

Introducción a Clúster, Cloud y DevOps

Remo Suppi Boldrito

PID_00238614

Índice

Introducción	5
Objetivos	8
1. Clusterización	9
1.1. Virtualización	9
1.1.1. Plataformas de virtualización	10
1.2. Beowulf	12
1.2.1. ¿Cómo configurar los nodos?	14
1.3. Beneficios del cómputo distribuido	15
1.3.1. ¿Cómo hay que programar para aprovechar la conurrencia?	17
1.4. Memoria compartida. Modelos de hilos (<i>threading</i>)	19
1.4.1. Multihilos (<i>multithreading</i>)	20
1.5. OpenMP	23
1.6. MPI, <i>Message Passing Interface</i>	27
1.6.1. Configuración de un conjunto de máquinas para hacer un clúster adaptado a OpenMPI	29
1.7. Rocks Cluster	32
1.7.1. Guía rápida de instalación	33
1.8. Despliegue automático	35
1.8.1. FAI (Fully Automatic Installation)	35
1.9. Guía de instalación de FAI	36
1.9.1. Instalación del servidor	37
1.9.2. Configuración del nodo	38
1.10. <i>Logs</i>	39
1.11. Visualización de logs	41
2. Cloud	45
2.1. Opennebula	48
3. DevOps	55
3.1. Linux Containers, LXC	56
3.2. Docker	61
3.3. Instalación de un servidor Apache sobre un contenedor Docker	63
3.4. Puppet	64
3.4.1. Instalación	64
3.5. <i>Main Manifest File</i>	65
3.6. Chef	69
3.7. Estructura de Chef	70

3.8. Ansible	72
3.9. Vagrant	75
Actividades	78
Bibliografía	79

Introducción

Los avances en la tecnología han llevado, por un lado, al desarrollo de procesadores más rápidos, con más de un elemento de cómputo (núcleos o *cores*) en cada uno de ellos, de bajo coste y con soporte a la virtualización *hardware* y por otro, al desarrollo de redes altamente eficientes. Esto, junto con el desarrollo de sistemas operativos como GNU/Linux, ha favorecido un cambio radical en la utilización de sistemas de procesadores de múltiple-*cores* y altamente interconectados, en lugar de un único sistema de alta velocidad (como los sistemas vectoriales o los sistemas de procesamiento simétrico SMP). Además, estos sistemas han tenido una curva de despliegue muy rápida, ya que la relación precio/prestaciones ha sido, y lo es cada día más, muy favorable tanto para el responsable de los sistemas TIC de una organización como para los usuarios finales en diferentes aspectos: prestaciones, utilidad, fiabilidad, facilidad y eficiencia. Por otro lado, las crecientes necesidades de cómputo (y de almacenamiento) se han convertido en un elemento tractor de esta tecnología y su desarrollo, vinculados a la provisión dinámica de servicios, contratación de infraestructura y utilización por uso -incluso en minutos- (y no por compra o por alquiler), lo cual ha generado una evolución total e importante en la forma de diseñar, gestionar y administrar estos centros de cómputo. No menos importante han sido los desarrollos, además de los sistemas operativos, para soportar todas estas tecnologías, en lenguajes de desarrollo, API y entornos (*frameworks*) para que los desarrolladores puedan utilizar la potencialidad de la arquitectura subyacente para el desarrollo de sus aplicaciones, a los administradores y gestores para que puedan desplegar, ofrecer y gestionar servicios contratados y valorados por minutos de uso o bytes de E/S, por ejemplo, y a los usuarios finales para que puedan hacer un uso rápido y eficiente de los sistemas de cómputo sin preocuparse de nada de lo que existe por debajo, dónde está ubicado y cómo se ejecuta.

Es por ello por lo que, en el presente capítulo, se desarrollan tres aspectos básicos que son parte de la misma idea, a diferentes niveles, cuando se desea ofrecer la infraestructura de cómputo de altas prestaciones. O una plataforma como servicio o su utilización por un usuario final en el desarrollo de una aplicación, que permita utilizar todas estas infraestructuras de cómputo distribuidas y de altas prestaciones. En sistemas de cómputo de altas prestaciones (*High Performance Computing -HPC*), podemos distinguir dos grandes configuraciones:

1) Sistemas fuertemente acoplados (*tightly coupled systems*): son sistemas donde la memoria es compartida por todos los procesadores (*shared memory systems*) y la memoria de todos ellos “se ve” (por parte del programador) como una única memoria.

2) Sistemas débilmente acoplados (*loosely coupled systems*): no comparten memoria (cada procesador posee la suya) y se comunican mediante mensajes pasados a través de una red (*message passing systems*).

En el primer caso, son conocidos como *sistemas paralelos de cómputo (parallel processing system)* y en el segundo, como *sistemas distribuidos de cómputo (distributed computing systems)*.

En la actualidad la gran mayoría de los sistemas (básicamente por su relación precio-prestaciones) son del segundo tipo y se conocen como clústers donde los sistemas de cómputo están interconectados por una red de gran ancho de banda y todos ellos trabajan en estrecha colaboración, viéndose para el usuario final como un único equipo. Es importante indicar que cada uno de los sistemas que integra el clúster a su vez puede tener más de un procesador con más de un *core* en cada uno de ellos y por lo cual, el programador deberá tener en cuenta estas características cuando desarrolle sus aplicaciones (lenguaje, memoria compartida|distribuida, niveles de caché, etc) y así poder aprovechar toda la potencialidad de la arquitectura agregada. El tipo más común de clúster es el llamado Beowulf, que es un clúster implementado con múltiples sistemas de cómputo (generalmente similares pero pueden ser heterogéneos) e interconectados por una red de área local (generalmente Ethernet, pero existen redes más eficientes como Infiniband o Myrinet). Algunos autores denominan *MPP (massively parallel processing)* cuando es un clúster, pero cuentan con redes especializadas de interconexión (mientras que los clústers utilizan hardware estándar en sus redes) y están formados por un alto número de recursos de cómputo (1.000 procesadores, no más). Hoy en día la mayoría de los sistemas publicados en el TOP500 (<https://www.top500.org/statistics/list/>) son clústers y en la última lista publicada (junio 2016) ocupaba el primer puesto el ordenador Sunway TaihuLight (China) con 10,6 millones de cores, 1.300 Terabytes de memoria (1,3 Petabytes) y un consumo de 15,3 MW, y que puede ejecutar 93.014 TFlops (es decir, aproximadamente 93.000.000.000.000.000 instrucciones de punto flotante por segundo).

El otro concepto vinculado a los desarrollos tecnológicos mencionados a las nuevas formas de entender el mundo de las TIC en la actualidad es el *Cloud Computing* (o cómputo|servicios en la nube) que surge como concepto de ofrecer servicios de cómputo a través de Internet y donde el usuario final no tiene conocimiento de dónde se están ejecutando sus aplicaciones y tampoco necesita ser un experto para contactar, subir y ejecutar sus aplicaciones en minutos. Los proveedores de este tipo de servicio tienen recursos (generalmente distribuidos en todo el mundo) y permiten que sus usuarios contraten y gestionen estos recursos en forma *on-line* y sin la mayor intervención y en forma casi automática, lo cual permite reducir los costes teniendo ventajas (además de que no se tiene que instalar-mantener la infraestructura física –ni la obra civil–) como la fiabilidad, flexibilidad, rapidez en el aprovisionamiento, facilidad de uso, pago por uso, contratación de lo que se necesita, etc. Obviamente, también cuenta con sus desventajas, como lo son la dependencia de un proveedor

y la centralización de las aplicaciones|almacenamiento de datos, servicio vinculado a la disponibilidad de acceso a Internet, cuestiones de seguridad, que los datos 'sensibles' del negocio no residen en las instalaciones de las empresas y pueden generar riesgos de sustracción/robo de información, confiabilidad de los servicios prestados por el proveedor (siempre es necesario firmar una SLA –contrato de calidad de servicio–), tecnología susceptible al monopolio, servicios estándar (solo los ofrecidos por el proveedor), escalabilidad o privacidad.

Un tercer concepto vinculado a estas tecnologías, pero probablemente más del lado de los desarrolladores de aplicaciones y *frameworks*, es el de DevOps. DevOps surge de la unión de las palabras *Development* (desarrollo) y *Operations* (operaciones), y se refiere a una metodología de desarrollo de software que se centra en la comunicación, colaboración e integración entre desarrolladores de software y los profesionales de operaciones|administradores en las tecnologías de la información (IT). DevOps se presenta como una metodología que da respuesta a la relación existente entre el desarrollo de software y las operaciones IT, teniendo como objetivo que los productos y servicios software desarrollados por una entidad se puedan hacer más eficientemente, tengan una alta calidad y además sean seguros y fáciles de mantener. El término DevOps es relativamente nuevo y fue popularizado a través de *DevOps Open Days* (Bélgica, 2009) y como metodología de desarrollo ha ido ganando adeptos desde grande a pequeñas compañías, para hacer más y mejor sus productos y servicios y, sobre todo, debido a diferentes factores con una fuerte presencia hoy en día como: el uso de los procesos y metodologías de desarrollo ágil, necesidad de una mayor tasa de versiones, amplia disponibilidad de entornos virtualizados|*cloud*, mayor automatización de centros de datos y aumento de las herramientas de gestión de configuración.

En este módulo se verán diferentes formas de crear y programar un sistema de cómputo distribuido (clúster|*cloud*), las herramientas y librerías más importantes para cumplir este objetivo, y los conceptos y herramientas vinculadas a las metodológicas DevOps.

Objetivos

En los materiales didácticos de este módulo encontraréis los contenidos y las herramientas procedimentales para conseguir los objetivos siguientes:

- 1.** Analizar las diferentes infraestructuras y herramientas para el cómputo de altas prestaciones (incluida la virtualización) (HPC).
- 2.** Configurar e instalar un clúster de HPC y las herramientas de monitorización correspondientes.
- 3.** Instalar y desarrollar programas de ejemplos en las principales API de programación: Posix Threads, OpenMPI, y OpenMP.
- 4.** Instalar un clúster específico basado en una distribución ad hoc (Rocks).
- 5.** Instalar una infraestructura para prestar servicios en *cloud* (IaaS).
- 6.** Herramientas DevOps para automatizar un centros de datos y gestiones de la configuración.

1. Clusterización

La historia de los sistemas informáticos es muy reciente (se puede decir que comienza en la década de 1960). En un principio, eran sistemas grandes, pesados, caros, de pocos usuarios expertos, no accesibles y lentos. En la década de 1970, la evolución permitió mejoras sustanciales llevadas a cabo por tareas interactivas (*interactive jobs*), tiempo compartido (*time sharing*), terminales y con una considerable reducción del tamaño. La década de 1980 se caracteriza por un aumento notable de las prestaciones (hasta hoy en día) y una reducción del tamaño en los llamados microordenadores. Su evolución ha sido a través de las estaciones de trabajo (*workstations*) y los avances en redes (LAN de 10 Mb/s y WAN de 56 kB/s en 1973 a LAN de 1/10 Gb/s y WAN con ATM, *asynchronous transfer mode* de 1,2 Gb/s en la actualidad o redes de alto rendimiento como Infiniband –96Gb/s– o Myrinet –10Gb/s–), que es un factor fundamental en las aplicaciones multimedia actuales y de un futuro próximo. Los sistemas distribuidos, por su parte, comenzaron su historia en la década de 1970 (sistemas de 4 u 8 ordenadores) y su salto a la popularidad lo hicieron en la década de 1990. Si bien su administración, instalación y mantenimiento pueden tener una cierta complejidad (cada vez menos) porque continúan creciendo en tamaño, las razones básicas de su popularidad son el incremento de prestaciones que presentan en aplicaciones intrínsecamente distribuidas (aplicaciones que por su naturaleza son distribuidas), la información compartida por un conjunto de usuarios, la compartición de recursos, la alta tolerancia a los fallos y la posibilidad de expansión incremental (capacidad de agregar más nodos para aumentar las prestaciones y de forma incremental). Otro aspecto muy importante en la actualidad es la posibilidad, en esta evolución, de la virtualización. Las arquitecturas cada vez más eficientes, con sistemas *multicores*, han permitido que la virtualización de sistemas se transforme en una realidad con todas las ventajas (y posibles desventajas) que ello comporta.

1.1. Virtualización

La virtualización es una técnica que está basada en la abstracción de los recursos de un ordenador, llamada *Hypervisor* o VMM (*Virtual Machine Monitor*) que crea una capa de separación entre el hardware de la máquina física (*host*) y el sistema operativo de la máquina virtual (*virtual machine, guest*), y es un medio para crear una “versión virtual” de un dispositivo o recurso, como un servidor, un dispositivo de almacenamiento, una red o incluso un sistema operativo, donde se divide el recurso en uno o más entornos de ejecución. Esta capa de software (VMM) maneja, gestiona y administra los cuatro recur-

Los principales de un ordenador (CPU, memoria, red y almacenamiento) y los reparte de forma dinámica entre todas las máquinas virtuales definidas en el computador central. De este modo nos permite tener varios ordenadores virtuales ejecutándose sobre el mismo ordenador físico.

La máquina virtual, en general, es un sistema operativo completo que se ejecuta como si estuviera instalado en una plataforma de hardware autónoma.

Enlace de interés

Para saber más sobre virtualización podéis visitar: <http://en.wikipedia.org/wiki/Virtualization>. Se puede consultar una lista completa de programas de virtualización en: http://en.wikipedia.org/wiki/Comparison_of_platform_virtual_machines.

1.1.1. Plataformas de virtualización

Los ejemplos más comunes de plataforma de virtualización con licencias GPLs o similares son Xen, KVM, Qemu (emulación), OpenVz, VirtualBox, Oracle VM (solo el servidor), y entre las plataformas propietarias (algunas como *free to use*) VMware ESX/i, Virtual PC/Hyper-V, Parallels, Virtuozzo. Existen diferentes taxonomías para clasificar la virtualización (y algunas plataformas pueden ejecutarse en más de una categoría) siendo la más conocida la siguiente:

1) Virtualización por HW: considerando como *Full Virtualization* cuando la máquina virtual simula todo el HW para permitir a un sistema *guest* ejecutarse sin modificaciones (siempre y cuando esté diseñado para el mismo set de instrucciones). Fue la 1.ª generación de VM para procesadores x86 y lo que hace es una captura de instrucciones que acceden en modo prioritario a la CPU y las transforma en una llamada a la VM para que sean emuladas por software. En este tipo de virtualización pueden ejecutarse Parallels, VirtualBox, Virtual Iron, Oracle VM, Virtual PC, Hyper-V, VMware por ejemplo.

2) Virtualización asistida por hardware, donde el hardware provee soporte que facilita la construcción y trabajo del VM y mejora notablemente las prestaciones (a partir de 2005/6 Intel y AMD proveen este soporte como VT-x y AMD-V respectivamente). Las principales ventajas, respecto a otras formas de virtualización, es que no se debe tocar el sistema operativo *guest* (como en paravirtualization) y se obtienen mejores prestaciones, pero como contrapartida, se necesita soporte explícito en la CPU, lo cual no está disponible en todos los procesadores x86/86_64. En este tipo de virtualización pueden ejecutarse KVM, VMware Fusion, Hyper-V, Virtual PC, Xen, Parallels, Oracle VM, VirtualBox.

3) Paravirtualization: es una técnica de virtualización donde la VM no necesariamente simula el HW, sino que presenta una API que puede ser utilizada por el sistema *guest* (por lo cual, se debe tener acceso al código fuente para reemplazar las llamadas de un tipo por otro). Este tipo de llamadas a esta API se denominan *hypercall* y Xen puede ejecutarse en esta forma.

4) Virtualización parcial: es lo que se denomina *Address Space Virtualization*. La máquina virtual simula múltiples instancias del entorno subyacente del hardware (pero no de todo), particularmente el *address space*. Este tipo de virtualización acepta compartir recursos y alojar procesos, pero no permite instancias separadas de sistemas operativos *guest* y se encuentra en desuso actualmente.

5) Virtualización a nivel del sistema operativo: en este caso la virtualización permite tener múltiples instancias aisladas y seguras del servidor, todas ejecutándose sobre el mismo servidor físico y donde el sistema operativo *guest* coincidirá con el sistema operativo base del servidor (se utiliza el mismo kernel). Plataformas dentro de este tipo de virtualización son FreeBSD jails (el pionero), OpenVZ, Linux-VServer, LXC, Virtuozzo.

La diferencia entre instalar dos sistemas operativos y virtualizar dos sistemas operativos es que en el primer caso todos los sistemas operativos que tengamos instalados funcionarán de la misma manera que si estuvieran instalados en distintos ordenadores, y necesitaremos un gestor de arranque que al encender el ordenador nos permita elegir qué sistema operativo queremos utilizar, pero solo podremos tener funcionando simultáneamente uno de ellos. En cambio, la virtualización permite ejecutar muchas máquinas virtuales con sus sistemas operativos y cambiar de sistema operativo como si se tratase de cualquier otro programa; sin embargo, se deben valorar muy bien las cuestiones relacionadas con las prestaciones, ya que si el HW subyacente no es el adecuado, podremos notar muchas diferencias en las prestaciones entre el sistema operativo instalado en base o el virtualizado.

Entre las principales ventajas de la virtualización podemos contemplar:

- 1) Consolidación de servidores y mejora de la eficiencia de la inversión en HW con reutilización de la infraestructura existente.
- 2) Rápido despliegue de nuevos servicios con balanceo dinámico de carga y reducción de los sobredimensionamientos de la infraestructura.
- 3) Incremento de *Uptime* (tiempo que el sistema está al 100 % en la prestación de servicios), incremento de la tolerancia a fallos (siempre y cuando exista redundancia física) y eliminación del tiempo de parada por mantenimiento del sistema físico (migración de las máquinas virtuales).
- 4) Mantenimiento a coste aceptables de entornos software obsoletos pero necesarios para el negocio.
- 5) Facilidad de diseño y test de nuevas infraestructuras y entornos de desarrollo con un bajo impacto en los sistemas de producción y rápida puesta en marcha.

- 6) Mejora de TCO (*Total Cost of Ownership*) y ROI (*Return on Investment*).
- 7) Menor consumo de energía que en servidores físicos equivalentes.

Como desventajas podemos mencionar:

- 1) Aumenta la probabilidad de fallos si no se considera redundancia/alta disponibilidad (si se consolidan 10 servidores físicos en uno potente equivalente, con servidores virtualizados, y dejan de funcionar todos los servidores en él, dejarán de prestar servicio).
- 2) Rendimiento inferior (posible) en función de la técnica de virtualización utilizada y recursos disponibles.
- 3) Proliferación de servicios y máquinas que incrementan los gastos de administración/gestión (básicamente por el efecto derivado de la 'facilidad de despliegue' se tiende a tener más de lo necesarios).
- 4) Infraestructura desaprovechada (posible) ya que es habitual comprar una infraestructura mayor que la necesaria en ese momento para el posible crecimiento futuro inmediato.
- 5) Pueden existir problemas de portabilidad, hardware específico no soportado, y compromiso a largo término con la infraestructura adquirida.
- 6) Tomas de decisiones en la selección del sistema anfitrión puede ser complicada o condicionante.

Como es posible observar, las desventajas se pueden resolver con una planificación y toma de decisiones adecuada, y esta tecnología es habitual en la prestación de servicios y totalmente imprescindible en entornos de *Cloud Computing* (que veremos en este capítulo también).

1.2. Beowulf

Beowulf [Rad, Beo, Beo1, Kur] es una arquitectura multiordenador que puede ser utilizada para aplicaciones paralelas/distribuidas. El sistema consiste básicamente en un servidor y uno o más clientes conectados (generalmente) a través de Ethernet y sin la utilización de ningún hardware específico. Para explotar esta capacidad de cómputo, es necesario que los programadores tengan un modelo de programación distribuido que, si bien es posible mediante UNIX (Sockets, RPC), puede implicar un esfuerzo considerable, ya que son modelos de programación a nivel de *systems calls* y lenguaje C, por ejemplo; pero este modo de trabajo puede ser considerado de bajo nivel. Un modelo

más avanzado en esta línea de trabajo son los **Posix Threads**, que permiten explotar sistemas de memoria compartida y *multicores* de forma simple y fácil. La capa de software (interfaz de programación de aplicaciones, API) aportada por sistemas tales como **Parallel Virtual Machine** (PVM) y **Message Passing Interface** (MPI) facilita notablemente la abstracción del sistema y permite programar aplicaciones paralelas/distribuidas de modo sencillo y simple. Una de las formas básicas de trabajo es la de maestro-trabajadores (*master-workers*), en que existe un servidor (maestro) que distribuye la tarea que realizarán los trabajadores. En grandes sistemas (por ejemplo, de 1.024 nodos) existe más de un maestro y nodos dedicados a tareas especiales, como por ejemplo entrada/salida o monitorización. Otra opción no menos interesante, sobre todo con el auge de procesadores *multicores*, es **OpenMP** que es una API para la programación multiproceso de memoria compartida. Esta capa de software permite añadir concurrencia a los programas sobre la base del modelo de ejecución *fork-join* y se compone de un conjunto de directivas de compilador, rutinas de biblioteca y variables de entorno, que influyen el comportamiento en tiempo de ejecución y proporciona a los programadores una interfaz simple y flexible para el desarrollo de aplicaciones paralelas y arquitecturas de CPU que puedan ejecutar más de un hilo de ejecución simultáneo (*hyperthreading*) o que dispongan de más de un núcleo por procesador accediendo a la misma memoria compartida.

Existen dos conceptos que pueden dar lugar a dudas y que son Cluster Beowulf y COW (*Cluster of Workstations*). Una de las principales diferencias es que Beowulf “se ve” como una única máquina donde se accede a los nodos remotamente, ya que no disponen de terminal (ni de teclado), mientras que un COW es una agrupación de ordenadores que pueden ser utilizados tanto por los usuarios de la COW como por otros usuarios en forma interactiva, a través de su pantalla y teclado. Hay que considerar que Beowulf no es un software que transforma el código del usuario en distribuido ni afecta al núcleo del sistema operativo (como por ejemplo, Mosix). Simplemente, es una forma de agrupación (un clúster) de máquinas que ejecutan GNU/Linux y actúan como un super-ordenador. Obviamente, existe una gran cantidad de herramientas que permiten obtener una configuración más fácil, bibliotecas o modificaciones al núcleo para obtener mejores prestaciones, pero es posible construir un clúster Beowulf a partir de un GNU/Linux estándar y de software convencional. La construcción de un clúster Beowulf de dos nodos, por ejemplo, se puede llevar a cabo simplemente con las dos máquinas conectadas por Ethernet mediante un concentrador (*hub*), una distribución de GNU/Linux estándar (Debian), el sistema de archivos compartido (NFS) y tener habilitados los servicios de red, como por ejemplo `ssh`. En estas condiciones, se puede argumentar que se dispone de un clúster simple de dos nodos. En cambio, en un COW no necesitamos tener conocimientos sobre la arquitectura subyacente (CPU, red) necesaria para hacer una aplicación distribuida en Beowulf pero sí que es necesario tener un sistema operativo (por ejemplo, Mosix) que permita este tipo de compartición de recursos y distribución de tareas (además de un API específica para programar la aplicación que es necesaria en los dos

sistemas). El primer Beowulf se construyó en 1994 y en 1998 comienzan a aparecer sistemas Beowulf/linux en las listas del Top500 (Avalon en la posición 314 con 68 cores y 19,3 Gflops, junio'98).

1.2.1. ¿Cómo configurar los nodos?

Consideramos que el servidor tiene dos interfaces de red: *eth0* conectada a Internet y *eth1* conectada a una red interna a la cual están conectados los nodos. Primero se debe modificar el archivo `/etc/hosts` para que contenga la línea de `localhost` (con 127.0.0.1) y el resto de los nombres a las IP internas de los restantes nodos (este archivo deberá ser igual en todos los nodos). Por ejemplo:

```
127.0.0.1 localhost
192.168.0.254 srv.nteum.org srv
```

Y añadir las IP de los nodos (y para todos los nodos), por ejemplo:

```
192.168.0.1 lucix1.nteum.org lucix1
192.168.0.2 lucix2.nteum.org lucix2
...
```

Se debe crear un usuario (`vteum`) en todos los nodos, crear un grupo y añadir este usuario al grupo:

```
groupadd beowulf
adduser vteum beowulf
echo umask 007 >> /home/vteum/.bash_profile
```

Así, cualquier archivo creado por el usuario `vteum` o cualquiera dentro del grupo será modificable por el grupo `beowulf`. Se debe crear un servidor de NFS (y los demás nodos serán clientes de este NFS) y exportar el directorio `/home` así todos los clientes verán el directorio `$HOME` de los usuarios. Para ello sobre el servidor se hace la edición del archivo `/etc/exports` y se agrega una línea como `/home 192.168.0.0/24(rw, sync, no_root_squash)`. Sobre cada nodo cliente se monta agregando en `/etc/fstab` para que en el arranque el nodo monte el directorio como

```
192.168.0.254:/home /home nfs defaults 0 0
```

también es aconsejable tener en el servidor un directorio `/soft` (donde se instalará todo el software para que lo vean todos los nodos, de tal manera que nos evitamos de tener que instalar este en cada nodo), y agregarlos al `export` como `/soft 192.16.0.0/24(rw, async, no_root_squash)` y en cada nodo agregamos en el `fstab`

```
192.168.0.254:/soft /soft nfs defaults 0 0
```

También deberemos hacer la configuración del archivo `/etc/resolv.conf` y el `iptables` (para hacer un NAT) sobre el servidor (por ejemplo, `iptables -t nat -A POSTROUTING -o eth0 -j MASQUERADE`) si queremos que los nodos también tengan acceso (es muy útil sobre todo a la hora de actualizar el sistema operativo).

