



MESOS.

Administrar clústers compartint recursos de forma dinàmica.

Joan Carles Martínez Rodríguez

Treball de Final Grau en Enginyeria Informàtica

Itinerari d'Arquitectura de Computadors i Sistemes Operatius

Consultor: Francesc Guim Bernat

Lliurament: 30/12/2014

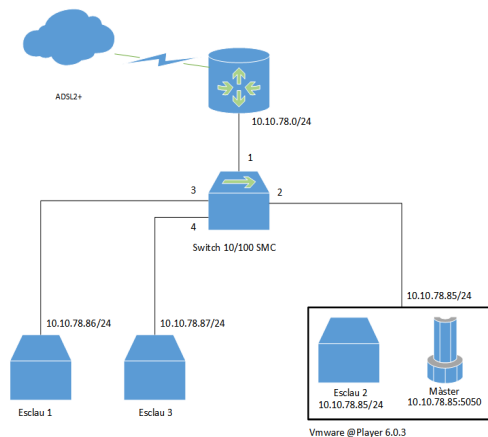
Motivació del Treball



Conèixer en detall el funcionament intern dels grans clústers de computació respecte a:

1. Conèixer problemàtiques i mecanismes d'eficiència en el tractament de grans volums de dades i computació.
2. Estudi de solucions actuals dedicades a la gestió eficient de recursos.
3. Implementació d'una solució real de gestió de clústers que pugui ser escalable.
4. Aplicar de forma pràctica els coneixements adquirits al llarg dels estudis.

Introducció al Treball



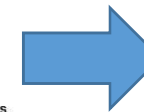
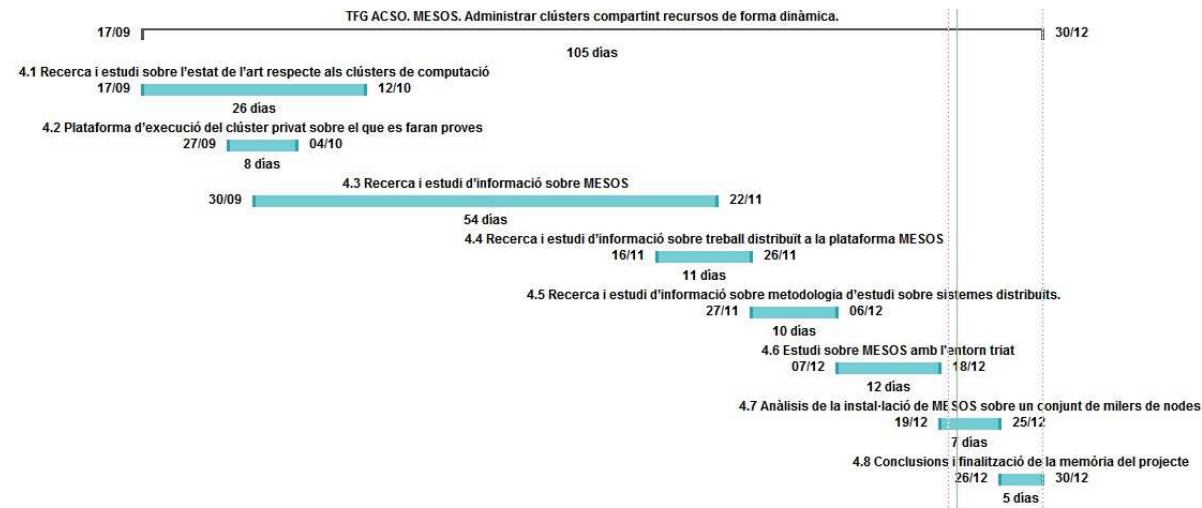
Parts principals:

- Estudi de l'estat de l'art en solucions de gestió eficient de recursos als clústers.
- Estudi, instal·lació i gestió d'un clúster administrat per MESOS.
- Obtenció de mètriques del clúster implementat i conclusions de la gestió de MESOS.
- Conclusions finals del treball desenvolupat.

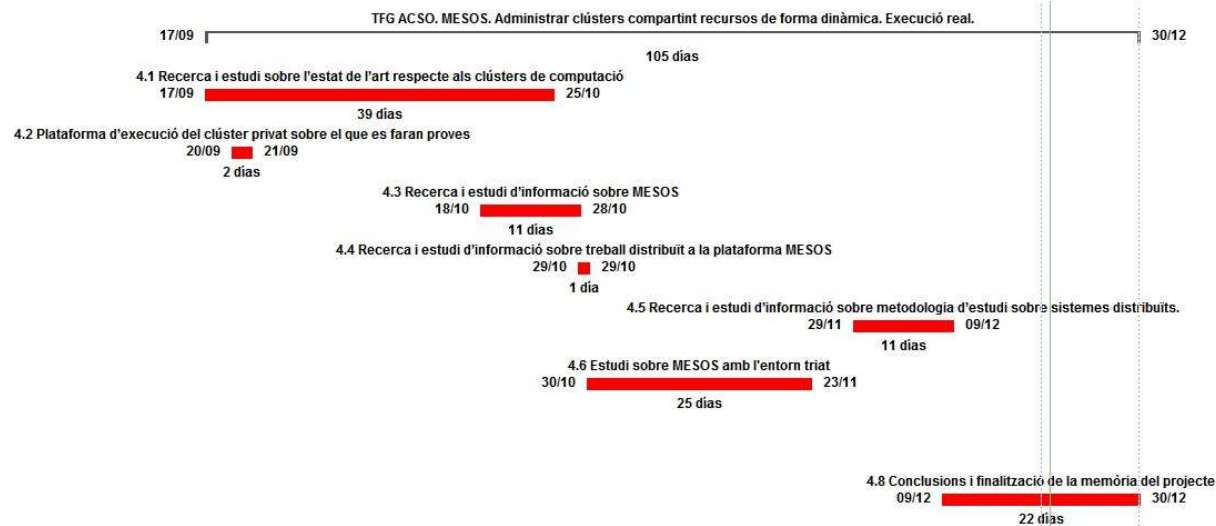
Objectius del Treball

1. Conèixer les bases del disseny de clústers de computació i les millores d'eficiència.
2. Estudi del gestor de clústers MESOS, en les següents funcions:
 - 2.1 Creació d'un clúster totalment operatiu.
 - 2.2 Administració del clúster amb la plataforma MESOS.
 - 2.3 Execució activa del clúster utilitzant un dels entorns d'aplicació.
 - 2.4 Projecció en el control de milers de nodes per a treball distribuït.
3. Extracció de conclusions i comparació amb altres gestors actuals.

