01/home/maxime.kermarguer/hello world/./file output 1284207.err StdIn=/dev/null StdOut=/network/lustre/iss01/home/maxime.kermarguer/hello world/./file output 1284207.out Power=



Monitoring jobs

To cancel jobs in queue or running jobs, you can use scancel command.

Can be filtered on:

- > scancel 13082016
- ➢ scancel −u \$USER
- > scancel --name=JobName
- > scancel --state=PENDING
- > scancel --state=RUNNING
- > scancel --nodelist=lmb024,lmb025

- # To delete the job 13082016
- # To delete all your jobs
- # To delete job with the name JobName
- # To delete all jobs in the queue
- # To delete all RUNNING jobs
- # To delete all your jobs running on nodes Imb024 and Imb025



Monitoring jobs sacct

Sacct command give information for past and running job

Submit	Partition	User	JobID	JobNa	Timelimit	Elapsed	ReqMem	MaxRSS	AveRSS	State
2018-11-27T16:48:22	normal	gizem.temiz	1278336_54	bedp+	1-00:00:00	08:27:14	1Gn			COMPLETED
2018-11-28T02:35:33			1278336_54.batch	batch		08:27:14	1Gn	0.17G	0.17G	COMPLETED
2018-11-27T16:48:22	normal	gizem.temiz	1278336_56	bedp+	1-00:00:00	08:27:24	1Gn			COMPLETED
2018-11-28T02:40:04	_		1278336_56.batch	batch		08:27:24	1Gn	0.17G	0.17G	COMPLETED
2018-11-27T16:48:22	normal	gizem.temiz	1278336_57	bedp+	1-00:00:00	08:29:17	1Gn			COMPLETED
2018-11-28T02:40:16			1278336_57.batch	batch		08:29:17	1Gn	0.17G	0.17G	COMPLETED
2018-11-2/116:48:22	normal	gizem.temiz	12/8336_58	beap+	1-00:00:00	08:28:44	1GN	0 170	0 170	
2018-11-28102:40:28			1278336_58.Datch	batch		08:28:44	IGN	0.1/6	0.1/6	COMPLETED

- > Can be filtered on:
 - Time (beginning and end)
 - Partition / accounting
 - Job ID
- > You can display various information:
 - Number of processor / node / tasks
 - Memory size / Wall-time
 - •••



Jobs array Options

> Job arrays offer a mechanism for submitting and managing collections of similar jobs quickly and easily.

All jobs must have the same initial options (e.g. size, time limit, ...)

- > **Option:** --array:
 - sbatch --array=0-50
 - sbatch --array=1,3,5,7
 - sbatch --array=1-50%2
 - sbatch --array=1-50:10
- > Environment variable that can be use in our script:
 - SLURM_ARRAY_JOB_ID
 - SLURM_ARRAY_TASK_ID



Jobs array Example

maxime.kerm	arquer@login01	~/job_array \$	tree					
outputs run_arra scripts script script	ay.sh t_1.py t_2.py t_3.py							
2 directories,	4 files							
maxime.kerm	arquer@login01	~/job_array \$	sbatch -	array=	=1-3	run_array.sh		
Submitted ba	tch job 1284209			-				
maxime.kerm	arquer@login01	~/job_array \$	squeuer	ne				
JOBID PA	RTITION	NAME	ST	TIME	CPU	MIN_MEMO	NODES	NODELIST(REASON)
1284218_1	bigmem	run_array.sh	R	0:04	1	2G 1	hmb001	
1284218_2	bigmem	run_array.sh	R	0:04	1	2G 1	hmb001	
1284218_3	bigmem	run_array.sh	R	0:04	1	2G 1	hmb001	



#!/bin/bash
#SBATCH --partition=bigmem
#SBATCH --time=02:00:00
#SBATCH --mem=2G
#SBATCH --ntasks=1
#SBATCH --ntasks=1
#SBATCH --error=outputs/job_%N_%A_%a.log
#SBATCH --output=outputs/job_%N_%A_%a.log

echo "SLURM_JOBID: " \$SLURM_JOBID echo "SLURM_ARRAY_JOB_ID: " \$SLURM_ARRAY_JOB_ID echo "SLURM_ARRAY_TASK_ID: " \$SLURM_ARRAY_TASK_ID

module load python/2.7 python scripts/script_\${SLURM_ARRAY_TASK_ID}.py

- Coordinate: SLURM_ARRAY_JOB_ID + SLURM_ARRAY_TASK_ID
- Cancel a job array, e.g.:
 - scancel 1284218
- Cancel only some jobs, e.g.:
 - scancel 1284218_1
 - scancel 42_[1-2]

Jobs array Example



Jobs array Example



02/06/2020

Jobs array Example

> Job array run 10 MATLAB scripts with 10 different paths by using sed shell command and getenv MATLAB function.