A continuación se debe verificar que los servicios están activos (con la orden `systemctl status sshd`; `systemctl status nfs` o también con la utilidad `chkconfig`, aunque en Debian por ejemplo no está disponible y se deberá instalar). En caso de que estos servicios no estén habilitados deberemos habilitarlos (por ejemplo, `systemctl enable sshd`). Lo siguiente es podernos conectar desde un determinado usuario a todas la máquinas mediante `passwd` por lo que utilizaremos una característica de `ssh` que usa un método de clave pública-privada para la autenticación. Para ello verificamos que están habilitadas en el archivo `/etc/ssh/sshd_config` las siguientes líneas (sin # delante):

```
RSAAuthentication yes
AuthorizedKeysFile .ssh/authorized_keys
```

Reiniciamos el servicio (con la orden `service sshd restart` o también con `systemctl restart ssh`) y nos conectamos con el usuario que se desea generar las llaves (por ejemplo, `vteum`). A continuación ejecutamos la instrucción `ssh-keygen -t rsa`, que creará en `$HOME/.ssh` las llaves `id_rsa` y `id_rsa.pub`, en referencia a la llave privada y pública respectivamente. Finalmente ejecutamos `ssh-copy-id lucix1.nteum.org` para copiar al primer nodo la llave pública del usuario `vteum`. Como tendremos el directorio `$HOME` montado para todos los nodos por NFS solo la deberemos hacer la copia para el primer nodo y ya nos podremos conectar a los restantes nodos, ya que la llave pública es la misma para todos ellos. A continuación debemos comprobar que la conexión a `lucix1` desde y hacia el usuario `vteum` no requiere `passwd`. Es conveniente también instalar NIS sobre el clúster así evitamos tener que definir cada usuario en los nodos y un servidor DHCP para definir todos los parámetros de red de cada nodo que cada nodo solicite estos durante la etapa de `boot`, para ello consultad en el capítulo de red de la asignatura de Administración GNU/Linux donde se explica con detalle estas configuraciones.

A partir de aquí ya tenemos un clúster Beowulf para ejecutar aplicaciones con interfaz a MPI para aplicaciones distribuidas (también puede ser aplicaciones con interfaz a PVM, pero es recomendable MPI por cuestiones de eficiencia y prestaciones). Existen sobre diferentes distribuciones (Debian, FC incluidas) la aplicación `system-config-cluster`, que permite configurar un clúster en base a una herramienta gráfica o `clusterssh`, `parallel`, `pssh` o `dish`, que permiten gestionar y ejecutar comandos en un conjunto de máquinas de forma simultánea (comandos útiles cuando necesitamos hacer operaciones sobre todos los nodos del clúster, por ejemplo, apagarlos, reiniciarlos o hacer copias de un archivo a todos ellos).

Enlace de interés

Información Adicional:
<https://wiki.debian.org/HighPerformanceComputing>

1.3. Beneficios del cómputo distribuido

¿Cuáles son los beneficios del cómputo en paralelo? Veremos esto con un ejemplo [Rad]. Consideremos un programa para sumar números (por ejemplo, `4 + 5 + 6 + ...`) llamado `sumdis.c`:

```
#include <stdio.h>

int main (int argc, char** argv){
long inicial, final, resultado, tmp;
    if (argc < 2) {
        printf (" Uso: %s N° inicial N° final\n",argv[0]);
        return (4); }
    else {
        inicial = atoll(argv[1]);
        final = atoll(argv[2]);
        resultado = 0;}
    for (tmp = inicial; tmp <= final; tmp++){resultado += tmp; };
    printf("%lu\n", resultado);
    return 0;
}
```

Lo compilamos con `gcc -o sumdis sumdis.c` y si miramos la ejecución de este programa con, por ejemplo,

```
time ./sumdis 1 1000000
500000500000

real 0m0.005s
user 0m0.000s
sys 0m0.004s
```

se podrá observar que el tiempo en una máquina Debian sobre una máquina virtual (Virtualbox) con procesador i7 es 5 milésimas de segundo. Si, en cambio, hacemos desde 1 a 16 millones, el tiempo real sube hasta 0,050s, es decir, 10 veces más, lo cual, si se considera 1600 millones, el tiempo será del orden de unos 4 segundos.

La idea básica del cómputo distribuido es repartir el trabajo. Así, si disponemos de un clúster de 4 máquinas (lucix1–lucix4) con un servidor y donde el archivo ejecutable se comparte por NFS, sería interesante dividir la ejecución mediante ssh de modo que el primero sume de 1 a 400.000.000, el segundo, de 400.000.001 a 800.000.000, el tercero, de 800.000.001 a 1.200.000.000 y el cuarto, de 1.200.000.001 a 1.600.000.000. Los siguientes comandos muestran una posibilidad. Consideramos que el sistema tiene el directorio `/home` compartido por NFS y el usuario **adminp**, que tiene adecuadamente configurado las llaves para ejecutar el código sin *password* sobre los nodos, ejecutará:

```
mkfifo out1      Crea una cola fifo en /home/adminp
./distr.sh & time cat out1 | awk '{total += $1 } END {printf "%lf", total}'
```

Se ejecuta el comando `distr.sh`; se recolectan los resultados y se suman mientras se mide el tiempo de ejecución. El *shell script* `distr.sh` puede ser algo como:


```
ssh lucix1 /home/adminp/sumdis 1 400000000 > /home/adminp/out1 < /dev/null &
ssh lucix2 /home/adminp/sumdis 400000001 800000000 > /home/adminp/out1 < /dev/null &
ssh lucix3 /home/adminp/sumdis 800000001 1200000000 > /home/adminp/out1 < /dev/null &
ssh lucix4 /home/adminp/sumdis 1200000001 1600000000 > /home/adminp/out1 < /dev/null &
```

Podremos observar que el tiempo se reduce notablemente (aproximadamente en un valor cercano a 4) y no exactamente de forma lineal, pero muy próxima. Obviamente, este ejemplo es muy simple y solo válido para fines demostrativos. Los programadores utilizan bibliotecas que les permiten realizar el tiempo de ejecución, la creación y comunicación de procesos en un sistema distribuido (por ejemplo MPI u OpenMP). Una forma más eficiente de ejecutar en forma paralela múltiples comandos es utilizando `pssh`, `parallel` o `taktuk` (todos en el repositorio de Debian). Por ejemplo, después de instalado `parallel` ejecutamos:

```
parallel -S lucix1 -S lucix2 -S lucix3 -verbose ./sumdis 1 :::
10000 100000 100000000
```

En este caso, si medimos el tiempo con el comando `time` (anteponiéndolo al comando `parallel`) veremos que el tiempo de ejecución es mejor que en el caso anterior (le podemos agregar igual que en el comando anterior el `awk` para que sume todos los resultados también). El comando `parallel` admite una gran cantidad de parámetros y opciones (igual que `pssh` y `taktuk`), que se pueden consultar en el tutorial [Par].

1.3.1. ¿Cómo hay que programar para aprovechar la concurrencia?

Existen diversas maneras de expresar la concurrencia en un programa. Las tres más comunes son:

- 1) Utilizando hilos (o procesos) en el mismo procesador (multiprogramación con solapamiento del cómputo y la E/S).
- 2) Utilizando hilos (o procesos) en sistemas *multicore*.
- 3) Utilizando procesos en diferentes procesadores que se comunican por medio de mensajes (MPS, *Message Passing System*).

Estos métodos pueden ser implementados sobre diferentes configuraciones de hardware (memoria compartida o mensajes) y, si bien ambos métodos tienen sus ventajas y desventajas, los principales problemas de la memoria compartida son las limitaciones en la escalabilidad (ya que todos los *cores*/procesadores utilizan la misma memoria y el número de estos en el sistema está limitado por el ancho de banda de la memoria) y, en los los sistemas de paso de mensajes, la latencia y velocidad de los mensajes en la red. El programador deberá evaluar

qué tipo de prestaciones necesita, las características de la aplicación subyacente y el problema que se desea solucionar. No obstante, con los avances de las tecnologías de *multicores* y de red, estos sistemas han crecido en popularidad (y en cantidad). Las API más comunes hoy en día son Posix Threads y OpenMP para memoria compartida y MPI (en sus versiones OpenMPI o Mpich) para paso de mensajes. Como hemos mencionado anteriormente, existe otra biblioteca muy difundida para pasos de mensajes, llamada PVM, pero que la versatilidad y prestaciones que se obtienen con MPI ha dejado relegada a aplicaciones pequeñas o para aprender a programar en sistemas distribuidos. Estas bibliotecas, además, no limitan la posibilidad de utilizar hilos (aunque a nivel local) y tener concurrencia entre procesamiento y entrada/salida.

Para realizar una aplicación paralela/distribuida, se puede partir de la versión serie o mirando la estructura física del problema y determinar qué partes pueden ser concurrentes (independientes). Las partes concurrentes serán candidatas a re-escribirse como código paralelo. Además, se debe considerar si es posible reemplazar las funciones algebraicas por sus versiones paralelizadas (por ejemplo, ScaLapack *Scalable Linear Algebra Package* (se puede probar en Debian los diferentes programas de prueba que se encuentran en el paquete *scalapack-mpi-test*, por ejemplo y directamente, instalar las librerías *libscalapack-mpi-dev*). También es conveniente averiguar si hay alguna aplicación similar paralela que pueda orientarnos sobre el modo de construcción de la aplicación paralela*.

*<http://www.mpich.org/>

Paralelizar un programa no es una tarea fácil, ya que se debe tener en cuenta la **ley de Amdahl**, que afirma que el incremento de velocidad (*speedup*) está limitado por la fracción de código (f), que puede ser paralelizado de la siguiente manera:

$$\text{speedup} = \frac{1}{1-f}$$

Esta ley implica que con una aplicación secuencial $f = 0$ y el $\text{speedup} = 1$, mientras que con todo el código paralelo $f = 1$ y el speedup se hace infinito. Si consideramos valores posibles, un 90% ($f = 0,9$) del código paralelo significa un speedup igual a 10, pero con $f = 0,99$ el speedup es igual a 100. Esta limitación se puede evitar con algoritmos escalables y diferentes modelos de programación de aplicación (paradigmas):

- 1) Maestro-trabajador: el maestro inicia a todos los trabajadores y coordina su trabajo y el de entrada/salida.
- 2) *Single Process Multiple Data* (SPMD): mismo programa que se ejecuta con diferentes conjuntos de datos.

3) Funcional: varios programas que realizan una función diferente en la aplicación.

En resumen, podemos concluir:

- 1) Proliferación de máquinas multitarea (multiusuario) conectadas por red con servicios distribuidos (NFS y NIS).
- 2) Son sistemas heterogéneos con sistemas operativos de tipo NOS (*Networked Operating System*), que ofrecen una serie de servicios distribuidos y remotos.
- 3) La programación de aplicaciones distribuidas se puede efectuar a diferentes niveles:
 - a) Utilizando un modelo cliente-servidor y programando a bajo nivel (*sockets*) o utilizando memoria compartida a bajo nivel (Posix Threads).
 - b) El mismo modelo, pero con API de “alto” nivel (OpenMP, MPI).
 - c) Utilizando otros modelos de programación como, por ejemplo, programación orientada a objetos distribuidos (RMI, CORBA, Agents, etc.).

1.4. Memoria compartida. Modelos de hilos (*threading*)

Normalmente, en una arquitectura cliente-servidor, los clientes solicitan a los servidores determinados servicios y esperan que estos les contesten con la mayor eficacia posible. Para sistema distribuidos con servidores con una carga muy alta (por ejemplo, sistemas de archivos de red, bases de datos centralizadas o distribuidas), el diseño del servidor se convierte en una cuestión crítica para determinar el rendimiento general del sistema distribuido. Un aspecto crucial en este sentido es encontrar la manera óptima de manejar la E/S, teniendo en cuenta el tipo de servicio que ofrece, el tiempo de respuesta esperado y la carga de clientes. No existe un diseño predeterminado para cada servicio, y escoger el correcto dependerá de los objetivos y restricciones del servicio y de las necesidades de los clientes.

Las preguntas que debemos contestar antes de elegir un determinado diseño son: ¿Cuánto tiempo se tarda en un proceso de solicitud del cliente? ¿Cuántas de esas solicitudes es probable que lleguen durante ese tiempo? ¿Cuánto tiempo puede esperar el cliente? ¿Cuánto afecta esta carga del servidor a las prestaciones del sistema distribuido? Además, con el avance de la tecnología de procesadores nos encontramos con que disponemos de sistemas *multicore* (múltiples núcleos de ejecución) que pueden ejecutar secciones de código independientes. Si se diseñan los programas en forma de múltiples secuencias de ejecución y el sistema operativo lo soporta (y GNU/Linux es uno de ellos), la ejecución de los programas se reducirá notablemente y se incrementarán en forma (casi) lineal las prestaciones en función de los *cores* de la arquitectura.

1.4.1. Multihilos (*multithreading*)

Las últimas tecnologías en programación para este tipo aplicaciones (y así lo demuestra la experiencia) es que los diseños más adecuados son aquellos que utilizan modelos de multi-hilos (*multithreading models*), en los cuales el servidor tiene una organización interna de procesos paralelos o hilos cooperarantes y concurrentes.

Un hilo (*thread*) es una secuencia de ejecución (hilo de ejecución) de un programa, es decir, diferentes partes o rutinas de un programa que se ejecutan concurrentemente en un único procesador y accederán a los datos compartidos al mismo tiempo.

¿Qué ventajas aporta esto respecto a un programa secuencial? Consideremos que un programa tiene tres rutinas A, B y C. En un programa secuencial, la rutina C no se ejecutará hasta que se hayan ejecutado A y B. Si, en cambio, A, B y C son hilos, las tres rutinas se ejecutarán concurrentemente y, si en ellas hay E/S, tendremos concurrencia de ejecución con E/S del mismo programa (proceso), cosa que mejorará notablemente las prestaciones de dicho programa. Generalmente, los hilos están contenidos dentro de un proceso y diferentes hilos de un mismo proceso pueden compartir algunos recursos, mientras que diferentes procesos no. La ejecución de múltiples hilos en paralelo necesita el soporte del sistema operativo y en los procesadores modernos existen optimizaciones del procesador para soportar modelos multihilo (*multithreading*) además de las arquitectura multicores donde existe múltiples núcleo y donde cada uno de ellos puede ejecutar un thread.

Generalmente, existen cuatro modelos de diseño por hilos (en orden de complejidad creciente):

- 1) Un hilo y un cliente:** en este caso el servidor entra en un bucle sin fin escuchando por un puerto y ante la petición de un cliente se ejecutan los servicios en el mismo hilo. Otros clientes deberán esperar a que termine el primero. Es fácil de implementar pero solo atiende a un cliente a la vez.
- 2) Un hilo y varios clientes con selección:** en este caso el servidor utiliza un solo hilo, pero puede aceptar múltiples clientes y multiplexar el tiempo de CPU entre ellos. Se necesita una gestión más compleja de los puntos de comunicación (*sockets*), pero permite crear servicios más eficientes, aunque presenta problemas cuando los servicios necesitan una alta carga de CPU.
- 3) Un hilo por cliente:** es, probablemente, el modelo más popular. El servidor espera por peticiones y crea un hilo de servicio para atender a cada nueva petición de los clientes. Esto genera simplicidad en el servicio y una alta disponibilidad, pero el sistema no escala con el número de clientes y puede

saturar el sistema muy rápidamente, ya que el tiempo de CPU dedicado ante una gran carga de clientes se reduce notablemente y la gestión del sistema operativo puede ser muy compleja.

4) Servidor con hilos en granja (*worker threads*): este método es más complejo pero mejora la escalabilidad de los anteriores. Existe un número fijo de hilos trabajadores (*workers*) a los cuales el hilo principal distribuye el trabajo de los clientes. El problema de este método es la elección del número de trabajadores: con un número elevado, caerán las prestaciones del sistema por saturación; con un número demasiado bajo, el servicio será deficiente (los clientes deberán esperar). Normalmente, será necesario sintonizar la aplicación para trabajar con un determinado entorno distribuido.

Existen diferentes formas de expresar a nivel de programación con hilos: paralelismo a nivel de tareas o paralelismo a través de los datos. Elegir el modelo adecuado minimiza el tiempo necesario para modificar, depurar y sintonizar el código. La solución a esta disyuntiva es describir la aplicación en términos de dos modelos basados en un trabajo en concreto:

- Tareas paralelas con hilos independientes que pueden atender tareas independientes de la aplicación. Estas tareas independientes serán encapsuladas en hilos que se ejecutarán asincrónicamente y se deberán utilizar bibliotecas como Posix Threads (Linux/Unix) o Win32 Thread API (Windows), que han sido diseñadas para soportar concurrencia a nivel de tarea.
- Modelo de datos paralelos para calcular lazos intensivos; es decir, la misma operación debe repetirse un número elevado de veces (por ejemplo comparar una palabra frente a las palabras de un diccionario). Para este caso es posible encargar la tarea al compilador de la aplicación o, si no es posible, que el programador describa el paralelismo utilizando el entorno OpenMP, que es una API que permite escribir aplicaciones eficientes bajo este tipo de modelos.

Una aplicación de información personal (*Personal Information Manager*) es un buen ejemplo de una aplicación que contiene concurrencia a nivel de tareas (por ejemplo, acceso a la base de datos, libreta de direcciones, calendario, etc.). Esto podría ser en pseudocódigo:

```
Function addressBook;  
Function inBox;  
Function calendar;  
Program PIM {  
    CreateThread (addressBook);  
    CreateThread (inBox);  
    CreateThread (calendar); }
```

Podemos observar que existen tres ejecuciones concurrentes sin relación entre ellas. Otro ejemplo de operaciones con paralelismo de datos podría ser un

corrector de ortografía, que en pseudocódigo sería: `Function SpellCheck {loop (word = 1, words_in_file) compareToDictionary (word);}`

Se debe tener en cuenta que ambos modelos (hilos paralelos y datos paralelos) pueden existir en una misma aplicación. A continuación se mostrará el código de un productor de datos y un consumidor de datos basado en Posix Threads. Para compilar sobre Linux, por ejemplo, se debe utilizar `gcc -o pc pc.c -lpthread`.

```
#include <pthread.h>
#include <stdio.h>
#include <unistd.h>
#include <stdlib.h>
#define QUEUESIZE 10
#define LOOP 20

void *producer (void *args);
void *consumer (void *args);
typedef struct { /* Estructura del buffer compartido y descriptores de threads */
int buf[QUEUESIZE]; long head, tail; int full, empty;
pthread_mutex_t *mut; pthread_cond_t *notFull, *notEmpty;
} queue;
queue *queueInit (void); /* Prototipo de función: inicialización del buffer */
void queueDelete (queue *q); /* Prototipo de función: Borrado del buffer*/
void queueAdd (queue *q, int in); /* Prototipo de función: insertar elemento en el buffer */
void queueDel (queue *q, int *out); /* Prototipo de función: quitar elemento del buffer */

int main () {
queue *fifo; pthread_t pro, con; fifo = queueInit ();
if (fifo == NULL) { fprintf (stderr, " Error al crear buffer.\n"); exit (1); }
pthread_create (&pro, NULL, producer, fifo); /* Creación del thread productor */
pthread_create (&con, NULL, consumer, fifo); /* Creación del thread consumidor*/
pthread_join (pro, NULL); /* main () espera hasta que terminen ambos threads */
pthread_join (con, NULL);
queueDelete (fifo); /* Eliminación del buffer compartido */
return 0; } /* Fin */

void *producer (void *q) { /*Función del productor */
queue *fifo; int i;
fifo = (queue *)q;
for (i = 0; i < LOOP; i++) { /* Inserto en el buffer elementos=LOOP*/
pthread_mutex_lock (fifo->mut); /* Semáforo para entrar a insertar */
while (fifo->full) {
printf ("Productor: queue FULL.\n");
pthread_cond_wait (fifo->notFull, fifo->mut); }
/* Bloqueo del productor si el buffer está lleno, liberando el semáforo mut */
/*para que pueda entrar el consumidor. Continuará cuando el consumidor ejecute*/
pthread_cond_signal (fifo->notEmpty);*/
queueAdd (fifo, i); /* Inserto elemento en el buffer */
pthread_mutex_unlock (fifo->mut); /* Libero el semáforo */
pthread_cond_signal (fifo->notFull); /*Desbloqueo consumidor si está bloqueado*/
usleep (100000); /* Duermo 100 mseg para permitir que el consumidor se active */
}
return (NULL); }

void *consumer (void *q) { /*Función del consumidor */
queue *fifo; int i, d;
fifo = (queue *)q;
for (i = 0; i < LOOP; i++) { /* Quito del buffer elementos=LOOP*/
pthread_mutex_lock (fifo->mut); /* Semáforo para entrar a quitar */
while (fifo->empty) {
printf (" Consumidor: queue EMPTY.\n");
pthread_cond_wait (fifo->notEmpty, fifo->mut); }
/* Bloqueo del consumidor si el buffer está vacío, liberando el semáforo mut */
/*para que pueda entrar el productor. Continuará cuando el consumidor ejecute */
pthread_cond_signal (fifo->notFull);*/
}
```

```

    queueDel (fifo, &d); /* Quito elemento del buffer */
    pthread_mutex_unlock (fifo->mut); /* Libero el semáforo */
    pthread_cond_signal (fifo->notFull); /*Desbloqueo productor si está bloqueado*/
    printf (" Consumidor: Recibido %d.\n", d);
    usleep(200000);/* Duermo 200 mseg para permitir que el productor se active */
}
return (NULL); }

queue *queueInit (void) {
    queue *q;
    q = (queue *)malloc (sizeof (queue)); /* Creación del buffer */
    if (q == NULL) return (NULL);
    q->empty = 1; q->full = 0; q->head = 0; q->tail = 0;
    q->mut = (pthread_mutex_t *) malloc (sizeof (pthread_mutex_t));
    pthread_mutex_init (q->mut, NULL); /* Creación del semáforo */
    q->notFull = (pthread_cond_t *) malloc (sizeof (pthread_cond_t));
    pthread_cond_init (q->notFull, NULL); /* Creación de la variable condicional notFull*/
    q->notEmpty = (pthread_cond_t *) malloc (sizeof (pthread_cond_t));
    pthread_cond_init (q->notEmpty, NULL); /* Creación de la variable condicional notEmpty*/
    return (q); }

void queueDelete (queue *q) {
    pthread_mutex_destroy (q->mut); free (q->mut);
    pthread_cond_destroy (q->notFull); free (q->notFull);
    pthread_cond_destroy (q->notEmpty); free (q->notEmpty);
    free (q); }

void queueAdd (queue *q, int in) {
    q->buf[q->tail] = in; q->tail++;
    if (q->tail == QUEUESIZE) q->tail = 0;
    if (q->tail == q->head) q->full = 1;
    q->empty = 0;
    return; }

void queueDel (queue *q, int *out){
    *out = q->buf[q->head]; q->head++;
    if (q->head == QUEUESIZE) q->head = 0;
    if (q->head == q->tail) q->empty = 1;
    q->full = 0;
    return; }

```

1.5. OpenMP

El **OpenMP** (*Open-Multi Processing*) es una interfaz de programación de aplicaciones (API) con soporte multiplataforma para la programación en C/C++ y Fortran de procesos con uso de memoria compartida sobre plataformas Linux/Unix (y también Windows). Esta infraestructura se compone de un conjunto de directivas del compilador, rutinas de la biblioteca y variables de entorno que permiten aprovechar recursos compartidos en memoria y en tiempo de ejecución. Definido conjuntamente por un grupo de los principales fabricantes de *hardware* y *software*, OpenMP permite utilizar un modelo escalable y portátil de programación, que proporciona a los usuarios un interfaz simple y flexible para el desarrollo, sobre plataformas paralelas, de aplicaciones de escritorio hasta aplicaciones de altas prestaciones sobre superordenadores. Una aplicación construida con el modelo híbrido de la programación paralela puede ejecutarse en un ordenador utilizando tanto OpenMP como Message Passing Interface (MPI) [Bla].