Planificació temporal

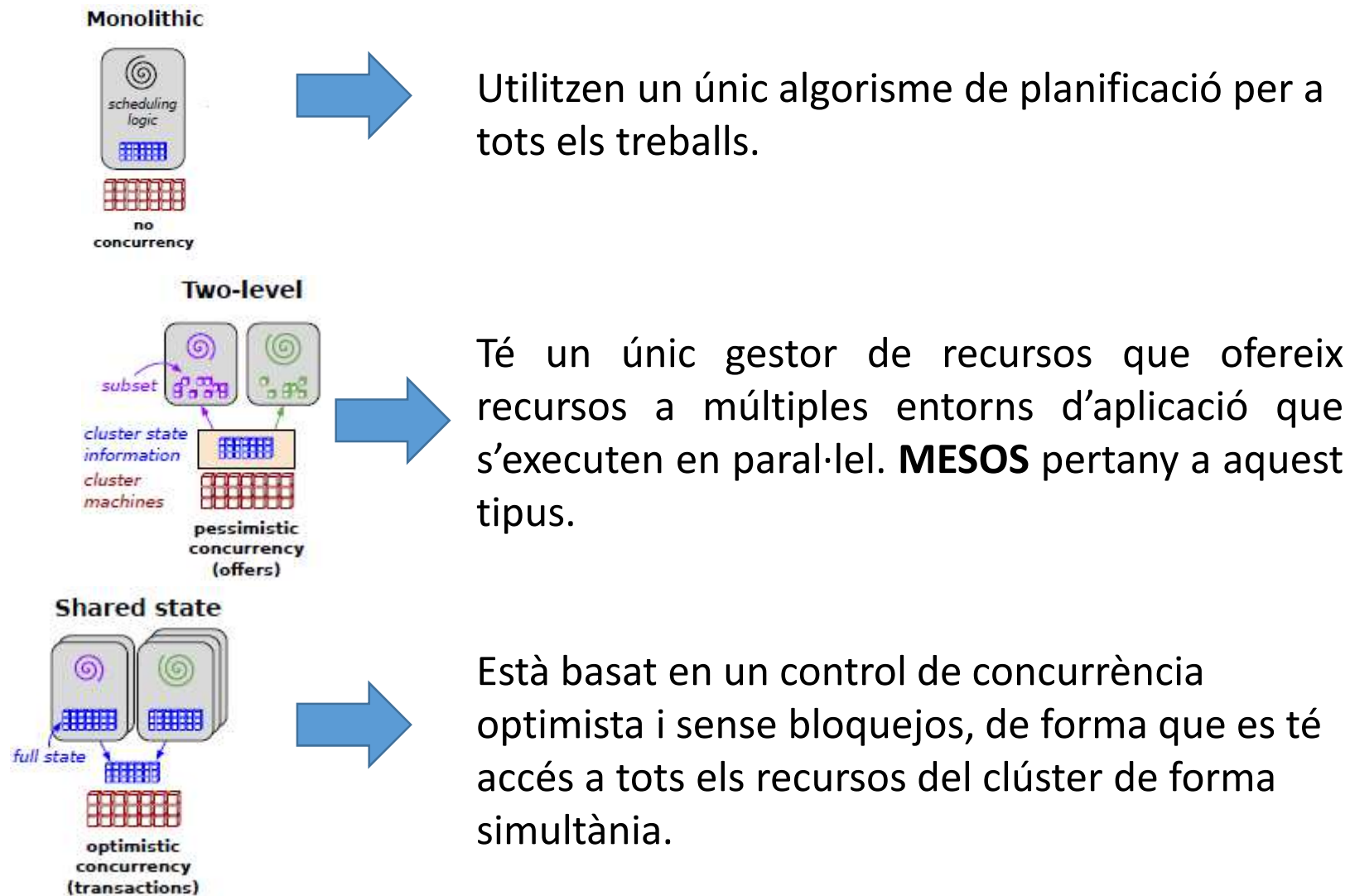


Planificació inicial
Previst: 408 hores



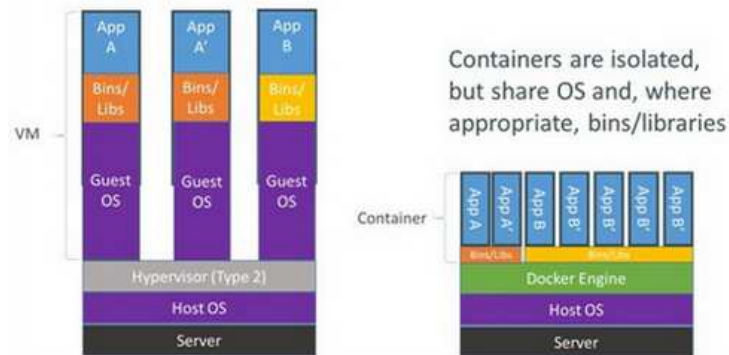
Planificació real
Cost: 509 hores

Estat de l'art. Exemples de gestió de clústers



Estat de l'art. Contenedors DOCKER

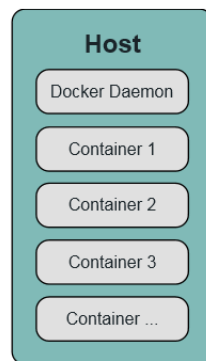
Containers vs. VMs



- Nou paradigma en virtualització.
- Menys consum de recursos que amb VM.
- Aïllament entre aplicacions.
- Ideal per aplicacions que utilitzen el mateix sistema operatiu.
- Arquitectura basada en client-servidor (dimonis).
- Utilitza plantilles d'aplicació (imatges).

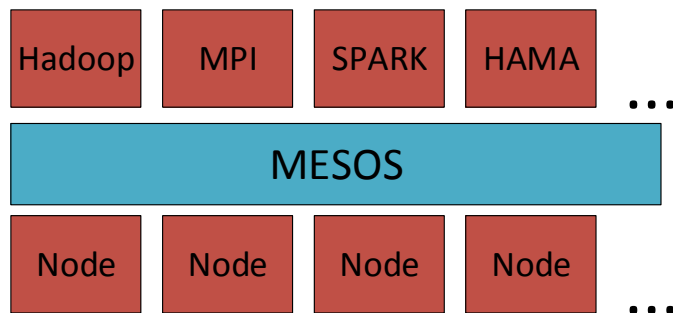
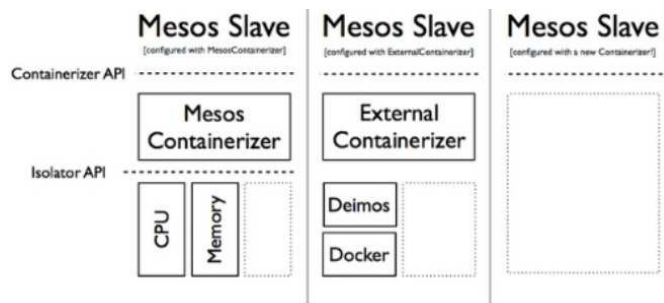
Docker Client

```
docker pull
docker run
docker ...
```



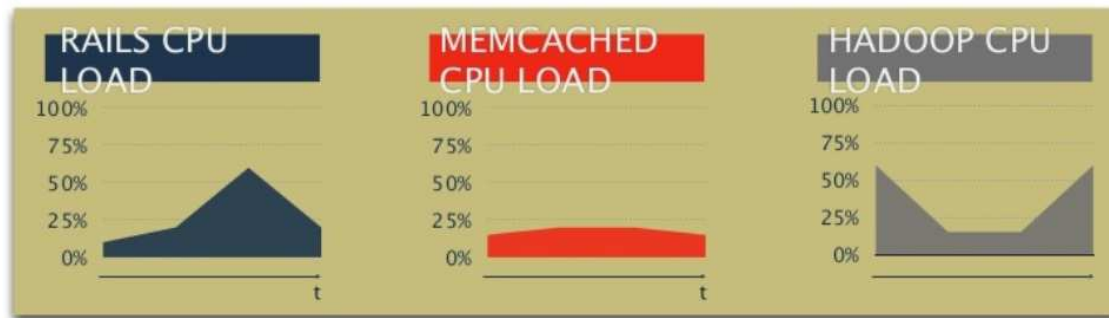
Docker Index

Què és MESOS?



- És un sistema de *kernel* distribuït.
 - Les aplicacions a executar utilitzen entorns d'aplicació (frameworks).
- Capacitats:
- Escalable fins 10.000 nodes.
 - Capacitat de recuperació d'errors.
 - Contenedors LINUX o DOCKER.
 - Planificació multirekurs (memòria, CPU, disc i ports).
 - API Java, Python i C++ APIs, per desenvolupar noves aplicacions paral·leles.
 - Web UI per visualitzar l'estat del clúster..

