# Introduction to SLURM

for users

Jordi Blasco
(jordi.blasco@nesi.org.nz)

Introduction
0000000

Running a job
0000000000000000000000

System Information
00000

## Outline

**Introduction**
○●○○○○○○

Running a job
○○○○○○○○○○○○○○○○○○○○○○○

System Information
○○○○○

# About SLURM

## About SLURM

**Introduction**
○○●○○○○

Running a job
○○○○○○○○○○○○○○○○○○○○○○○

System Information
○○○○○

## About SLURM

### About SLURM

- **SLURM** was an acronym for Simple Linux Utility for Resource Management.
- Development started in 2002 at Lawrence Livermore National Laboratory as a simple resource manager for Linux clusters
- Has evolved into a capable job scheduler
- About 500,000 lines of C code. Not Simple anymore.
- Now is called : Slurm Workload Manager
- Supports AIX, Linux, Solaris, other Unix variants
- Used on many of the world's largest computers
- Also used in the coolest facility in the world : NeSI
- Commercial support provided by SchedMD

**Introduction**
○○○●○○○

Running a job
○○○○○○○○○○○○○○○○○○○○○○○○

System Information
○○○○○

## New Features introduced with Slurm

### Already deployed features

- Full control over CPU and Memory usage
- Better Scheduling techniques & performance
- Job Array support
- Better integration with MPI
- Interactive sessions support
- High Availability (2 masters)
- Debugger friendly
- Topology aware (better MPI performance)
- Privacy Environment

**Introduction**
○○○○●○○

Running a job
○○○○○○○○○○○○○○○○○○○○○○○○○

System Information
○○○○○

## New Features introduced with Slurm

### Future features planed to be deployed

- Kernel Level Checkpointing & Restart
- Job Migration
- Shared FlexLM integration
- Job profiling (`srun --profile=All`)

**Introduction**
ooooo●o

Running a job
oooooooooooooooooooo

System Information
ooooo

Resource Management

## Node and Job States

- Nodes
  - state (up/down/idle/allocated/mix/drained)
- Jobs
  - queued/pending and running
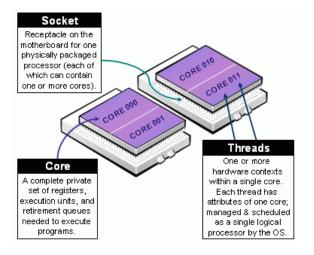  - suspended/preempted
  - cancelled/completed/failed

**Introduction**
○○○○○○●

Running a job
○○○○○○○○○○○○○○○○○○○○○○○

System Information
○○○○○

## Definitions of Socket, Core, & Thread



Figure : Definitions of Socket, Core, & Thread. Source SchedMD

Introduction
0000000

Running a job
●0000000000000000000000

System Information
00000

## Running a job

### SLURM Commands

- sbatch – submits a script job. (=llsubmit)
- scancel – cancels a running or pending job. (=llcancel)
- srun – runs a command across nodes.
- sbcast – Transfer file to a compute nodes allocated to a job.
- interactive – opens an interactive job session.
- sattach – Connect stdin/out/err for an existing job or job step.

Introduction
0000000

Running a job
0●0000000000000000000

System Information
00000

## Running a job

### srun : Simple way to manage MPI, OpenMP, pthreads & serial jobs

- Slurm provides a single command line to manage all the MPI flavours, OpenMP, Pthreads and serial applications
- Users don't need to worry about MPI flags and options for each MPI implementation mpirun/mpiexec/mpiexec.hydra
- The tool is called **srun** and it is **mandatory** for submitting jobs in the cluster.

