Submitting and Running Jobs





Scott French NERSC User Services Group

New User Training February 23, 2015







- Most jobs are parallel, using 10s to 100,000+ cores
- Production runs execute in batch mode
- Interactive and debug jobs are supported for up to 30 minutes
- Typically run times are a few to 10s of hours.
 - Each machine has different limits.
 - Limits are necessary because of MTBF and the need to accommodate 5,500 users' jobs
- Also a number of 'serial' jobs
 - Typically some kind of pleasantly parallel simulation





Login Nodes and Compute Nodes



- Each supercomputer has 3 types of nodes visible to users
 - Login nodes
 - Compute nodes
 - Application launcher or "MOM" nodes
- Login nodes
 - Edit files, compile codes, run UNIX commands
 - Submit batch jobs
 - Run short, small utilities and applications
- Compute nodes
 - Execute your application; dedicated to your job
 - No direct login access
- Application launcher or "MOM" nodes
 - Execute your batch script commands
 - Carver: "head" compute node
 - Edison / Hopper: shared "service" node (not a compute node)





Launching Parallel Jobs







- An "application launcher" executes your code
 - Distributes your executables to all your nodes
 - Starts concurrent execution of N instances of your program
 - Manages execution of your application
 - On Edison / Hopper: the launcher is called "aprun"
 - On Carver: "mpirun"
- Only the application launcher can start your application on compute nodes
- You can't run the application launcher from login nodes







- To run a job on the compute nodes you must write a "batch script" that contains
 - Batch directives to allow the system to schedule your job
 - An aprun or mpirun command that launches your parallel executable
- Submit the job to the queuing system with the qsub command
 - % qsub my_batch_script





Edison - Cray XC30





- 133,824 cores, 5,576 nodes
- "Aries" interconnect
- 2 x 12-core Intel 'Ivy Bridge'
 2.4 GHz processors per node
- 24 processor cores per node, 48 with hyperthreading
- 64 GB of memory per node
- 357 TB of aggregate memory

- 2.7 GB memory / core for applications
- /scratch disk quota of 10 TB
- 7.6 PB of /scratch disk
- Choice of full Linux operating system or optimized Linux OS (Cray Linux)
- Intel, Cray, and GNU compilers







```
#PBS -q debug
```

```
#PBS -1 mppwidth=96
```

```
#PBS -1 walltime=00:10:00
```

```
#PBS -N my_job
```

```
cd $PBS_O_WORKDIR
aprun -n 96 ./my_executable
```









Job directives: instructions for the batch system

- Submission queue
- How many compute cores to reserve for your job (/ 24 = # nodes)
- How long to reserve those nodes
- Optional: what to name STDOUT files, what account to charge, whether to notify you by email when your job finishes, etc.







Change from home directory to job submission directory

- Script is initially run from your home directory, which is not advisable (as we mention in the filesystem intro)
- You will see much better performance if your job reads / writes from one of the high-performance scratch filesystems









Launches parallel executable on the compute nodes

- Carries over (partial) login environment
- Controls how your executable:
 - maps to processors on the compute nodes (e.g. how many tasks?)
 - accesses the memory on each processor









mppwidth is number of compute cores requested for your job

- mppwidth = 24 x # of nodes on Edison (and Hopper)
- must be greater than or equal to the number of tasks requested (-n)





- #PBS -1 mppwidth=192
- #PBS -1 walltime=00:10:00
- #PBS -N my_job

cd \$PBS_O_WORKDIR aprun -n 96 -N 12 ./my_executable

> –N = number of tasks per node Might do this to get more memory / task Note that mppwidth has changed accordingly


```
#PBS -q regular
#PBS -1 mppwidth=96
#PBS -1 walltime=00:10:00
#PBS -N my job
cd $PBS O WORKDIR
export OMP NUM THREADS=6
aprun -n 16 -d 6 -N 4 -S 2 ./hybrid.x
```

A more complex example for mixing MPI and OpenMP:

16 tasks (**-n**), 4 on each node (**-N**), 6 OpenMP threads per task (**-d**), • assign 2 tasks to each NUMA node (**-S**)

Many more examples on www.nersc.gov

Carver - IBM iDataPlex

- 9,984 compute cores
- 1,202 compute nodes
- 2 quad-core Intel Nehalem 2.67 GHz processors per node
- 8 processor cores per node
- 24 GB of memory per node (48 GB on 80 "fat" nodes)
- 2.5 GB / core for applications (5.5 GB / core on "fat" nodes)
- Interconnect: InfiniBand 4X QDR

- NERSC global /scratch directory quota of 20 TB
- Full Linux operating system
- PGI, GNU, Intel compilers

Number of MPI tasks (similar to aprun's -n argument) Distribution across compute cores primarily controlled via the ppn value on the PBS -1 directive (other ways possible)

 You can run small parallel jobs interactively for up to 30 minutes (ex. is for Edison)

login% qsub -I -1 mppwidth=48
[wait for job to start]
mom% cd \$PBS_O_WORKDIR
mom% aprun -n 48 ./mycode.x

- Both Carver and the Crays have a special queue for running serial jobs
 - A single process running on a single core
 - Each serial node can run up to 12 jobs from different users on Carver and 24 jobs on Hopper or Edison

```
#PBS -q serial
#PBS -l walltime=00:10:00
#PBS -N my_job
cd $PBS_O_WORKDIR
./myexecutable
```


- Once your job is submitted, it will start when resources are available
- You can monitor it with:
 - qstat -a
 - qstat -u username
 - showq
 - qs
 - NERSC web site

https://www.nersc.gov/users/live-status/global-queue-look/

https://www.nersc.gov/users/job-logs-and-analytics/completed-jobs/

- MyNERSC

https://my.nersc.gov

I qsub'd my job, but it's not running!

• You are not alone!

That's what the batch queue is for!

source: itu.dk

- Your jobs will wait until the resources are available for them to run.
- Your job's place in the queue is a mix of time and priority, so line jumping is allowed, but it may cost more.

There are per user, per machine job limits. Here are the limits on Edison as of October 16, 2014.

Specify the #PBS –q	ese queues queue_nai	with M me	Not these!													
	Submit Queue	Execution Queue	Nodes	Physical Cores	Max Wallclock (hours)	Relative Priority	Run Limit	Queued Limit	Queue Charge Factor							
	debug	debug	1-512	1-12,288	30 mins	2	2	2	1							
	ccm_int ¹	ccm_int	1-512	1-12,288	30 mins	2	2	2	1							
		reg_small	1-682	1-16,368	48 hrs	3	24	24	1							
	regular	reg_med	683- 2048	16,369- 49,152	36 hrs	3	8	8	0.6							
		reg_big	2049- 4096	49,153- 98,304	36 hrs	2	2	2	0.6							
		reg_xbig	4097- 5462	98,305- 131,088	12 hrs	2	2	2	0.6							
	ccm_queue	ccm_queue	1-682	1-16,368	48 hrs	3	16	16	1							
	premium	premium	1-2048	1-49,152	36	1	1	1	2							
	low	low	1-682	1-16,368	24	4	8	8	0.5							
	killable ²	killable	1-682	1-16,368	48 hrs	3	8	8	1							
	serial ³	serial	1	1	48 hrs	-	50	50	1/24							

- Submit shorter jobs, they are easier to schedule
 - Checkpoint if possible to break up long jobs
 - Short jobs can take advantage of 'backfill' opportunities
 - Run short jobs just before maintenance
- Very important: make sure the wall clock time you request is accurate
 - As noted above, shorter jobs are easier to schedule
 - Many users unnecessarily enter the largest wall clock time possible as a default

- Your repository is charged for each node your job was allocated for the entire duration of your job.
 - The minimum allocatable unit is a node (*except for the serial queues*). Hopper and Edison have 24 cores / node, so your minimum charge is 24*walltime.