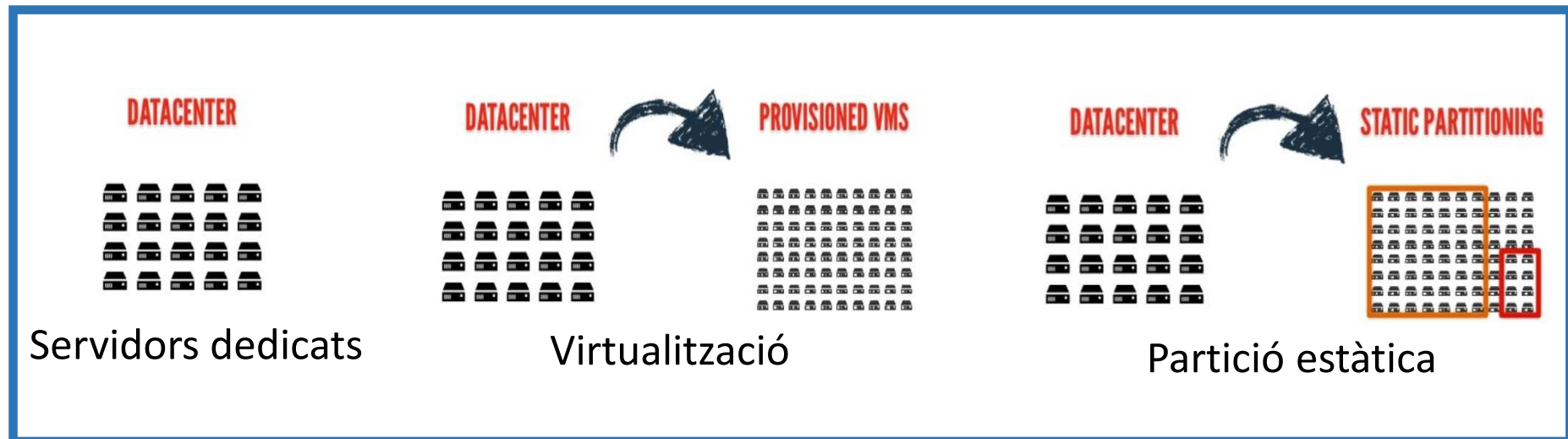
MESOS. Compartir recursos de forma dinàmica.



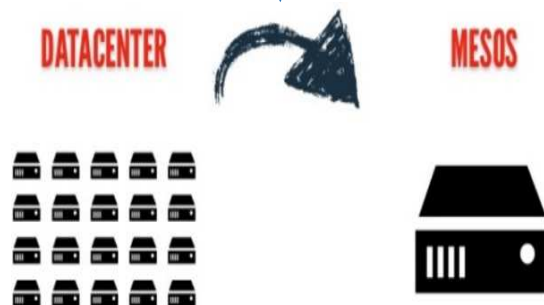
Entorn que s'adapti a les necessitats de **forma dinàmica i flexible** o recursos dedicats de forma estàtica?



MESOS. Evolució tecnològica



Solució: Visió simplificada del clúster a administrar (*pool de recursos*)



The screenshot shows the Mesos web interface with the following sections:

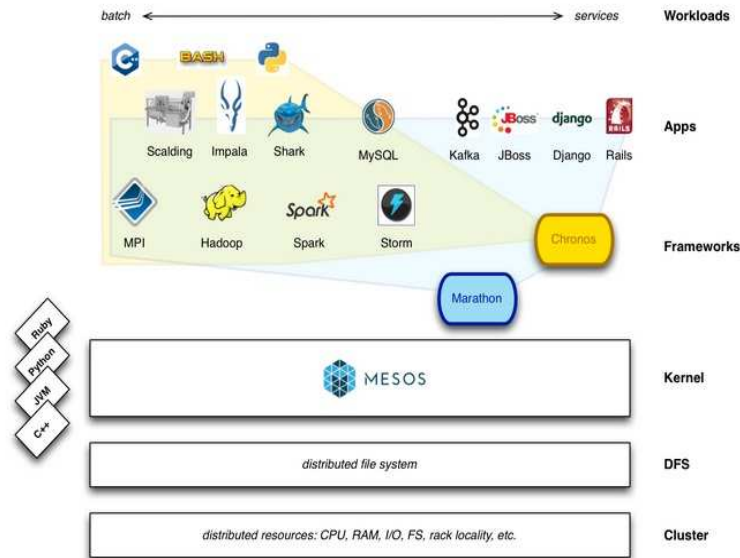
- Cluster Information:**
 - Cluster: (unnamed)
 - Server: 10.10.78.50:0050
 - Version: 0.26.0
 - Built: 4 days ago by master
 - Started: 1.5 minutes ago
 - Elected: 1.4 minutes ago
- Tasks:**
 - Staged: 0
 - Started: 0
 - Finished: 0
 - Killed: 0
 - Failed: 0
 - Lost: 0
- Resources:**

| | CPU% | Mem |
|---------|------|--------|
| Total | 6 | 2.9 GB |
| Used | 0 | 0 B |
| Offered | 0 | 0 B |
| Idle | 6 | 2.9 GB |
- Active Tasks:**

| ID | Name | State | Started | Host |
|------------------|------|-------|---------|------|
| No active tasks. | | | | |
- Completed Tasks:**

| ID | Name | State | Started | Stopped | Host |
|---------------------|------|-------|---------|---------|------|
| No completed tasks. | | | | | |

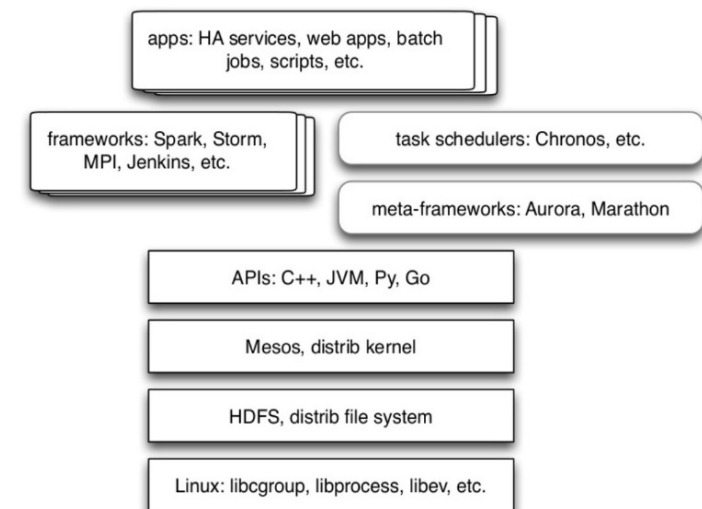
MESOS. Arquitectura



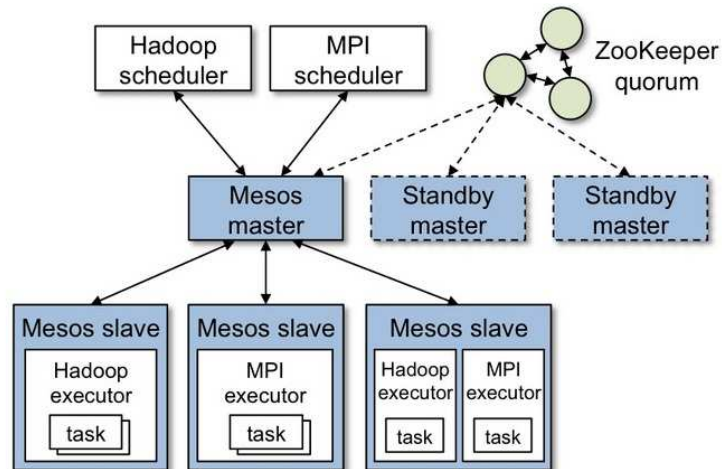
1. **MESOS** ofereix recursos als entorns d'aplicació (*frameworks*).

2. Els entorns d'aplicació decideixen si accepten els recursos.

MESOS resol els detalls tècnics de baix nivell de la interfície, i s'espera que llibreries d'alt nivell resolguin la resta de problemes.

