

Job Submission in Ibex Cluster

Passant Hafez
HPC Applications Specialist
Supercomputing Core Lab

- 1) Job Submission
- 2) Creating a Job Submission Script
- 3) Job Monitoring

1) Job Submission

1) Job Submission

- Resource Manager: SLURM v 19.05.2
 - SLURM Resource Allocation Methods:
 - srun
 - salloc
 - sbatch
- (most options are common between them)
- Don't SSH directly to the compute nodes.

1) Job Submission

General Notes:

- Requested resources vs waiting time (also --exclusive)
- 2 queues:
 - batch queue (default): up to 14 days (more than that job will be pending FOREVER)
 - debug queue: up to 2 hrs

1) Job Submission

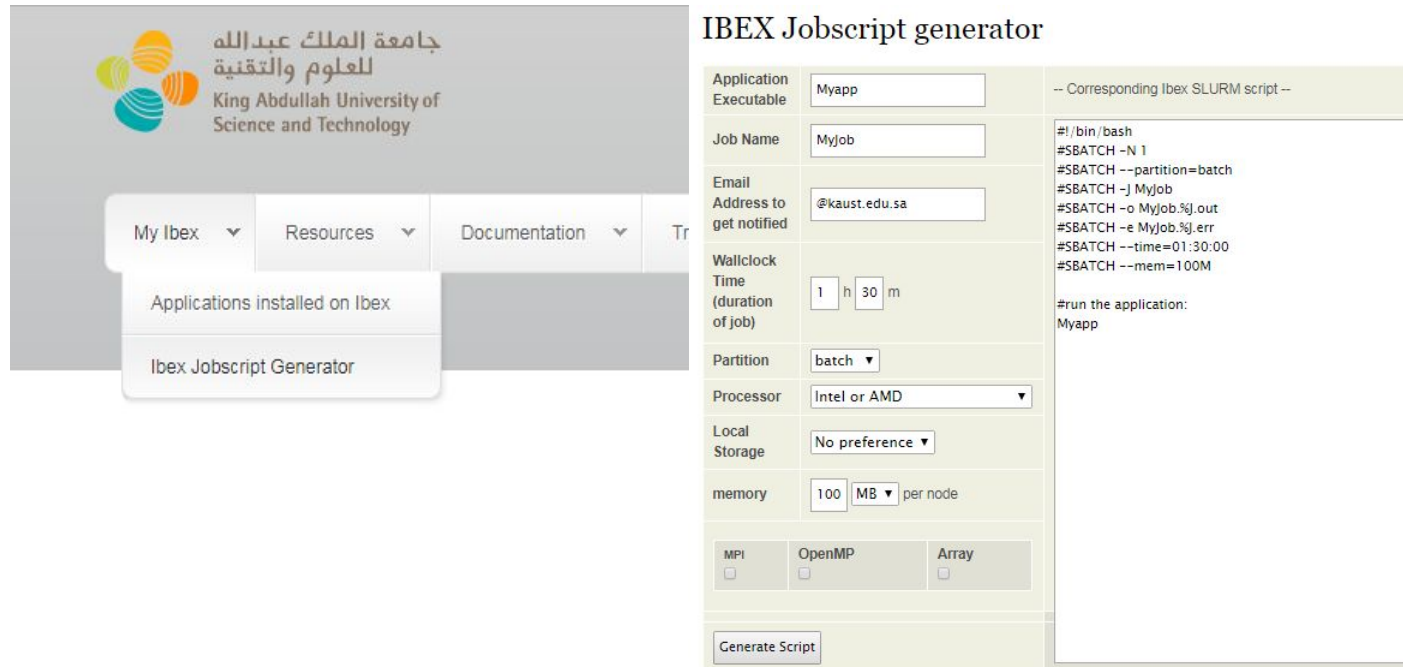
SLURM Options/Directives

-t or --time=(default it's no. of mins , format D-HH:MM)	--gres=gpu:p100:4	-c or --cpus-per-task= (default is 1)
-n or --ntasks=	or	--ntasks-per-core=
-N or --nodes=	--gres=gpu:4 --constraint=[p100]	--ntasks-per-socket=
--partition= or -p (default batch)	--mail-type=ALL (equivalent to BEGIN, END, FAIL, REQUEUE, and STAGE_OUT)	--ntasks-per-node=
-J or --job-name= (default is job script's name)	TIME_LIMIT	--threads-per-core=
-C or --constraint=	TIME_LIMIT_90 (reached 90% of time limit)	-a or --array=
--mem (default 2 GB/core , > ~362 GB job is marked as large-memory)	TIME_LIMIT_80 (reached 80% of time limit)	
--mem-per-cpu=	TIME_LIMIT_50 (reached 50% of time limit)	-d or --dependency=
-e or --error=	(useful for extensions)	--no-requeue
-o or --output=	--mail-user=	
(default both standard output and standard error are directed to a file of the name " slurm-%j.out " and Filename patterns (%J, %j, %A, %a, %x, %N)	-x or --exclude	
	-w or --nodelist=cn603-11-r,cn603-12-l or db202-02-[1-8]	