MPP hours = (# nodes) * (# cores / node) * (walltime) * (QCF) * (MCF)

- Example: 96 Edison cores for 1 hour in regular queue
 MPP hours = (4) * (24) * (1 hour) * (1) * (2) = 192 MPP hours
- Serial jobs are charged with: (walltime) * (MCF)
- If you have access to multiple repos, pick which one to charge in your batch script

#PBS -A repo_name

- Each machine has a "machine charge factor" (MCF) that multiplies the "raw hours" used
 - Edison MCF = 2.0
 - Hopper MCF = 1.0
 - Carver MCF = 1.5
- Each queue has a "queue charge factor" (QCF) and corresponding relative scheduling priorities
 - Premium QCF = 2.0
 - Low QCF = 0.5
 - Regular (and everything else) QCF = 1.0 (Hopper: 0.8)

• On Edison:

Jobs requesting more than 682 nodes (reg_med, reg_big, reg_xbig queues) get a 40% discount (QCF = 0.6)

NERSC Web pages

Hopper:

http://www.nersc.gov/users/computational-systems/hopper/running-jobs/

Edison:

http://www.nersc.gov/users/computational-systems/edison/running-jobs/

Carver:

http://www.nersc.gov/users/computational-systems/carver/running-jobs/

Contact NERSC Consulting:

- Toll-free 800-666-3772
- 510-486-8611, #3
- Email <u>consult@nersc.gov</u>.