OpenMP es una implementación multihilo, mediante la cual un hilo maestro divide la tareas sobre un conjunto de hilos trabajadores. Estos hilos se ejecutan simultáneamente y el entorno de ejecución realiza la asignación de estos a los diferentes procesadores de la arquitectura. La sección del código que está diseñada para funcionar en paralelo está marcada con una directiva de preprocesamiento que creará los hilos antes que la sección se ejecute. Cada hilo tendrá un identificador (*id*) que se obtendrá a partir de una función (`omp_get_thread_num()` en C/C++) y, después de la ejecución paralela, los hilos se unirán de nuevo en su ejecución sobre el hilo maestro, que continuará con la ejecución del programa. Por defecto, cada hilo ejecutará una sección paralela de código independiente pero se pueden declarar secciones de “trabajo compartido” para dividir una tarea entre los hilos, de manera que cada hilo ejecute parte del código asignado. De esta forma, es posible tener en un programa OpenMP paralelismo de datos y paralelismo de tareas conviviendo conjuntamente.

Los principales elementos de OpenMP son las sentencias para la creación de hilos, la distribución de carga de trabajo, la gestión de datos de entorno, la sincronización de hilos y las rutinas a nivel de usuario. OpenMP utiliza en C/C++ las directivas de preprocesamiento conocidas como *pragma* (`#pragma omp <resto del pragma>`) para diferentes construcciones. Así por ejemplo, `omp parallel` se utiliza para crear hilos adicionales para ejecutar el trabajo indicado en la sentencia paralela donde el proceso original es el hilo maestro (*id*=0). El conocido programa que imprime “Hello, world” utilizando OpenMP y multihilos es*:

```
#include <stdio.h>
#include <omp.h>

int main(int argc, char* argv[]){
    #pragma omp parallel
    printf("Hello, world.\n");
    return 0;}
```

```
*Compilad con gcc -fopenmp
-o hello hello.c (en
algunas distribuciones -debian
está instalada por defecto- se
debe tener instalada la
biblioteca GCC OpenMP
(GOMP) apt-get install
libgomp1
```

Donde ejecutará un thread por cada core disponible en la arquitectura. Para especificar *work-sharing constructs* se utiliza:

- **omp for** o **omp do**: reparte las iteraciones de un lazo en múltiples hilos.
- **sections**: asigna bloques de código independientes consecutivos a diferentes hilos.
- **single**: especifica que el bloque de código será ejecutado por un solo hilo con una sincronización (*barrier*) al final del mismo.
- **master**: similar a *single*, pero el código del bloque será ejecutado por el hilo maestro y no hay *barrier* implicado al final.

Por ejemplo, para inicializar los valores de un *array* en paralelo utilizando hilos para hacer una porción del trabajo (compilad, por ejemplo, con la orden `gcc -fopenmp -o init2 init2.c`):

```
#include <stdio.h>
#include <omp.h>
#define N 1000000
int main(int argc, char *argv[]) {
    float a[N]; long i;
    #pragma omp parallel for
    for (i=0;i<N;i++) a[i]= 2*i;
    return 0;
}
```

Si ejecutamos con el *pragma* y después lo comentamos y calculamos el tiempo de ejecución (`time ./init2`), vemos que el tiempo de ejecución pasa de 0.003 s a 0.007 s, lo que muestra la utilización de los múltiples *cores* (si se tienen disponibles) del procesador (si se está utilizando un máquina virtual, hay que verificar que esta tiene más de un *core* asignado, por ejemplo, en VBox en el apartado *System->Processor* de la MV).

Ya que OpenMP es un modelo de memoria compartida, muchas variables en el código son visibles para todos los hilos por defecto. Pero a veces es necesario tener variables privadas y pasar valores entre bloques secuenciales del código y bloques paralelos, por lo cual es necesario definir atributos a los datos (*data clauses*) que permitan diferentes situaciones:

- **shared**: los datos en la región paralela son compartidos y accesibles por todos los hilos simultáneamente.
- **private**: los datos en la región paralela son privados para cada hilo, y cada uno tiene una copia de ellos sobre una variable temporal.
- **default**: permite al programador definir cómo serán los datos dentro de la región paralela (`shared`, `private` o `none`).

Otro aspecto interesante de OpenMP son las directivas de sincronización:

- **critical section**: el código enmarcado será ejecutado por hilos, pero solo uno por vez (no habrá ejecución simultánea) y se mantiene la exclusión mutua.
- **atomic**: similar a `critical section`, pero avisa al compilador para que use instrucciones de hardware especiales de sincronización y así obtener mejores prestaciones.
- **ordered**: el bloque es ejecutado en el orden como si de un programa secuencial se tratara.
- **barrier**: cada hilo espera que los restantes hayan acabado su ejecución (implica sincronización de todos los hilos al final del código).

- **nowait**: especifica que los hilos que terminen el trabajo asignado pueden continuar.

Además, OpenMP provee de sentencias para la planificación (*scheduling*) del tipo `schedule(type, chunk)` (donde el tipo puede ser `static`, `dynamic` o `guided`) o proporciona control sobre las sentencias `if`, lo que permitirá definir si se paraleliza o no en función de si la expresión es verdadera o no. OpenMP también proporciona un conjunto de funciones de biblioteca, como por ejemplo:

- **omp_set_num_threads**: define el número de hilos a usar en la siguiente región paralela.
- **omp_get_num_threads**: obtiene el número de hilos que se están usando en una región paralela.
- **omp_get_max_threads**: obtiene la máxima cantidad posible de hilos.
- **omp_get_thread_num**: devuelve el número del hilo.
- **omp_get_num_procs**: devuelve el máximo número de procesadores que se pueden asignar al programa.
- **omp_in_parallel**: devuelve un valor distinto de cero si se ejecuta dentro de una región paralela.

Veremos a continuación algunos ejemplos simples (compilados con la instrucción `gcc -fopenmp -o out_file input_file.c`):

```
/* Programa simple multithreading con OpenMP */
#include <omp.h>
#include <stdio.h>

int main() {
    int iam = 0, np = 1;
    #pragma omp parallel private(iam, np)
    {
        #if defined (_OPENMP)
            np = omp_get_num_threads();
            iam = omp_get_thread_num();
        #endif
        printf("Hello from thread %d out of %d \n", iam, np);
    }
}

/* Programa simple con Threads anidados con OpenMP */
#include <omp.h>
#include <stdio.h>
main(){
    int x=0,nt,tid,ris;

    omp_set_nested(2);
    ris=omp_get_nested();
    if (ris) printf("Paralelismo anidado activo %d\n", ris);
    omp_set_num_threads(25);
    #pragma omp parallel private (nt,tid)
    {
        tid = omp_get_thread_num();
        printf("Thread %d\n",tid);
        nt = omp_get_num_threads();
        if (omp_get_thread_num()==1)
            printf("Número de Threads: %d\n",nt);
    }
}
```

```

/* Programa simple de integración con OpenMP */
#include <omp.h>
#include <stdio.h>
#define N 100
main() {
double local, pi=0.0, w; long i;
w = 1.0 / N;
#pragma omp parallel private(i, local)
{
#pragma omp single
pi = 0.0;
#pragma omp for reduction(+: pi)
for (i = 0; i < N; i++) {
local = (i + 0.5)*w;
pi = pi + 4.0/(1.0 + local*local);
printf ("Pi: %f\n",pi);
}
}
}

/* Programa simple de reducción con OpenMP */
#include <omp.h>
#include <stdio.h>
#define NUM_THREADS 2
main () {
int i; double ZZ, res=0.0;
omp_set_num_threads(NUM_THREADS);
#pragma omp parallel for reduction(+:res) private(ZZ)
for (i=0; i< 1000; i++) {
ZZ = i*i;
res = res + ZZ;
printf("ZZ: %f, res: %f\n", ZZ, res);}
}

```

Nota

Se recomienda ver también los ejemplos que se pueden encontrar en la referencia [opm].

1.6. MPI, *Message Passing Interface*

La definición de la API de MPI [Proc] ha sido el trabajo resultante del MPI Forum (MPIF), que es un consorcio de más de 40 organizaciones. MPI tiene influencias de diferentes arquitecturas, lenguajes y trabajos en el mundo del paralelismo, como por ejemplo: WRC (Ibm), Intel NX/2, Express, nCUBE, Vertex, p4, Parmac y contribuciones de ZipCode, Chimp, PVM, Chamaleon, PICL.

El principal objetivo de MPIF fue diseñar una API, sin relación particular con ningún compilador ni biblioteca, y que permitiera la comunicación eficiente (*memory-to-memory copy*), cómputo y comunicación concurrente y descarga de comunicación, siempre y cuando exista un coprocesador de comunicaciones. Además, se pedía que soportara el desarrollo en ambientes heterogéneos, con interfaz C y F77 (incluyendo C++, F90), donde la comunicación fuera fiable y los fallos, resueltos por el sistema. La API también debía tener interfaz para diferentes entornos, disponer una implementación adaptable a diferentes plataformas con cambios insignificantes y que no interfiera con el sistema operativo (*thread-safety*). Esta API fue diseñada especialmente para programadores que utilizaran el *Message Passing Paradigm* (MPP) en C y F77, para aprovechar su característica más relevante: la portabilidad. El MPP se puede ejecutar sobre máquinas multiprocesador, redes de estaciones de trabajo e incluso sobre

máquinas de memoria compartida. La primera versión del estándar fue MPI-1 (que si bien hay muchos desarrollos sobre esta versión, se considera en EOL), la versión MPI-2 incorporó un conjunto de mejoras, como creación de procesos dinámicos, *one-sided communication*, entrada/salida paralela entre otras, y finalmente la última versión, MPI-3 (considerada como una revisión mayor), incluye *nonblocking collective operations*, *one-sided operations* y soporte para Fortran 2008.[mpi3]

Muchos aspectos han sido diseñados para aprovechar las ventajas del hardware de comunicaciones sobre SPC (*scalable parallel computers*) y el estándar ha sido aceptado en forma mayoritaria por los fabricantes de hardware en paralelo y distribuido (SGI, SUN, Cray, HPConvex, IBM, etc.). Existen versiones libres (por ejemplo, Mpich, LAM/MPI y openMPI) que son totalmente compatibles con las implementaciones comerciales realizadas por los fabricantes de hardware e incluyen comunicaciones punto a punto, operaciones colectivas y grupos de procesos, contexto de comunicaciones y topología, soporte para F77 y C y un entorno de control, administración y *profiling*. [lam, ompi, Proc]

Pero existen también algunos puntos que pueden presentar algunos problemas en determinadas arquitecturas, como son la memoria compartida, la ejecución remota, las herramientas de construcción de programas, la depuración, el control de hilos, la administración de tareas y las funciones de entrada/salida concurrentes (la mayor parte de estos problemas de falta de herramientas están resueltos a partir de la versión 2 de la API -MPI2-). Una de los problemas de MPI1, al no tener creación dinámica de procesos, es que solo soporta modelos de programación MIMD (*Multiple Instruction Multiple Data*) y comunicándose vía llamadas MPI. A partir de MPI-2 y con la ventajas de la creación dinámica de procesos, ya se pueden implementar diferentes paradigmas de programación como *master-worker/farmer-tasks*, *divide & conquer*, paralelismo especulativo, etc. (o al menos sin tanta complejidad y mayor eficiencia en la utilización de los recursos).

Para la instalación de MPI se recomienda utilizar la distribución (en algunos casos la compilación puede ser compleja debido a las dependencias de otros paquetes que puede necesitar. Debian incluye la versión OpenMPI (sobre Debian 8.5 es versión 2, pero se pueden bajar los fuentes y compilarlos) y Mpich2 (Mpich3 disponible en <http://www.mpich.org/downloads/>). La mejor elección será OpenMPI, ya que combina las tecnologías y los recursos de varios otros proyectos (FT-MPI, LA-MPI, LAM/MPI y PACX-MPI) y soporta totalmente el estándar MPI-2 (y desde la versión 1.75 soporta la versión MPI3). Entre otras características de OpenMPI tenemos: es conforme a MPI-2/3, *thread safety & concurrency*, creación dinámica de procesos, alto rendimiento y gestión de trabajos tolerantes a fallos, instrumentación en tiempo de ejecución, *job schedulers*, etc. Para ello se deben instalar los paquetes `openmpi-dev`, `openmpi-bin`, `openmpi-common` y `openmpi-doc`. Además, Debian 8.5 incluye otra implementación de MPI llamada LAM (paquetes `lam*`). Se debe considerar que si bien las implementaciones son equivalentes desde el punto de vista de MPI,

tienen diferencias en cuanto a la gestión y procedimientos de compilación/ejecución/gestión.

1.6.1. Configuración de un conjunto de máquinas para hacer un clúster adaptado a OpenMPI

Para la configuración de un conjunto de máquinas para hacer un clúster adaptado a OpenMPI [tec], se han de seguir los siguientes pasos:

- 1) Hay que tener las máquinas “visibles” (por ejemplo a través de un `ping`) a través de TCP/IP (IP pública/privada).
- 2) Es recomendable que todas las máquinas tengan la misma arquitectura de procesador, así es más fácil distribuir el código, con versiones similares de Linux (a ser posible con el mismo tipo de distribución).
- 3) Se recomienda tener NIS o si no se debe generar un mismo usuario (por ejemplo, `mpiuser`) en todas las máquinas y el mismo directorio `$HOME` montado por NFS.
- 4) Nosotros llamaremos a las máquinas como `slave1`, `slave2`, etc., (ya que luego resultará más fácil hacer las configuraciones) pero se puede llamar a las máquinas como cada uno prefiera.
- 5) Una de las máquinas será el maestro y las restantes, `slaveX`.
- 6) Se debe instalar en todos los nodos (supongamos que tenemos Debian): `openmpi-bin`, `openmpi-common`, `libopenmpi-dev`. Hay que verificar que en todas las distribuciones se trabaja con la misma versión de OpenMPI.
- 7) En Debian los ejecutables están en `/usr/bin` pero si están en un *path* diferente, deberá agregarse a `mpiuser` y también verificar que `LD_LIBRARY_PATH` apunta a `/usr/lib`.
- 8) En cada nodo esclavo debe instalarse el SSH server (instalad el paquete `openssh-server`) y sobre el maestro, el cliente (paquete `openssh-client`).
- 9) Se crearán las claves públicas y privadas haciendo `ssh-keygen -t dsa` (tampoco hay problema si se usa `rsa`) y copiar a cada nodo con `ssh-copy-id` para este usuario (solo se debe hacer en un nodo, ya que, como tendremos el directorio `$HOME` compartido por NFS para todos los nodos, con una copia basta).
- 10) Si no se comparte el directorio hay que asegurar que cada esclavo conoce que el usuario `mpiuser` se puede conectar sin `passwd`, por ejemplo haciendo:
`ssh slave1`.
- 11) Se debe configurar la lista de las máquinas sobre las cuales se ejecutará el programa, por ejemplo `/home/mpiuser/.mpi_hostfile` y con el siguiente contenido:

```
# The Hostfile for Open MPI
# The master node, slots=2 is used because it is a dual-processor machine.
  localhost slots=2
# The following slave nodes are single processor machines:
  slave1
  slave2
  slave3
```

12) OpenMPI permite utilizar diferentes lenguajes, pero aquí utilizaremos C. Para ello hay que ejecutar sobre el maestro, por ejemplo,

```
mpicc testprogram.c -o testprogram
```

Si se desea ver que incorpora mpicc, se puede hacer `mpicc -showme`.

13) Para ejecutar en local podríamos hacer `mpirun -np 2 ./testprogram` y para ejecutar sobre los nodos remotos (por ejemplo 5 procesos)

```
mpirun -np 2 -hostfile ./mpi_hostfile ./testprogram
```

Es importante notar que `np` es el número de procesos o procesadores en que se ejecutará el programa y se puede poner el número que se desee, ya que OpenMPI intentará distribuir los procesos de forma equilibrada entre todas las máquinas. Si hay más procesos que procesadores, OpenMPI/Mpich utilizará las características de intercambio de tareas de GNU/Linux para simular la ejecución paralela. A continuación se verán dos ejemplos: `Srtest` es un programa simple para establecer comunicaciones entre procesos punto a punto, y `cpi` calcula el valor del número π de forma distribuida (por integración).

```
/* Srtest Program */
#include "mpi.h"
#include <stdio.h>
#include <string.h>
#define BUFLLEN 512

int main(int argc, char *argv[]){
  int myid, numprocs, next, namelen;
  char buffer[BUFLLEN], processor_name[MPI_MAX_PROCESSOR_NAME];
  MPI_Status status;
  MPI_Init(&argc,&argv); /* Debe ponerse antes de otras llamadas MPI, siempre */
  MPI_Comm_size(MPI_COMM_WORLD,&numprocs);
  MPI_Comm_rank(MPI_COMM_WORLD,&myid); /*Integra el proceso en un grupo de comunicaciones*/
  MPI_Get_processor_name(processor_name,&namelen); /*Obtiene el nombre del procesador*/

  fprintf(stderr,"Proceso %d sobre %s\n", myid, processor_name);
  fprintf(stderr,"Proceso %d de %d\n", myid, numprocs);
  strcpy(buffer,"hello there");
  if (myid == numprocs-1) next = 0;
  else next = myid+1;

  if (myid == 0) { /*Si es el inicial, envía string de buffer*/
    printf("%d sending '%s' \n",myid,buffer);fflush(stdout);
    MPI_Send(buffer, strlen(buffer)+1, MPI_CHAR, next, 99, MPI_COMM_WORLD);
    /*Blocking Send, 1:buffer, 2:size, 3:tipo, 4:destino, 5:tag, 6:contexto*/
    printf("%d receiving \n",myid);fflush(stdout);
    MPI_Recv(buffer, BUFLLEN, MPI_CHAR, MPI_ANY_SOURCE, 99, MPI_COMM_WORLD,&status);
    printf("%d received '%s' \n",myid,buffer);fflush(stdout);
    /* mpdprintf(001,"%d receiving \n",myid); */
  }
  else {
    printf("%d receiving \n",myid);fflush(stdout);
    MPI_Recv(buffer, BUFLLEN, MPI_CHAR, MPI_ANY_SOURCE, 99, MPI_COMM_WORLD,&status);
```

```

    /* Blocking Recv, 1:buffer, 2:size, 3:tipo, 4:fuente, 5:tag, 6:contexto, 7:status*/
    printf("%d received '%s' \n",myid,buffer);fflush(stdout);
    /* mpdprintf(001,"%d receiving \n",myid); */
    MPI_Send(buffer, strlen(buffer)+1, MPI_CHAR, next, 99, MPI_COMM_WORLD);
    printf("%d sent '%s' \n",myid,buffer);fflush(stdout);
    }
    MPI_Barrier(MPI_COMM_WORLD); /*Sincroniza todos los procesos*/
    MPI_Finalize(); /*Libera los recursos y termina*/
    return (0);
}

/* CPI Program */
#include "mpi.h"
#include <stdio.h>
#include <math.h>
double f( double );
double f( double a) { return (4.0 / (1.0 + a*a)); }
int main( int argc, char *argv[] ) {
    int done = 0, n, myid, numprocs, i;
    double PI25DT = 3.141592653589793238462643;
    double mypi, pi, h, sum, x; double startwtime = 0.0, endwtime;
    int namelen; char processor_name[MPI_MAX_PROCESSOR_NAME];
    MPI_Init(&argc,&argv);
    MPI_Comm_size(MPI_COMM_WORLD,&numprocs); /*Indica el número de procesos en el grupo*/
    MPI_Comm_rank(MPI_COMM_WORLD,&myid); /*Id del proceso*/
    MPI_Get_processor_name(processor_name,&namelen); /*Nombre del proceso*/
    fprintf(stderr, "Proceso %d sobre %s\n", myid, processor_name);
    n = 0;
    while (!done) {
        if (myid ==0) { /*Si es el primero...*/
            if (n ==0) n = 100; else n = 0;
            startwtime = MPI_Wtime();} /* Time Clock */
        MPI_Bcast(&n, 1, MPI_INT, 0, MPI_COMM_WORLD); /*Broadcast al resto*/
        /*Envía desde el 4 arg. a todos los procesos del grupo Los restantes que no son 0
        copiarán el buffer desde 4 o arg -proceso 0-*/
        /*1:buffer, 2:size, 3:tipo, 5:grupo */
        if (n == 0) done = 1;
        else {h = 1.0 / (double) n;
            sum = 0.0;
            for (i = myid + 1; i <= n; i += numprocs) {
                x = h * ((double)i - 0.5); sum += f(x); }
            mypi = h * sum;
            MPI_Reduce(&mypi, &pi, 1, MPI_DOUBLE, MPI_SUM, 0,MPI_COMM_WORLD);
            /* Combina los elementos del Send Buffer de cada proceso del grupo usando la
            operación MPI_SUM y retorna el resultado en el Recv Buffer. Debe ser llamada
            por todos los procesos del grupo usando los mismos argumentos*/
            /*1:sendbuffer, 2:recvbuffer, 3:size, 4:tipo, 5:oper, 6:root, 7:contexto*/
            if (myid == 0){ /*solo el P0 imprime el resultado*/
                printf("Pi es aproximadamente %.16f, el error es %.16f\n", pi,fabs(pi - PI25DT));
                endwtime = MPI_Wtime();
                printf("Tiempo de ejecución = %f\n", endwtime-startwtime); }
        }
    }
    MPI_Finalize(); /*Libera recursos y termina*/
    return 0;
}

```

Para visualizar la ejecución de un código paralelo/distribuido en MPI existe una aplicación llamada XMPI (en Debian `xmpi`) que permite 'ver', durante la ejecución, el estado de la ejecución de las aplicaciones sobre MPI, pero está vinculada al paquete LAM/MPI. Para visualizar y analizar código sobre OpenMPI, se debería (y es recomendable) bajar y compilar el código de MPE (<http://www.mcs.anl.gov/research/projects/perfvis/software/MPE/>) o TAU (<http://www.cs.uoregon.edu/research/tau/home.php>) que son herramientas de *profiling* muy potentes y que no requieren gran trabajo ni dedicación para ponerlas en marcha.

1.7. Rocks Cluster

Rocks Cluster es una distribución de Linux para clústers de computadores de alto rendimiento. Las versiones actuales de Rocks Cluster están basadas en CentOS (CentOS 6.6 a mayo 2015) y, como instalador, Anaconda con ciertas modificaciones, que simplifica la instalación 'en masa' en muchos ordenadores. Rocks Cluster incluye muchas herramientas (tales como MPI) que no forman parte de CentOS pero son los componentes que transforman un grupo de ordenadores en un clúster. Las instalaciones pueden personalizarse con paquetes de software adicionales llamados *rolls*. Los *rolls* extienden el sistema integrando automáticamente los mecanismos de gestión y empaquetamiento usados por el software básico, y simplifican ampliamente la instalación y configuración de un gran número de computadores. Se han creado una gran cantidad de *rolls*, como por ejemplo SGE *roll*, Cónдор *roll*, Xen *roll*, el Java *roll*, Ganglia *roll*, etc. (http://www.rocksclusters.org/wordpress/?page_id=4). Rocks Cluster es una distribución altamente empleada en el ámbito de clústers, por su facilidad de instalación e incorporación de nuevos nodos y por la gran cantidad de programas para el mantenimiento y monitorización del clúster.

Enlaces de interés

Para una lista detallada de las herramientas incluidas en Rocks Cluster, consultad: <http://www.rocksclusters.org/roll-documentation/base/5.5/>.

Las principales características de Rocks Cluster son:

- 1) Facilidad de instalación, ya que solo es necesario completar la instalación de un nodo llamado *nodo maestro (frontend)*, el resto se instala con Avalache, que es un programa P2P que lo hace de forma automática y evita tener que instalar los nodos uno a uno.
- 2) Disponibilidad de programas (conjunto muy amplio de programas que no es necesario compilar y transportar) y facilidad de mantenimiento (solo se debe mantener el nodo maestro).
- 3) Diseño modular y eficiente, pensado para minimizar el tráfico de red y utilizar el disco duro propio de cada nodo para solo compartir la información mínima e imprescindible.

La instalación se puede seguir paso a paso desde el sitio web de la distribución [Rock] y los autores garantizan que no se tarda más de una hora para una instalación básica. Rocks Cluster permite, en la etapa de instalación, diferentes módulos de software, los *rolls*, que contienen todo lo necesario para realizar la instalación y la configuración de sistema con estas nuevas 'adiciones' de forma automática y, además, decidir en el *frontend* cómo será la instalación en los nodos esclavos, qué *rolls* estarán activos y qué arquitectura se usará. Para el mantenimiento, incluye un sistema de copia de respaldo del estado, llamado *Roll Restore*. Este *roll* guarda los archivos de configuración y *scripts* (e incluso se pueden añadir los archivos que se desee).

