

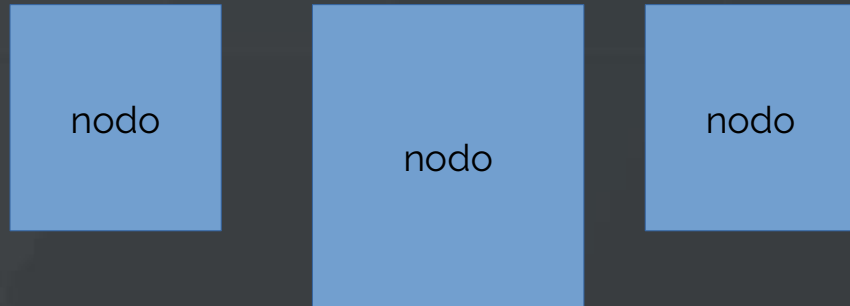
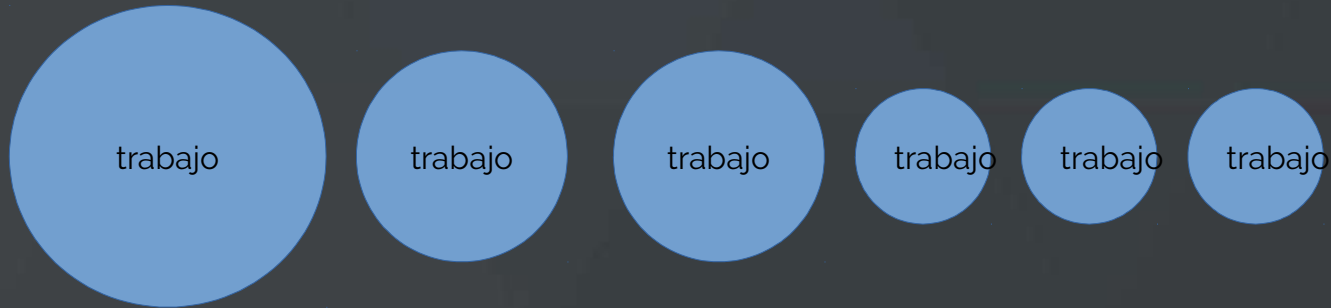
# Grid Engine y Slurm

HPC Admintech  
Valencia, Mayo 2018

Jesús Cuenca  
Senior HPC Consultant @ SIE

# Introducción sistema de colas

# Reto



# Solución

trabajo

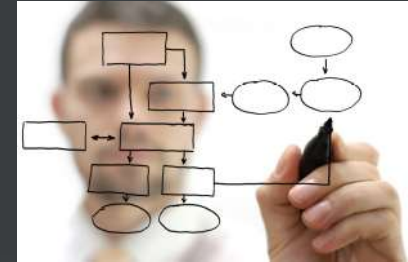
trabajo

trabajo

trabajo

trabajo

trabajo



nodo

nodo

nodo



# Grid Engine



# Historia

1992: Codine

2000: Sun

2001: Open Source

2010: Oracle. Closed source

2011: Univa Grid Engine (comercial, ~1 release/año – marzo 2017)

Son of Grid Engine (última versión: marzo 2016)

Open Grid Scheduler (abandonware)

[http://www.softpanorama.org/HPC/Grid\\_engine/history.shtml](http://www.softpanorama.org/HPC/Grid_engine/history.shtml)



# Arquitectura

Master host: monitoring + scheduling

Execution host: sge\_execd

Administration host

Submit host



# Arquitectura

“Parallel Environments” (PE)

Slots: máximo procesos

allocation rule: procesos/nodo (constante), fill\_up (todos los slots), round\_robin (de uno en uno)

start/stop scripts (MPI)



# SLURM



## Presentación

Simple Linux Utility for Resource Management

“Slurm es un sistema de gestión de clústers y planificación de trabajos escalable, tolerante a fallos y de código abierto (...)

Su uso es autocontenido (...)”



