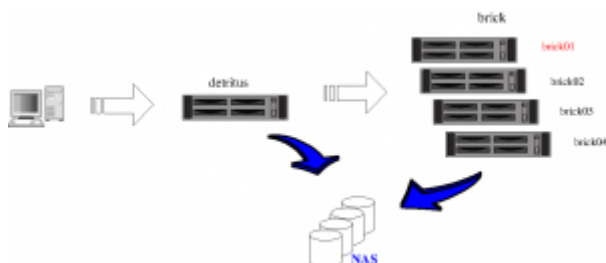


# Como usar el cluster sin morir en el intento

## Get it started

El cluster (*brick*) consta de cuatro nodos (*brick00*, *brick01*, *brick02* y *brick03*). Lo que sigue describe la operativa básica para ejecutar tareas en estas maquinas de una manera ordenada



El directorio */nas* de *detritus* se monta como */home* en cada uno de los *bricks* por lo que todo lo que se modifique en *detritus:/nas/user* cambiara automaticamente en *brick0X:/home/user*

La presentacion inicial (hecha con pinpoint y exportada a pdf) se puede bajar de aqui:  
<http://detritus.fundacioace.com/files/cluster.pdf>

## como hacer que los nodos no pidan password

### We don't talk anymore

#### Esto es básico para el funcionamiento correcto del resto de las herramientas

Empezamos creando una clave RSA en nuestro directorio de *detritus* y moviendo la clave publica a un archivo que podamos identificar facilmente

```
$ ssh-keygen
$ mv $HOME/.ssh/id_rsa.pub $HOME/.ssh/detritus.pub
```

Ahora añadimos esta clave a la lista de claves autorizadas en el HOME de los nodos. Como hasta ahora no hay ninguna se puede hacer con *cp*. En caso de haber alguna se podría hacer con *cat*. Y no podemos olvidarnos de dar los permisos correctos!

En mi caso:

```
$ mkdir /nas/osotolongo/.ssh
$ chmod 700 /nas/osotolongo/.ssh
$ cp ~/.ssh/detritus.pub /nas/osotolongo/.ssh/authorized_keys
$ chmod 600 /nas/osotolongo/.ssh/authorized_keys
```

A partir de ahora, cuando se entra de *detritus* a cualquiera de los nodos no deberia ser necesario introducir ningun password.

# como ejecutar la misma orden en varios nodos del cluster (pssh y mpssh)

## Nobody Wants To Be Lonely

### pssh

- primero crear un archivo con la lista de hosts:

```
[osotolongo@detritus ~]$ cat host.pssh
osotolongo@brick01
osotolongo@brick02
osotolongo@brick03
```

- despues ejecutar un comando en todas las maquinas

```
[osotolongo@detritus ~]$ pssh -h host.pssh 'ps ax | grep ssh'
[1] 13:11:46 [SUCCESS] osotolongo@brick03
[2] 13:11:46 [SUCCESS] osotolongo@brick02
[3] 13:11:46 [SUCCESS] osotolongo@brick01
```

### mpssh

- la syntaxis cambia ligeramente, y la salida tambien

```
[osotolongo@detritus ~]$ mpssh -f host.pssh 'ps ax | grep ssh'
MPSSH - Mass Parallel Ssh Ver.1.3.3
(c)2005-2013 Nikolay Denev <ndenev@gmail.com>

[*] read (3) hosts from the list
[*] executing "ps ax | grep ssh" as user "osotolongo"
[*] spawning 3 parallel ssh sessions

brick01 -> 1920 ?      Ss      0:00 /usr/sbin/sshd -D
brick01 -> 73087 ?     Ss      0:00 sshd: osotolongo [priv]
brick01 -> 73090 ?      S       0:00 sshd: osotolongo@notty
brick01 -> 73091 ?     Ss      0:00 bash -c ps ax | grep ssh
brick01 -> 73101 ?      S       0:00 grep ssh
brick02 -> 1887 ?      Ss      0:00 /usr/sbin/sshd -D
brick02 -> 21242 ?     Ss      0:00 sshd: osotolongo [priv]
brick02 -> 21245 ?      S       0:00 sshd: osotolongo@notty
brick02 -> 21246 ?     Ss      0:00 bash -c ps ax | grep ssh
brick02 -> 21256 ?      S       0:00 grep ssh
brick03 -> 1928 ?      Ss      0:00 /usr/sbin/sshd -D
brick03 -> 21270 ?     Ss      0:00 sshd: osotolongo [priv]
brick03 -> 21273 ?      S       0:00 sshd: osotolongo@notty
```

```
brick03 -> 21274 ?      Ss      0:00 bash -c ps ax | grep ssh
brick03 -> 21284 ?      S       0:00 grep ssh
```

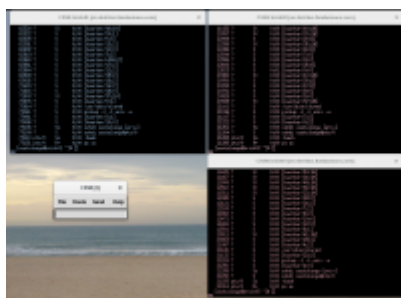
Done. 3 hosts processed.

