

# Future of Batch Processing at CERN

## HEPiX Fall 2013

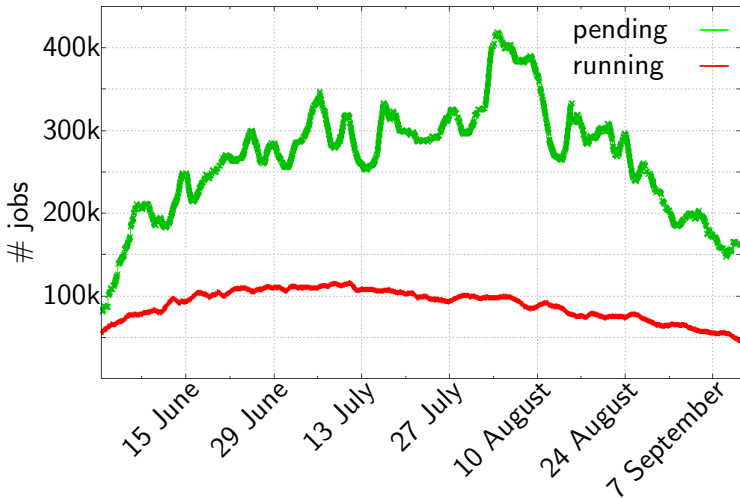
Jérôme Belleman   Daniel Pek  
CERN IT  
October 2013

- 1 The Future, LSF and its Shortcomings
- 2 Alternative Batch Systems
- 3 Initial Results

## Section 1

# The Future, LSF and its Shortcomings

- IBM LSF 7.0.6
- 4 000 nodes
  - SLC5 → SLC6
  - Physical → virtual ( $\approx$  1000 virtual so far)
- > 65 000 cores
- 400 000 jobs/day
- $\pm$ 70 000 running jobs



Goals	Concerns with LSF
30 000 to 50 000 nodes	6 500 nodes max
Cluster dynamism	Adding/Removing nodes requires reconfiguration
10 to 100 Hz dispatch rate	Transient dispatch problems
100 Hz query scaling	Slow query/submission response times

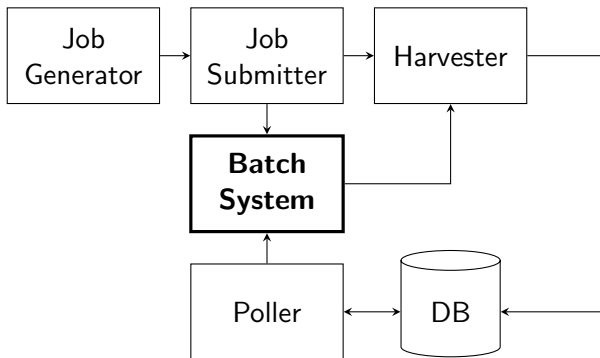
- Hired an LSF consultant
- Suggested a few minor enhancements:
  - Kernel parameters
  - Memory limits
  - CPU binding
  - Logging and accounting file handling
- No miraculous improvements
- Taught us a lot