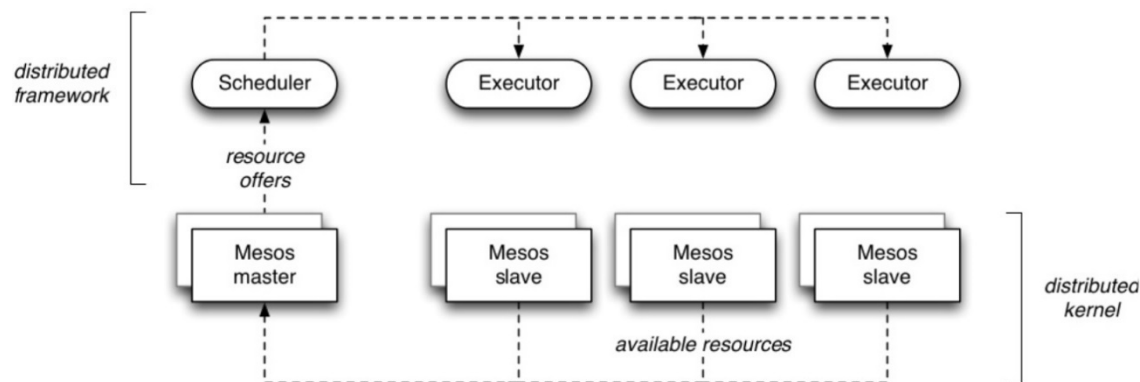


MESOS. Components i funcionament



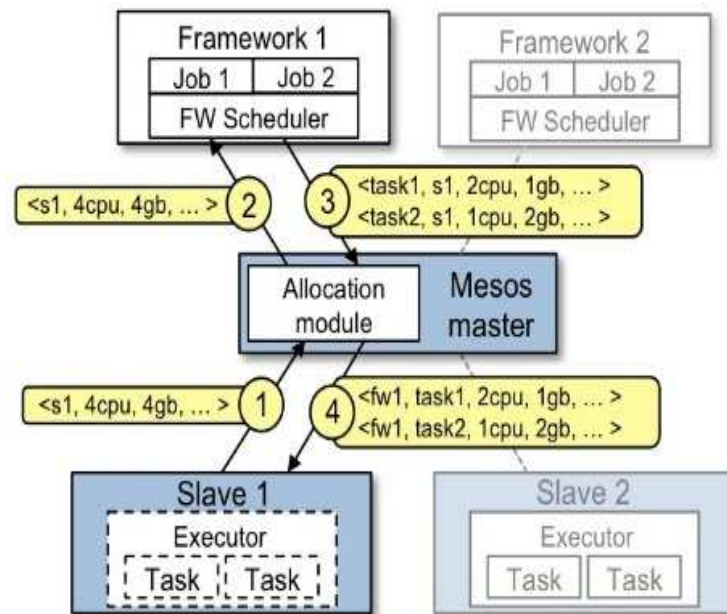
Els components bàsics de MESOS són el node màster i els nodes esclaus.

La resiliència existeix a nivell de node màster i autorecuperació dels nodes esclaus.



Gestió: el node màster ofereix els recursos lliures al planificador de cada entorn d'aplicació.

MESOS. Components i funcionament



Oferiment i assignació dinàmica de recursos segons les necessitats dels entorns d'aplicació:

1. Els esclaus indiquen els recursos lliures.
2. El màster informa als entorns d'aplicació dels recursos lliures
3. Els planificadors dels entorns d'aplicació accepten o rebutgen els recursos oferts.

MESOS. Entorns d'aplicació disponibles



Processar Big Data

(Big Data Processing)

Cray Chapel

Dpark

Exelixa

Hadoop

Hama

MPI

Spark

Storm

Serveis de llarga durada

(Long running)

Aurora

Marathon

Singularity

SSSP

Planificació per lots

(Batch scheduling)

Chronos

Jenkins

JobServer

Torque

Emmagatzematge de dades

(Data storage)

Cassandra

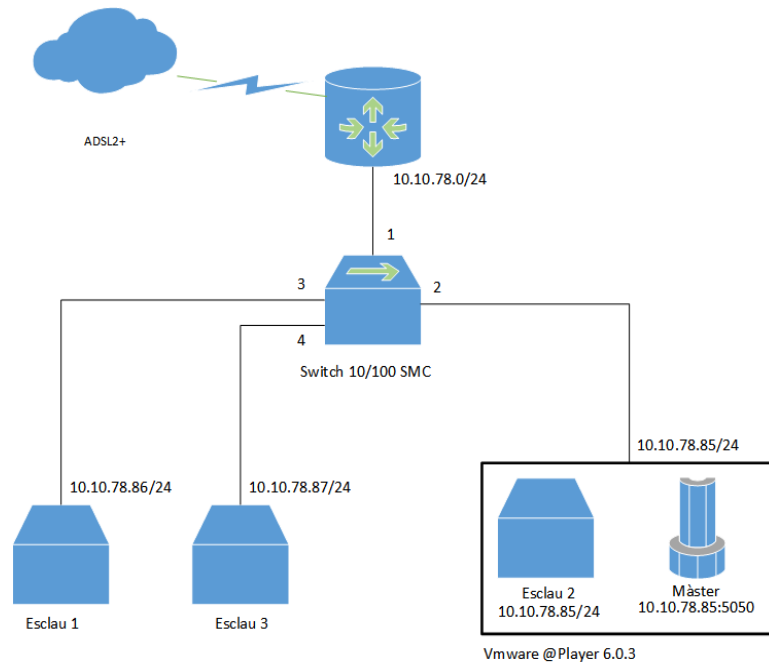
ElasticSearch

Hypertable

MESOS. Utilització actual?

| | | | |
|----------------|---------------|-----------------|-------------|
| Airbnb | GoCardless | OpenTable | UC Berkeley |
| Artirix | Groupon | Oscar Insurance | URX |
| Atigeo | HubSpot | PayPal | Viadeo |
| Atlassian | Ignidata | Pinkbike | Vimeo |
| Branding Brand | iQIYI | ProfitStars | Virdata |
| Categorize | LIFX | Qubit | Wizcorp |
| CloudPhysics | Localsensor | Revisely | WooRank |
| Conviva | Magine TV | Sailthru | Yieldbot |
| CorvisaCloud | Medidata | Sharethrough | Xogito |
| Coursera | Solutions | Sigmoid | |
| CRP-Gabriel | meemo | Analytics | |
| Lippmann | MediaCrossing | SiQueries | |
| Daemon | Mesosphere | Squarespace | |
| Devicescape | Netflix | The Factory | |
| DueDil | OakmoreLabs | Twitter | |
| eBay | OpenCredo | UCSF | |

Clúster físic privat



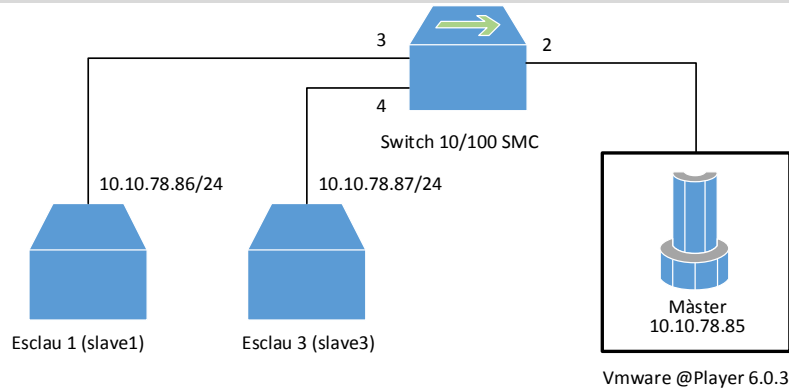
Maquinari utilitzat

- Un PC Intel (R) core (TM)2 Quad CPU Q6600@2.4 GHz, 3.25 GB DRAM.
- Dos PC AMD 64 3200+ 2 GHz, 2 GB DRAM.
- Un commutador (*switch*) de 8 ports a 10/100 Mbps marca SMC, model SMCFS8
- Un enrutador (*router*) ADSL2+ amb xarxa DHCP 10/100 Mbps

Programari utilitzat

- Sistema operatiu Ubuntu server versió 14.04
- MPICH2 de l'Argonne National Laboratory
- NPB 3.3, mètrica IS, de NAS Parallel Benchmark
- PERFTEST-1.5
- Gnuplot
- VMWare© Player, versió 6.0.3
- MESOS 0.20.0
- Java 1.7