Introduction
0000000

Running a job
00●0000000000000000000

System Information
00000

Running a job

### Commonly used SLURM variables

- $SLURM_JOBID
- $SLURM_JOB_NODELIST : (example sb[004,006])
- $SLURM_NNODES (Number of nodes)
- $SLURM_SUBMIT_DIR (Directory from which the job was submitted)

Introduction
ooooooo

Running a job
ooo●oooooooooooooooooo

System Information
ooooo

## Examples

### sbatch

```
login-01.uoa.nesi.org.nz ~ $ vim testjob.sl
login-01.uoa.nesi.org.nz ~ $ cat testjob.sl
#!/bin/bash
#SBATCH --nodes=10
srun echo "running on : $(hostname)"
srun echo "allocation : $SLURM_NODELIST"

login-01.uoa.nesi.org.nz ~ $ sbatch testjob.sl
Submitted batch job 11109

login-01.uoa.nesi.org.nz ~ $ cat slurm-11109.out
running on : wm001
allocation : wm[001-010]
```

Introduction
ooooooo

Running a job
ooo●oooooooooooooooooo

System Information
ooooo

## Examples

### sbatch

```
login-01.uoa.nesi.org.nz ~ $ vim testjob.sl
login-01.uoa.nesi.org.nz ~ $ cat testjob.sl
#!/bin/bash
#SBATCH --nodes=10
srun echo "running on : $(hostname)"
srun echo "allocation : $SLURM_NODELIST"

login-01.uoa.nesi.org.nz ~ $ sbatch testjob.sl
Submitted batch job 11109

login-01.uoa.nesi.org.nz ~ $ cat slurm-11109.out
running on : wm001
allocation : wm[001-010]
```

Introduction
ooooooo

Running a job
ooo●oooooooooooooooooooo

System Information
ooooo

## Examples

### sbatch

```
login-01.uoa.nesi.org.nz ~ $ vim testjob.sl
login-01.uoa.nesi.org.nz ~ $ cat testjob.sl
#!/bin/bash
#SBATCH --nodes=10
srun echo "running on : $(hostname)"
srun echo "allocation : $SLURM_NODELIST"

login-01.uoa.nesi.org.nz ~ $ sbatch testjob.sl
Submitted batch job 11109

login-01.uoa.nesi.org.nz ~ $ cat slurm-11109.out
running on : wm001
allocation : wm[001-010]
```

Introduction
○○○○○○○

Running a job
○○○○●○○○○○○○○○○○○○○○○○

System Information
○○○○○

## Submitting a Job

### Standard Job Script Directives

```
#!/bin/bash
#SBATCH -J JobName
#SBATCH -A uoa99999              # Project Account
#SBATCH --time=08:00:00          # Walltime
#SBATCH --mem-per-cpu=4096       # memory/cpu (in MB)
#SBATCH --ntasks=2               # 2 tasks
#SBATCH --cpus-per-task=4        # number of cores per tasks
#SBATCH --nodes=1                # number nodes
#SBATCH -C sb                    # sb=Sandybridge,wm=Westmere
```

Introduction
0000000

Running a job
00000●000000000000000

System Information
00000

## Submitting a Job

### Optional Job Script Directives

```
#SBATCH --mail-type=end
#SBATCH --mail-user=jordi.blasco@nesi.org.nz
#SBATCH -D /path_to_working_directory/
```

Introduction
0000000

Running a job
000000●00000000000000

System Information
00000

## Submitting a Serial Job

### Job Description Example : Serial

```
#!/bin/bash
#SBATCH -J Serial_JOB
#SBATCH -A uoa99999        # Project Account
#SBATCH --time=01:00:00    # Walltime
#SBATCH --mem-per-cpu=8132 # memory/core (in MB)

srun my_serial_binary
```

Introduction
0000000

Running a job
00000000●000000000000

System Information
00000

Submitting a OpenMP Job

### Job Description Example : SMP

```
#!/bin/bash
#SBATCH -J OpenMP_JOB
#SBATCH -A uoa99999        # Project Account
#SBATCH --time=01:00:00    # Walltime
#SBATCH --mem-per-cpu=8132 # memory/core (in MB)
#SBATCH --cpus-per-task=8  # 8 OpenMP Threads

srun my_openmp_binary
```

