

High Performance Roundtable

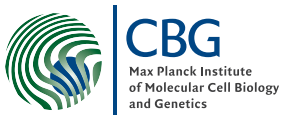
- Introducing Furiosa -

Òscar González¹ Peter Steinbach²

¹Computer Department

²Scientific Computing Facility

March 17th, 2016



1. Benchmarks
2. SLURM
3. Cluster Expansion
4. HPC Consolidation

1. Benchmarks
2. SLURM
3. Cluster Expansion
4. HPC Consolidation



- required to double check what was promised by vendor
- expected to be concluded until easter
- benchmark performance parts of it or entire cluster

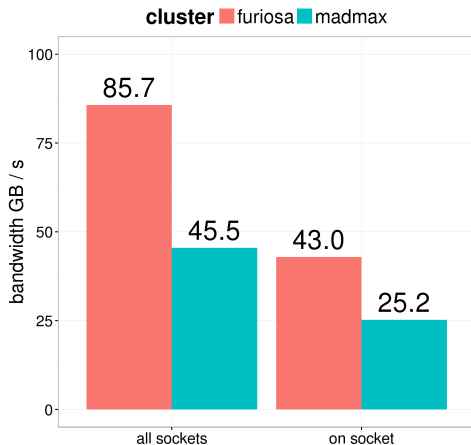
FINISHED memory
PRELIMINARY CPU+network
PRELIMINARY file system

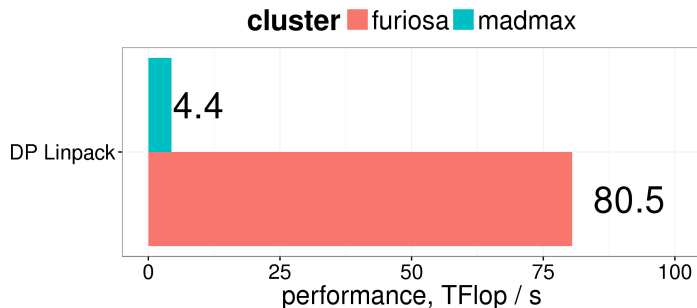
- stream benchmark 5.10
(simple operations on large arrays in memory)

- Triad

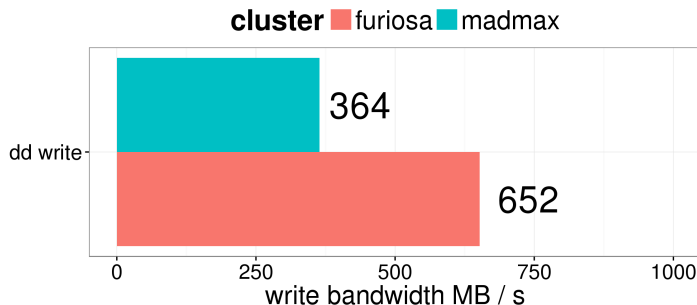
```
for(size_t i = 0; i < len; ++i)  
    a[i] = b[i] + c*d[i];
```

- double precision items
- array size > L3 cache size





- High Performance Linpack 2.1
- double precision workload
- huge matrix inversion through LU factorization stretching all nodes
- furiosa: **No. 1 in Top500 in 2004**
- enclosure (16 nodes): 11.5 TFlop/s



- simple dd call (sequential write of 48 GB file)
- single node i/o
- 2× bandwidth than madmax
- 10× storage volume than madmax
- expect > 20 GB/s cluster-wide

1. Benchmarks
- 2. SLURM**
3. Cluster Expansion
4. HPC Consolidation



- **Free and Open Source Job Scheduler** (GPL)
- 6 out of 10 Top500 clusters use it (June 2015)
- TUD uses it on taurus
- energy aware (built in profiling)
- scalable (islands of nodes)
- flexible (plugin system)
- fault tolerant
- ...



- **Free and Open Source Job Scheduler** (GPL)
- **6 out of 10 Top500 clusters** use it (June 2015)
- TUD uses it on taurus
- energy aware (built in profiling)
- scalable (islands of nodes)
- flexible (plugin system)
- fault tolerant
- ...

LSF on madmax

```
$ bsub -W 00:20 -o my.log -n 12 -R "span[hosts=1]" ./a.out
```

LSF on madmax

```
$ bsub -W 00:20 -o my.log -n 12 -R "span[hosts=1]" ./a.out
```

SLURM on furiosa

```
$ srun --time=00:20:00 -o my.log -N 1 -n 1 -c 24 ./a.out &
```

- srun does wait for job to complete
- how many nodes
-N 1
- how many tasks per node
-n 1
- how many cores per task
-c 24
- provide timing information to get better placement

lsf.script

```
#!/bin/bash
#BSUB -W 00:20
#BSUB -o my.log
#BSUB -n 12
#BSUB -R "span[hosts=1]"
```

```
./a.out
```

LSF on madmax

```
$ bsub < lsf.script
```

lsf.script

```
#!/bin/bash
#BSUB -W 00:20
#BSUB -o my.log
#BSUB -n 12
#BSUB -R "span[hosts=1]"

./a.out
```