1.7.1. Guía rápida de instalación

Este es un resumen breve de la instalación propuesta en [Rock] y se parte de la base que el *frontend* tiene (mínimo) 30 GB de disco duro, 1 GB de RAM, arquitectura x86-64 y 2 interfaces de red (eth0 para la comunicación con internet y eth1 para la red interna); para los nodos 30 GB de disco duro, 512 MB de RAM y 1 interfaz de red (red interna). Después de obtener los discos *Kernel/Boot Roll*, *Base Roll*, *OS Roll CD1/2* (o en su defecto DVD equivalente), insertamos *kernel boot*, seguimos los siguientes pasos:

- 1) Arrancamos el *frontend* y veremos una pantalla en la cual introducimos `build` y nos preguntará la configuración de la red (IPV4 o IPV6).
- 2) El siguiente paso es seleccionar los `rolls` (por ejemplo seleccionando los “CD/DVD-based Roll” y vamos introduciendo los siguientes *rolls*) y marcar en las sucesivas pantallas cuáles son los que deseamos.
- 3) La siguiente pantalla nos pedirá información sobre el clúster (es importante definir bien el nombre de la máquina, *Fully-Qualified Host Name*, ya que en caso contrario fallará la conexión con diferentes servicios) y también información para la red privada que conectará el *frontend* con los nodos y la red pública (por ejemplo, la que conectará el *frontend* con Internet) así como DNS y pasarelas.
- 4) A continuación se solicitará la contraseña para el *root*, la configuración del servicio de tiempo y el particionado del disco (se recomienda escoger “auto”).
- 5) Después de formatear los discos, solicitará los CD de *rolls* indicados e instalará los paquetes y hará un *reboot* del *frontend*.
- 6) Para instalar los nodos se debe entrar como *root* en el *frontend* y ejecutar `insert-ethers`, que capturarán las peticiones de DHCP de los nodos y los agregará a la base de datos del *frontend*, y seleccionar la opción *Compute* (por defecto, consultad la documentación para las otras opciones).
- 7) Encendemos el primer nodo y en el *boot order* de la BIOS generalmente se tendrá CD, PXE (Network Boot), Hard Disk, (si el ordenador no soporta PXE, entonces haced el *boot* del nodo con el *Kernel Roll CD*). En el *frontend* se verá la petición y el sistema lo agregará como `compute-0-0` y comenzará la descarga e instalación de los archivos. Si la instalación falla, se deberán reiniciar los servicios `httpd`, `mysqld` y `autofs` en el *frontend*.
- 8) A partir de este punto se puede monitorizar la instalación ejecutando la instrucción `rocks-console`, por ejemplo con `rocks-console compute-0-0`. Después de que se haya instalado todo, se puede salir de `insert-ethers` pulsando la tecla F8. Si se dispone de dos *racks*, después de instalado el primero se puede comenzar el segundo haciendo `insert-ethers -cabinet=1` los cuales recibirán los nombres `compute-1-0`, `compute-1-1`, etc.

Enlaces de interés

Se puede descargar los discos desde:
http://www.rocksclusters.org/wordpress/?page_id=80

9) A partir de la información generada en consola por `rocks list host`, generaremos la información para el archivo `machines.conf`, que tendrá un aspecto como:

```
nfeum slot=2
compute-0-0 slots=2
compute-0-1 slots=2
compute-0-2 slots=2
compute-0-3 slots=2
```

y que luego deberemos ejecutar con

```
mpirun -np 10 -hostfile ./machines.conf ./mpi_program_to_execute.
```

Un ejercicio interesante para hacer es instalar un clúster funcional virtualizado que obviamente no tendrá grandes prestaciones, pero que es útil para aprender todos los conceptos que intervienen y como prueba de concepto. Para ello utilizaremos VirtualBox sobre un procesador de 64 bits (las últimas versiones de Rocks solo están para 64 bits) y que tengan las VTx/AMD-V del procesador. El procedimiento sería:

1) Descargar la iso del DVD desde <http://www.rocksclusters.org>.

2) Crear una máquina virtual (por ejemplo, RocksFE) con 28 GB de disco y dos adaptadores de red (uno conectado a la red externa `-eth0-` y otro a una red interna `-eth1-` que la llamaremos *rocksnet*). Esta máquina será el *Frontend* de nuestro clúster. Luego insertamos el DVD sobre el DVD virtual y arrancamos la máquina virtual.

3) Introducir en la primera página,

```
build ksdevice=eth0 ip=IP-red-externa gateway=gateway-externo netmask=255.255.255.0.
```

4) Seguir los pasos antes indicados, seleccionando la media desde el DVD y los *rolls* (mínimo: os, kernel, base) que se deseen y contestando a las preguntas de la red y DNS.

5) Después de la instalación podremos entrar y ver el estado de la máquina.

6) Crear una segunda máquina (por ejemplo, RocksCN) similar a la primera pero con solo una tarjeta de red (que configuraremos como red interna y la conectaremos a *rocksnet*) y configurar el arranque (en *System*) como *Network*.

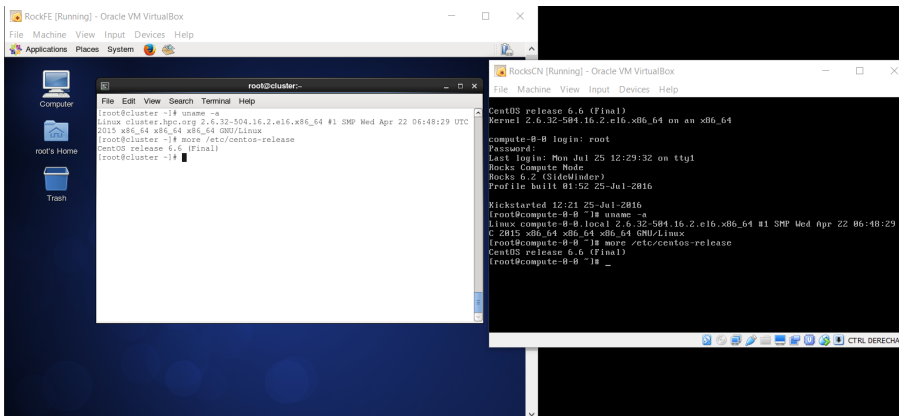
7) Sobre el *Frontend* y en una terminal ejecutar `insert-ethers` seleccionando *Compute* y poner en marcha la máquina virtual. Veremos después de un momento que la reconoce dentro de la ventana como `compute-0-0` y arranca y se instala el SO sobre RocksCN.

8) Cuando finalice la instalación sobre RocksCN debemos apagar la máquina y activar el *boot* desde el disco duro y volver a arrancar la máquina.

9) Para finalizar veremos que está dada de alta en `/etc/hosts` sobre el *Frontend* y que podemos acceder a ella a través de un `ssh compute-0-0`.

La figura 1 muestra la ejecución sobre VirtualBox de la prueba de concepto del *Frontend* (con GUI) y la de un nodo (en modo texto).

Figura 1



1.8. Despliegue automático

El despliegue automático (*Automatic Provisioning*) es una de las grandes preocupaciones y supone un consumo de tiempo para los administradores, sobre todo cuando se tiene que desplegar un conjunto importante de máquinas. Esto explica el esfuerzo de algunos proyectos para simplificar esta tarea, por ejemplo:

- **Kickstart [KS]**, para máquinas RH, Fedora y CentOS, básicamente.
- **FAI [FG]**, para sistemas Debian pero que pueden instalar otro tipo de distribuciones.
- **Cobbler [Co]**, que si bien es para sistemas de la rama de RH hay experiencias sobre otros sistemas (internamente utiliza KickStart).
- **Spacewalk [SW]**, para servidores RH pero que pueden instalar otros tipos de distribuciones.
- **OpenQRM [OQRM]**, para diferentes servidores (*versión community*).

En nuestro caso, para seguir con la misma distribución, veremos el caso de FAI sobre un sistema Debian.

1.8.1. FAI (Fully Automatic Installation)

FAI es una herramienta de instalación automatizada para desplegar Linux en un clúster o en un conjunto de máquinas en forma desatendida. Es equiva-

lente a *kickstart* de RH, o Alice de SuSE. FAI puede instalar Debian, Ubuntu y RPMs de distribuciones Linux. Sus ventajas radican en que mediante esta herramienta se puede desplegar Linux sobre un nodo (físico o virtual) simplemente haciendo que arranque por PXE y quede preparado para trabajar y totalmente configurado (cuando se dice uno, se podría decir 100 con el mismo esfuerzo por parte del administrador) y sin ninguna interacción por medio. Por lo tanto, es un método escalable para la instalación y actualización de un clúster Beowulf o una red de estaciones de trabajo sin supervisión y con poco esfuerzo. FAI utiliza la distribución Debian, una colección de scripts (su mayoría en Perl) y *cfengine* y *Preseeding d-i* para el proceso de instalación y cambios en los archivos de configuración.

Es una herramienta pensada para administradores de sistemas que deben instalar Debian en decenas o cientos de ordenadores y puede utilizarse además como herramienta de instalación de propósito general para instalar (o actualizar) un clúster de cómputo HPC, de servidores web, o un pool de bases de datos y configurar la gestión de la red y los recursos en forma desatendida. Esta herramienta permite tratar (a través de un concepto que incorpora llamado clases) con un hardware diferente y diferentes requisitos de instalación en función de la máquina y características que se disponga en su preconfiguración, permitiendo hacer despliegues masivos eficientes y sin intervención de administrador (una vez que la herramienta haya sido configurada adecuadamente). [fai1, fai2, fai3]

1.9. Guía de instalación de FAI

En este apartado veremos una guía de la instalación básica sobre máquinas virtuales (en Virtualbox) para desplegar otras máquinas virtuales, pero con mínimas intervenciones se puede adaptar/ampliar a las necesidades del entorno. FAI es una herramienta que no cuenta con *daemons/DB*; solo consiste en *scripts* (consultar las referencias indicadas).

La forma más fácil de hacer una prueba de concepto es a través de las instalaciones modulares que provee el proyecto (donde se puede usar una máquina virtual, tal como haremos nosotros). Para ello se utiliza una ISO en un USB o en un CD descargable desde <http://fai-project.org/fai-cd/>. Las imágenes auto-instalables disponibles son:

- 1) **FAI ISO large** (953 MB): imagen por defecto que incluye XFCE y GNOME (pero solo la base para CentOS 7 y Ubuntu 16.04).
- 2) **FAI ISO small** (585 MB): imagen que incluye solo XFCE. Los otros paquetes son descargados desde Internet.
- 3) **FAI ISO Ubuntu** (1.1 GB): todo Ubuntu 16.04 LTS.

4) Autodiscover CD (19 MB): imagen que permite solo el arranque y que busca un servidor FAI.

Todas las imágenes anteriores utilizan Debian Jessie, FAI 5.1.2, kernel 3.16, arquitectura amd64. Además está disponible la imagen **FAI ISO small** (stretch) (585 MB), que es la versión FAI ISO small pero utilizando Debian stretch.

Para poner las imágenes sobre un USB se puede utilizar el comando

```
dd if=faicd64-large_5.1.2.iso of=/dev/sdX
```

donde se debe reemplazar */dev/sdX* con el dispositivo del USB (por ejemplo, */dev/sdg*).

Con estas imágenes se puede instalar una Debian 8 (XFCE/GNOME), Ubuntu 16.04 o CentOS 7 sin necesidad de tener un servidor (tanto el usuario *fai*, *demo* como *root* tiene por *passwd fai*). Se debe tener cuidado, ya que el sistema será instalado sobre el primer disco y todos los datos serán borrados. Para instalarlo sobre una máquina Virtualbox solo se debe crear una máquina con un disco de, por ejemplo, 12 GB y cargar la imagen en el CDRom virtual y arrancar la máquina seleccionando solo lo que se desea al inicio. Luego la instalación es desatendida y tendremos la máquina instalada sin intervención.

1.9.1. Instalación del servidor

Si se desea instalar un servidor para aprovisionar otras máquinas y adaptar las configuraciones a las necesidades de la instalación, se puede seguir la guía del proyecto*. No obstante, sobre Jessie presenta algunos problemas con algunos paquetes y se debe invertir gran cantidad de tiempo para compatibilizar todos los requerimientos. Es por ello más fácil utilizar la imagen FAI ISO Large, ya que permite instalar el servidor en un tiempo breve y luego adaptarlo a las necesidades de la instalación y de los clientes. Para ello en esta prueba de concepto haremos:

*<http://fai-project.org/fai-guide/>

- 1) Descargar la imagen: http://fai-project.org/fai-cd/faicd64-large_5.1.2.iso.
- 2) Crear en VirtualBox una máquina virtual con 16 Gb de disco y una tarjeta de red con NAT (luego lo cambiaremos) y asociar esta imagen al DVD virtual.
- 3) Arrancar la máquina virtual y seleccionar *FAI server Installation using external DHCP*. Introducir como usuario *fai* e *install* como *passwd* y comenzará la instalación.
- 4) Cuando finalice y se reinicie la máquina, seleccionar *Boot OS of First Partition (1st disc)* (o en su defecto quitar la imagen del DVD) y comenzará la instalación del servidor durante unos minutos.

5) Cuando finalice, iniciar la sesión como **root** y *passwd fai* y hacer los siguientes ajustes:

- a) Cambiar el teclado en */etc/default/keyboard* en *XKBLAYOUT="es"* y ejecutar `service keyboard-setup restart`.
- b) Editar */etc/hosts* y cambiar la IP del *faiserver* a *192.168.33.1*, ya que queremos que el aprovisionamiento lo realice por la otra red que agregaremos.
- c) Editar el archivo */etc/network/interfaces* y agregar la configuración para *eth1* como:

```
auto lo eth0 eth1
...
iface eth1 inet static
    address 192.168.33.1
    netmask 255.255.255.0
```

- d) Modificar */etc/exports* y donde dice *10.0.2.0/24* cambiar por *192.168.33.0/24* para que los clientes puedan montar todos los recursos del servidor.
- e) Cambiar en */etc/dhcp/dhcpd.conf* la IP de la opción de *domain-name-servers* (por ejemplo a *192.168.33.1*).
- f) Reiniciar la máquina y antes que arranque pararla y agregar un segundo adaptador (*eth1*) conectado a la red interna (será la misma red a la que luego conectaremos los clientes). A continuación quitar el DVD (si no lo hemos hecho antes) para que arranque directamente del disco duro. Iniciar la sesión como *root* y verificar que todo está configurado correctamente y que tenemos los servicios *tftp* y *dhcpd* funcionando (por ejemplo, con `systemctl status`).
- g) **Consideraciones que debemos tener en cuenta.** Cuando se crea la MV para los clientes hay que escoger el tipo *Linux64*, ya que sino durante el arranque del cliente dará un error porque intentará instalar un *Linux32*. Cuando se ha instalado satisfactoriamente una máquina, queda registro en el servidor y si queremos hacer otra prueba sobre la misma MV deberemos cambiarle la MAC para reutilizar la misma máquina ya que si no dará un error porque esta consta como instalada.

1.9.2. Configuración del nodo

Nuestra prueba de concepto la haremos para un nodo en Virtualbox. Primero es necesario instalar las expansiones de la versión que tenemos instalada. Para ello deberemos ir a <https://www.virtualbox.org> y para la distribución instalada descargar e instalar el paquete *VirtualBox Extension Pack All supported platforms* (esto nos permitirá tener un dispositivo que haga *boot* por PXE). Podremos verificar que están instaladas en el menú de VirtualBox ->File ->Preferences ->Extensions. Luego deberemos crear una nueva máquina virtual con un disco vacío (por ejemplo, 8 GB si es un cliente texto o 14-16 GB si es un cliente XFE/Gnome) y una tarjeta de red conectada a red interna y a la misma red en la que está el servidor en su tarjeta *eth1*. Luego se deberá seleccionar en *System* de esta máquina la opción de ***boot order network*** en primer lugar.

Finalmente arrancaremos la máquina y veremos que recibe IP (verificar los posibles errores en `/var/log/messages|syslog`) y que comienza bajándose el `initrd-img` y luego el kernel. Nos solicitará qué tipo de cliente queremos instalar y tras unos minutos el sistema estará instalado. Cuando finalice la instalación deberemos apagar la máquina y cambiar el orden de `boot` en `System` para que vuelva a arrancar en primer lugar desde el disco duro.

La figura 2 muestra las pruebas realizadas con la selección de un cliente simple y la figura 3 con un cliente Gnome (a la izquierda se visualiza el servidor de aprovisionamiento).

Figura 2

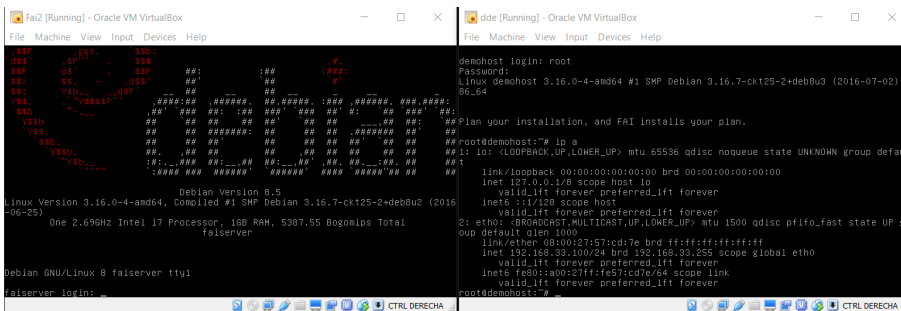
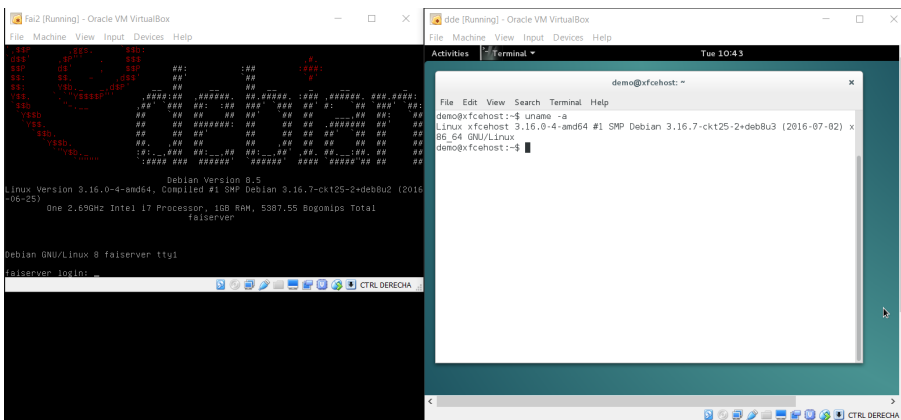


Figura 3



1.10. Logs

Linux mantiene un conjunto de registros llamados `system logs` o `logs` simplemente, que permiten analizar qué ha pasado y cuándo en el sistema, a través de los eventos que recoge del propio `kernel` o de las diferentes aplicaciones y casi todas las distribuciones se encuentran en `/var/log`. Probablemente los dos archivos más importantes (y que en mayor o menor presencia están en todas las distribuciones) son `/var/log/messages` y `/var/log/syslog` los que tienen registros (Ascii) de eventos tales como errores del sistema, (re)inicios/apagados, errores de aplicaciones, advertencias, etc. Existe un comando, `dmesg`, que permite además ver los mensajes de inicio (en algunas distribuciones son los que aparecen en la consola o con la tecla Esc en el modo gráfico de arranque, o Crtl+F8 en otras distribuciones) para visualizar los pasos seguidos durante

el arranque (teniendo en cuenta que puede ser muy extenso, se recomienda ejecutar `dmesg | more`).

Por lo cual, el sistema nos permitirá analizar qué ha pasado y obtener las causas que han producido este registro y dónde/cuándo. El sistema incluido en Debian es `rsyslog` (<http://www.rsyslog.com/>) con un servicio (a través del *daemon* `rsyslogd` y que reemplaza al original `syslog`). Es un sistema muy potente y versátil y su configuración es través del archivo `/etc/rsyslog.conf` o archivos en el directorio `/etc/rsyslog.d`. Mirar la documentación sobre su configuración (`man rsyslog.conf` o en la página indicada) pero en resumen incluye una serie de directivas, filtros, templates y reglas que permiten cambiar el comportamiento de los logs del sistema. Por ejemplo, las reglas son de tipo **recurso.nivel acción** pero se puede combinar más de un recurso separado por `'`, o diferentes niveles, o todos `'*` o negarlo `'!`, y si ponemos `'='` delante del nivel, indica solo ese nivel y no todos los superiores, que es lo que ocurre cuando solo se indica un nivel (los niveles pueden ser de menor a mayor importancia `debug`, `info`, `notice`, `warning`, `err`, `crit`, `alert`, `emerg`). Por ejemplo `*.warning /var/log/messages` indicará que todos los mensajes de `warning, err, crit, alert, emerg` de cualquier recurso vayan a parar a este archivo. Se puede incluir también un `'-'` delante del nombre del fichero, que indica que no se sincronice el fichero después de cada escritura para mejorar las prestaciones del sistema y reducir la carga.

Un aspecto interesante de los logs es que los podemos centralizar desde las diferentes máquinas de nuestra infraestructura en un determinado servidor para no tener que conectarnos si queremos 'ver' qué está ocurriendo a nivel de logs en esas máquinas (sobre todo si el número de máquinas es elevado). Para ello, debemos en cada cliente modificar el archivo `/etc/rsyslog.conf` (donde la IP será la del servidor que recogerá los logs):

```
# Provides TCP forwarding.
*. * @192.168.168.254:514
```

En el servidor deberemos quitar el comentario (#) de los módulos de recepción por TCP en `/etc/rsyslog.conf`:

```
# Provides TCP syslog reception
$ModLoad imtcp
$InputTCPServerRun 514
```

Después de hacer las modificaciones de cada archivo no hay que olvidar hacer un `/etc/init.d/rsyslog restart`. A partir de ese momento, podemos mirar el `/var/log/syslog` del servidor, por ejemplo, y veremos los registros marcados con el nombre (o IP) de la máquina donde se ha generado el log.

Existe un paquete llamado `syslog-ng` que presenta características avanzadas (por ejemplo, cifrado o filtros basados en contenidos) u otras equivalentes con `rsyslog`.*

*Más información en la siguiente página web:
<http://www.balabit.com/network-security/syslog-ng>.

1.11. Visualización de logs

Un aspecto importante cuando se tienen los *logs* –registros– (incluidos de otras máquinas) es cómo visualizarlos y hacer una búsqueda selectiva para prevenir situaciones o saber el estado del sistema, así como para generar informes o visualizar de forma integrada diferentes fuentes de información. Existen herramientas que nos permiten visualizar (*log browsing*) y buscar de forma simple como por ejemplo:

- 1) **glogg** (<http://glogg.bonnefon.org/>),
- 2) **kssystemlog** (<https://www.kde.org/applications/system/kssystemlog/>),
- 3) **Inav** (<http://lnav.org/>),
- 4) **multitail** (<https://www.vanheusden.com/multitail/>), que se encuentran en el repositorio de Debian, o
- 5) **logio** (<http://logio.org/>), cuando la visualización se debe hacer remotamente.

Cuando se desea realizar una supervisión y una gestión de forma eficiente y escalable, con capacidades de pre-procesado, filtrado, indexación y almacenamiento de gran cantidad de eventos (para su análisis *on-line* o posterior), es necesario contar con herramientas adecuadas a estas funciones. En esta categoría se pueden mencionar herramientas como:

- 1) **LogAnalyzer** (<http://loganalyzer.adiscon.com/>). Es una herramienta fácil de configurar/usar con una interfaz versátil para buscar y analizar datos de eventos de un sistema de diferentes fuentes de red (*syslog/rsyslog*). Por ejemplo, esta aplicación recoge los datos producidos por *rsyslog* (estándar en Debian) y los visualiza a través de un servidor Web y un navegador. La aplicación está escrita en PHP y permite la conexión a MySQL (u otras bases de datos SQL-NoSQL) cuando son necesarias opciones avanzadas o un procesamiento eficiente.
- 2) **Logstash** (<https://www.elastic.co/products/logstash>). Esta herramienta *open source* permite recolectar, analizar y almacenar los archivos de *logs* para su análisis posterior. Incorpora un conjunto de *plugins* para facilitar diferentes formatos, decodificación y reglas de filtrado (por ejemplo, *logs* de S3, RabbitMQ, Syslog, collectd, sockets TCP/UDP).
- 3) **Collectd** (<https://collectd.org/>). Si bien es más una herramienta de monitorización, permite a través de diversos *plugins* de entrada/salida procesar diferentes archivos de *logs* (locales o remotos). Está incluida en el repositorio de Debian.
- 4) **Fluentd** (<http://www.fluentd.org/> <https://github.com/fluent/fluentd>). Esta es una herramienta que presenta una capa unificada de agregación de *logs* y que permite el procesamiento *in-stream* para una variedad de flujo de datos y

archivos de *logs*. Su arquitectura es extensible y cuenta con más de 300 *plugins* para diversas fuentes de entrada/salida de interfaces diferentes (por ejemplo, Apache, syslog|rsyslog, bases de datos SQL-NoSQL, S3...).

5) Flume (<https://flume.apache.org>). Es el proyecto desarrollado por Apache como un servicio distribuido y eficiente para recoger, agregar y mover grandes cantidades de datos de *logs*. Esta herramienta presenta una arquitectura flexible que funciona sobre flujos de datos (*streaming data flows*). Es robusta y tolerante a fallos y utiliza un modelo simple de datos extensible lo cual permite que pueda ser utilizada como herramienta de gestión analítica y toma de decisiones.

Como ejemplo veremos la instalación y configuración de LogAnalyzer, que destaca por su simplicidad y prestaciones:

1) Verificamos que se dispone de Apache2 instalado y descargamos el paquete desde la web del desarrollador: <http://logalyzer.adiscon.com/downloads/>.