Introduction
0000000

Running a job
0000000●000000000000

System Information
00000

Submitting a MPI Job

### Job Description Example : MPI

```
#!/bin/bash
#SBATCH -J MPI_JOB
#SBATCH -A uoa99999         # Project Account
#SBATCH --time=01:00:00     # Walltime
#SBATCH --mem-per-cpu=8132  # memory/core (in MB)
#SBATCH --ntasks=2          # number of tasks

srun my_mpi_binary
```

Introduction
0000000

Running a job
0000000000●000000000000

System Information
00000

## Submitting a Hybrid(MPI+OpenMP) Job

### Job Description Example : Hybrid(MPI+OpenMP)

```bash
#!/bin/bash
#SBATCH -J Hybrid_JOB
#SBATCH -A uoa99999          # Project Account
#SBATCH --time=01:00:00      # Walltime
#SBATCH --mem-per-cpu=8132   # memory/core (in MB)
#SBATCH --ntasks=4           # number of tasks
#SBATCH --cpus-per-task=8    # 8 OpenMP Threads
#SBATCH --nodes=1            # Can be range eg --nodes=2-4

srun my_binary_hybrid
```

Introduction
0000000

Running a job
0000000000●0000000000000

System Information
00000

## Submitting an Array Job

### Array Job

- Slurm job arrays offer a mechanism for submitting and managing collections of similar jobs quickly and easily.

- In general, array jobs are useful for applying the same processing routine to a collection of multiple input data files.

- Array jobs offer a very simple way to submit a large number of independent processing jobs.

Introduction
0000000

**Running a job**
00000000000●0000000000

System Information
00000

## Submitting an Array Job

### Array Job Syntax

- Job array with index values between 1 and 1000
  `--array=1-1000`

- Job array with index values of 1, 3, 5 and 7
  `--array=1,3,5,7`

- Job array with index values between 1 and 7 with a step size
  of 2 (i.e. 1, 3, 5 and 7)
  `--array=1-7:2`

Introduction
0000000

Running a job
0000000000**000**000000000

System Information
00000

## Submitting an Array Job

### Array Job example

```
#!/bin/bash
#SBATCH -J JobArray
#SBATCH --time=01:00:00      # Walltime
#SBATCH -A uoa99999          # Project Account
#SBATCH --mem-per-cpu=8132   # memory/core (in MB)
#SBATCH --cpus-per-task=4    # 4 OpenMP Threads
#SBATCH --array=1-1000       # Array definition

srun my_binary_array $SLURM_ARRAY_TASK_ID
```

## Submitting an Array Job

### Array Job

To submit 1,000 element job array sbatch blast_array.sl
Submit time < 1 second
Environment variable with array index: SLURM_ARRAY_TASK_ID

Introduction
0000000

Running a job
00000000000000●0000000

System Information
00000

Submitting an Array Job

### Array Job

The management is really easy:

```
$ squeue -u sbae335
         JOBID PARTITION     NAME    USER ST       TIME NODES NODELIST
  28317_[1-1000]    high  SungSHM sbae335 PD       0:00     1 (Priority)
27817_[196-1000]    high Sung_BLA sbae335 PD       0:00     1 (Resources)
     27817_184      high Sung_BLA sbae335 R     4:15:27     1 wm001
     27817_185      high Sung_BLA sbae335 R     4:15:27     1 wm001
     27817_186      high Sung_BLA sbae335 R     4:15:27     1 wm001
     ...
     ...

$ scancel 28317_[900-1000]
```

## GRES subsystem

Generic Resource System to request special hardware like GPUs or Intel Phis

### Requesting GPUs

Add the following line in your submit script:

```
#SBATCH --gres=gpu:1       # GPUs per node
```

### Fine tuning

```
#SBATCH -C kepler # ask only for NVIDIA K20X
#SBATCH -C fermi  # ask only for NVIDIA Tesla M2090
```