JOBID	PARTITION	NAME	USER	ST	TIME	CPU	MIN_MEMO	START_TIME	NODES	NODELIST(REASON)
126567_3	bigmem	run_array_parameters-path.sh	maxime.kermarquer	CG	0:11	1	2G	2019-10-04T18:11:01	1	hmb002
126567 9	bigmem	run_array_parameters-path.sh	maxime.kermarquer	CG	0:11	1	2G	2019-10-04T18:11:01	1	hmb002
126567_10	bigmem	run_array_parameters-path.sh	maxime.kermarquer	CG	0:10	1	2G	2019-10-04T18:11:01	1	hmb002
126567_1	bigmem	run_array_parameters-path.sh	maxime.kermarquer	R	0:11	1	2G	2019-10-04T18:11:01	1	hmb002
126567_2	bigmem	run_array_parameters-path.sh	maxime.kermarquer	R	0:11	1	2G	2019-10-04T18:11:01	1	hmb002
126567_4	bigmem	run_array_parameters-path.sh	maxime.kermarquer	R	0:11	1	2G	2019-10-04T18:11:01	1	hmb002
126567 5	bigmem	run array parameters-path.sh	maxime.kermarquer	R	0:11	1	2G	2019-10-04T18:11:01	1	hmb002
126567 6	bigmem	run array parameters-path.sh	maxime.kermarquer	R	0:11	1	2G	2019-10-04T18:11:01	1	hmb002
126567 7	bigmem	run array parameters-path.sh	maxime.kermarquer	R	0:11	1	2G	2019-10-04T18:11:01	1	hmb002
126567 8	bigmem	run array parameters-path.sh	maxime.kermarquer	R	0:11	1	2G	2019-10-04T18:11:01	1	hmb002
100500				2	0.54	20	2 4 9 9	2010 10 01710 01 10		1 1 0 0 1

Allow for example to run same MATLAB script on different patient.

Use parameter :

/network/lustre/iss01/home/maxime.kermarquer/demo/job_array/parameters/patients/01

< M A T L A B (R) > Copyright 1984-2017 The MathWorks, Inc. R2017b (9.3.0.713579) 64-bit (glnxa64) September 14, 2017

For online documentation, see http://www.mathworks.com/support For product information, visit www.mathworks.com.

>> >> >> >> Hello! I'm a MATLAB script and I get the parameter : /network/lustre/iss01/home/maxime.kermarquer/demo/job_array/parameters/patients/**01** >> >> >> >> >> Use parameter :

/network/lustre/iss01/home/maxime.kermarquer/demo/job_array/parameters/patients/02

< M A T L A B (R) > Copyright 1984-2017 The MathWorks, Inc. R2017b (9.3.0.713579) 64-bit (glnxa64) September 14, 2017

For online documentation, see http://www.mathworks.com/support For product information, visit www.mathworks.com.

>> >> >> >> Hello! I'm a MATLAB script and I get the parameter : /network/lustre/iss01/home/maxime.kermarquer/demo/job_array/parameters/patients/**02** >> >> >> >> >>



02/06/2020

Jobs dependency

Options

> SLURM gives the possibility to set up pipeline of jobs with dependencies between them

- > sbatch command with --dependency option, a type of dependency and a jobid
- > Type of dependency
 - Begin execution after the specified jobs have begun execution.
 > after
 - Begin execution after the specified jobs have terminated.
 Afterany
 - Begin execution after the specified jobs have successfully executed afterok
 - Begin execution after the specified jobs have terminated with failed state afternotok



Jobs dependency Dependency graph

ICM Institut du Cerveau et de la Moelle épinièr



02/06/2020

Jobs dependency Example





58 Using the ICM cluster

02/06/2020

Scheduling policy Priority

Priority are determined with the following expression:

```
Priority =
(PriorityWeightAge) * (age_factor) +
(PriorityWeightFairshare) * (fair-share_factor) +
(PriorityWeightJobSize) * (job_size_factor)
```

- Slurm's fair-share is a floating point number between 0,0 and 1,0.
- > Reflects the amount of computing resources the user's jobs have consumed.
- Higher is the value, higher is the placement in the queue of jobs waiting to be scheduled.
- > The limits per user are **300 cores and 1200G** but we can raise them for special requests.



Scheduling policy

Classification of jobs

> To share equitably the use of the cluster, jobs are classify for normal partition

	Class 1	Class 2	Class 3	Class 4	Class 5
% Disponible	100,00%	80,00%	50,00%	30,00%	20,00%
Nb cores	700	560	350	210	140
Memory (Mb)	3194880	2555904	1597440	958464	638976
Memory (Gb)	3120	2496	1560	936	624
Nodes	25	20	12	7	5

	≤ 4h	≤ 1 semaine	≤ 2 semaine	≤ 3 mois
4G				
8G				
16G				
32G				
64G				
122G				



Deep learning Which framework are available ?

5



CHERCHER, TROUVER, GUÉRIR, POUR VOUS & AVEC VOUS.

Deep learning Resources

> To perform Deep Learning ICM, you can use the GPU node with NVIDIA Tesla P100 cards

> In python modules are installed the DL frameworks







Conclusion



CHERCHER, TROUVER, GUÉRIR, POUR VOUS & AVEC VOUS.

Conclusion Contact and materials

\succ The slides will be available



The wiki (only accessible from the ICM network) <u>https://dokuwiki.icm-institute.org/doku.php?id=cluster:wiki</u>

> If you need help, you can request

- > Open an issue in *Calculateurs-07* category
- > By mattermost, direct message to Maxime Kermarquer or computer cluster channel
- By email : <u>maxime.kermarquer@icm-institute.org</u>



02/06/2020



Thank you for your attention !



65 Using the ICM cluster

02/06/2020