2) Lo descomprimos `tar xzvf logalyzer-x.y.z.tar.gz`, donde la versión correspondiente es *x.y.z*, y movemos la carpeta interna *src* al directorio *DocumentRoot* de Apache, por ejemplo,

```
mv logalyzer-x.y.z/src /var/www/html/log2
```

y también el archivo *contrib/configure.sh*, por ejemplo,

```
cp logalyzer-x.y.z/contrib/configure.sh /var/www/html/log2/
```

Ejecutamos

```
chmod +x /var/www/html/log2/configure.sh
```

y después `/var/www/html/log2/configure.sh`.

3) Con esto y conectándonos al servidor/*logs2* podremos realizar la configuración restante y cuando hayamos finalizado tendremos la visualización de los *logs* pero no podremos acceder a cierta funcionalidad y configuración reservada para el administrador. Para ello es necesario disponer de una base de datos (MySQL por defecto). Para hacer la instalación de la base de datos haremos `apt-get install mysqlserver php5-mysql` (debemos recordar el *passwd* para el usuario *root* de la BD, ya que luego será necesario). Creamos una base de datos llamada *logalyzer* (p. ej., ejecutando `mysql -u root -p` y luego `create database logalyzer;`).

4) Desde un navegador conectamos al servidor/aplicación. En nuestro caso <http://srv.nteum.org/log2/> y seguimos las ventanas de configuración. No se debe tocar nada hasta el paso 3 (*frontend options*), en cuya ventana inferior *user database options* se debe seleccionar a *Enable User Database = yes* y configurar como usuario *root* y el *passwd* introducido cuando se creó la base de datos. En la ventana siguiente se crearán las tablas y se pedirá un usuario y *passwd* para acceder a la gestión. Con ello ya podremos ver los *log* desde *syslog*, agregar nuevas fuentes de *log*, seleccionarlos, buscar, o entrar en el *Admin Center*

para cambiar las fuentes de *logs*, aspecto, informes, usuarios, campos, vistas, estadísticas y gráficos y una gran cantidad de opciones adicionales. Podemos encontrar información adicional sobre LogAnalyzer en [LogA] (o a través del botón *Help* de la aplicación). También podemos encontrar información sobre cómo generar informes semanalmente o diariamente (el ejemplo está para W pero es equivalente en Linux) en [LogAR].

5) Una vez que finalizada configuración, es conveniente cambiar las protecciones del archivo *config.php*, por ejemplo,

```
chmod 644 /var/www/html/log2/config.php
```

Hay que tener en cuenta las recomendaciones de *Securing LogAnalyzer* de la documentación indicada.

Las figuras 4 y 5 muestran dos imágenes de LogAnalyzer (la primera del *syslog* y la segunda de *auth.log*, como se ve en el selector en la parte derecha-superior).

Figura 4

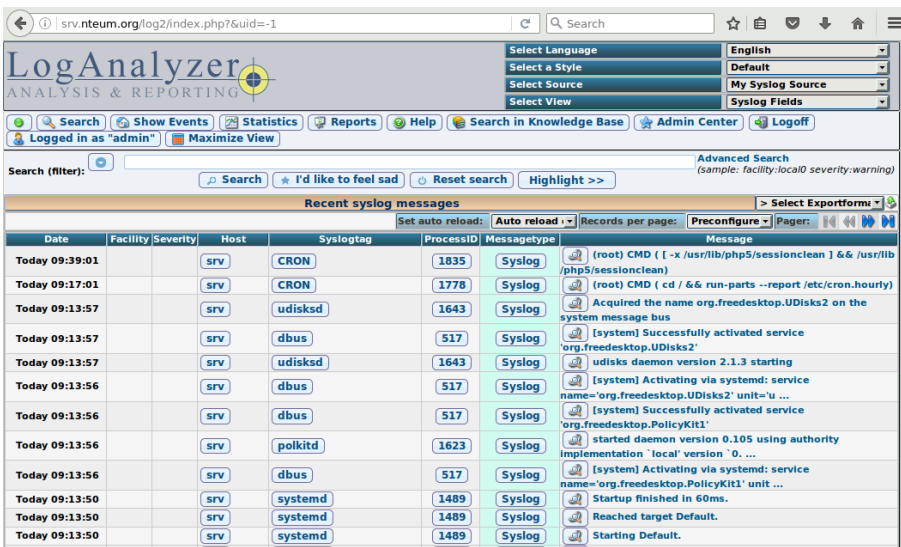
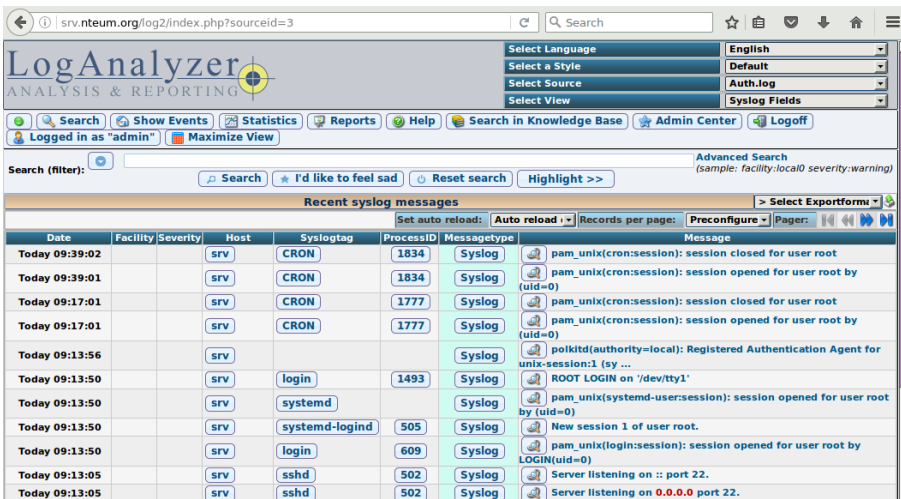


Figura 5



2. *Cloud*

Las infraestructuras *cloud* se pueden clasificar en 3+1 grandes grupos en función de los servicios que prestan y a quién:

1) Públicas: los servicios se encuentran en servidores externos y las aplicaciones de los clientes se mezclan en los servidores y con otras infraestructuras. La ventaja más clara es la capacidad de procesamiento y almacenamiento sin instalar máquinas localmente (no hay inversión inicial ni mantenimiento) y se paga por uso. Tiene un retorno de la inversión rápido y puede resultar difícil integrar estos servicios con otros propios.

2) Privadas: las plataformas se encuentran dentro de las instalaciones de la empresa/institución y son una buena opción para quienes necesitan alta protección de datos (estos continúan dentro de la propia empresa), es más fácil integrar estos servicios con otros propios, pero existe inversión inicial en infraestructura física, sistemas de virtualización, ancho de banda, seguridad y gasto de mantenimiento, lo cual supone un retorno más lento de la inversión. No obstante de tener infraestructura a tener infraestructura *cloud* y virtualizada, esta última es mejor por las ventajas de mayor eficiencia, mejor gestión y control, y mayor aislamiento entre proyectos y rapidez en la provisión.

3) Híbridas: combinan los modelos de nubes públicas y privadas y permite mantener el control de las aplicaciones principales al tiempo de aprovechar el *Cloud Computing* en los lugares donde tenga sentido con una inversión inicial moderada y a la vez contar los servicios que se necesiten bajo demanda.

4) Comunidad: Canalizan necesidades agrupadas y las sinergias de un conjunto o sector de empresas para ofrecer servicios de *cloud* a este grupos de usuarios.

Existen además diferentes capas bajo las cuales se pueden contratar estos servicios:

1) *Infrastructure as a Service* (IaaS): se contrata capacidad de proceso y almacenamiento que permiten desplegar aplicaciones propias que por motivos de inversión, infraestructura, coste o falta de conocimientos no se quiere instalar en la propia empresa (ejemplos de este tipo de EC2/S3 de Amazon y Azure de Microsoft).

2) *Platform as a Service* (PaaS): se proporciona además un servidor de aplicaciones (donde se ejecutarán nuestras aplicaciones) y una base de datos, donde se podrán instalar las aplicaciones y ejecutarlas -las cuales se deberán desarrollar de acuerdo a unas indicaciones del proveedor- (por ejemplo Google App Engine).

3) *Software as a Service* (SaaS): comúnmente se identifica con '*cloud*', donde el usuario final paga un alquiler por el uso de software sin adquirirlo en propiedad, instalarlo, configurarlo y mantenerlo (ejemplos Adobe Creative Cloud, Google Docs, o Office365).

4) *Business Process as a Service* (BPaaS): es la capa más nueva (y se sustenta encima de las otras 3), donde el modelo vertical (u horizontal) de un proceso de negocio puede ser ofrecido sobre una infraestructura *cloud*.

Si bien las ventajas son evidentes y muchos actores de esta tecnología la basan principalmente en el abaratamiento de los costes de servicio, comienzan a existir opiniones en contra (<http://deepvalue.net/ec2-is-380-more-expensive-than-internal-cluster/>) que consideran que un *cloud* público puede no ser la opción más adecuada para determinado tipo de servicios/infraestructura. Entre otros aspectos, los negativos pueden ser la fiabilidad en la prestación de servicio, seguridad y privacidad de los datos/información en los servidores externos, *lock-in* de datos en la infraestructura sin posibilidad de extraerlos, estabilidad del proveedor (como empresa/negocio) y posición de fuerza en el mercado no sujeto a la variabilidad del mercado, cuellos de botellas en la transferencias (empresa/proveedor), rendimiento no predecible, tiempo de respuesta a incidentes/accidentes, problemas derivados de la falta de madurez de la tecnología/infraestructura/gestión, acuerdos de servicios (SLA) pensados más para el proveedor que para el usuario, etc. No obstante es una tecnología que tiene sus ventajas y que con la adecuada planificación y toma de decisiones, valorando todos los factores que influyen en el negocio y escapando a conceptos superficiales (todos lo tienen, todos lo utilizan, bajo costo, expansión ilimitada, etc.), puede ser una elección adecuada para los objetivos de la empresa, su negocio y la prestación de servicios en IT que necesita o que forma parte de sus fines empresariales.

Es muy amplia la lista de proveedores de servicios *cloud* en las diferentes capas, pero de acuerdo a la información de Synergy Research Group en 2015*, los líderes en el mercado de infraestructura *cloud* son:

- 1) **Público IaaS/PaaS:** mercado dominado por Amazon y Microsoft (aproximadamente el 50%)
- 2) **Privado e Híbrido:** IBM y Amazon (aproximadamente el 45%)
- 3) **SaaS:** Salesforce y Microsoft (aproximadamente el 30%)

*<https://www.srgresearch.com/articles/2015-review-shows-110-billion-cloud-market-growing-28-annually>

4) **UCaaS (Unified Communications as a Service)**: Cisco y Citrix (aproximadamente el 18%)

5) **Infraestructura (hw & sw)**: Público: Cisco-HPE (aproximadamente 30%) y Privado: HPE-Cisco (aproximadamente el 18%)

El volumen de negocio (según este informe) crece un 28% anualmente y alcanzó los 110.000 millones de dólares. Esto significa un cambio apreciable con relación a los datos de 2014, donde se puede observar una variabilidad de la oferta muy alta*.

*<https://www.srgresearch.com/articles/amazon-salesforce-and-ibm-lead-cloud-infrastructure-service-segments>

En cuanto a plataformas para desplegar *clouds* con licencias GPL, Apache, BSD (o similares), podemos contar entre las más referenciadas (y por orden alfabético):

1) **AppScale**: es una plataforma que permite a los usuarios desarrollar y ejecutar/almacenar sus propias aplicaciones basadas en Google AppEngine y puede funcionar como servicio o en local. Se puede ejecutar sobre AWS EC2, Rackspace, Google Compute Engine, Eucalyptus, Openstack, CloudStack, así como sobre KVM y VirtualBox y soporta Python, Java, Go, y plataformas PHP Google AppEngine. <http://www.appscale.com/>

2) **Cloud Foundry**: es plataforma *open source* PaaS desarrollada por VMware escrita básicamente en Ruby and Go. Puede funcionar como servicio y también en local. <http://cloudfoundry.org/>

3) **Apache CloudStack**: diseñada para gestionar grandes redes de máquinas virtuales como IaaS. Esta plataforma incluye todas las características necesarias para el despliegue de un IaaS: CO (*compute orchestration*), *Network-as-a-Service*, gestión de usuarios/cuentas, API nativa, *resource accounting*, y UI mejorada (*User Interface*). Soporta los *hypervisores* más comunes, gestión del *cloud* vía web o CLI y una API compatible con AWS EC2/S3 que permiten desarrollar *clouds* híbridos. <http://cloudstack.apache.org/>

4) **Eucalyptus**: plataforma que permite construir *clouds* privados compatibles con AWS. Este software permite aprovechar los recursos de cómputo, red y almacenamiento para ofrecer autoservicio y despliegue de recursos de *cloud* privado. Se caracteriza por la simplicidad en su instalación y estabilidad en el entorno con una alta eficiencia en la utilización de los recursos. <https://www.eucalyptus.com/>

5) **Nimbus**: es una plataforma que una vez instalada sobre un clúster, provee IaaS para la construcción de *clouds* privados o de comunidad y puede ser configurada para soportar diferentes virtualizaciones, sistemas de colas, o Amazon EC2. <http://www.nimbusproject.org/>

6) OpenNebula: es una plataforma para gestionar todos los recursos de un centro de datos permitiendo la construcción de IaaS privados, públicos e híbridos. Provee una gran cantidad de servicios, prestaciones y adaptaciones que han permitido que sea una de las plataformas más difundidas en la actualidad. <http://opennebula.org/>

7) OpenQRM: es una plataforma para el despliegue de *clouds* sobre un centro de datos heterogéneos. Permite la construcción de *clouds* privados, públicos e híbridos con IaaS. Combina la gestión del la CPU/almacenamiento/red para ofrecer servicios sobre máquinas virtualizadas y permite la integración con recursos remotos o otros *clouds*. <http://www.openqrm-enterprise.com/index-2.html>

8) OpenShift: apuesta importante de la compañía RH y es una plataforma (version Origin) que permite prestar servicio *cloud* en modalidad PasS. OpenShift soporta la ejecución de binarios que son aplicaciones web tal y como se ejecutan en RHEL por lo cual permite un amplio número de lenguajes y *frameworks*. <https://www.openshift.com/products/origin>

9) OpenStack es una arquitectura software que permite el despliegue de *cloud* en la modalidad de IaaS. Se gestiona a través de una consola de control vía web que permite la provisión y control de todos los subsistemas y el aprovisionamiento de los recursos. El proyecto iniciado (2010) por Rackspace y NASA es actualmente gestionado por OpenStack Foundation y hay más de 200 compañías adheridas al proyecto, entre las cuales se encuentran las grandes proveedoras de servicios *clouds* públicos y desarrolladoras de SW/HW (ATT, AMD, Canonical, Cisco, Dell, EMC, Ericsson, HP, IBM, Intel, NEC, Oracle, RH, SUSE Linux, VMware, Yahoo entre otras). Es otra de las plataformas más referenciadas. <http://www.openstack.org/>

10) PetiteCloud: es una plataforma software que permite el despliegue de *clouds* privados (pequeños) y no orientado a datos. Esta plataforma puede ser utilizada sola o en unión a otras plataformas *clouds* y se caracteriza por su estabilidad/fiabilidad y facilidad de instalación. <http://www.petitecloud.org>

11) oVirt: si bien no puede considerarse una plataforma *cloud*, oVirt es una aplicación de gestión de entornos virtualizados. Esto significa que se puede utilizar para gestionar los nodos HW, el almacenamiento o la red y desplegar y monitorizar las máquinas virtuales que se están ejecutando en el centro de datos. Forma parte de RH Enterprise Virtualization y es desarrollado por esta compañía con licencia Apache. <http://www.ovirt.org/>

2.1. Opennebula

Considerando las opiniones en [il1, il2], nuestra prueba de concepto sobre infraestructura *cloud* la realizaremos sobre OpenNebula. Probablemente una de

las formas más fáciles de probar las funcionalidades de OpenNebula y hacer un despliegue inmediato es utilizar el Sandbox virtualizado y preconfigurado que provee el desarrollador. Para ello en nuestro caso utilizaremos la versión de VirtualBox que está sobre una máquina virtual CentOS 7 con OpenNebula 5.0.0 utilizando QEMU para ejecutar las máquinas virtuales. Además, pueden agregarse cualquier otro nodo físico, con algunos de los *hipervisores* soportados por OpenNebula para construir un *cloud* a pequeña escala. En esta prueba el usuario podrá conectarse a su OpenNebula *cloud* privado, mirar los recursos gestionados y lanzar nuevas instancias de máquinas virtuales sin las dificultades de configurar la infraestructura física. Los requerimientos serán 1 GB RAM libre, 10 GB de espacio libre de disco, un procesador/SO de 64 bits y las *Virtualization Extensions* de Virtualbox instaladas.

1) Descargamos el *virtual appliance* en formato OVA desde la siguiente página web: <http://openebula.org/tryout/sandboxvirtualbox/>.

2) Importamos el *appliance* desde VirtualBox: *File-> Import Appliance*. Ponemos en marcha la máquina virtual y saldrá *one-sandbox login*: y nos conectamos como *root* y *passwd* "openebula".

3) Cambiamos el teclado con `localectl set-keymap es` y recreamos las caché de los repositorios `yum makecache fast`.

4) Como utilizaremos la misma máquina virtual para la prueba, instalaremos un entorno gráfico mínimo haciendo:

```
yum groupinstall "X Window System"
yum groupinstall "Fonts"
yum install kde-workspace
yum install gdm firefox
unlink /etc/systemd/system/default.target
ln -sf /lib/systemd/system/graphical.target \
/etc/systemd/system/default.target
reboot (o también systemctl isolate graphical.target)
```

5) Accedemos al escritorio (se puede configurar el teclado desde las opciones del entorno de KDE), luego abrimos Firefox y nos conectamos a `127.0.0.1:9869` y saldrá la pantalla de OpenNebula. Introducimos como usuario *oneadmin* y *passwd* "openebula" y accederemos al escritorio de OpenNebula.

6) A partir de este punto veremos el *DashBoard* de OpenNebula y podemos crear una máquina virtual (signo + en VM) y crearla con el *template* que hay

disponible después de revisar los parámetros y ajustarlos. Una vez creada podremos acceder a ella y realizar todas las acciones necesarias con esta máquina virtual, acceder a ella o volver al *Dashboard* y ver los recursos consumidos.

Las figuras 6, 7 y 8 muestran un ejemplo del *Dashboard* y la ejecución de la máquina virtual en el mismo servidor. En la figura 8 se pueden ver en CLI los comandos útiles (en el recuadro amarillo) para obtener información de los recursos y los que está funcionando (hay que tener en cuenta que se deben hacer como el usuario *oneadmin*).

Figura 6

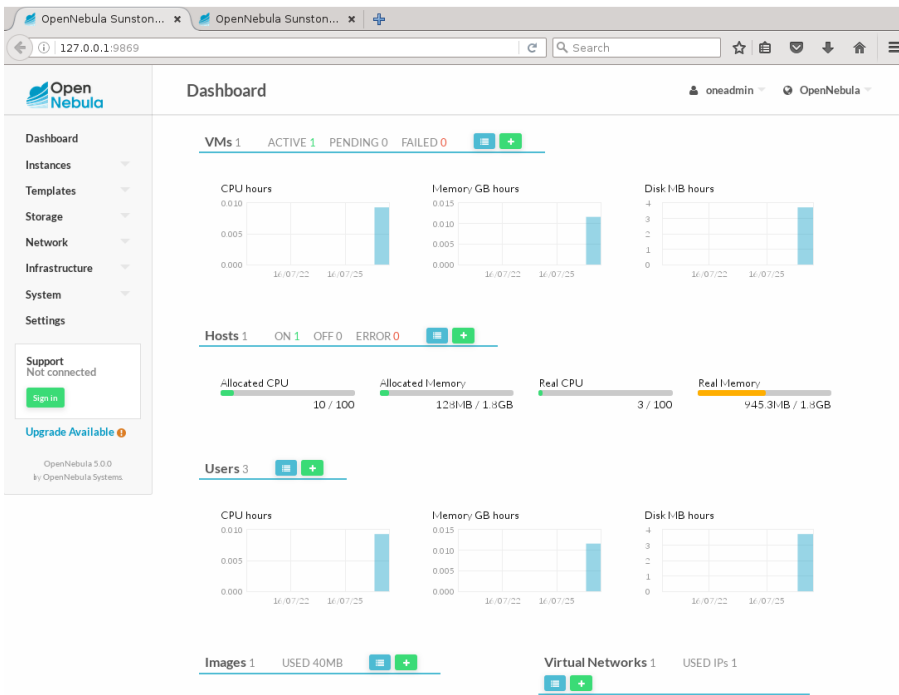
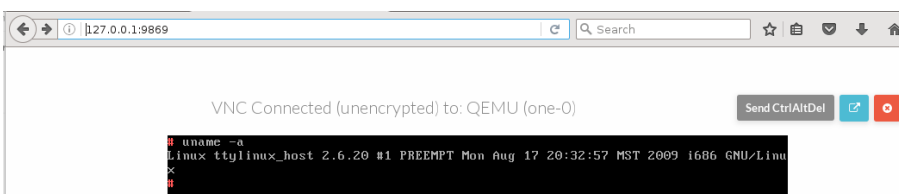


Figura 7



La instalación sobre Debian, por ejemplo, del *Frontend* se puede hacer fácilmente con (ver [ONDe]):

- 1) `wget -q -O- http://downloads.opennebula.org/repo/Debian/repo.key | apt-key add -`
- 2) `echo "deb http://downloads.opennebula.org/repo/5.0/Debian/8 stable opennebula"> /etc/apt/sources.list.d/opennebula.list`
- 3) `apt-get update`

Figura 8

```
[root@one-sandbox ~]# su - oneadmin
Last login: Wed Jun 15 10:42:19 BST 2016
[oneadmin@one-sandbox ~]$ onehost list


| ID | NAME        | CLUSTER | RVH | ALLOCATED CPU  | ALLOCATED MEM    | STAT |
|----|-------------|---------|-----|----------------|------------------|------|
| 0  | one-sandbox | default | 1   | 10 / 100 (10%) | 128M / 1.8G (6%) | on   |


[oneadmin@one-sandbox ~]$ onenetwork list


| ID | USER     | GROUP    | NAME  | CLUSTERS | BRIDGE | LEASES |
|----|----------|----------|-------|----------|--------|--------|
| 0  | oneadmin | oneadmin | cloud | 0        | br0    | 1      |


[oneadmin@one-sandbox ~]$ oneimage list


| ID | USER     | GROUP    | NAME    | DATASTORE | SIZE | TYPE | PER | STAT | RVMS |
|----|----------|----------|---------|-----------|------|------|-----|------|------|
| 0  | oneadmin | oneadmin | ttlinux | default   | 40M  | OS   | No  | used | 1    |


[oneadmin@one-sandbox ~]$ onetemplate list


| ID | USER     | GROUP    | NAME    | REGTIME        |
|----|----------|----------|---------|----------------|
| 0  | oneadmin | oneadmin | ttlinux | 06/15 10:46:50 |


[oneadmin@one-sandbox ~]$ onetemplate show 0
TEMPLATE 0 INFORMATION
ID : 0
NAME : ttlinux
USER : oneadmin
GROUP : oneadmin
REGISTER TIME : 06/15 10:46:50

PERMISSIONS
OWNER : um-
GROUP : ---
OTHER : ---

TEMPLATE CONTENTS
CONTEXT={
  NETWORK="YES",
  SSH_PUBLIC_KEY="root{SSH_PUBLIC_KEY}"
}
CPU="0.1"
DESCRIPTION="A small GNU/Linux system for testing"
DISK={
  IMAGE="ttlinux"
}
FEATURES={
  ACPI="no",
  APIC="no"
}
GRAPHICS={
  LISTEN="0.0.0.0",
  TYPE="vnc"
}
MEMORY="128"
NIC={
  NETWORK="cloud"
}
[oneadmin@one-sandbox ~]$
```

4) `apt-get install opennebula opennebula-sunstone opennebula-gate opennebula-flow`

Los paquetes que se instalarán son: *opennebula-common* (archivos comunes), *ruby-opennebula* (API Ruby), *libopennebula-java* (API Java), *libopennebula-java-doc* (documentación), *opennebula-node* (prepara un nodo como opennebula-node), *opennebula-sunstone* (GUI), *opennebula-tools* (CLI), *opennebula-gate* (comunicación entre VMs y OpenNebula), *opennebula-flow* (gestiona los servicios y la elasticidad) y *opennebula* (Daemon).

5) Ejecutamos `/usr/share/one/install_gems` (instala paquetes adicionales).

6) Si el entorno es de producción es recomendable substituir SQLite por MySQL (ver [ONDe]).

7) Hacemos `su - oneadmin` y luego `echo "oneadmin:mypassword" > .one/one_auth` para cambiar el *passwd* generado automáticamente.