Introduction
0000000

Running a job
0000000000000000●00000

System Information
00000

## GRES subsystem

### Requesting GPUs

```bash
#!/bin/bash
#SBATCH -J GPU_JOB
#SBATCH --time=01:00:00      # Walltime
#SBATCH -A uoa99999          # Project Account
#SBATCH --ntasks=4           # number of tasks
#SBATCH --ntasks-per-node=2  # number of tasks per node
#SBATCH --mem-per-cpu=8132   # memory/core (in MB)
#SBATCH --cpus-per-task=4    # 4 OpenMP Threads
#SBATCH --array=1-1000       # Array definition
#SBATCH --gres=gpu:2         # GPUs per node
#SBATCH -C kepler

srun my_binary_cuda_mpi
```

Introduction
○○○○○○○

Running a job
○○○○○○○○○○○○○○●○○○○○

System Information
○○○○○

## GRES subsystem

### Requesting Intel Phi (MIC)

Add the following line in your submit script:

```
#SBATCH --gres=mic:1    # Intel Phi per node
```

Introduction
0000000

Running a job
00000000000000000000000000000

System Information
00000

Job dependencies

### Job dependencies

Add the following line in your submit script:

`--dependency=afterok:$SLURM_JOB_ID`

## Running a job

### Interactive Job Session

```
[4845] login-01.uoa.nesi.org.nz ~ $interactive -h
Usage: interactive [-A] [-a] [-c] [-m] [-J]

Mandatory arguments:
 -A: account
Optional arguments:
 -a: architecture (default: wm, values sb=SandyBridge wm=Westmere)
 -c: number of CPU cores (default: 1)
 -m: amount of memory (GB) per core (default: 1 [GB])
 -J: job name
example : interactive -A nesi99999 -a wm -c 4 -J MyInteractiveJob

Written by: Alan Orth <a.orth@cgiar.org>
Modified by: Jordi Blasco <jordi.blasco@nesi.org.nz>
```

Introduction
○○○○○○○

Running a job
○○○○○○○○○○○○○○○○○○○○○○●○

System Information
○○○○○

# Limits

### Current limits in the cluster

- Max array size : 1000
- Max number of submitted jobs : 10,000
- Max allocatable memory per node : 92GB (Westmere), 124GB (SandyBridge), 508GB LargeMemory nodes)
- Number of cores per node : 12 (Westmere), 16 (SandyBridge), 40 LargeMemory nodes)

Introduction
0000000

Running a job
0000000000000000000000000

System Information
00000

Temporary File Systems

## Temporary File Systems

- $TMP_DIR (local filesystem)
- $SCRATCH_DIR (shared filesystem)
- $SHM_DIR (local RAM filesystem)

Introduction
0000000

Running a job
00000000000000000000000

System Information
●0000

## System Information

### System Information

- squeue – shows the status of jobs. (=llq)
- sinfo – provides information on partitions and nodes. (=llstatus)
- sview – GUI to view job, node and partition information.
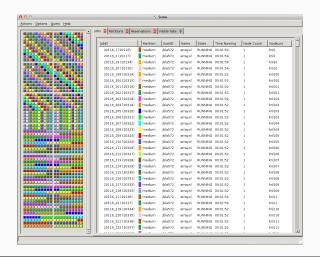- smap – CLI to view job, node and partition information.

Introduction
0000000

Running a job
0000000000000000000000

System Information
0●000

## System Information

### squeue

Show jobid and allocated nodes for running jobs of the user jblasco:

```
[4925] login-01.uoa.nesi.org.nz ~$ squeue
JOBID PARTITION     NAME     USER ST     TIME  NODES NODELIST(REASON)
24258     high Migrate- jbla572 PD     0:00      4 (Resources)
24259     high Migrate- jbla572 PD     0:00      4 (Priority)
24257     high Migrate- jbla572  R     0:27    512 sb[1-512]
```

# System Information

## sview

# Questions & Answers