Mètriques. Configuració dels clústers



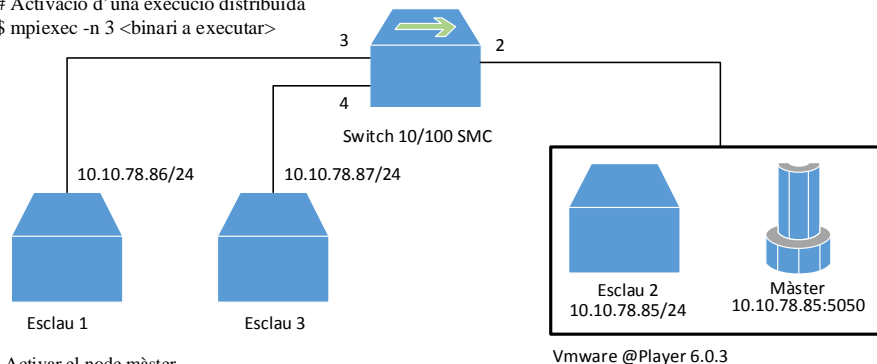
Clúster **MPD**. Permet executar aplicacions MPI de forma directa.



```
# Activació d'anell de dimonis MPD
$ mpdboot -n 3 -f mpd.hosts
```

```
# Contingut del fitxer mpd.hosts
slave1
slave3
```

```
# Activació d'una execució distribuïda
$ mpiexec -n 3 <binari a executar>
```



Clúster **MESOS** amb entorn d'aplicació MPI.



MESOS

```
# Activar el node màster
$ ./bin/mesos-master.sh --ip= 10.10.78.85 --work_dir=/var/lib/mesos
```

```
# Activar node esclau.
$ ./bin/mesos-slave.sh --ip= <ip node esclau> --master=10.10.78.85:5050 --resources: "mem(*):1512;cpus(*):1" --no-switch_user
```

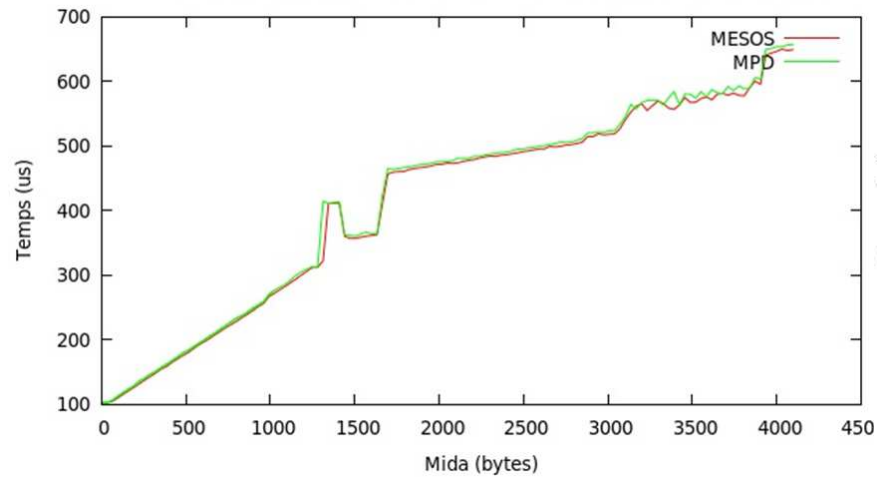
```
# Activar treball MPI a node màster
$ mpiexec-mesos -n 3 -m 1512 10.10.78.85:5050 ./<binari a executar>
```

```
# WEBUI
10.10.78.85:5050
```

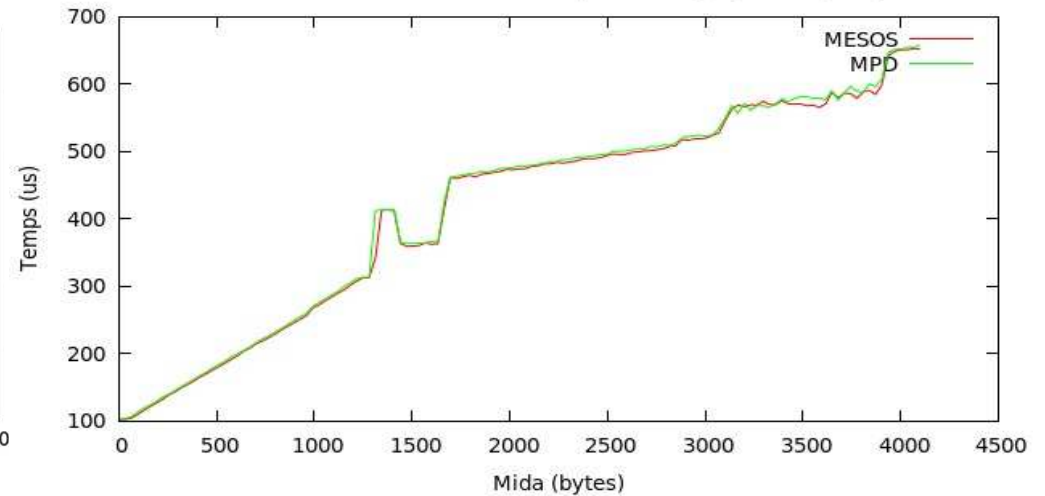
Mètriques. Resultats MPI PERFTEST-1.5



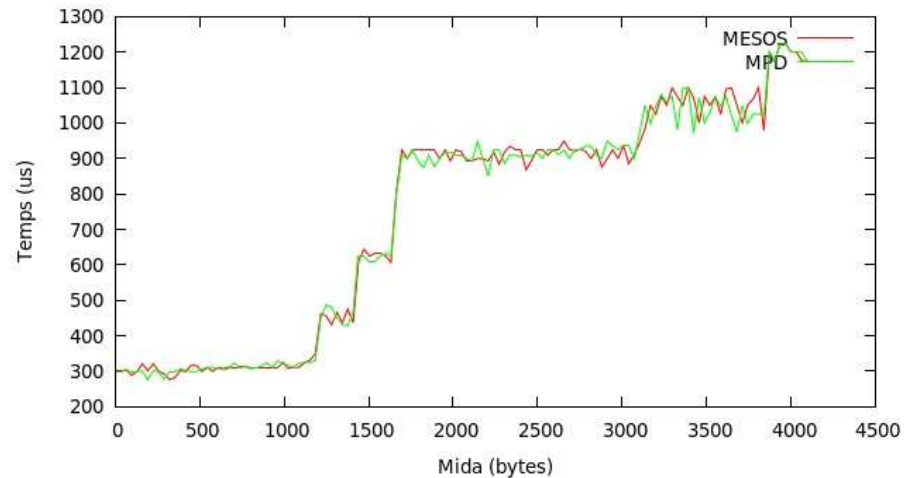
Cluster. Rendiment de Comunicacio de tipus bloquejant (mpptest) per a MPI



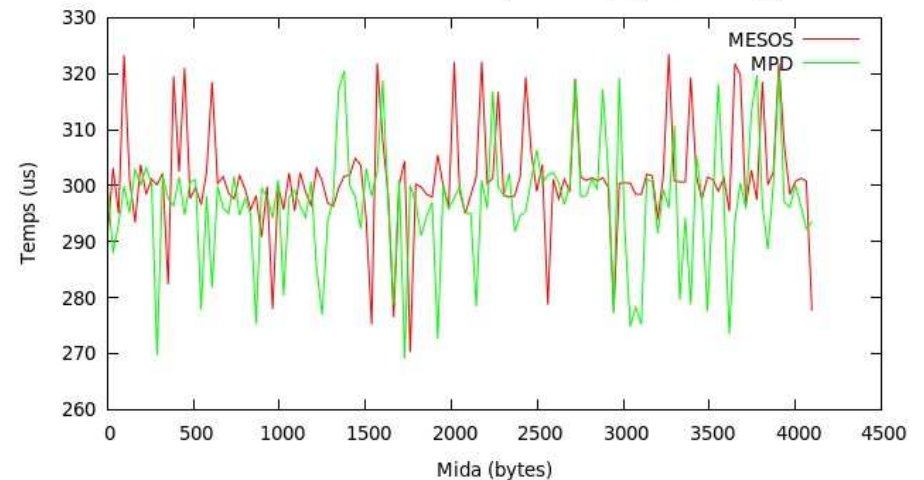
Cluster. Rendiment de Comunicacio de tipus NO bloquejant (-async) per a MPI