## o como hacerlo mas interactivo (cluster-ssh)

### Dance again

Una imagen vale mas que mil palabras

```
$ cssh brick
```



Esta herramienta permite, cuando se escribe en la ventana comun, escribir lo mismo en todas las terminales abiertas o en una sola si se selecciona expresamente.

## Slurm

### This is what you came for

#### **srun:**

Este comando es capaz de lanzar un proceso en varios nodos del cluster al mismo tiempo

```
[osotolongo@detritus ~]$ srun -N3 -l /bin/hostname
0: brick01
2: brick03
1: brick02
```

#### **sbatch:**

Este comando lanza un script (tipicamente una lista de comandos *srun*) en el cluster. Es posible escoger el numero de procesos a lanzar e incluso los nodos en los cuales hacerlo

```
[osotolongo@detritus ~]$ cat test_script.sh
#!/bin/sh
#SBATCH --time=1
/bin/hostname
```

```
srun -l /bin/hostname
srun -l /bin/pwd
[osotolongo@detritus ~]$ sbatch -n6 -w "brick0[1-3]" -o test_stdout
test_script.sh
Submitted batch job 30
[osotolongo@detritus ~]$ cat /nas/osotolongo/test_stdout
brick01
1: brick01
2: brick02
4: brick03
3: brick02
5: brick03
0: brick01
0: /home/osotolongo
5: /home/osotolongo
3: /home/osotolongo
1: /home/osotolongo
4: /home/osotolongo
2: /home/osotolongo
```

### squeue:

Este comando permite ver los trabajos que se han lanzado e informacion sobre ellos

```
[osotolongo@detritus ~]$ cat test_script.sh
#!/bin/sh
#SBATCH --time=1
/bin/hostname
srun -l /bin/hostname
sleep 20
srun -l /bin/pwd
[osotolongo@detritus ~]$ squeue
      JOBID PARTITION      NAME      USER ST      TIME  NODES
NODELIST(REASON)
[osotolongo@detritus ~]$ sbatch -n6 -w "brick0[1-3]" -o test_stdout
test_script.sh
Submitted batch job 35
[osotolongo@detritus ~]$ squeue
      JOBID PARTITION      NAME      USER ST      TIME  NODES
NODELIST(REASON)
      35      debug test_scr osotolon R      0:01      3
brick[01-03]
```

### scancel:

OK, y cuando la caguemos podemos usar *scancel* para parar los trabajos. El kill de toda la vida pero a traves de un workload manager

```
[osotolongo@detritus ~]$ cat test_script.sh
#!/bin/sh
/bin/hostname
```

```

srun -l /bin/hostname
sleep 200
srun -l /bin/pwd
[osotolongo@detritus ~]$ sbatch -n6 -w "brick0[1-3]" -o test_stdout
test_script.sh
Submitted batch job 36
[osotolongo@detritus ~]$ squeue
          JOBID PARTITION      NAME      USER ST       TIME  NODES
NODELIST(REASON)
          36      debug test_scr osotolon  R       0:03     3
brick[01-03]
[osotolongo@detritus ~]$ scancel 36
[osotolongo@detritus ~]$ squeue
          JOBID PARTITION      NAME      USER ST       TIME  NODES
NODELIST(REASON)

```

<https://slurm.schedmd.com/quickstart.html>

mas: <https://wiki.fysik.dtu.dk/niflheim/SLURM#configure-slurm>

**Nota:** Todo lo que estaba hecho en Perl, usando HPC, ha dejado de funcionar. Los modulos dan error. NO obstante, los docs estan aqui: [HPC::Runner::Slurm \(AKA the hat trick\)](https://wiki.fysik.dtu.dk/niflheim/SLURM#configure-slurm)

## The real life (--multi-prog)

### Fighter

Vamos a intentar usar todo esto para algo util. Digamos que queremos ejecutar un programa con diferentes argumentos en los nodos del cluster. Para ello usamos **sbatch**

```
$ sbatch test.sh
```

El script *test.sh* no es un programa real sino que es una orden de ejecucion de las tareas.

### test.sh

```

#!/bin/bash
#SBATCH --time=100
#SBATCH --ntasks-per-node=20
#SBATCH -n 50
#SBATCH --mail-type=ALL
#SBATCH --mail-user=osotolongo
srun --exclusive --multi-prog test.conf

```

Este script lo unico que hace es configurar minimamente **sbatch** y lanzar el comando **srun**.