## Section 2

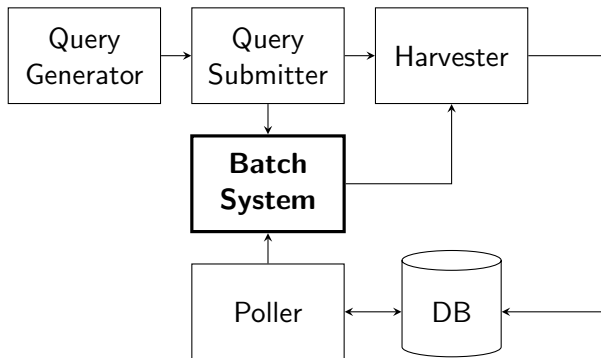
# Alternative Batch Systems



- SLURM 2.5.7
- HTCondor 8.1.0
- Son of Grid Engine 8.1.3
- LSF 8/9



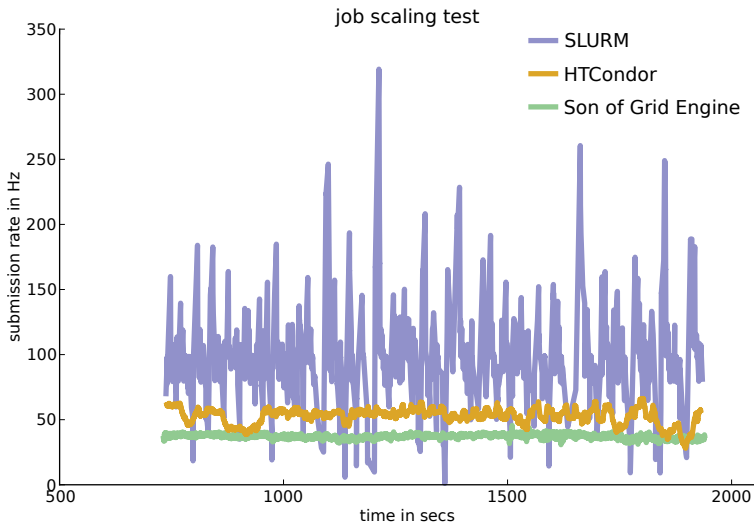
```
[  
  { "cmd":      "tests.d/psleep.py 120",  
    "count":   1000,  
    "instance": "batch",  
    "user":    "atlas",  
    "queue":   "debug" }  
]
```



```
[  
  { "cmd":      "tests.d/bjobs.py",  
    "count":   1,  
    "instance": "batch" }  
]
```

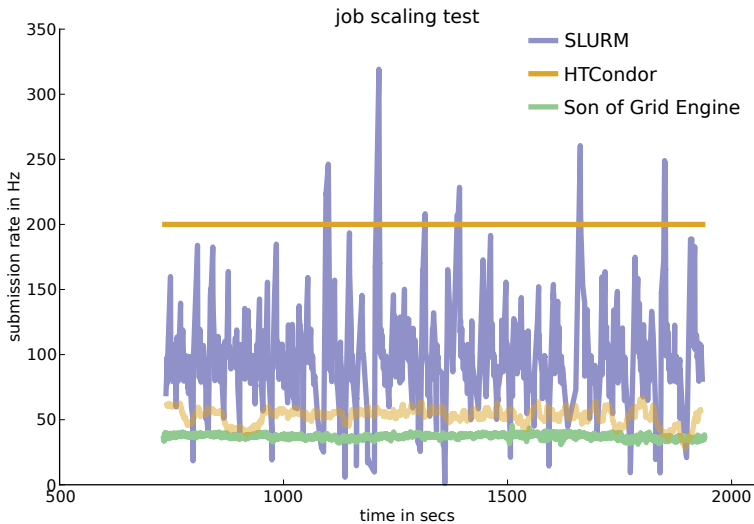
## Section 3

# Initial Results













### Requirements for measurement:

- Piggyback on existing resources
- Need to simulate many queries

### Experience:

- Slow SLURM startup with too many hosts?
- HTCondor made to scale out
- CPU load on the master(s)

	SLURM	HTCondor	Son of Grid Engine	Current LSF
Config (Puppet)	Struggle	Comfortable 1 <sup>st</sup> time, updates harder	Needs shared filesystem	Complex
Maturity	Easy to freeze	Trustworthy	Rough around the edges	Trustworthy
Dynamic	Not really	Yes	Yes	Not really
Doc	Poor	Solid	Solid	Solid
Community support	Rather uninterested	Very enthusiastic	Enthusiastic	Commercial

- Grid support
- Kerberos/AFS
- Accounting
- Host normalisation
- Fairshare scheduling
- Support for commercial applications
- IPv6?

### Replacement candidates:

- SLURM feels too young
- HTCondor mature and promising
- Son of Grid Engine fast, a bit rough

### What's next:

- Host scalability
- Query load
- Features



# Questions?