8) Ejecutamos `service opennebula start`

9) Ejecutamos `service opennebula-sunstone start`

10) Probamos el servicio ejecutando `oneuser show` y nos dará algo como:

```
USER 0 INFORMATION
ID                : 0
NAME              : oneadmin
GROUP             : oneadmin
PASSWORD         : e6b47f32bdbcdb96aca2df6c8ccc190df3020828
AUTH_DRIVER      : core
ENABLED          : Yes
USER TEMPLATE
TOKEN_PASSWORD="f605ba14857d6e4ee28bce2735c832225727bb7a"
RESOURCE USAGE & QUOTAS
```

11) Luego nos conectamos a la máquina y al puerto 9869, en nuestro caso `http://srv.nteum.org:9869` e ingresamos como usuario `oneadmin` y como `passwd` el cambiado anteriormente. Así, volveremos a ver una imagen similar a la que se mostró anteriormente del Dashboard.

12) El siguiente paso es configurar un nodo (donde se aprovisionarán las máquinas virtuales). No tiene mayores dificultades y se puede hacer siguiendo la guía desde [ONNo].

13) Para hacer una prueba de concepto sobre la misma máquina, y dado que es una máquina virtual sobre VirtualBox y no podemos utilizar KVM, utilizaremos QEMU. Para ello instalamos

```
install libvirt-bin libvirt0 libvirt-daemon libvirt-clients
apt-get install qemu qemu-kvm qemu-system-x86 qemu-utils
```

14) Luego deberemos modificar los archivos de configuración para agregar:

```
/etc/libvirt/qemu.conf
user = "oneadmin"
group = "oneadmin"
dynamic_ownership = 0
/etc/libvirt/libvirtd.conf
listen_tls = 0
listen_tcp = 1
mdns_adv = 0
unix_sock_group = "oneadmin"
unix_sock_rw_perms = "0777"
auth_unix_ro = "none"
auth_unix_rw = "none"
/etc/one/oned.conf
    VM_MAD = [
        NAME           = "kvm",
        SUNSTONE_NAME  = "KVMQemu",
    ...
        TYPE = "qemu",
    ...
```

15) Reiniciamos la máquina y arrancamos los *daemons* de `opennebula` y `opennebula-sunstone`. Ingresamos en la página web (`http://srv.nteum.org:9869/`) y desde `Storage->App` descargamos una imagen de prueba (por ejemplo, la ID=0 `tty-linux - kvm`) indicándole que queremos que pase a OpenNebula desde el Marketplace. Cuando haya finalizado la descarga la veremos en `Storage->Images` como `Ready` y también como `template`.

16) Después en Infraestructura crearemos un *Host* seleccionando como driver el nombre que le hemos puesto (por ejemplo, nosotros le hemos dado como

nombre KVMQemu) y como máquina *localhost*. La deberemos ver en el estatus *ON* y con recursos asignados.

17) Finalmente podremos crear una máquina virtual (desde *Instances* o desde el *DashBoard*) con el *template* cargado y sobre la máquina creada. Esta la veremos en estado *Running* y podremos acceder a una consola para verificar su funcionamiento. Se debe tener en cuenta que no le hemos asignado una red con lo cual la máquina no tendrá dispositivo de red, pero como prueba sí que podremos ver el despliegue y acceder a la consola.

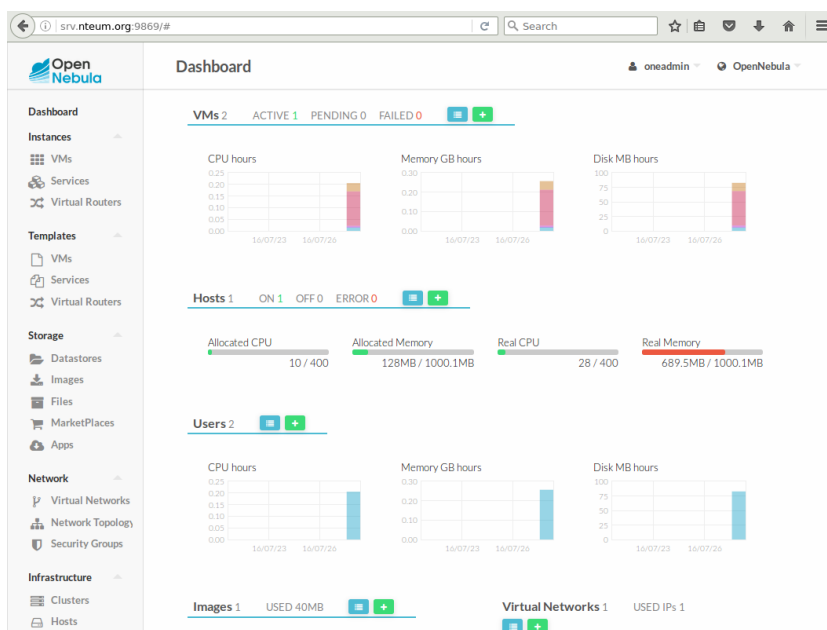
18) Para agregar una red deberemos configurar el sistema como *bridge* modificando el archivo */etc/network/interfaces*:

```
auto lo br0
iface lo inet loopback
allow-hotplug eth0 eth1
iface eth0 inet dhcp
iface br0 inet static
    address 172.16.1.1
    network 172.16.1.0
    netmask 255.255.255.0
    broadcast 172.16.1.255
    #gateway 192.168.0.1
    bridge_ports eth1
    bridge_fd 9
    bridge_hello 2
    bridge_maxage 12
    bridge_stp off
```

Luego seleccionamos Network->VirtualNetwork en el *DashBoard* y creamos la red asociada a *br0* y con los parámetros correspondientes. Luego se puede crear una nueva máquina virtual y asignar la red verificando que cuando arranque tendrá IP y el resto de parámetros de red asignados.

La figura 9 muestra la ejecución de esta última máquina virtual sobre el servidor configurado (<http://srv.nteum.org:9869>).

Figura 9



Como se ha podido comprobar, la instalación está bastante automatizada pero se debe tener en cuenta que es una infraestructura compleja y que se debe dedicar tiempo y análisis a determinar las causas por las cuales se producen los errores y solucionarlos. Existe una gran cantidad de documentación y sitios en Internet pero recomiendo comenzar por las fuentes [on1] y [on2].

En este apartado se ha visto una prueba de concepto funcional pero OpenNebula es muy extenso y flexible, y permite gran cantidad de opciones/extensiones. A través de OpenNebula MarketPlace* se podrán encontrar y descargar imágenes de máquinas virtuales creadas para OpenNebula y listas para ponerlas en funcionamiento (son las que se ven desde el MarketPlace de la interfaz web), pero al descargarlas sobre nuestro servidor la podremos utilizar siempre cuando la necesitemos y no se deberá descargar cada vez (no se debe descomprimir ni hacer nada, utilizar tal y como está).

[*http://marketplace.opennebula.org/systems/appliance](http://marketplace.opennebula.org/systems/appliance)

3. DevOps

Tal y como hemos comentado al inicio de este capítulo, Devops se puede considerar, en palabras de los expertos, un movimiento tanto en lo profesional como en lo cultural del mundo IT. Si bien no existen todas las respuestas todavía, un administrador se enfrentará a diferentes 'comunidades' (y él mismo formará parte de una de ellas), que tendrán nuevas necesidades de desarrollo y producción de servicios/productos dentro del mundo IT y con las premisas de 'más rápido', 'más eficiente', 'de mayor calidad' y totalmente adaptable a los 'diferentes entornos'. Es decir, el grupo de trabajo deberá romper las barreras existentes entre los departamentos de una empresa permitiendo que un producto pase rápidamente desde el departamento de investigación al de diseño, luego al de desarrollo + producción, luego al de test + calidad y por último, a ventas, y con las necesidades de herramientas que permitan desplegar todas estas actividades en cada una de las fases y llevar el control de todos por los responsables de cada uno de los ámbitos, incluidos los funcionales y los de la organización/directivos. Es por ello por lo que un administrador necesitará herramientas de gestión, control, despliegue, automatización y configuración. Una lista (corta) de las herramientas que podríamos considerar de acuerdo a nuestros objetivos de licencia GPL-BSD-Apache o similares (es interesante el sitio <http://devs.info/>, ya que permite acceder a la mayor parte de las herramientas/entornos/documentación para programadores y una lista más detallada -que incluye SW propietario- se puede encontrar en [NR]):

- 1) **Linux:** Ubuntu|Debian^v, Fedora^v|CentOS|SL
- 2) **IaaS:** Cloud Foundry, OpenNebula^v, OpenStack
- 3) **Virtualización:** KVM^v, Xen, VirtualBox^v, Vagrant^v
- 4) **Contenedores:** LXC^v, Docker^v
- 5) **Instalación SO:** Kickstart y Cobbler (rh), Preseed|Fai^v (deb), Rocks^v
- 6) **Gestión de la configuración:** Puppet^v, Chef^v, CFEngine, SaltStack, Juju, bcfg2, mcollective, fpm (effing), Ansible^v
- 7) **Servidores Web y aceleradores:** Apache^v, nginx^v, varnish, squid^v
- 8) **BD:** MySQL^v|MariaDB^v, PostgreSQL^v, OpenLDAP^v, MongoDB^v, Redis
- 9) **Entornos:** Lamp, Lamr, AppServ, Xampp, Mamp
- 10) **Gestión de versiones:** Git^v, Subversion^v, Mercurial^v
- 11) **Monitorización/supervisión:** Nagios^v, Icinga, Ganglia^v, Cacti^v, Monin^v, MRTG^v, XYmon^v

- 12) Misc: pdsh, ppsch, pssh, GNUparallel^v, nfsroot, Multihost SSH Wrapper, lldpd, Benchmarks^v, Librerías^v
- 13) Supervisión: Monit^v, runit, Supervisor^v, Godrb, BluePill-rb, Upstart, Systemd^v
- 14) Security: OpenVas^v, Tripwire^v, Snort^v
- 15) Desarrollo y Test: Jenkins, Maven, Ant, Gradle, CruiseControl, Hudson
- 16) Despliegue y Workflow: Capistrano
- 17) Servidores de aplicaciones: JBoss, Tomcat, Jetty, Glassfish,
- 18) Gestión de Logs: Rsyslog^v, Octopussy^v, Logstash

Excepto las que son muy orientadas a desarrollo de aplicaciones y servicios, en las dos asignaturas se han visto (o se verán dentro de este capítulo) gran parte de ellas (marcadas con ^v) y sin ninguna duda, con los conocimientos obtenidos, el alumno podrá rápidamente desplegar todas aquellas que necesite y que no se han tratado en estos cursos.

A continuación veremos algunas herramientas (muy útiles en entornos DevOps) más orientadas a generar automatizaciones en las instalaciones o generar entornos aislados de desarrollos/test/ejecución que permiten de forma simple y fácil (y sin pérdidas de prestaciones/rendimiento) disponer de herramientas y entornos adecuados a nuestras necesidades.

3.1. Linux Containers, LXC

LXC (LinuX Containers) es un método de virtualización a nivel del sistema operativo para ejecutar múltiples sistemas Linux aislados (llamados contenedores) sobre un único *host*. El *kernel* de Linux utiliza `cgroups` para poder aislar los recursos (CPU, memoria, E/S, network, etc.) lo cual no requiere iniciar ninguna máquina virtual. `cgroups` también provee aislamiento de los espacios de nombres para aislar por completo la aplicación del sistema operativo, incluido árbol de procesos, red, id de usuarios y sistemas de archivos montados. A través de una API muy potente y herramientas simples, permite crear y gestionar contenedores de sistema o aplicaciones. LXC utiliza diferentes módulos del *kernel* (`ipc`, `uts`, `mount`, `pid`, `network`, `user`) y de aplicaciones (`Apparmor`, `SELinux profiles`, `Seccomp policies`, `Chroots -pivot_root-` y `Control groups -cgroups-`) para crear y gestionar los contenedores. Se puede considerar que LXC está a medio camino de un 'potente' `chroot` y una máquina virtual, ofreciendo un entorno muy cercano a un Linux estándar pero sin necesidad de tener un kernel separado. Esto es más eficiente que utilizar virtualización con un *hypervisor* (KVM, Virtualbox) y más rápido de (re)iniciar, sobre todo si se están haciendo desarrollos y es necesario hacerlo frecuentemente, y su impacto en el rendimiento es muy bajo (el contenedor no ocupa recursos) y todos se dedicarán a los procesos que se estén ejecutando.

La instalación se realiza a través de `apt-get install lxc lxcctl` y se pueden instalar otros paquetes adicionales que son opcionales (`bridge-utils`, `libvirt-bin`, `debootstrap`), no obstante, si queremos que los contenedores tengan acceso a la red, con Ip propia es conveniente instalar el paquete `bridge-utils`: `apt-get install bridge-utils`. En Debian Jessie no se debe hacer nada más, pero en versiones anteriores se deben hacer ajustes con el directorio `cgroups` (consultar la documentación en [LXC1, LXC2, LXC4]). Después podemos verificar la instalación con `lxc-checkconfig` que dará una salida similar a:

```
Kernel configuration not found at /proc/config.gz; searching...
Kernel configuration found at /boot/config-3.16.0-4-amd64
--- Namespaces ---
Namespaces: enabled
Utsname namespace: enabled
Ipc namespace: enabled
Pid namespace: enabled
User namespace: enabled
Network namespace: enabled
Multiple /dev/pts instances: enabled

--- Control groups ---
Cgroup: enabled
Cgroup clone_children flag: enabled
Cgroup device: enabled
Cgroup sched: enabled
Cgroup cpu account: enabled
Cgroup memory controller: enabled
Cgroup cpuset: enabled

--- Misc ---
Veth pair device: enabled
Macvlan: enabled
Vlan: enabled
File capabilities: enabled

Note : Before booting a new kernel, you can check its configuration
usage : CONFIG=/path/to/config /usr/bin/lxc-checkconfig
```

Si surgen opciones deshabilitadas, se debe mirar la causa y solucionarlo (si bien en algunas configuraciones algunas opciones pueden estarlo).

En el directorio `/usr/share/lxc/templates` existe un conjunto de plantillas para crear diferentes contenedores, si bien alguna puede tener necesidades específicas. Para crear un contenedor Debian por ejemplo, ejecutamos:

```
lxc-create -n debian8 -t debian -- -r jessie
```

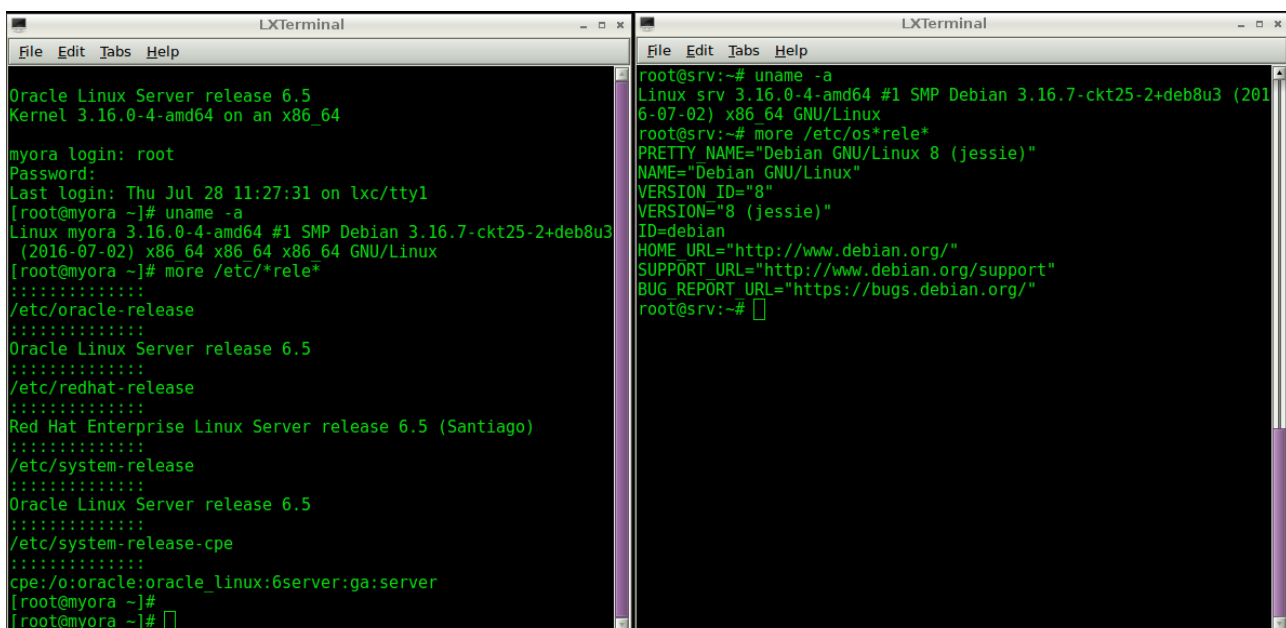
Veremos, luego de unos instantes de instalación, que en `/var/lib/lxc` tenemos un directorio que se llama `debian8` (el nombre del contenedor y de tamaño 287 MB). Es importante estar atento a la información que da la creación del contenedor, ya que allí nos informará del usuario y `passwd` (especialmente el de `root`). Los comandos más útiles para trabajar con el contenedor serán:

1) `lxc-ls` para listar los contenedores y `lxc-info` para obtener información sobre los contenedores disponibles y los que se están ejecutando (con el parámetro `--active`).

- 2) `lxc-start -n debian8 -d` para poner en marcha un contenedor en *background* (`-d=daemon`).
- 3) `lxc-console -n debian8` para conectarse a la consola del contenedor (con el *passwd* generado para el *root* durante la creación). Es recomendable verificar `/etc/issue.net` (o cualquier otro de los archivos que indique la versión). Usaremos `uname -a` para verificar el kernel y cambiar el *passwd* del *root*. Para volver al sistema operativo original deberemos presionar las teclas `Ctrl+a` y luego `'q'`.
- 4) `lxc-stop -n debian8` para la ejecución del contenedor.
- 5) `lxc-destroy -n debian8` para eliminar un contenedor.
- 6) `lxc-clone debian8 debian8V2` para clonar un contenedor.
- 7) Para montar un directorio del *host* en el contenedor se puede añadir en `/var/lib/lxc/containername/config` la sentencia
`lxc.mount.entry = /usr/bin mnt none ro,bind 0.0`
 que montará el `/usr/bin` de *host* en el directorio `/mnt` del container.
- 8) Si deseamos instalar un contenedor de otra rama de la del *host* (p. ej., Oracle basado en Red Hat), deberemos instalar una serie de dependencias como p. ej. los paquetes *rpm* y *yum* para Debian (`apt-get install rpm yum`). Luego podremos crear el contenedor con `lxc-create -n myoracle -t oracle` y veremos que se descarga una gran cantidad de paquetes y que el directorio (luego de finalizada la instalación) ocupa algo más que la distribución Debian previamente instalada (448 MB).

La figura 10 muestra el *host* a la derecha (Debian 8.5) y el contenedor a la izquierda (Oracle 6.5) funcionando sobre el mismo kernel, como se puede ver de la ejecución del comando `uname -a` desde cada uno de ellos.

Figura 10



Es interesante instalar el paquete *lxctl* (`apt-get install lxctl`) para gestionar de forma más simple los contenedores LXC.

Si se desea desde un *host* Debian crear un Ubuntu con el *template* proporcionado antes, se deberá instalar las *keys* de Ubuntu con

```
apt-get install ubuntu-archive-keyring
```

La configuración de la red es simple gracias al paquete *bridge-utils* (se pueden seguir los pasos indicados en [23-25]), pero simplificando, sobre el *host* creamos un dispositivo de *Bridge*, en */etc/network/interfaces* (en este caso el contenedor está sobre una máquina virtualizada en VBox donde el primer adaptador está como *Bridged*):

```
auto lo br0
iface lo inet loopback

iface br0 inet static
    address ip_dentro de la red
    netmask 255.255.255.0
    gateway gw_dentro_de_red
    bridge_ports eth0
    bridge_maxwait 0
    bridge_fd 0
```

Luego en */var/lib/lxc/container_name/config* modificamos (en nuestro caso es un contenedor Ubuntu creado anteriormente):

```
lxc.network.type = veth
lxc.network.flags = up
lxc.network.link = br0
#lxc.network.ipv4 = 158.109.74.24/24
lxc.network.hwaddr = 00:1E:1E:1a:00:00
#lxc.network.ipv4.gateway = 158.109.64.1
```

En este caso (y para evitar el tiempo de espera del dhcp) no le indicamos la IP ni el GW y se lo asignaremos en forma estática dentro del contenedor. Arrancamos el contenedor y modificamos */etc/network/interfaces* con:

```
auto lo
iface lo inet loopback
auto eth0
iface eth0 inet static
address ip_dentro_red_del_host
netmask 255.255.255.0
gateway ip_gw_igual_host
```

Reiniciamos el servicio de red y verificamos la conexión entre el contenedor y el *host* (probar un *ping* o una conexión *ssh*).

No obstante, si se desea una red **ALL-Nat**, es decir VirtualBox configurado en NAT en el primer adaptador (y Debian con dhcp sobre *eth0*) y que el contenedor tenga salida hacia afuera y NAT, lo mejor es crear una red virtual. Para ello instalamos `apt-get install libvirt-bin ebttables dnsmasq`. Re-

iniciamos y a continuación podremos ver el dispositivo por defecto con la orden `virsh net-info default` (en nuestro caso `vrbr0`); los parámetros se pueden modificar en `/etc/libvirt/qemu/networks/default.xml`. Luego iniciamos la red con `virsh net-start default` y veremos el dispositivo en marcha con `ip address`. Luego editamos `/var/lib/lxc/container_name/config` y configuramos:

```
[...]
lxc.network.type = veth
lxc.network.flags = up
lxc.network.link = virbr0
lxc.network.hwaddr = 00:FF:AA:00:00:01
lxc.network.ipv4 = 0.0.0.0/24
[...]
```

Ponemos en marcha el contenedor y verificamos la conexión externa con un `ping 8.8.8.8`

Si se desea dar a los contenedores la misma IP (mapeados por la MAC) se puede hacer en el archivo anterior:

```
<dhcp>
  <range start="192.168.122.100" end="192.168.122.254" />
  <host mac="00:FF:AA:00:00:01" name="foo.example.com" ip="192.168.122.101" />
  <host mac="00:FF:AA:00:00:02" name="bar.example.com" ip="192.168.122.102" />
</dhcp>
```

Si además se desea crear el dispositivo virtual cuando el `host` arranca, se puede hacer

```
virsh net-autostart default; virsh net-info default
```

Se debe tener cuidado cuando se inicia el contenedor sin `'-d'`, ya que no hay forma de salir (en la versión actual el `'Ctrl+a q'` no funciona). Iniciar siempre los contenedores en `background` (con `'-d'`) a no ser que necesite depurar porque el contenedor no arranca. Otra consideración a tener es que el comando `lxc-halt` ejecutará el `telinit` sobre el archivo `/run/initctl` si existe y apagará el `host`, por lo cual hay que apagarlo con `lxc-stop` y tampoco hacer un `shutdown -h now` dentro del contenedor, ya que también apagará el `host`. También existe un panel gráfico (si bien el desarrollo no ha sido actualizado últimamente) para gestionar los contenedores vía web <http://lxc-webpanel.github.io/install.html> (en Debian hay problemas con el apartado de red; algunos de ellos los soluciona <https://github.com/vaytess/LXC-Web-Panel/tree/lwp-backup>). Para ello clonar el sitio

```
git clone https://github.com/vaytess/LXC-Web-Panel.git
```

y luego reemplazar el directorio obtenido por `/srv/lwp` (renombrar este antes y hacer un `/etc/init.d/lwp restart`).

<https://linuxcontainers.org/downloads/>

3.2. Docker

Docker es una plataforma abierta para el desarrollo, empaquetado y ejecución de aplicaciones de forma que se puedan poner en producción o compartir más rápido separando estas de la infraestructura, de forma que sea menos costoso en recursos (básicamente espacio de disco, CPU y puesta en marcha) y que esté todo preparado para el siguiente desarrollador con todo lo que Ud. puso pero nada de la parte de infraestructura. Esto significará menor tiempo para para probar y acelerará el despliegue acortando en forma significativa el ciclo entre que se escribe el código y pasa a producción. Docker emplea una plataforma (contenedor) de virtualización ligera con flujos de trabajo y herramientas que le ayudan a administrar e implementar las aplicaciones proporcionando una forma de ejecutar casi cualquier aplicación en forma segura en un contenedor aislado. Este aislamiento y seguridad permiten ejecutar muchos contenedores de forma simultánea en el *host* y dada la naturaleza (ligera) de los contenedores todo ello se ejecuta sin la carga adicional de un *hypervisor* (que sería la otra forma de gestionar estas necesidades compartiendo la VM) lo cual significa que podemos obtener mejores prestaciones y utilización de los recursos. Es decir, con una VM cada aplicación virtualizada incluye no solo la aplicación, que pueden ser decenas de MBytes -binarios + librerías- sino también el sistema operativo 'guest', que pueden ser varios Gbytes; en cambio, en Docker lo que se denomina DE (Docker Engine) solo comprende la aplicación y sus dependencias, que se ejecuta como un proceso aislado en el espacio de usuario del SO 'host', compartiendo el *kernel* con otros contenedores. Por lo tanto, tiene el beneficio del aislamiento y la asignación de recursos de las máquinas virtuales, pero es mucho más portátil y eficiente transformándose en el entorno perfecto para dar soporte al ciclo de vida del desarrollo de software, test de plataformas, entornos, etc.[d1]

El entorno está formado por dos grandes componentes: **Docker** [d2] (es la plataforma de virtualización -contenedor-) y **Docker Hub** [d4] (una plataforma SasS que permite obtener y publicar/gestionar contenedores ya configurados). En su arquitectura, Docker utiliza una estructura cliente-servidor donde el cliente (CLI **docker**) interactúa con el *daemon* Docker, que hace el trabajo de la construcción, ejecución y distribución de los contenedores de Docker. Tanto el cliente como el *daemon* se pueden ejecutar en el mismo sistema, o se puede conectar un cliente a un *daemon* de Docker remoto. El cliente Docker y el servicio se comunican a través de sockets o una API RESTful*.