2) Creating a Job Submission Script

2) Creating a Job Submission Script

- Use Ibex Job Script Generator: <https://www.hpc.kaust.edu.sa/ibex/job>



IBEX Jobscrip generator

Application Executable	<input type="text" value="Myapp"/>	-- Corresponding Ibex SLURM script --
Job Name	<input type="text" value="Myjob"/>	<pre>#!/bin/bash #SBATCH -N 1 #SBATCH --partition=batch #SBATCH -j Myjob #SBATCH -o Myjob.%j.out #SBATCH -e Myjob.%j.err #SBATCH --time=01:30:00 #SBATCH --mem=100M #run the application: Myapp</pre>
Email Address to get notified	<input type="text" value="@kaust.edu.sa"/>	
Wallclock Time (duration of job)	<input type="text" value="1"/> h <input type="text" value="30"/> m	
Partition	<input type="text" value="batch"/>	
Processor	<input type="text" value="Intel or AMD"/>	
Local Storage	<input type="text" value="No preference"/>	
memory	<input type="text" value="100"/> MB <input type="text" value="per node"/>	
MPI	<input type="checkbox"/>	
OpenMP	<input type="checkbox"/>	
Array	<input type="checkbox"/>	
<input type="button" value="Generate Script"/>		

2) Creating a Job Submission Script

Notes:

- Make sure you're not loading any conflicting modules before submitting the job, or just add **module purge** before loading other modules in the script.
- Use **#SBATCH** with any of the previously mentioned options

For example: `#SBATCH --time=00:30:00`

SLURM Output Environment Variables:

SLURM controller sets some variables in the environment of the batch script, for example:

- **SLURM_JOB_ID**
- **SLURM_JOB_NODELIST**
- **SLURM_ARRAY_JOB_ID** → Array's master job ID number
- **SLURM_ARRAY_TASK_ID** → Job array ID (index) number

2) Creating a Job Submission Script

Filename Pattern:

sbatch allows for a filename pattern to contain one or more replacement symbols, which are a percent sign "%" followed by a letter, for example:

%j represents jobid of the running job.

%x Job name.

%A Job array's master job allocation number.

%a Job array ID (index) number.

For more patterns check: <https://slurm.schedmd.com/sbatch.html#lBAH>



جامعة الملك عبد الله
للعلوم والتقنية
King Abdullah University of
Science and Technology

When is my job starting?

Is my job running?

3) Job Monitoring

Where's the output
of my job?

3) Job Monitoring

- **queue** (for running and pending jobs)
- **scontrol show job <job_id>**
- **sacct -j <job_id>** (for any job state)
- **seff <job_id>** (for finished jobs)

3) Job Monitoring

sacct --format

Fields available:

Account	AdminComment	AllocCPUS	AllocGRES
AllocNodes	AllocTRES	AssocID	AveCPU
AveCPUFreq	AveDiskRead	AveDiskWrite	AvePages
AveRSS	AveVMSize	BlockID	Cluster
Comment	Constraints	ConsumedEnergy	ConsumedEnergyRaw
CPUTime	CPUTimeRAW	DerivedExitCode	Elapsed
ElapsedRaw	Eligible	End	ExitCode
Flags	GID	Group	JobID
JobIDRaw	JobName	Layout	MaxDiskRead
MaxDiskReadNode	MaxDiskReadTask	MaxDiskWrite	MaxDiskWriteNode
MaxDiskWriteTask	MaxPages	MaxPagesNode	MaxPagesTask
MaxRSS	MaxRSSNode	MaxRSSTask	MaxVMSize
MaxVMSizeNode	MaxVMSizeTask	McsLabel	MinCPU
MinCPUNode	MinCPUTask	NCPUS	NNodes
NodeList	NTasks	Priority	Partition
QOS	QOSRAW	Reason	ReqCPUFreq
ReqCPUFreqMin	ReqCPUFreqMax	ReqCPUFreqGov	ReqCPUS
ReqGRES	ReqMem	ReqNodes	ReqTRES
Reservation	ReservationId	Reserved	ResvCPU
ResvCPURAW	Start	State	Submit
Suspended	SystemCPU	SystemComment	Timelimit
TimelimitRaw	TotalCPU	TRESUsageInAve	TRESUsageInMax
TRESUsageInMaxNode	TRESUsageInMaxTask	TRESUsageInMin	TRESUsageInMinNode
TRESUsageInMinTask	TRESUsageInTot	TRESUsageOutAve	TRESUsageOutMax
TRESUsageOutMaxNode	TRESUsageOutMaxTask	TRESUsageOutMin	TRESUsageOutMinNode
TRESUsageOutMinTask	TRESUsageOutTot	UID	User
UserCPU	WCKey	WCKeyID	WorkDir