Thank You

Average Queue Wait Time

	Hours	Req	ues	sted																																												
Nodes	<1	1	2	3	4	5	6	7	8		9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27 2	8 2	9 3	0 31	32	33	34	35	36	37 3	8 3	9 40	41	42	43	44	45	46 4	7 4	848+
1	1.6	5.3	1.1	5.6	5.6	8.5	13	38	58		24	9.8	17	18	10	67	15	31	16	42	6.2	17	10	26	26	29	11	19	20	7 2	6 4	2 13	12	27	34	21	46	. 1	8 2	6 26	17	20	. 0	28	20	21 2	4 4	3 29
2	1.7	7.58	3.0	9.6	22	15	25	8.6	21		14	15	13	23	18	22	23	15	14	30	34	34	31	25	30	41	21	35	33 2	20 1	9 2	50.0	41	0.0	36	29	47	400.	1 3	3 45	6 O . C	35	17	58	31	33 3	3 5	8 76
3	1.4	17	1.8	5.1	4.4	8.0	48	22	11	2,59	0.3	15	5.0	19	28	29	19	26	14	23	23	51	0.0	35	37	45	66	89	.0 3	32 4	0 5	10.0	42	0.0	1.0	501	12	44 3	10	0 35	6 O . C	0.0	8	.3	22	39 3	7 9	3 24
4	3.6	11	5.3	141	10.0	12	14	99	16		27	27	29	38	16	24	26	39	20	71	28	48	0.0	14	29	52	60	440	. 0 3	85 0.	11	2	54	0.0		69	540	.0 2	5 (5 148	0.0	0.0).00	.00		9	3 6	5 44
5	2.2	1.1	7.2	6.0	5.9	9.7	20	13	16		6.5	49	8.6	19	13	38	25	12	28	35	33	58	43	25	50	76	61	. 0 0	.00.	0 2	3 4	30.0	47	0.0	1.0	43	49	.0 4	4 0	6 44	0.0	45).00	.00	. 00	.0 2	5 7	4 32
6	1.1	6.0	3.5	10	6.4	13	11	13	41		11	40	21	19	24	21	23	34	445	32		43	0.0	35	84	80		27	7	8	5	8 34	56	49	28	53	59	.0 4	2 0	0 53	80.0	0.0).00	.00		56 9	310	2 34
7	0.4	5.4	2.5	7.5	9.9	6.5	14	3.2	15		7.8	38	31	24	. 0	0.0	0.0	0.0	0.0	0.0	0.0	47	0.0	83	16	57).00	. 0 0	.00.	.00.	0 3	70.0	47	0.0	1.0	55	19	.00.	00	00.0	0.0	0.0	1	.04	. 00	.00.	6	4 35
8	1.8	16	3.4	36	9.8	16	19	17	16		16	59	41	224	.7	27	36	37	0.0	33	57	64	0.0	0.0	73	68	18	18	5	6	2	50.0	75	44	28		56	.00.	00.	00.0	0.0	0.0		86	17	7	7 6	7 35
9	0.4	2.12	2.5	5.8	8.7	15	11	12	14			32	18	18	. 0 0		36	23	0.0	0.0	0.0	130	0.0	0.0	0.0	45) . 0 0	. 0 0		.00.	23	6	0.0	0.0).0		63	.00.	00	0 51	0.0	0.0).0	90	49	72 5	9 10	4 16
10	1.8	7.28	3.6	4.8	13	11	12	17	21		18	28	26	23	50	81	72	40	0.0	13	25	80	0.0	72	50	110	60	. 0 0	.00.	.00.	6	2	0.0	0.0	1.0	36	47	.00.	00.	51	0.0	0.0	1	.13	91	5	7 8	7 76
11	10	26	7.0	6.1	16	13	30	15	15		11	22	27	31	38	27	31	25	39	45	61	75	0.0	0.0	125	56		410		.00.	5	9	0.0	0.0		47	77	410.	00	0 52	20.0	0.0).00		48	58 8	110	2 45
12	2.4	9.6	10	5.3	23	47	96	16	36		31	254	18	32 9	.5	0.0	45	40	0.0	22	53	60	0.0	0.0	113	57	61	. 0 0	3	34 O .	3	60.0	0.0	0.0) . 0	1	38	.00.	00	0 38	63	0.0		68	37	63 4	7 5	9 42
13	0.4	1.4	7.7	5.0	10	19	13	6.48	3.8		10	80		18	. 0	0.0	0.0	0.0	0.0	0.0	0.0	48	0.0	0.0	0.0	57).00	. 0 0	.00.	00.	5	0.0	0.0	0.0) . 0		15	.00.	00	64	0.0	0.0).00	.00		70	6	6
14	1.6	1.1	18	12	8.7	22	17	14	27		27	27	1.0	29		0.0	0.0	0.0	0.0	105	0.0	97	0.0	0.0	37	54).00	. 0 0	.00.	00.	0 3	80.0	0.0	0.0) . O (37	.00.	00	00.0	0.0	53	1	51	. 0 0	.00.	10	9 120
15	12	6.7	15	16	7.7	9.9	62	81	14		47	53	37	25	47	86	38	17	38	0.0	0.0	72	0.0	0.0	65	204	54		480.	00.	0 4	40.0	0.0	0.0	1.0	11	93	.00.	00	00.0	0.0	0.0).00		51 0	.00.	7	7 87