## Historia

2002: LLNL Slurm

2010: SchedMD

2018: Slurm 17.11.05 (varias releases por año)

[https://computing.llnl.gov/tutorials/linux\\_clusters/LinuxProjectReport.2002.08.18.pdf](https://computing.llnl.gov/tutorials/linux_clusters/LinuxProjectReport.2002.08.18.pdf)



PC / Linux

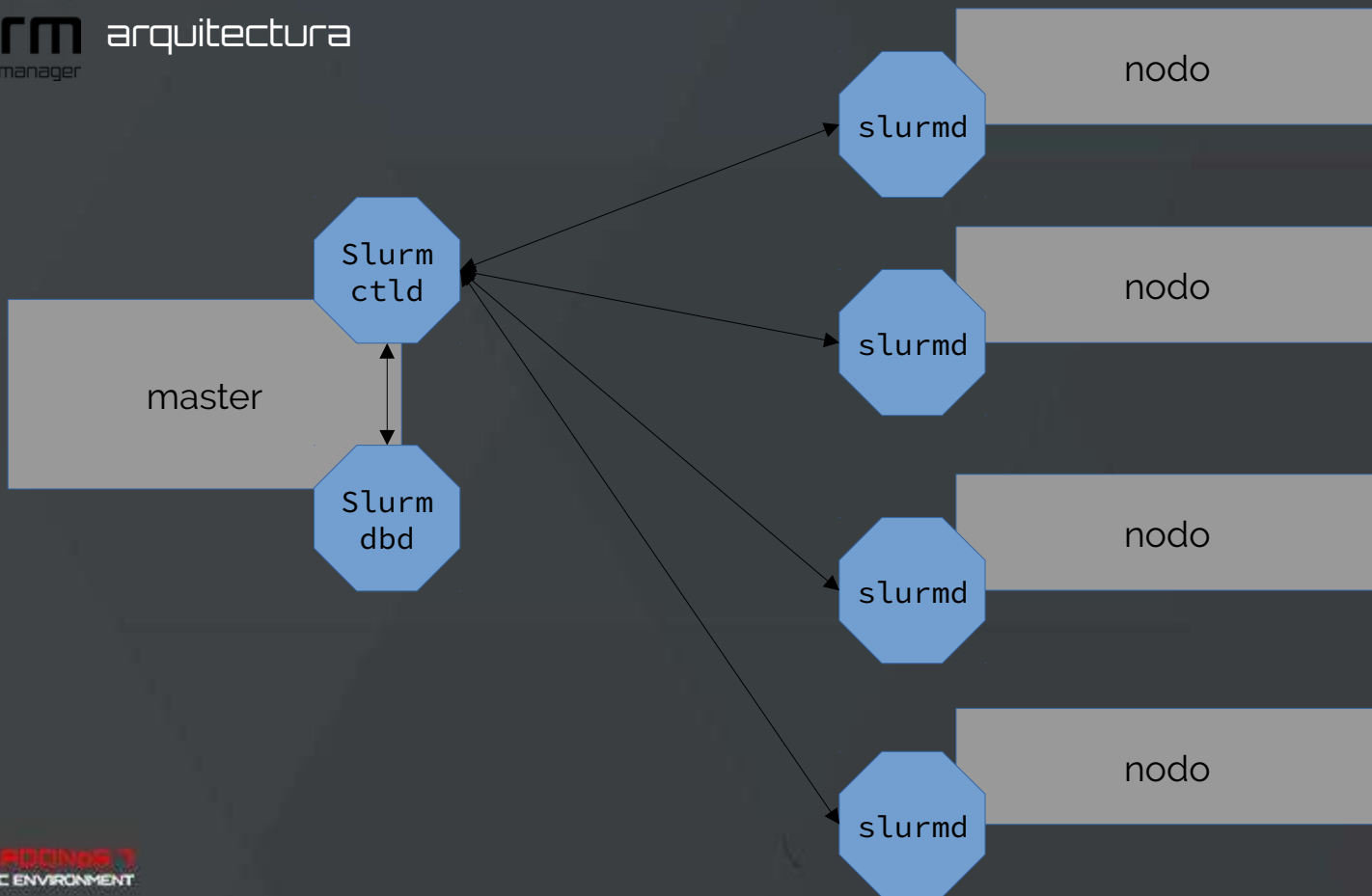
IBM BlueGene

Cray/XT

IBM Parallel Environment



Rank	System	Cores	Rmax (TFlop/s)	Rpeak (TFlop/s)	Power (kW)
1	<b>Sunway TaihuLight</b> - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway , <b>NRCPC</b> National Supercomputing Center in Wuxi China	10,649,600	93,014.6	125,435.9	15,371
2	<b>Tianhe-2 (MilkyWay-2)</b> - TH-IVB-FEP Cluster, Intel Xeon E5-2692 12C 2.200GHz, TH Express-2, Intel Xeon Phi 31S1P , <b>NUDT</b> National Super Computer Center in Guangzhou China	3,120,000	33,862.7	54,902.4	17,808
3	<b>Piz Daint</b> - Cray XC50, Xeon E5-2690v3 12C 2.6GHz, Aries interconnect , <b>NVIDIA Tesla P100</b> , <b>Cray Inc.</b> Swiss National Supercomputing Centre (CSCS) Switzerland	361,760	19,590.0	25,326.3	2,272
4	<b>Gyokou</b> - ZettaScaler-2.2 HPC system, Xeon D-1571 16C 1.3GHz, Infiniband EDR, PEZY-SC2 700Mhz , <b>ExaScaler</b> Japan Agency for Marine-Earth Science and Technology Japan	19,860,000	19,135.8	28,192.0	1,350