LSF on madmax

```
$ bsub < lsf.script
```

slurm.script

```
#!/bin/bash
#SBATCH --time=00:20:00
#SBATCH -N 1
#SBATCH -n 1
#SBATCH -c 24
#SBATCH -o my.log

./a.out
```

SLURM on furiosa

```
$ sbatch slurm.script
```

from the command-line (very customizable)

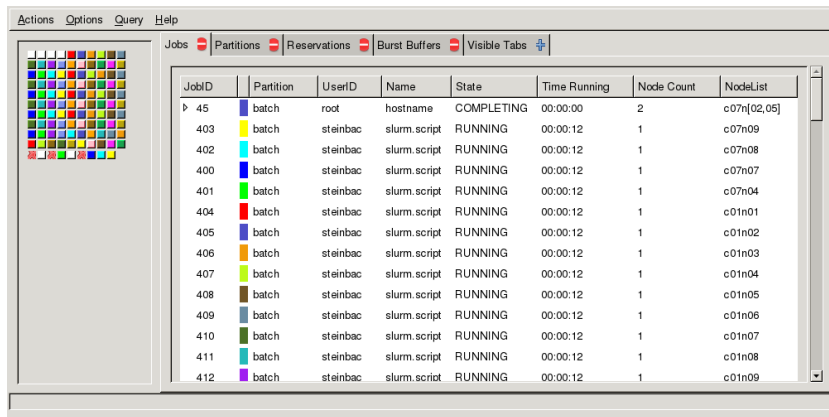
```
$ squeue -u ${USER}
JOBID PARTITION NAME USER ST TIME NODES NODELIST(REASON)
180 batch slurm.sc steinbac R 0:11 1 c07n04
181 batch slurm.sc steinbac R 0:11 1 c07n07
182 batch slurm.sc steinbac R 0:11 1 c07n08
...
$ squeue -t RUNNING -o "%.6i %p"
JOBID PRIORITY
180 0.00000081583858
181 0.00000081583858
182 0.00000081583858
...
```

- as every SLURM command, squeue very versatile
- -o supports 54 different columns types

Or you can have something more fancy ...

GTK2 based GUI (Xorg session required)

```
$ sview
```



The screenshot shows the sview GUI interface. At the top, there are menu items: Actions, Options, Query, and Help. Below the menu is a toolbar with buttons for Jobs, Partitions, Reservations, Burst Buffers, and Visible Tabs. On the left side, there is a color palette with various colored squares. The main area displays a table of jobs with the following columns: JobID, Partition, UserID, Name, State, Time Running, Node Count, and NodeList.

JobID	Partition	UserID	Name	State	Time Running	Node Count	NodeList
45	batch	root	hostname	COMPLETING	00:00:00	2	c07n(02,05)
403	batch	steinbac	slurm.script	RUNNING	00:00:12	1	c07n09
402	batch	steinbac	slurm.script	RUNNING	00:00:12	1	c07n08
400	batch	steinbac	slurm.script	RUNNING	00:00:12	1	c07n07
401	batch	steinbac	slurm.script	RUNNING	00:00:12	1	c07n04
404	batch	steinbac	slurm.script	RUNNING	00:00:12	1	c01n01
405	batch	steinbac	slurm.script	RUNNING	00:00:12	1	c01n02
406	batch	steinbac	slurm.script	RUNNING	00:00:12	1	c01n03
407	batch	steinbac	slurm.script	RUNNING	00:00:12	1	c01n04
408	batch	steinbac	slurm.script	RUNNING	00:00:12	1	c01n05
409	batch	steinbac	slurm.script	RUNNING	00:00:12	1	c01n06
410	batch	steinbac	slurm.script	RUNNING	00:00:12	1	c01n07
411	batch	steinbac	slurm.script	RUNNING	00:00:12	1	c01n08
412	batch	steinbac	slurm.script	RUNNING	00:00:12	1	c01n09

In case something goes wrong

```
$ scancel <job-id>
```

In case everything goes wrong

```
$ scancel -u $USER
```

- SLURM does not require anything like queues
- some clusters can live without them
(but require walltime limit, memory limits, ...)
- **We suggest:**
To clone the queues we have in madmax and set them up in furiosa!

- SLURM does not require anything like queues
- some clusters can live without them
(but require walltime limit, memory limits, ...)
- **We suggest:**
To clone the queues we have in madmax and set them up in furiosa!

What do you think?

1. Benchmarks
2. SLURM
- 3. Cluster Expansion**
4. HPC Consolidation

- budget available for cluster expansion in 2017
- question: What to use it for?

- budget available for cluster expansion in 2017
- question: What to use it for?

More CPUs



More Disk



GPUs*



What do our users think?

1. Benchmarks
2. SLURM
3. Cluster Expansion
- 4. HPC Consolidation**



- keep madmax operational
- reinstall OS to match furiosa's and upgrade to use SLURM
- ...

What do our users think?