*Más detalles sobre la arquitectura en <https://docs.docker.com/engine/understanding-docker/>.

No se debe confundir una vieja aplicación llamada *docker* en Debian con la **Plataforma Docker** (para evitar confusiones en Debian se llama *docker.io*). Para agregar el repositorio en Jessie hacemos:

- 1) `apt-get install apt-transport-https ca-certificates`
- 2) `apt-key adv --keyserver hkp://p80.pool.sks-keyserver.net:80 --recv-keys 58118E89F3A912897C070ADB76221572C52609D`

- 3) `echo "deb https://apt.dockerproject.org/repo debian-jessie main" >\n/etc/apt/sources.list.d/docker.list`
- 4) `apt-get update`
- 5) Se puede verificar que vemos el repositorio con `apt-cache policy docker-engine`.
- 6) Para instalarlo `apt-get install docker-engine`.
- 7) Iniciar el *daemon* `service docker start`
- 8) Verificar que está instalado correctamente `docker run hello-world`
- 9) Instalar un contenedor Ubuntu `docker run -it ubuntu bash` (como el contenedor no existe lo descargará y lo ejecutará por lo cual veremos el *prompt* del *bash* de Ubuntu). Si hacemos `cat /etc/os-release`, veremos que estamos en Ubuntu. 16.04 LTS (Xenial Xerus). Con `exit` (o `Crtl-D`) salimos al *host* nuevamente.

Los comandos para iniciarse (además de los que hemos visto) son:

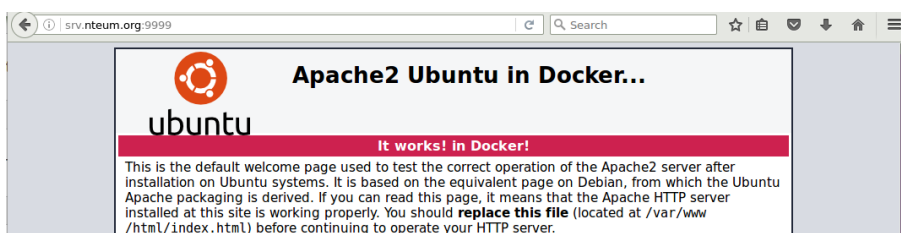
- 1) `docker`: muestra todas las opciones.
- 2) `docker images`: lista las imágenes localmente.
- 3) `docker search patrón`: busca contenedores/imágenes que tengan ese patrón.
- 4) `docker pull nombre`: obtiene imágenes del Hub. nombre puede ser `<user-name>/<repository>` para la particulares.
- 5) `docker run name cmd`: ejecutará el contenedor y dentro el comando indicado por `cmd`
- 6) `docker ps -l`: permite obtener el ID e información de una imagen.
- 7) `docker commit ID name`: actualiza la imagen con lo que se ha instalado hasta ese momento (con 3-4 números del ID es suficiente).
- 8) `docker inspect`: muestra las imágenes corriendo -similar a `docker ps`-.
9) `docker push`: salva la imagen en el Hub (debemos registrarnos primero) y está disponible para otros usuarios/host que se la quieran instalar con todo lo que hemos configurado dentro.
- 10) `docker cp`: copia archivos/folders desde el contenedor al host.
- 11) `docker export/import`: exporta/importa el contenedor en un archivo `tar`.
- 12) `docker history`: muestra la historia de una imagen.
- 13) `docker info`: muestra información general (imágenes, directorios, ...).
- 14) `docker kill`: para la ejecución de un contenedor.
- 15) `docker restart`: reinicia un contenedor.
- 16) `docker rm`: elimina un contenedor.
- 17) `docker rmi`: elimina imágenes.
- 18) `docker start/stop`: inicia/para un contenedor.

3.3. Instalación de un servidor Apache sobre un contenedor Docker

Para este objetivo comenzaremos con un contenedor base, instalaremos Apache y luego accederemos desde el *host* mapeando los puertos del contenedor. Para ello:

- 1) Ejecutamos `docker run -it ubuntu /bin/bash` para acceder al contenedor. Si no está instalado se bajará y se instalará.
- 2) Dentro del contenedor ejecutamos `apt-get update` y luego ejecutamos `apt-get install apache2 vim`. Podemos poner a prueba Apache2 con la instrucción `service apache2 start` y ver que arranca sin inconvenientes (solo mostrará una advertencia que no tiene el *ServerName* configurado). Se puede ver con `ps -edaf`. Modificamos la página inicial para poner alguna cosa relativa a Docker y así poder identificarla `vi /var/www/html/index.html`.
- 3) Es importante no salir del contenedor, sino hacer Ctrl-P y luego Ctrl-Q para salir al *host*, ya que si hacemos Ctrl-D apagaremos el contenedor y perderemos los cambios. También se puede hacer desde otro terminal.
- 4) Con el contenedor en marcha verificamos que está en *UP* con `docker ps` y hacemos `docker commit ID_contenedor apache`, donde *ID_contenedor* es el número que nos da `docker ps` y solo debemos poner los 3-4 primeros números. Esto salvará una nueva imagen llamada *apache* y que podemos verificar con `docker images`.
- 5) Ahora salimos de la sesión del contenedor (Ctrl-D o *exit*) o lo apagamos con `docker stop ID_contenedor`; `docker ps` nos deberá mostrar que no hay ninguna imagen en ejecución.
- 6) Ponemos en marcha el contenedor con
`docker run -d -p 9999:80 apache /usr/sbin/apache2ctl -D FOREGROUND`, donde le indicamos con `-d` que se ejecute en modo *daemon*, con `-p` que mapeemos el puerto 80 del contenedor en el 9999 del *host*, y ejecutamos el comando `apache2ctl` indicándole que ponga en marcha Apache2.
- 7) Luego desde el *host* podemos abrir un navegador con *localhost:9999* o el nombre del *host* mapeado en una IP del *host* como *srv.nteum.org:9999* y veremos la página modificada del contenedor como se muestra en la figura 11.

Figura 11



Lectura complementaria

Para más información de cómo utilizar los contenedores ir a [DoUs]

3.4. Puppet

Puppet es una herramienta de gestión de configuración y automatización en sistemas IT (licencia Apache, -antes GPL-). Su funcionamiento es gestionado a través de archivos (llamados *manifests*) de descripciones de los recursos del sistema y sus estados, utilizando un lenguaje declarativo propio de la herramienta. El lenguaje Puppet puede ser aplicado directamente al sistema, o compilado en un catálogo y distribuido al sistema de destino utilizando un modelo cliente-servidor (mediante una API REST), donde un agente interpela a los proveedores específicos del sistema para aplicar el recurso especificado en los *manifests*. Este lenguaje permite una gran abstracción que habilita a los administradores describir la configuración en términos de alto nivel, tales como usuarios, servicios y paquetes sin necesidad de especificar los comandos específicos del sistema operativo (apt, dpkg, etc.).[p1, p2, p12] El entorno Puppet tiene dos versiones Puppet (*Open Source*) y Puppet Enterprise (producto comercial) cuyas diferencias se pueden ver en <http://puppetlabs.com/puppet/enterprise-vs-open-source> y además se integran con estos entornos otras herramientas como MCollective (*orchestration framework*), Puppet Dashboard (consola web pero abandonada en su desarrollo), PuppetDB (*datawarehouse* para Puppet), Hiera (herramienta de búsqueda de datos de configuración), Factor (herramienta para la creación de catálogos) y Geppetto (IDE *–integrated development environment–* para Puppet).

3.4.1. Instalación

Es importante comenzar con dos máquinas que tengan sus *hostname* y definición en */etc/hosts* correctamente y que sean accesibles a través de la red con los datos obtenidos del */etc/hosts* (es importante que el nombre de la máquina -sin dominio- de este archivo coincida con el nombre especificado en *hostname* ya que si no habrá problemas cuando se generen los certificados).

Sobre el servidor instalamos *puppetmaster* (consideramos que tenemos una Debian Jessie ya que los paquetes están en el repositorio; para versiones anteriores podéis consultar la documentación):

```
apt-get install puppetmaster
```

Si está sobre una máquina virtual, reducimos los requerimientos de memoria editando */etc/default/puppetserver* y poniendo: `JAVA_ARGS="-Xms1g -Xmx1g"` si se desea dar 1 GB o `JAVA_ARGS="-Xms512m -Xmx512m"` para 512Mb. A continuación ejecutamos:

```
service puppetmaster restart
```

Sobre el cliente, instalamos *puppet*:

```
apt-get install puppet
```

Sobre el cliente modificamos `/etc/puppet/puppet.conf` para agregar en la sección `[main]` el servidor (`srv.nteum.org`, en nuestro caso) y reiniciamos:

```
[main] server=srv.nteum.org
```

```
puppet resource service puppet ensure=running enable=true service
puppet restart
```

Luego ejecutamos sobre el *master*: `puppet cert list`. Veremos algo como `"node.nteum.org"(SHA256): 7A:B7:CE:...:01:6D`

Firmamos el certificado:

```
puppet cert sign node.nteum.org.
```

La salida será algo como:

```
Notice: Signed certificate request for node.nteum.org
Notice: Removing file Puppet::SSL::CertificateRequest node,teum.orgat
'/var/lib/puppet/ssl/ca/requests/node.nteum.org'
```

Esto significa que `node.nteum.org` es aceptada por el *master* (`srv.nteum.org`) y se puede verificar con `puppet cert list -all`. La salida será algo como:

```
+ "node.nteum.org" (SHA256) 7A:B7:CE:...:01:6D
+ "srv.nteum.org" (SHA256) FF:E2:49:...:62:5F (alt names: "DNS:puppet",
"DNS:puppet.nteum.org", "DNS:srv.nteum.org")
```

En esta salida el signo "+" indica que los certificados están correctos.

Sobre el cliente podemos hacer `puppet agent --fingerprint`. Y veremos la clave SHA256 correspondiente.

```
puppet agent --test
```

Veremos que obtiene todos los *pluginfacts*, *puglins*, *catalogs* y muestra la versión de la configuración (esta información saldrá sobre el terminal en color verde; si hay algo que no funciona o un *warning*, saldrá en rojo y se deberá solucionar).

3.5. Main Manifest File

Puppet utiliza un lenguaje de dominio específico para describir las configuraciones del sistema. Estas descripciones se guardan en archivos denominados *manifestos*, que tienen una extensión *name.pp*. El archivo de manifiesto principal por defecto se encuentra en `/etc/puppet/manifests/site.pp` (para crearlo podemos hacer `touch /etc/puppet/manifests/site.pp`).

El comando `puppet apply` permite ejecutar bajo demanda manifiestos que no son relativos al principal. En este caso, se aplicará el manifiesto sobre el nodo que se haya indicado, por ejemplo:

```
puppet apply /etc/puppet/modules/test/init.pp.
```

Ejecutar el manifiesto de esta forma solo es útil si se desea probar un nuevo manifiesto sobre un agente o si solo se busca ejecutarlo una vez para inicializar el agente a un estado determinado.

Crearemos un manifiesto simple editando `/etc/puppet/manifests/site.pp` y agregamos:

```
file {'/tmp/example-ip':
  ensure => present,          # resource type file and filename
  mode   => 0644,             # make sure it exists
  content => "Here is my Public IP Address: ${ipaddress_eth0}.\n", # file permissions
}                             # note the ipaddress_eth0 fact
```

En este fragmento se asegura que el archivo exista y que además tenga permisos `rw-r - -r - -` e inserta un texto con la IP pública. Se puede esperar hasta que el agente verifique automáticamente; o se puede ejecutar sobre el cliente `puppet agent --test` y luego verificar `cat /tmp/example-ip`, lo que dará como resultado *Here is my Public IP Address: 172.16.1.2*. Se puede agregar en `site.pp` una especificación concreta a un nodo, por ejemplo. Si solo queremos que se instale sobre el nodo llamado `nodo.nteum.org`, se pone el contenido `file` dentro de una sentencia `node 'nodo.nteum.org' { file {..} }` y veremos que solo se aplica a este nodo. Hay que tener en cuenta que si no definimos un recurso, Puppet hará lo necesario para no tocarlo, con lo cual aunque se borran estos recursos del manifiesto Puppet no se borrarán los recursos creados. Si, por ejemplo, se desea borrar estos archivos, se debe cambiar en `ensure` la etiqueta `present` por `absent`.

Otra posibilidad que incluye Puppet, y que ahora utilizaremos, son los *modules*. Los módulos son útiles para agrupar tareas y están disponibles para la comunidad Puppet, si bien uno puede crear sus propios módulos. Para instalar Apache utilizaremos el `puppetlabs-apache module` desde `forgeapi` ejecutando `sudo puppet module install puppetlabs-apache` (la instalación de este módulo eliminará todas las configuraciones anteriores). A continuación editamos `vi /etc/puppet/manifest/site.pp` e insertamos:

```
node 'node' {
class { 'apache': } # use apache module
  apache::vhost { 'node.nteum.org': # define vhost resource
    port   => '80',
    docroot => '/var/www/html'
  }
}
```

Si ejecutamos sobre nodo `puppet agent -test`, veremos que luego de unos cuantos mensajes habrá instalado y configurado Apache y podremos conectarnos, por ejemplo, desde `srv.nteum.org` a `http://nodo.nteum.org`. Instalar un paquete, por ejemplo `nmap`, sobre el nodo es tan fácil como insertar el siguiente código dentro del `node 'node' { ... }`:

```
package {'nmap':
  ensure => installed,
}
```

Finalmente, para crear un módulo que pueda crear un usuario se puede hacer:

1) `mkdir -p /etc/puppet/modules/accounts/manifests`

2) `vi /etc/puppet/modules/accounts/manifests/groups.pp`

```
class accounts::groups {
  group { 'username':
    ensure => present,
  }
}
```

3) `vi /etc/puppet/modules/accounts/manifests/init.pp`

```
class accounts {
  include groups
  user { 'usernteum':
    ensure      => present,
    home        => '/home/usernteum',
    shell        => '/bin/bash',
    managehome  => true,
    gid          => 'username',
    password     => '$1$v9ESs.OZ$p.Y39eCnQ1T8A.BsHGGvO1'
  }
}
```

Para generar el `passwd` utilizamos `openssl passwd -1`

4) Incluimos como primera línea en `/etc/puppet/manifests/sites.pp`

```
node 'node' {
  include accounts
  ...
}
```

5) Ejecutamos sobre el cliente `puppet agent - -test` y luego verificamos que nos podemos conectar a `ssh usernteum@localhost`.

Una herramienta interesante que se complementa con Puppet es Foreman [p3, p4]. Esta herramienta *open source* permite gestionar los servidores en todo su ciclo de vida, desde el aprovisionamiento y la configuración de la orquestación y el seguimiento. El uso de Puppet, Chef y Ansible y la arquitectura de proxy inteligente de Foreman permite automatizar fácilmente tareas repetitivas, desplegar aplicaciones y gestionar de forma proactiva el cambio, tanto en las instalaciones con máquinas virtuales y *baremetal* o en la nube (*cloud*). Presenta una interfaz *web*, una CLI y una API REST para los diferentes niveles

Lectura recomendada

Véase la configuración y los detalles de la potencialidad del lenguaje Puppet y la arquitectura en [PuDoc]. Hay que tener en cuenta la versión instalada (`puppet -v`).

de interacción y permite gestionar desde decenas a decenas de miles de servidores. Su instalación no es compleja aunque sí tardará unos minutos y son importantes los requerimientos de memoria (por ejemplo para esta prueba de concepto se ha utilizado una MV en VirtualBox con 2,5 GB de memoria RAM).

La instalación se ha realizado para Debian Jessie (recomendado) aunque se puede hacer para otras versiones (véase [FOR]). Los pasos que se han de seguir son los siguientes:

1) Activamos los repositorios y certificados:

```
apt-get -y install ca-certificates
wget https://apt.puppetlabs.com/puppetlabs-release-pcl-jessie.deb
dpkg -i puppetlabs-release-pcl-jessie.deb
echo "deb http://deb.theforeman.org/ jessie 1.12" > /etc/apt/sources.list.d/foreman.list
echo "deb http://deb.theforeman.org/ plugins 1.12" >> /etc/apt/sources.list.d/foreman.list
apt-get -y install ca-certificates
wget -q https://deb.theforeman.org/pubkey.gpg -O- | apt-key add -
```

2) Descargamos el instalador:

```
apt-get update && apt-get -y install foreman-installer
```

3) Ejecutamos el instalador. La instalación es no-interactiva pero si se desean ajustar algunos parámetros se puede pasar al comando `foreman-installer` el parámetro `-i`.

Después de unos minutos tendremos una información que hay que verificar que ha terminado sin errores y en el caso de que los haya se deberán solucionar y repetir la instalación.

```
* Foreman is running at https://srv.nteum.org
  Initial credentials are admin / 4eftk...
* Foreman Proxy is running at https://srv.nteum.org:8443
* Puppetmaster is running at port 8140
```

El log queda registrado en `/var/log/foreman-installer/foreman-installer.log`.

4) Reiniciamos la máquina y verificamos con `service puppetmaster status` que nos da una información similar a:

```
puppetserver.service - LSB: puppetserver
  Loaded: loaded (/etc/init.d/puppetserver)
  Active: active (running) since Sun 2016-07-31 10:53:44 BST; 22min ago
  CGroup: /system.slice/puppetserver.service
          +-602 /usr/bin/java -XX:OnOutOfMemoryError=kill -9 %p -Djava.secur.
```

En el caso de que nos dé errores, tendremos que verificar que son de memoria en `var/log/puppetlabs/puppetserver/` y modificaremos `/etc/default/puppetserver` con, por ejemplo, `JAVA_ARGS="-Xms1g -Xmx1g -XX:MaxPermSize=256m"`.

Lo siguiente es acceder a la dirección `https://srv.nteum.org` y aceptar el certificado y conectarnos a `admin` con el `passwd` indicado antes. Para hacer una

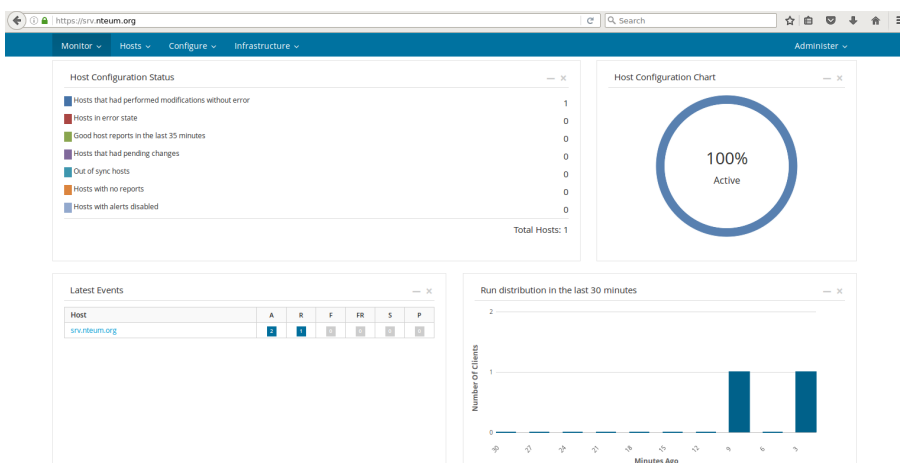
prueba se puede ejecutar desde la misma máquina `puppet agent --test` (dará unos mensajes de advertencia) y veremos que en *Tab->Host->All Hosts* de Foreman queda registrada la actividad y con el estatus de OK.

A continuación instalaremos un módulo de Puppet para manejar *module* con la intención de manejar el servicio NTP ejecutando

```
puppet module install puppetlabs/ntp
```

En Foreman, vamos a *Configure->Classes->Import from hostname* (arriba a la derecha) y la clase “ntp” aparecerá instalada correctamente. Haciendo un clic sobre “ntp” y en el tab *Smart Class Parameters seleccionar* marcamos *Override*; así, Foreman puede manejar esta clase. Sobre el tab *Hosts*, accedemos al servidor *->tab Puppet Classes ->ntp module* y agregamos ‘+’ para habilitar esta clase sobre *host*, y finalmente salvamos. Si desde un terminal ejecutamos la orden `puppet agent --test`, veremos cómo se instala el servicio automáticamente y se pone en marcha (se puede comprobar con `ntpq -p`) [FOR]. La figura 12 muestra el *Dashboard* de Foreman después de la acción.

Figura 12



3.6. Chef

Chef es otra de la grandes herramientas de configuración con funcionalidades similares a Puppet (existe un buen artículo donde realiza una comparación sobre el dilema de Puppet o Chef [ch1]). Esta herramienta utiliza un lenguaje (basado en Ruby) DSL (*domain-specific language*) para escribir las 'recetas' de configuración que serán utilizadas para configurar y administrar los servidores de la compañía/institución y que además puede integrarse con diferentes plataformas (como Rackspace, Amazon EC2, Google y Microsoft Azure) para, automáticamente, gestionar la provisión de recursos de nuevas máquinas. El administrador comienza escribiendo las 'recetas' que describen cómo maneja Chef las aplicaciones del servidor (tales como Apache, MySQL, or Hadoop) y cómo serán configuradas indicando qué paquetes deberán ser instalados, los servicios que deberán ejecutarse y los archivos que deberán modificarse informando además de todo ello en el servidor para que el administrador tenga el

control del despliegue. Chef junto con Puppet, CFEngine, Ansible, Bcfg2 es una de las herramientas más utilizadas para este tipo de funcionalidad y que ya forma parte de las media-grandes instalaciones actuales de GNU/Linux.

3.7. Estructura de Chef

La estructura de Chef se basa en una **estación de trabajo** (*Workstation*) desde donde se gestiona toda la infraestructura incluyendo Chef DK, las “recetas” y el uso de herramientas como *kichen*, *chef-cero*, las herramientas de línea de comandos como Knife (para interactuar con el servidor de Chef) y el propio Chef (para interactuar con su chef-repo local), y recursos como el entorno básico de Chef (para la creación de recetas) e InSpec (para la construcción de los controles de seguridad y flujos de trabajo). Esta estación de trabajo está configurada para permitir a los usuarios crear, probar y mantener los *cookbooks* que se cargan en el servidor y que pueden ser propios para la organización o de otros disponibles en el *Chef Supermarket*.

Los **nodos** son las máquinas-físicos, virtuales, nubes, etc., que están bajo la gestión de Chef. El chef-cliente se instala en cada nodo y es lo que lleva a cabo la automatización en esa máquina.

El **servidor Chef** actúa como un centro de información. Los *cookbooks* y las políticas de configuración se cargan en el servidor Chef desde las estaciones de trabajo, aunque también pueden ser mantenidas desde el propio servidor Chef, a través de la interfaz web de la consola de gestión de Chef.

El **chef-cliente** se ejecuta en el nodo y tiene acceso al servidor de Chef para obtener los datos de configuración, realizar búsquedas de datos históricos, y una vez finalizada la ejecución del Chef y del cliente este actualiza los datos de ejecución en el servidor como nodo actualizado.

Chef Supermarket es el lugar en el que se almacena los *cookbooks* de la comunidad y pueden ser utilizados por cualquier usuario Chef.

Una forma de probar las recetas es utilizar el programa **Chef-Solo** para después pasarlas al servidor. No obstante, esto es obsoleto ya que el cliente actual puede funcionar en modo local (y emula toda la funcionalidad de Chef-Solo). Un programa complementario es **Chef-Zero** (que surgió como una experiencia de laboratorio), que es un servidor Chef en memoria y para fines de desarrollo, sin datos en el disco ni tampoco autenticación o autorización y que está incluido en Chef cliente (a partir de la versión 11.8.0) y que permite que pasarle el parámetro `-local-mode` al chef-client. El modo local pone en marcha el servidor local Chef limitado solo al *localhost*, se cargan todas las recetas locales, se ejecuta el Chef cliente, se realizan las acciones y luego termina el servidor Chef-Zero. La experiencia del usuario final es la misma que si contara con la

Chef-cero

Chef-cero es una herramienta de línea de comandos que se ejecuta localmente como si estuviera conectado a un servidor auténtico Chef.