5	<b>Titan</b> - Cray XK7, Opteron 6274 16C 2.200GHz, Cray Gemini interconnect, <b>NVIDIA K20x</b> , <b>Cray Inc.</b> DOE/SC/Oak Ridge National Laboratory United States	560,640	17,590.0	27,112.5	8,209
6	<b>Sequoia</b> - BlueGene/Q, Power BQC 16C 1.60 GHz, Custom , <b>IBM</b> DOE/NNSA/LLNL United States	1,572,864	17,173.2	20,132.7	7,890
7	<b>Trinity</b> - Cray XC40, Intel Xeon Phi 7250 68C 1.4GHz, Aries interconnect , <b>Cray Inc.</b> DOE/NNSA/LANL/SNL United States	979,968	14,137.3	43,902.6	3,844
8	<b>Cori</b> - Cray XC40, Intel Xeon Phi 7250 68C 1.4GHz, Aries interconnect , <b>Cray Inc.</b> DOE/SC/LBNL/NERSC United States	622,336	14,014.7	27,880.7	3,939
9	<b>Oakforest-PACS</b> - PRIMERGY CX1640 M1, Intel Xeon Phi 7250 68C 1.4GHz, Intel Omni-Path , <b>Fujitsu</b> Joint Center for Advanced High Performance Computing Japan	556,104	13,554.6	24,913.5	2,719
10	<b>K computer</b> , SPARC64 VIIIfx 2.0GHz, Tofu interconnect , <b>Fujitsu</b> RIKEN Advanced Institute for Computational Science (AICS) Japan	705,024	10,510.0	11,280.4	12,660
16	<b>MareNostrum</b> - Lenovo SD530, Xeon Platinum 8160 24C 2.1GHz, Intel Omni-Path , <b>Lenovo</b> Barcelona Supercomputing Center Spain	153,216	6,470.8	10,296.1	1,632



/etc/slurm/slurm.conf  
(idéntica en todos los nodos)

### Parámetros:

Valores por defecto  
Recursos genéricos  
Re-encolado automático  
Características nodos cálculo  
Particiones

...