Cluster. Rendiment de Comunicacio de tipus bloquejant (-bisect) per a MPI



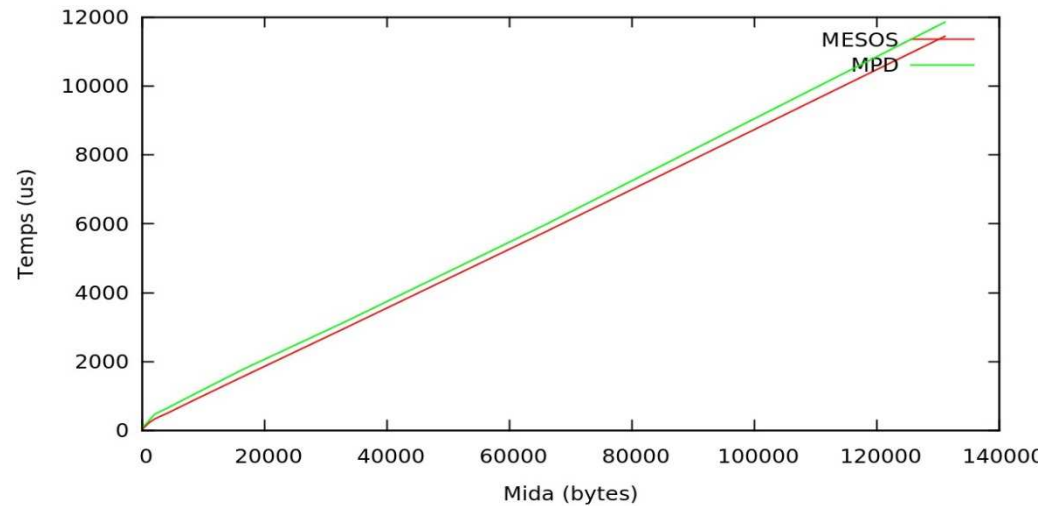
Cluster. Rendiment de Comunicacio de tipus NO bloquejant (-overlap) per a MPI



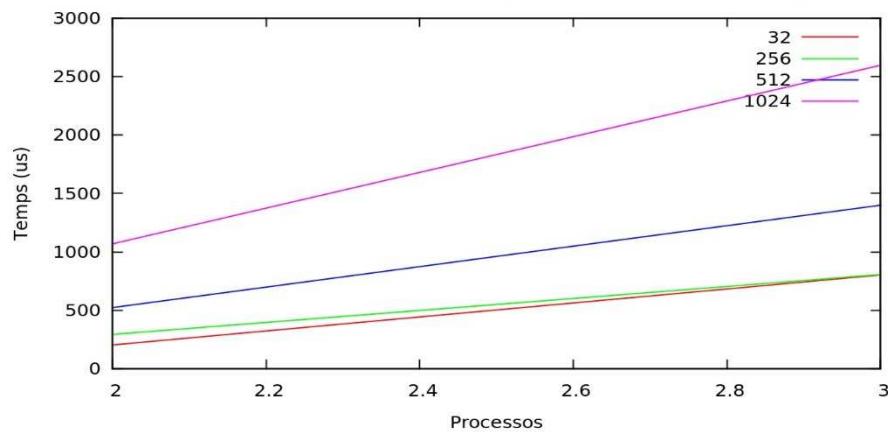
Mètriques. Resultats MPI PERFTEST-1.5



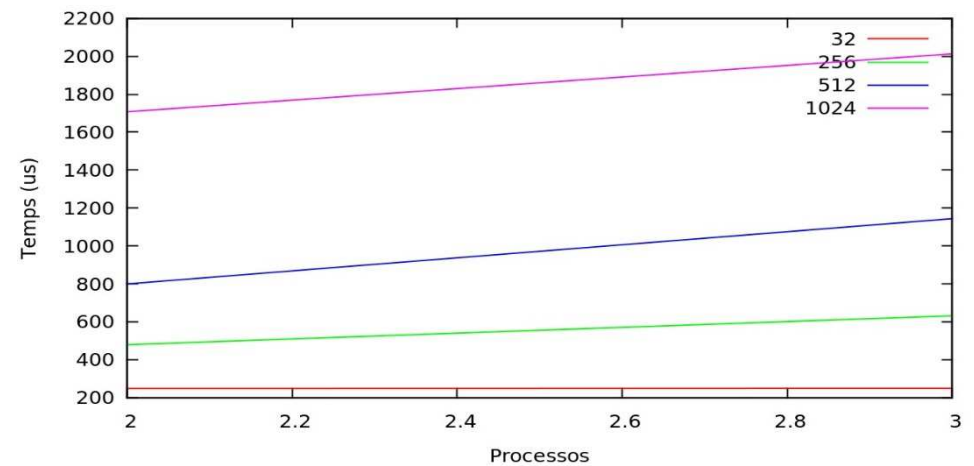
Cluster. Rendiment de Comunicacio de tipus bloquejant (-logscale) per a MPI



Cluster. Rendiment de Comunicacio de tipus isum (2 i 3 nodes, MESOS) per a MPI



Cluster. Rendiment de Comunicacio de tipus isum (2 i 3 nodes MPD) per a MPI



Mètriques. Consum de recursos als clústers

```

slave3@slave3: ~
top - 22:03:09 up 2 min, 1 user, load average: 0.52, 0.21, 0.08
Tasks: 95 total, 3 running, 92 sleeping, 0 stopped, 0 zombie
%Cpu(s): 81.4 us, 17.3 sy, 0.0 ni, 0.0 id, 0.0 wa, 1.3 hi, 0.0 si, 0.0 st
KiB Mem: 1984184 total, 143564 used, 1840620 free, 20492 buffers
KiB Swap: 2031612 total, 0 used, 2031612 free. 54928 cached Mem

  PID USER      PR  NI  VIRT  RES  SHR S %CPU %MEM    TIME+  COMMAND
 1175 master    20   0 17784   3180  816 R 99.4  0.2    0:34.13 mptest
     1 root      20   0 33644  2896 1404 S  0.0  0.1    0:03.21 init
     2 root      20   0   0     0   0 S  0.0  0.0    0:00.00 kthreadd
     3 root      20   0   0     0   0 S  0.0  0.0    0:00.00 ksoftirqd/0
     4 root      20   0   0     0   0 S  0.0  0.0    0:00.00 kworker/0:0
     5 root      0 -20   0     0   0 S  0.0  0.0    0:00.00 kworker/0:0H
     6 root      20   0   0     0   0 S  0.0  0.0    0:00.11 kworker/u4:0
     7 root      20   0   0     0   0 S  0.0  0.0    0:00.63 rcu_sched
     8 root      20   0   0     0   0 R  0.0  0.0    0:00.06 rcuos/0
     9 root      20   0   0     0   0 S  0.0  0.0    0:00.00 rcuos/1
    10 root      20   0   0     0   0 S  0.0  0.0    0:00.00 rcu_bh
    11 root      20   0   0     0   0 S  0.0  0.0    0:00.00 rcuob/0
    12 root      20   0   0     0   0 S  0.0  0.0    0:00.00 rcuob/1
    13 root      rt   0   0     0   0 S  0.0  0.0    0:00.00 migration/0
    14 root      rt   0   0     0   0 S  0.0  0.0    0:00.00 watchdog/0
    15 root      0 -20   0     0   0 S  0.0  0.0    0:00.00 khelper
    16 root      20   0   0     0   0 S  0.0  0.0    0:00.00 kdevtmpfs

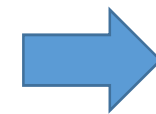
```

```

slave3@slave3: ~
top - 22:34:49 up 34 min, 1 user, load average: 1.03, 0.90, 0.68
Tasks: 96 total, 3 running, 93 sleeping, 0 stopped, 0 zombie
%Cpu(s): 82.7 us, 17.3 sy, 0.0 ni, 0.0 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0 st
KiB Mem: 1984184 total, 173840 used, 1810344 free, 21248 buffers
KiB Swap: 2031612 total, 0 used, 2031612 free. 76384 cached Mem

  PID USER      PR  NI  VIRT  RES  SHR S %CPU %MEM    TIME+  COMMAND
 1535 slave3    20   0 17784   5172  816 R 98.0  0.3    0:46.15 mptest
 1271 slave3    20   0 650760 13348 10744 S  2.0  0.7    0:10.40 lt-mesos-sl+
     1 root      20   0 33644  2896 1404 S  0.0  0.1    0:03.21 init
     2 root      20   0   0     0   0 S  0.0  0.0    0:00.00 kthreadd
     3 root      20   0   0     0   0 S  0.0  0.0    0:00.04 ksoftirqd/0
     5 root      0 -20   0     0   0 S  0.0  0.0    0:00.00 kworker/0:0H
     7 root      20   0   0     0   0 S  0.0  0.0    0:00.69 rcu_sched
     8 root      20   0   0     0   0 R  0.0  0.0    0:00.15 rcuos/0
     9 root      20   0   0     0   0 S  0.0  0.0    0:00.00 rcuos/1
    10 root      20   0   0     0   0 S  0.0  0.0    0:00.00 rcu_bh
    11 root      20   0   0     0   0 S  0.0  0.0    0:00.00 rcuob/0
    12 root      20   0   0     0   0 S  0.0  0.0    0:00.00 rcuob/1
    13 root      rt   0   0     0   0 S  0.0  0.0    0:00.00 migration/0
    14 root      rt   0   0     0   0 S  0.0  0.0    0:00.01 watchdog/0
    15 root      0 -20   0     0   0 S  0.0  0.0    0:00.00 khelper
    16 root      20   0   0     0   0 S  0.0  0.0    0:00.00 kdevtmpfs
    17 root      0 -20   0     0   0 S  0.0  0.0    0:00.00 netns