Enlace de interés

La mejor forma de ver la potencialidad de Chef es utilizar la máquina virtual Ubuntu que provee el desarrollador (<https://learn.chef.io/learn-the-basics/ubuntu/get-setup/>) y que se ejecuta para el usuario durante 12 horas. En un paso-a-paso enseña las potencialidades de Chef sin grandes inversiones en instalación y configuración.

infraestructura completa pero solo en una máquina local. En Debian se puede instalar Chef (`apt-get install chef`), que instalará Chef-Zero y que nos permitirá crear las recetas y aplicarlas en el mismo nodo.

Como prueba de concepto probemos algunas recetas:

1) Creamos un directorio de recetas: `mkdir $HOME/repo; cd $HOME/repo`

2) Creamos la primera receta `vi hello.rb`, la cual crea el archivo `/tmp/motd` con el contenido `'hello world'`:

```
file '/tmp/motd' do
  content 'hello world'
end
```

3) Ejecutamos `chef-client -local-mode hello.rb` que nos dará algo como:

```
...
Starting Chef Client, version 11.12.8
resolving cookbooks for run list: []
Synchronizing Cookbooks:
Compiling Cookbooks...
[2016-07-31T23:43:58+01:00] WARN: Node srv.nteum.org has an empty run list.
Converging 1 resources
Recipe: @recipe_files::/root/recetas/hello.rb
 * file [/tmp/motd] action create
   - create new file /tmp/motd
   - update content in file /tmp/motd from none to c38c60
     --- /tmp/motd 2016-07-31 23:43:58.287615557 +0100
     +++ /tmp/.motd20160731-19247-1wil5o5 2016-07-31 23:43:58.287615557 +0100
     @@ -1,2 @@
     +hello chef
Running handlers complete
Chef Client finished, 1/1 resources updated in 1.531396599 second
```

4) Podemos verificar que existe el archivo con el contenido correcto y se puede probar de ejecutar el comando por segunda vez. Veremos que no hay cambios (igual si cambiamos el texto del mensaje y volvemos a ejecutar veremos que Chef actualiza el contenido o si cambiamos el contenido del archivo Chef restituirá el contenido original). Si se cambia `content 'hello world'` por `action :delete` lo que haremos es borrar el archivo.

5) Para instalar un paquete primero deberemos hacer un `update` de los paquetes y luego instalar por ejemplo `apache2`: `vi webserver.rb`

```
execute 'apt-update' do
  command 'apt-get update'
end

package 'apache2' package 'apache2'
```

Se puede comprobar que luego tenemos Apache instalado. Si quisiéramos cambiar el archivo inicial solo deberíamos agregar luego de `package`:

```
file '/var/www/html/index.html' do
  content '<html>
```

```
<body>
  <h1>hello world</h1>
</body>
</html>'
end
```

Como se puede ver, la potencialidad y facilidad de las recetas es grande y es fácil hacer una descripción de lo que se desea para configurar un nodo. Una forma de organizar las recetas y de adaptarlas a diferentes escenarios es construir un *cookbook* que sea una estructura que luego se pueda ligar para tomar diferentes partes de ella y mantener todo organizado [ChLB].

Los software (*ChefDK*, *Server*, *Client*, *Automate*...) se pueden descargar desde <https://downloads.chef.io/> (hay que tener en cuenta que algunos paquetes tienen limitaciones de nodos –hasta 25–) y todos están para una arquitectura Ubuntu. Hay diversos tutoriales, por ejemplo [DO], sobre cómo instalar paso a paso el servidor y la infraestructura si el usuario está interesado en probar toda la infraestructura, pero requieren infraestructura específica o máquinas virtuales con altas prestaciones que se salen del objetivo de este subapartado.

3.8. Ansible

Ansible es una plataforma *open source* para configurar y administrar muy fácilmente sistemas informáticos e implementar aplicaciones de software en los nodos. Y todo ello utilizando solo SSH. Como hemos visto, las herramientas de despliegue consideradas anteriormente necesitan un agente en el *host* remoto. En cambio, Ansible solo necesita una conexión SSH y Python (2.4 o posterior) para llevar a cabo una acción en los nodos que se han de configurar. Además, integra módulos que trabajan sobre formato JSON y utiliza YAML para describir configuraciones reusables. Los desarrolladores son los mismos que los del *software* de aprovisionamiento Cobbler y tiene una arquitectura que distingue entre el **controlador** y los **nodos**. El controlador es dónde se inicia la orquestación y este gestiona por `ssh` los nodos que conoce a través de un inventario.

Ansible inserta módulos a los nodos (mediante `ssh`) que se guardan temporalmente y se comunican con el controlador a través del protocolo JSON sobre una salida estándar.

Su diseño está basado en una arquitectura minimalista (sin imponer dependencias adicionales), consistente, segura y de alta confiabilidad. Como cada módulo puede ser escrito en cualquier lenguaje estándar (Python, Perl, Ruby, Bash, etc.), la curva de aprendizaje es suave y permite estar operativo rápidamente incluso para infraestructuras complejas.

Su instalación y comprobación es sumamente fácil:

- 1) Instalamos el paquete sobre el controlador: `apt-get install ansible`
- 2) Verificamos su versión: `ansible --version`
- 3) Generamos las llaves SSH:

```
ssh-keygen -t rsa -b 4096 -C "admin@srv.nteum.org"
```

- 4) Las copiamos a los clientes: `ssh-copy-id root@172.16.1.2` y lo mismo para todos los clientes. Tendremos que verificar primero que en el cliente tengamos el parámetro `PermitRootLogin yes` en `/etc/ssh/sshd_config`. Si no es así habrá que cambiarlo y reiniciar el servicio.

- 5) Probamos que nos podemos conectar a todos los nodos sin `passwd`: `ssh root@172.16.1.2`

- 6) Editamos `vi /etc/ansible/hosts` y agregamos todos los `hosts` clientes (bajo una misma etiqueta, por ejemplo), en nuestro caso `webservers`.

- 7) Comprobamos `ansible -m ping webservers`, que nos dará una respuesta como

```
172.16.1.2 | success >> {
  "changed": false,
  "ping": "pong"
}
```

- 8) Ejecutamos el comando `ansible -m command -a "df -h" webservers` que nos dará como respuesta:

```
172.16.1.2 | success | rc=0 >>
Filesystem      Size  Used Avail Use% Mounted on
/dev/sda1        3.7G  1.6G  1.9G  46% /
udev            10M    0    10M   0% /dev
tmpfs           201M   4.5M  196M   3% /run
tmpfs           501M    0    501M   0% /dev/shm
tmpfs           5.0M    0    5.0M   0% /run/lock
tmpfs           501M    0    501M   0% /sys/fs/cgroup
tmpfs          101M    0    101M   0% /run/user/0
```

- 9) Ahora podemos ejecutar otros comandos como

```
ansible -m command -a "free -mt" web-servers
ansible -m command -a "free -mt" webservers
ansible -m command -a "uptime" webservers
ansible -m command -a "arch" webservers
ansible -m shell -a "hostname" webservers
ansible webservers -m copy -a "src=/etc/hosts dest=/tmp/hosts"
ansible all -m user -a 'password=$6$n... n2Mn1 name=testing'
para crear un usuario testing, el passwd es generado con mkpasswd --method=SHA-512.
```

Como se puede observar es muy simple de poner en marcha y muy potente y sin las particularidades/complejidades de instalación/configuración de Puppet, Capistrano, Salt o Chef. Al utilizar como agente SSH (generalmente disponible por defecto en todos los clientes Linux) su despliegue es muy fácil. Nuestro paso siguiente es configurar *playbooks* y trabajar con ellos. [Ans]

Los *playbooks* son descripciones en texto plano (en formato YAML) que permite organizar tareas complejas a través de ítems y pares *key: values*.

Dentro de un *playbook* podemos encontrar uno o más grupos de *hosts* (cada uno de estos es llamado *play*) donde las tareas serán llevadas a cabo. Por ejemplo:

```
---
- hosts: webservers
  remote_user: root
  vars:
    variable1: value1
    variable2: value2
  remote_user: root
  tasks:
    - name: description for task1
      task1: parameter1=value_for_parameter1 parameter2=value_for_parameter2
    - name: description for task2
      task2: parameter1=value_for_parameter1 parameter2=value_for_parameter2
  handlers:
    - name: description for handler 1
      service: name=name_of_service state=service_status
- hosts: dbbservers
  remote_user: root
  vars:
    variable1: value1
    variable2: value2
...
```

Los *handlers* son acciones que se ejecutarán al final de una tarea como por ejemplo al reiniciar un servicio o la máquina donde se han instalado unos paquetes. Un ejemplo para desplegar Apache sería:

1) `mkdir /etc/ansible/playbooks; vi /etc/ansible/playbooks/apache.yml` con el siguiente contenido:

```
---
- hosts: webservers
  tasks:
    - name: install apache2
      apt: name=apache2 update_cache=yes state=latest
    - name: copy index.html
      copy: src=/tmp/index.html dest=/var/www/html/ mode=0644
  handlers:
    - name: restart apache2
      service: name=apache2 state=restarted
```

2) Creamos el archivo `vi /tmp/index.html`

```
<!DOCTYPE html>
<html lang="en">
  <head>
    <meta charset="utf-8"/>
  </head>
  <body>
    <h1>Apache was started in this host via Ansible</h1><br>
    <h2>Ansible is Easy!!</h2>
  </body>
</html>
```

3) Ejecutamos el *playbook*

```
ansible-playbook /etc/ansible/playbooks/apache.yml
```

y accedemos al servidor poniendo en la URL 172.16.1.2 (la IP o el nombre del nodo desplegado).

Figura 13



Apache was started in this host via Ansible

Ansible is Easy!!!

Hay una gran cantidad de tutoriales (por ejemplo, [An4H]) que muestran cómo hacer *playbook* más complejos y cómo gestionar diferentes recursos. De esta manera el administrador estará en condiciones de implementar los suyos propios en forma muy simple y con un rápido despliegue.

Ansible dispone de una GUI llamada Ansible Tower basada en web que permite todas las opciones de la CLI además de características especiales como monitorización de los nodos en tiempo real, interfaz REST API para Ansible, *job scheduling*, y una interfaz gráfica para la gestión del inventario e integración con el *cloud* (por ejemplo, EC2, Rackspace, Azure) pero es un producto comercial (hasta 100 nodos U\$S 5.000/año).

Como alternativas se pueden encontrar RunDeck* que permite ejecutar *playbooks* de Ansible (mediante un *plugin* que se debe descargar**) o Semaphore*** que es una UI para Ansible *Open Source*.

*<http://rundeck.org/>
**<https://github.com/Batix/rundeck-ansible-plugin>
***<https://github.com/ansible-semaphore/semaphore>

3.9. Vagrant

Esta herramienta es útil en entornos DevOps y juega un papel diferente al de Docker pero orientado hacia los mismos objetivos: proporcionar entornos fáciles de configurar, reproducibles y portátiles con un único flujo de trabajo que ayudará a maximizar la productividad y la flexibilidad en el desarrollo de aplicaciones/servicios. Puede aprovisionar máquinas de diferentes proveedores (VirtualBox, VMware, AWS, u otros) utilizando *scripts*, Chef o Puppet para instalar y configurar automáticamente el software de la VM. Para los desarrolladores Vagrant aislará las dependencias y configuraciones dentro de un único entorno disponible, consistente, sin sacrificar ninguna de las herramientas que el desarrollador utiliza habitualmente y teniendo en cuenta un archivo, llamado *Vagrantfile*, el resto de desarrolladores tendrá el mismo entorno aun cuando trabajen desde otros SO o entornos, logrando que todos los miembros de un equipo estén ejecutando código en el mismo entorno y con las mismas

dependencias. Además, en el ámbito IT permite tener entornos desechables con un flujo de trabajo coherente para desarrollar y probar scripts de administración de la infraestructura ya que rápidamente se pueden hacer pruebas de scripts, 'cookbooks' de Chef, módulos de Puppets y otros que utilizan la virtualización local, tal como VirtualBox o VMware. Luego, con la misma configuración, puede probar estos *scripts* en el *cloud* (p. ej., AWS o Rackspace) con el mismo flujo de trabajo.

En primer lugar es necesario indicar que deberemos trabajar sobre un sistema Gnu/Linux base (lo que se define como *Bare-Metal*, es decir, sin virtualizar, ya que necesitaremos las extensiones de HW visibles y si hemos virtualizado con Virtualbox, por ejemplo, esto no es posible). Es recomendable instalar la versión de Virtualbox (puede ser la del repositorio Debian), verificar que todo funciona, y luego descargar la última versión de Vagrant desde <http://www.vagrantup.com/downloads.html> e instalarla con la instrucción `dpkg -i vagrant_x.y.z_x86_64.deb` (donde x.y.z será la versión que hemos descargado).

A partir de este punto es realmente muy simple ejecutando [v1] la instrucción `vagrant init hashicorp/precise32`, que inicializará/generará un archivo llamado **Vagrantfile** con las definiciones de la VM y cuando ejecutemos **vagrant up** descargará del repositorio *cloud* de Vagrant una imagen de Ubuntu 12.04-32b y la pondrá en marcha. Para acceder a ella simplemente debemos hacer `vagrant ssh` y si queremos desecharla, `vagrant destroy`. En el *cloud* de Vagrant [v2] podemos acceder a diferentes imágenes preconfiguradas* que podremos cargar haciendo `vagrant box add nombre`, por ejemplo `vagrant box add chef/centos-6.5`. Cuando se ejecute el comando 'up' veremos que nos da una serie de mensajes, entre los cuales nos indicará el puerto para conectarse a esa máquina virtual directamente (por ejemplo, `ssh vagrant@localhost -p 2222` y con *passwd vagrant*). Para utilizar una VM como base, podemos modificar el archivo Vagrantfile y cambiar el contenido:

```
Vagrant.configure("2") do |config|
  config.vm.box = "hashicorp/precise32"
end
```

Si deseamos trabajar con dos máquinas en forma simultánea deberemos modificar el archivo Vagrantfile con lo siguiente:

```
Vagrant.configure("2") do |config|
  config.vm.define "centos" do |centos|
    centos.vm.box = "chef/centos-6.5"
  end
  config.vm.define "ubu" do |ubu|
    ubu.vm.box = "hashicorp/precise32"
  end
end
```

*<https://vagrantcloud.com/discover/featured>

Podremos ver que les asigna un puerto SSH a cada una para evitar colisiones cuando hacemos el `vagrant up`. Desde la VM podríamos instalar el software como se hace habitualmente, pero para evitar que cada persona haga lo mismo existe una forma de aprovisionamiento que se ejecutará cuando se haga el 'up' de la VM. Para ello, escribimos un *script* `init.sh` con el siguiente contenido:

```
#!/usr/bin/env bash
apt-get update
apt-get install -y apache2
rm -rf /var/www
ln -fs /tmp /var/www
```

En este script (para no tener problemas de permisos para esta prueba) hemos apuntando el directorio `/var/www` a `/tmp`. A continuación modificamos el Vagrantfile (solo hemos dejado una VM para acelerar los procesos de carga):

```
Vagrant.configure("2") do |config|
  config.vm.define "ubu" do |ubu|
    ubu.vm.box = "hashicorp/precise32"
    ubu.vm.provision :shell, path: "init.sh"
  end
end
```

Luego deberemos hacer `vagrant reload -provision`. Veremos cómo se aprovisiona y carga Apache y luego, si en entramos en la máquina y creamos un `/tmp/index.html` y hacemos `wget 127.0.0.1` (o la ip interna de la máquina), veremos que accedemos al archivo y lo bajamos (igualmente si hacemos `ps -ef | grep apache2` veremos que está funcionando). Por último y para verla desde el *host* debemos redireccionar el puerto 80 por ejemplo al 4444. Para ello debemos agregar en el mismo archivo (debajo de donde pone `ubu.vm.provision`) la línea:

```
config.vm.network :forwarded_port, host: 4444, guest: 80
```

Volvemos a hacer `vagrant reload -provision` y desde el *host* nos conectamos a la url: `127.0.0.1:4444` y veremos el contenido del `index.html`.

En estas pruebas de concepto hemos mostrado algunos aspectos interesantes pero es una herramienta muy potente que se debe analizar con cuidado para configurar las opciones necesarias para nuestro entorno, como por ejemplo, aspectos del *Vagrant Share* o cuestiones avanzadas del *Provisioning* que aquí solo hemos tratado superficialmente.[v1]

Actividades

1. Instalad y configurad OpenMPI sobre un nodo; compilad y ejecutad el programa `cpi.c` y observad su comportamiento.
2. Instalad y configurad OpenMP; compilad y ejecutad el programa de multiplicación de matrices (<https://computing.llnl.gov/tutorials/openMP/exercise.html>) en 2 *cores* y obtened pruebas de la mejora en la ejecución.
3. Utilizando Rocks y VirtualBox, instalad dos máquinas para simular un clúster.
4. Instalad Docker y cread 4 entornos diferentes.
5. Instalad Puppet y configurad un máquina cliente.
6. Ídem punto anterior con Chef.
7. Ídem con Ansible.
8. Con Vagrant cread dos VM, una con Centos y otra con Ubuntu y aprovisionar la primera con Apache y la segunda con Mysql. Las máquinas se deben poder acceder desde el host y comunicar entre ellas.

Bibliografía

Todos los enlaces visitados en Julio de 2016.

- [An4H] Ansible Tutorial. <<https://serversforhackers.com/an-ansible-tutorial>>
- [Ans] Ansible: Getting Started <http://docs.ansible.com/ansible/intro_getting_started.html>
- [Beo] *Beowulf cluster*. <http://en.wikipedia.org/wiki/Beowulf_cluster>
- [Beo1] **S. Pereira** *Building a simple Beowulf cluster with Ubuntu*.
<https://www-users.cs.york.ac.uk/mjf/pi_cluster/src/Building_a_simple_Beowulf_cluster.html>
- [Bla] **Barney, B.** *OpenMP*. Lawrence Livermore National Laboratory.
<<https://computing.llnl.gov/tutorials/openMP/>>
- [ch1] *Puppet or Chef: The configuration management dilemma*.
<<http://www.infoworld.com/d/data-center/puppet-or-chef-the-configuration-management-dilemma-215279?page=0,0>>
- [ch2] **A. Gale**. *Getting started with Chef*. <<http://gettingstartedwithchef.com/>>
- [ch3] *Download Chef: Server and Client*. <<http://www.getchef.com/chef/install/>>
- [ch4] *Chef Cookbooks*. <<https://supermarket.getchef.com/cookbooks-directory>>
- [Co] Cobbler. <<http://cobbler.github.io/>>
- [ChLB] Learn the Chef basics on Ubuntu <<https://learn.chef.io/learn-the-basics/ubuntu/>>
- [ChSer12] How To Set Up a Chef 12 Configuration Management System on Ubuntu 14.04 Servers. <https://www.digitalocean.com/community/tutorial_series/getting-started-managing-your-infrastructure-using-chef>
- [d1] *Understanding Docker*. <<https://docs.docker.com/engine/understanding-docker/>>
- [d2] *Docker Docs*. <<https://docs.docker.com/>>
- [d4] *Docker HUB*. <<https://registry.hub.docker.com/>>
- [DO] Getting Started Managing Your Infrastructure Using Chef.
<https://www.digitalocean.com/community/tutorial_series/getting-started-managing-your-infrastructure-using-chef>
- [DoUs] Working with containers. <<https://docs.docker.com/v1.8/userguide/usingdocker/>>
- [empi] *MPI Examples*.
<<http://www.mcs.anl.gov/research/projects/mpi/usingmpi2/examples/starting/main.htm>>
- [fai1] *FAI (Fully Automatic Installation) for Debian GNU/Linux*. <<https://wiki.debian.org/FAI>>
- [fai2] *FAI (Fully Automatic Installation) project and documentation*. <<http://fai-project.org/>>
- [fai3] *El libro del administrador de Debian. Instalación automatizada*.
<<http://debian-handbook.info/browse/es-ES/stable/sect.automated-installation.html>>
- [faqmpi] *FAQ OpenMPI*. <<http://www.open-mpi.org/faq/>>
- [FG] FAI. <<http://fai-project.org/fai-guide>>
- [FOR] Foreman QuickStart. <<https://theforeman.org/manuals/1.12/index.html#2.Quickstart>>
- [KG] Kickstart
<https://access.redhat.com/documentation/en-US/Red_Hat_Enterprise_Linux/6/html/Installation_Guide/ch-kickstart2.html>
- [Kur] **Swendson, K.** *Beowulf HOWTO* (tldp). <<http://www.tldp.org/HOWTO/Beowulf-HOWTO/>>
- [kvm] *KVM*. <<https://wiki.debian.org/es/KVM>>
- [il1] *Eucalyptus, CloudStack, OpenStack and OpenNebula: A Tale of Two Cloud Models*.
<<http://opennebula.org/eucalyptus-cloudstack-openstack-and-opennebula-a-tale-of-two-cloud-models/>>

- [il2] *OpenNebula vs. OpenStack: User Needs vs. Vendor Driven.*
<<http://opennebula.org/opennebula-vs-openstack-user-needs-vs-vendor-driven/>>
- [lam] *LAM/MPI.* <<http://lam.fries.net/>>
- [LogA] *LogAnalyzer.* <<http://loganalyzer.adiscon.com/doc/>>
- [LogAR] *How to create automated daily-weekly reports.*
<<http://loganalyzer.adiscon.com/articles/how-to-create-automated-dailyweekly-reports/>>
- [LXC1] *LXC.* <<https://wiki.debian.org/LXC>>
- [LXC2] *LXC en Ubuntu.* <<https://help.ubuntu.com/lts/serverguide/lxc.html>>
- [LXC3] **S. Graber.** *LXC 1.0.* <<https://wiki.debian.org/BridgeNetworkConnections>>
- [LXC4] *LXC: Step-by-Step Guide.* <<https://www.stgraber.org/2013/12/20/lxc-1-0-blog-post-series/>>
- [mpi3] *MPI 3.*
<<http://www.mpi-forum.org/docs/mpi-3.0/mpi30-report.pdf>>
- [NR] *A Short List of DevOps Tools.* <<http://newrelic.com/devops/toolset>>
- [on1] *OpenNebula Documentation.* <<http://opennebula.org/documentation/>>
- [on2] *Quickstart: OpenNebula on Ubuntu 14.04 and KVM.*
<http://docs.opennebula.org/4.6/design_and_installation/quick_starts/qs_ubuntu_kvm.html>
- [ONDe] *OpenNebula Deployment - Front-end Installation.*
<http://docs.opennebula.org/5.0/deployment/opennebula_installation/frontend_installation.html>
- [ONNo] *OpenNebula. KVM Node Installation.*
http://docs.opennebula.org/5.0/deployment/node_installation/kvm_node_installation.html#kvm-node
- [ompi] *OpenMPI.* <<http://www.open-mpi.org/>>
- [opm] *OpenMP Exercise.* <<https://computing.llnl.gov/tutorials/openMP/exercise.html>>
- [OQRM] *OpenQRM.* <<https://sourceforge.net/projects/openqrm/>>
- [p1] *Puppet Labs Documentation.* <<http://docs.puppetlabs.com/>>
- [p2] *Puppet - Configuration Management Tool.* <<http://puppet-cmt.blogspot.com.es/>>
- [p3] *Foreman - A complete lifecycle management tool.* <<http://theforeman.org/>>
- [p4] *Foreman - Quick Start Guide.* <https://theforeman.org/manuals/1.12/quickstart_guide.html>
- [pr] *Preseeding d-i en Debian.* <<https://wiki.debian.org/DebianInstaller/Preseed>>
- [p12] *Learning Puppet.* <<http://docs.puppetlabs.com/learning/ral.html>>
- [Par] *GNU Parallel Tutorial* <https://www.gnu.org/software/parallel/parallel_tutorial.html>
- [Proc] *MPICH, High-Performance Portable MPI .* <<http://www.mpich.org/>>
- [PuDoc] *The Puppet Language.* <https://docs.puppet.com/puppet/3.7/reference/lang_visual_index.html>
- [Qe] *Qemu.* <http://wiki.qemu.org/Main_Page>
- [que] *QEMU. Guía rápida de instalación e integración con KVM y KQemu.*
<<http://wiki.debian.org/QEMU>>
- [Rad] **Radajewski, J.; Eadline, D.** *Beowulf: Installation and Administration.* TLDP.
<<http://www2.ic.uff.br/~vefr/research/clcomp/Beowulf-Installation-and-Administration-HOWTO.html>>
- [Rock] *Rocks Cluster. Base Users Guide.* <<http://central6.rocksclusters.org/roll-documentation/base/6.1.1/>>
- [tec] **Woodman, L.** *Setting up a Beowulf cluster Using Open MPI on Linux.*
<<http://techtinkering.com/2009/12/02/setting-up-a-beowulf-cluster-using-open-mpi-on-linux/>>
- [v1] *Vagrant: Getting Started.* <<https://www.vagrantup.com/docs/getting-started/>>

[v2] *Vagrant Cloud*. <<https://vagrantcloud.com/>>

[SW] *Spacewalk*. <<http://spacewalk.redhat.com/>>

[xen] *The Xen hypervisor*. <<http://www.xen.org/>>

[xeni] *Guía rápida de instalación de Xen*. <<http://wiki.debian.org/Xen>>