# Comandos

sbatch

scancel

squeue

scontrol show nodes

sinfo

sview

qsub

qdel

qstat

qhost

qhost -q

qmon





# Parámetros trabajo

Duración y programación horaria:

```
--time={DD-HH / HH:MM:SS}  
--begin=YYYY-MM-DDTHH:MM:SS  
--deadline=YYYY-MM-DDTHH:MM:SS
```

```
-l h_rt=[segundos]  
-a [YYMMDDhhmm]
```

Entorno

```
--workdir=  
--output=  
--error=  
--export  
--job-name=
```

```
-wd  
-o  
-e  
-V  
-N
```

# Parámetros trabajo

Cola

--partition=

-q

Tareas

--ntasks=

-pe [PE] [N]

Distribución de tareas

--tasks-per-node=

--nodelist=n,m

allocation\_rule en PE

-q [cola]@[nodo]



# Parámetros trabajo

## Recursos

### CPU

--cpus-per-task=

-pe [PE] [N]

### Genéricos

--gres=

-l

### Memoria

--mem=<size><unit>

-l h\_vmem=

--mem-per-cpu=<size><unit>

-l h\_vmem=



# Parámetros trabajo

Aviso

--mail-user=

-M

Dependencias

--dependency=

after:JOB\_ID

afterok:JOB\_ID

afternotok:JOB\_ID

singleton

-hold\_jid JOB\_ID



# Parámetros trabajo

## Variables informativas

\$SLURM\_JOB\_ID  
\$SLURM\_JOB\_NAME  
\$SLURM\_JOB\_NODELIST  
\$SLURM\_JOB\_PARTITION  
\$SLURM\_SUBMIT\_DIR  
\$SLURM\_SUBMIT\_HOST  
(...)



\$JOB\_ID  
\$JOB\_NAME  
\$PE\_HOSTFILE  
\$SLURM\_JOB\_PARTITION  
\$SGE\_O\_WORKDIR  
\$SGE\_O\_HOST  
(...)



**Partition:** {MaxCPUsPerNode, MaxMemPerCPU, MaxMemPerNode, MaxNodes, MaxTime}

**Account:** {GrpTres, MaxJobs...}

GrpTRES: “The total count of TRES able to be used at any given time from jobs running from an association”

TRES: Trackable RESources {cpu,mem,node,gres}

**QoS** (Quality of Service): {MinTRESPerJob, ...}

Age: tiempo de espera

Job size: nodos o CPUs solicitadas

TRES

QoS

Partition

Fair-share: reparto por cuentas / grupos

```
Job_priority =  
  (PriorityWeightAge) * (age_factor) +  
  (PriorityWeightFairshare) * (fair-share_factor) +  
  (PriorityWeightJobSize) * (job_size_factor) +  
  (PriorityWeightPartition) * (partition_factor) +  
  (PriorityWeightQOS) * (QOS_factor) +  
  SUM(TRES_weight_cpu * TRES_factor_cpu,  
      TRES_weight_<type> * TRES_factor_<type>,  
      ...)
```



Slurm.conf

SuspendTime

SuspendProgram

ResumeProgram

SuspendExcNodes

# SLURM demo

## MASTER

Munge

Slurm user

Paquetes slurm

Permisos

MySQL DB

Configuración Slurm (slurm.conf, slurmdbd.conf)

Slurmdbd

Slurmctld

sacctmgr

## BACKUP

Munge

Slurm user

Paquetes slurm

Permisos

(MySQL DB)

Replicar slurm.conf, crear slurmdb

Slurmdbd

Slurmctld

## NODOS

Munge

Slurm user

Paquetes slurm

Permisos

(Replicar configuración Slurm)

Slurmd

sinfo

scontrol

sacct

sreport

sview



Alta disponibilidad

Parar nodo `slurmctld` principal

Arrancar nodo, `munge`, `slurmctld`

# Gracias



# Contacto



Sistemas Informáticos Europeos

Calle Marqués de Mondejar nº 29

913 61 10 02

[www.sie.es](http://www.sie.es)



[/HPCSIE](https://www.facebook.com/HPCSIE)



[soporte@sie.es](mailto:soporte@sie.es)



[@HPCSIE](https://twitter.com/HPCSIE)



[+SistemasInformaticosEuropeosSLMadrid](https://plus.google.com/+SistemasInformaticosEuropeosSLMadrid)