16	1.6	1.9	10	27	11	22	39	18	18		30	27	28	22		124	42	40	43	29	66	58	77	0.0	65	65	79	33	450.	9	6 6	8140	0.0	53		20	530	.00.	00	105	0.0	0.0).00	.00	. 0 0	8	0 5	5 32
17-19	1.3	26	19	5.3	10	8.0	20	13	46		21	30	11	43		104	27	30	0.0	75	0.0	76	0.0	0.0	144	112	132	. 0 0	15	52 6	3 3	60.0	116	0.0) . 0	64	98	.0 4	90	69	0.0	0.0		67	. 0 0	.00.	7	7 53
20-23	8.9	8.3	9.3	19	20	19	23	24	29		46	32	25	29	68	53	57	154	0.0	0.0	67	57	76	41	61	66	59	.00	.00.	3	9 7	5	202	0.0) . 0	59	61	.00.	00.	6	0.0	27).00	.00		44 4	1 7	5 36
24-31	2.6	2.6	1.8	11	15	24	48	15	46		30	19	11	25	66	47	43	47	69	57	67	53	0.0	28	66	125	78	.00	8	36 3	9 6	6	132	41		102	75	.00.	00	94	0.0	0.0).00	1	00	40 4	0 8	7 104
32-47	1.5	9.2	9.5	26	18	27	58	18	36		51	45	28	37	82	51	54	49	65	86	44	60	128	47	66	101	46	.00	.00.	0 2	1 8	9	134	43) . 0	67	67	8	4 4	9 88	47	0.0).00		62	65 7	9 8	8 55
48-63	6.0	12	9.3	15	18	18	32	54	83		43	31	65	40	1.0	47	64	32	0.0	74	74	52	0.0	0.0	71	114	62	.00	8	16	6	40.0	69	0.0	1.0	L08	72	.0 3	40	58	0.0	0.0).00	1	06	95 7	7 8	0 98
64-127	1.6	23	20	21	31	29	36	45	53		33	39	43	52	77	67	59	73	60	86	61	72	45	86	94	78	63	61	10)5 11	9 5	9 0 0	64	0.0		127	64	.00.	0 4	2 65	0.0	0.0).00	. 0	59	58	7	0101
128-255	4.8-17	/8.1	21	38	40	50	54	45	77		51	33	149	48	66	130	58	73	73	78	0.0	63	0.0	55	79	101).00	.00	.00.	6 4	1 7	30.0	34	0.0		175	64	.00.	00	0 45	5 50	0.0		98	_ 1	0910	29	1 63
256-511	7.8	41	47	38	83	47	131	30	67		51	56	32	44	43	71	51	79	2.5	105	58	95	0.0	0.0	98	863	349	.00	.00.	7	79	50.0	0.0	0.0).0() . 0	540	.00.	OD,	68	0.0	0.0).00	.00	. 0 0	. 11	210	5
512-1023	12	53	100	91	48	48	65	30	50		94	44	60	56	62	28	66	54	44	36	41	53	0.0	82	32	42	16	. 0	32 4	20.	8	7 0	144	0.0	36	82	76	.00.	00	00.0	0.0	0.0).00	.00			12	8
1024-1535	21	31	32	26	301	29	77	35	56		55	90	39	89	39	41	0.0	33	65	93	0.0	45	0.0	0.0	0.0	45).00	.00	.00.	00.	7	2 0	0.0	0.0		1392	08	.00.	00	00.0	0.0	0.0).00	.00	. 0 0	.00.	00.	00.0
1536-2047	24	43	21	11	35	32	66	0.0	48		0.0	47).0	35		0.0	0.0	39	0.0	31	0.0	142	0.0	19	0.0	64).00	.00	.00.	00.	00.	0.0	0.0	0.0).0(1	10	.00.	00.	00.0	0.0	0.0).00	.00	. 0 0	.00.	00.	00.0
2048-3071	30	32	59	41	32	45	50	49	45		0.0	106	1.0	21	46	0.0	70	203	78	144	101	33	0.0	0.0	67	53	31	.00	5	i 6	00.	0.0	0.0	0.0().0(500	.00.	00.	00.0	0.0	0.0).00	.00	. 00	.00,	00.	00.0
3072-4095	25	57	102	0.0	0.0	0.0	0.0	0.00	. 0		130	0.0		193	. 0	0.0	0.0	0.0	0.0	72	0.0	0.0	0.0	0.0	0.0	31).00	.00	.00.	00.	00.	0.0	0.0	0.0	1.0	0.0	. 0 0	.00.	00.	00.0	0.0	0.0).00	.00	.00	.00.	00.	00.0
4096-6143	133	75	75	50	50	0.0	129	60			14	0.0	.	0.00	. 0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0).00	. 0 0	.00.	00.	00.	0.0	0.0	0.0	. 0	0.0	. 0 0	.00.	OD.	00.0	0.0	0.0		.00	.00	.00.	00.	00.0
6144-9531	0.0	55		29	0.0	0.0	0.0	0.00).0		0.0	0.0).O	0.00		0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0).00	.00	.00.	00.	00.	00.0	0.0	0.0	0.0).00	. 0 0	.00.	00	00.0	0.0	0.0).00	.00	. 0 0	. 0 0 .	00.	00.0