```



Consum de recursos en esclau clúster **MPD**.



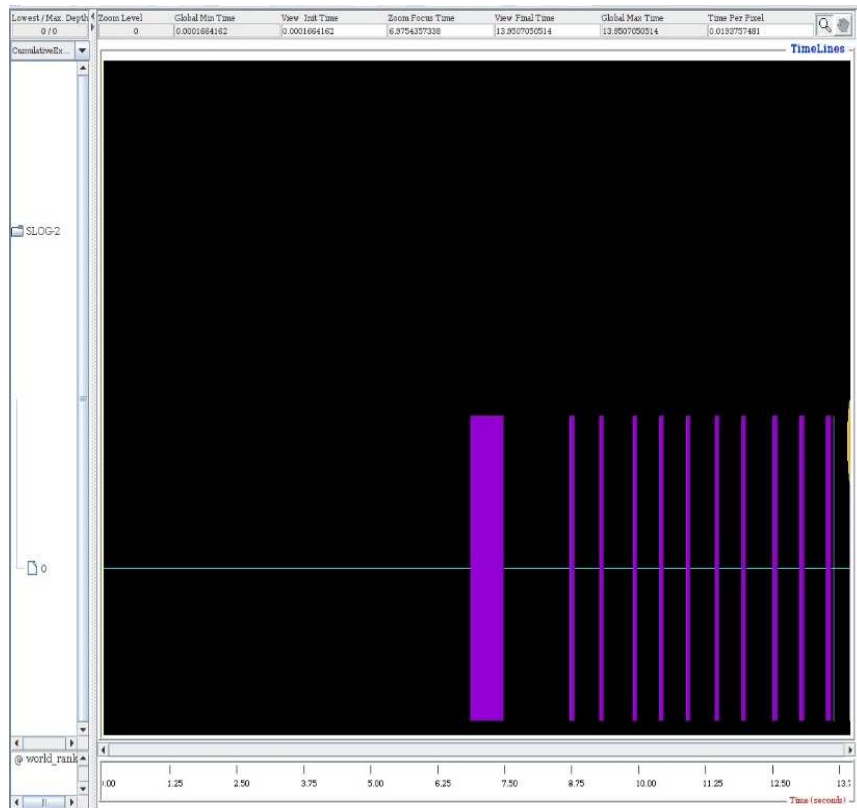
Consum de recursos en esclau clúster **MESOS**.