- Las directivas importantes aqui son `-ntasks-per-node=20` y `-n 50`. La primera dice a **sbatch** que lance 20 tareas en cada nodo, la segunda que hay un total de 50 tareas. para que todo vaya bien  $n/ntask-per-node < 3$ . Por ejemplo si ponemos `-ntasks-per-node=10`, entonces dara error porque  $50/10=5$ . *Esto se puede quitar pero si lanzamos 50 tareas las metera todas en brick01 y el ejemplo no servira de nada.*
- Las directivas `-mail-type=ALL` y `-mail-user=osotolongo` dicen que envíe por email todo lo que ocurra al usuario `osotolongo`. Si el usuario de cada uno esta bien configurado, esto deberia funcionar para cada uno con solo cambiar el nombre.
- El comando **srun** lee el archivo `test.conf` y ejecuta las ordenes que hay dentro de este archivo.

El archivo `test.conf` tiene la siguiente estructura

```
0      test_runner.sh test/blah1.txt
1      test_runner.sh test/blah2.txt
2      test_runner.sh test/blah3.txt
3      test_runner.sh test/blah4.txt
4      test_runner.sh test/blah5.txt
5      test_runner.sh test/blah6.txt
.
.
.
45     test_runner.sh test/blah46.txt
46     test_runner.sh test/blah47.txt
47     test_runner.sh test/blah48.txt
48     test_runner.sh test/blah49.txt
49     test_runner.sh test/blah50.txt
```

La primer columna debe numerar las tareas, comenzando en **0** y terminando en **n-1**. La segunda columna no es mas que la orden a ejecutar (en el ejemplo `test_runner.sh`), con todos los argumentos correspondientes (en este caso un nombre de archivo).

## Probando

```
[osotolongo@detritus cluster]$ sbatch test.sh
Submitted batch job 842
[osotolongo@detritus cluster]$ squeue
          JOBID PARTITION     NAME     USER ST       TIME  NODES
NODELIST(REASON)
          842      debug  test.sh osotolon  R           0:03     3
brick[01-03]
```

El programa a ejecutar es un script sencillo que recibe como argumento un nombre de archivo,

```
[osotolongo@detritus cluster]$ cat test_runner.sh
#!/bin/sh
file=$1
count=0
while [ $count -le 10 ]
do
```

```
    echo $count >> $file
    ((count++))
    sleep 50
done
echo Nice! >> $file
```

y tambien me he escrito un script para que me escriba el *test.conf*,

```
[osotolongo@detritus cluster]$ cat test_generator.sh
#!/bin/sh
count=0
while [ $count -lt 50 ]
do
    cn=$((count+1))
    echo "$count    test_runner.sh test/blah${cn}.txt" >> test.conf
    ((count++))
done
```

Luego, corro primero,

```
$ ./test_generator.sh
```

que me hace el *test.conf* y luego,

```
$ sbatch test.sh
```

Los archivos se generan correctamente,

```
[osotolongo@detritus cluster]$ ls test
blah10.txt  blah15.txt  blah1.txt   blah24.txt  blah29.txt  blah33.txt
blah38.txt  blah42.txt  blah47.txt  blah5.txt   blah2.txt   blah34.txt
blah11.txt  blah16.txt  blah20.txt  blah25.txt  blah6.txt   blah35.txt
blah39.txt  blah43.txt  blah48.txt  blah26.txt  blah30.txt  blah36.txt
blah12.txt  blah17.txt  blah21.txt  blah7.txt   blah31.txt  blah37.txt
blah3.txt   blah44.txt  blah49.txt  blah8.txt   blah32.txt  blah38.txt
blah13.txt  blah18.txt  blah22.txt  blah27.txt  blah4.txt   blah39.txt
blah40.txt  blah45.txt  blah4.txt   blah28.txt  blah9.txt   blah40.txt
blah14.txt  blah19.txt  blah23.txt  blah9.txt   blah41.txt  blah41.txt
blah41.txt  blah46.txt  blah50.txt  blah9.txt
```

y el contenido es el que debe ser,

```
[osotolongo@detritus cluster]$ cat test/blah1.txt
0
1
2
3
4
5
6
```

7  
8  
9  
10  
Nice!

Los emails también se envían correctamente,

---

<input type="checkbox"/> ☆ > Trolls United Co.	SLURM Job_id=842 Name=test.sh Ended, Run time 00:09:11, COMPLETED, ExitCode 0
<input type="checkbox"/> ☆ > Trolls United Co.	SLURM Job_id=842 Name=test.sh Began, Queued time 00:00:00

From:  
<http://detritus.fundacioace.com/wiki/> - **Detritus Wiki**

Permanent link:  
<http://detritus.fundacioace.com/wiki/doku.php?id=cluster&rev=1537439582>

Last update: **2020/08/04 10:47**