Font: programa TOP. Ubuntu 14.04

Mètriques. Execució d'aplicació MPI NPB-IS

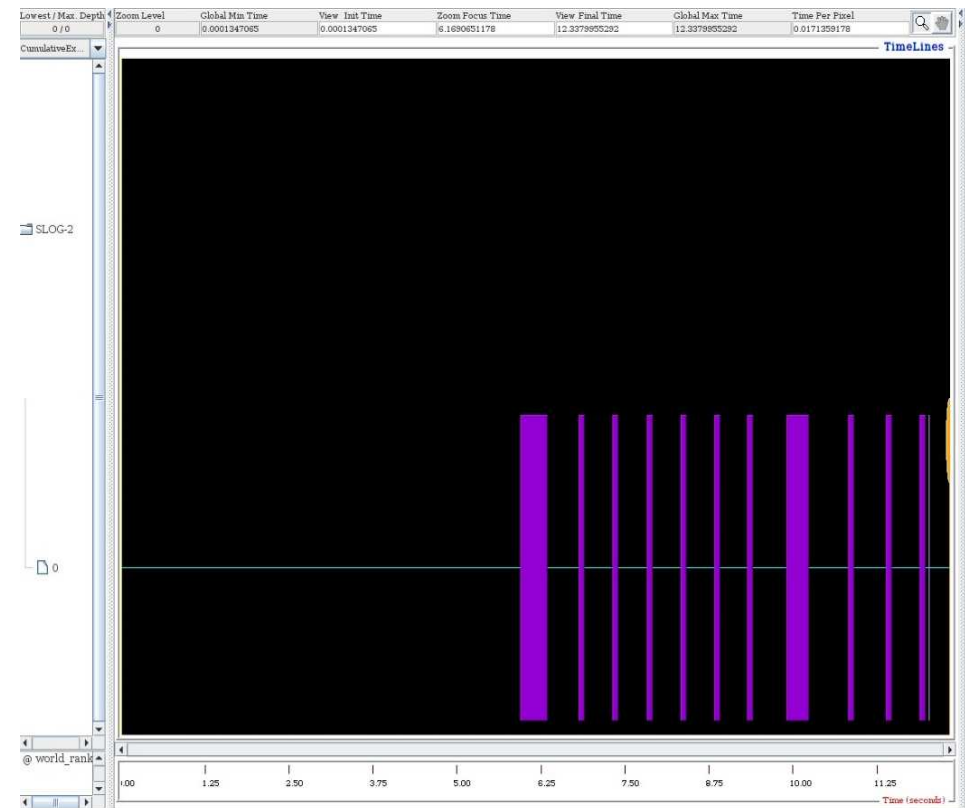
Aplicació MPI executada: is.B.1



Temps d'execució clúster **MESOS: 13,95 S**



MESOS



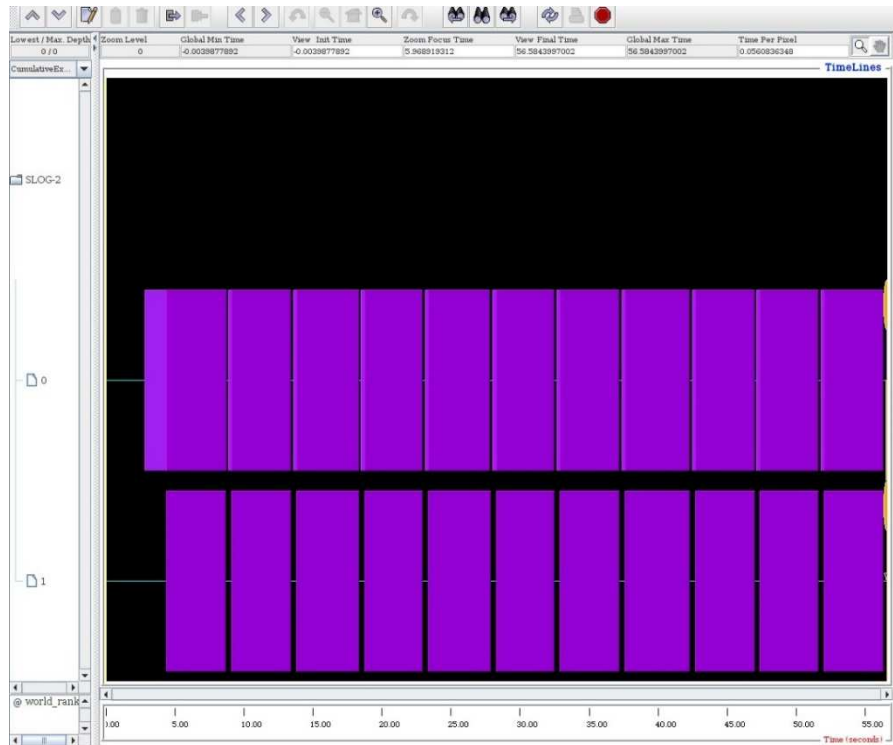
Temps d'execució clúster **MPD: 12,33 S**



Font: programa JUMPSHOT (MPICH2).

Mètriques. Execució d'aplicació MPI NPB-IS

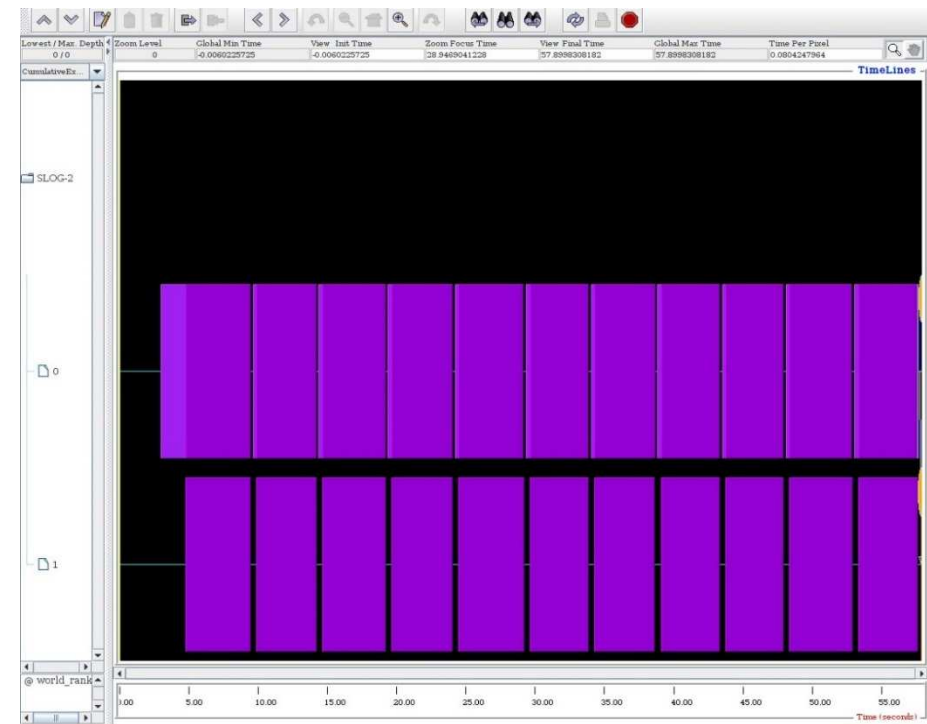
Aplicació MPI executada: **is.B.2**



Temps d'execució clúster **MESOS: 56,58 S**



MESOS



Temps d'execució clúster **MPD: 57,89 S**



Font: programa JUMPSHOT (MPICH2).

Conclusions de les mètriques



Rendiment del clúster: proves de **traspàs de missatges MPI** amb la suite *PERFTEST-1.5* de l'*Argonne National Laboratory*.

Rendiment en execució d'aplicació MPI: prova amb 1 i 2 processadors de l'aplicació **IS.B** del *NAS Parallel Benchmark*.

Consum de recursos en execució de rendiment del clúster: prova **MPPTTEST**, suite *PERFTEST-1.5* de l'*Argonne National Laboratory*.

L'execució sota la gestió de MESOS **no penalitza ni afavoreix**, respecte a no fer servir MESOS.

Conclusions de les mètriques



MESOS



Després d'observar els resultats de les mètriques, és útil **MESOS**?

- Sí, reutilitza recursos lliures de forma dinàmica per guanyar eficiència.
- Sí, podem executar en paral·lel diverses aplicacions, i que poden ser diferents.
- Sí, aporta mecanismes d'alta disponibilitat davant d'errors.
- Sí, aporta mecanismes de planificació dels treballs de forma nativa.
- Sí, podem controlar amb *WebUI* que passa al clúster en temps real.

Conclusions del treball

- Millor coneixement de les problemàtiques de gestió de clústers de computació i les seves solucions.
- Respecte al coneixement adquirit del gestor MESOS:
 - Apreciació del model tecnològic de dos nivells en que està basat.
 - Gestor estable i robust, i no necessita grans recursos operatius.
 - Observada la gestió de recursos de forma dinàmica per oferiment.
 - Permet conviure maquinaris no homogenis en cas necessari.
 - El desplegament és senzill i ràpid.
 - Conèixer el significat d'entorn d'aplicació (*framework*)
 - Estudi de mètriques del clúster.
 - Observació en temps real del funcionament del clúster.
- Respecte al coneixement general adquirit:
 - Observació de treball de computació paral·lela distribuïda.
 - Configuració de sistemes operatius, xarxes i aplicacions.

Conclusions del treball. Estudi de futur

Desplegar **MESOS** en un clúster de desenes de servidors, per observar.

- **Apache Zookeeper** i les opcions de clúster d'alta disponibilitat.
- Eines de control i programació del clúster (**Chronos**).
- El comportament de la xarxa analitzant les comunicacions entre nodes.
- El funcionament dels scripts d'iniciació automàtica dels nodes del clúster.
- Altres entorns d'aplicació diferents d'MPI (**Spark, Hadoop, ...**)

Altres temàtiques interessants:

- Estudi de realització d'un entorn d'aplicació concret.
- En un futur pròxim: **MESOSPHERE DCOS**.

Referències de les imatges utilitzades

- [0][2][7][14][19][20][21][22][23] Font: <http://mesos.apache.org/>
- [1] Font: http://upload.wikimedia.org/wikipedia/commons/2/29/Virginia_tech_xserve_cluster.jpg
- [5] Font: <http://static.googleusercontent.com/media/research.google.com/es//pubs/archive/41684.pdf>
- [6] Font: <http://www.zdnet.com/article/what-is-docker-and-why-is-it-so-darn-popular/>
- [6] Font: <https://docs.docker.com/introduction/understanding-docker/>
- [7] Font: <https://amplab.cs.berkeley.edu/projects/mesos-dynamic-resource-sharing-for-clusters/>
- [7][8][9][10][11] Font: <http://www.slideshare.net/pacoid/datacenter-computing-with-apache-mesos>
- [11][12] Font: http://people.csail.mit.edu/matei/papers/2011/nsdi_mesos.pdf
- [13] Font: mesosphere.com/docs/frameworks
- [16][19][20][21][22][23] Font: <http://www.anl.gov/> <http://www.nas.nasa.gov/publications/npb.html>
- [17][18] Captures programa **GNUPLLOT** sobre mètriques **PERFTEST-1.5**
- [19] Captura d'execució programa **TOP**, Ubuntu Server 14.04.
- [20][21] Captures de l'aplicació **JUMPSHOT** sobre execucions MPI amb llibreria MPE.