

Aplicación de técnicas de virtualización ligera para la evaluación de redes de comunicaciones

Trabajo Final de Estudios
Ingeniería Telemática

Enrique Fernández Sánchez
Universidad Politécnica de Cartagena

Revisión Septiembre 2021

Índice

Índice de figuras	4
Índice de tablas	4
Listado de códigos	4
Agradecimientos	5
1 Introducción	6
1.1 Contexto del trabajo	6
1.2 Objetivos	7
1.3 Descripción de los capítulos restantes de la memoria	7
2 NFV	8
3 Interfaces de red virtuales en <i>Linux</i>	9
3.1 Nombres predecibles interfaces de red. systemd networkd v197	9
3.2 enp2s0:{0,1,2...}	10
3.3 VLAN. enp2s0.{0,1,2...}	10
3.4 VLAN 802.1ad. enp2s0.{0,1,2...}	11
3.5 Pares VETH.	11
3.6 TUN/TAP	12
4 Espacio de nombres en <i>Linux</i>	13
4.1 ¿Qué es un <i>espacio de nombres</i> ?	13
4.2 ¿Cuántos <i>namespaces</i> hay?	13
4.2.1 UTS namespace	14
4.2.2 Mount namespace	15
4.2.3 Process ID namespace	16
4.2.4 Network namespace	17
4.2.4.1 Ejemplo práctico	18
4.2.5 User ID (user)	20
4.2.6 Interprocess Communication namespace (IPC)	21
4.2.7 Control group (cgroup)	21
4.2.8 Time	22
4.3 Ejemplo de uso de 'netns' usando comando ip	23
4.4 Ejemplo de uso de 'netns' usando comando unshare	25
5 Virtualización	26
5.1 Contenedores usando unshare	26
5.2 Contenedores LXC	26
5.3 Contenedores Docker	26

6	Caso práctico: Virtualización para simulación de redes	27
6.1	Evaluación de prestaciones	27
6.2	Interconexión física de diferentes red virtuales	27
6.3	Openflow OKO	27
7	Conclusiones	28
Anexos		29
Anexo 1.	Instalación ansible	29
Anexo 2.	Configuración de <i>Guest Network</i> para comunicar con la VM	31
Anexo 3.	Playbooks ansible	32
Glosario de términos		33
Bibliografía		34
Enlaces y referencias		34
Imágenes		35

Índice de figuras

1	Diagrama conexión VLANs	10
2	Diagrama VETH y netns	11
3	Diferentes namespaces en Linux y su API de acceso. (1)	21
4	Como funcionan los contenedores. (2)	22

Índice de tablas

Listado de códigos

1	Example usage UTS namespace	14
2	Example of usage mount namespace	15
3	Example of usage process id namespace	16
4	Creation persistent network namespace	17
5	Comprobar network namespaces existentes	17
6	Asignar interfaz loopback a un namespace	17
7	Ejecutar cualquier programa con un network namespace	17
8	Ejemplo configuración de NAT entre eth0 y veth	18
9	Configuración de NAT para dar Internet a un network namespace	19
10	Ejemplo de uso UID namespace	20

Agradecimientos

1 Introducción

Con el fin de concluir los estudios de grado en ingeniería telemática, es necesaria la investigación y el posterior desarrollo del *Trabajo Fin de Estudios* (TFE). Dicho trabajo, tiene como objetivo enfrentar al alumno a un proceso de investigación en el que pueda aplicar los conceptos que ha ido aprendiendo durante su paso por el grado, además pudiendo añadir puntos de innovación, y aportar soluciones nuevas a un proyecto en específico.

En este documento, recojo lo que sería mi memoria en relación al TFE. En dicho documento comentaremos los detallaremos las diferentes investigaciones realizadas sobre virtualización en sistemas Linux, como funcionan los contenedores y como utilizar virtualización ligera para la evaluación de redes de comunicaciones, en nuestro caso, de conmutación de paquetes.

1.1 Contexto del trabajo

Este proyecto nace con el objetivo de profundizar en conceptos tan importantes como pueden ser la virtualización, o los bien conocidos *contenedores*. Además, particularizamos estas tecnologías y las acercamos al campo de la telemática para utilizarlas con el fin de evaluar y simular redes de conmutación de tipo IP.

Partiendo de supuesto de que las herramientas más sencillas de simulación de redes se nos pueden quedar cortas en cuanto queremos escalar el sistema, nos topamos en que otras soluciones pueden ser mucho más *resource hungry* de lo que podríamos imaginar. Además, tampoco podemos contar con realizar dichas pruebas de manera física, ya que el presupuesto del proyecto escalaría de nivel exponencial. Esto sucede ya que deberíamos configurar cada dispositivo de manera específica, y después proceder a interconectarlos con una tecnología de red adecuada, que en nuestro caso, son tecnologías en fase de desarrollo o evaluación de prestaciones.

Conociendo las limitaciones, establecemos una serie de objetivos de cara a evaluar el viaje que realicemos durante ese proyecto, y poder evaluar posteriormente los avances que hemos podido conseguir, y además las conclusiones vinculadas a estos. Estos objetivos podrían ser:

- Análisis de diferentes opciones de simulación de redes. Tanto basadas en virtualización, como las que no.
- Comprender la diferencia entre virtualización *ligera* y virtualización *dura*.
- Profundizar en el concepto de *interfaces virtuales* en Linux.
- Definir que es el *espacio de nombres* y ejemplificar cada uno de ellos.
- Aproximación al concepto de contenedor, relacionarlo con los *espacio de nombres*.
- Aportar una solución válida para la simulación y evaluación de un sistema, utilizando la virtualización ligera.

1.2 Objetivos

Como bien hemos adelantado en el apartado anterior, el objetivo principal de este proyecto es el de analizar la situación en el ámbito de simulación y evaluación de redes de comunicación, además, concretar en aquellas soluciones que estén basadas en virtualización ligera. Por lo tanto, se pretende:

- Aprender los conceptos básicos de NFV.
- Estudiar, dentro del sistema operativo Linux, las diferentes tecnologías que nos permiten adoptar soluciones NFV.
- Definir que es *espacio de nombres* y como podemos aplicarlo para virtualizar redes.
- Desgranar el concepto amplio de contenedor.
- Desarrollar una aplicación conceptual para la evaluación de un sistema en concreto, utilizando virtualización ligera.

1.3 Descripción de los capítulos restantes de la memoria

2 NFV

NFV surge a partir de las necesidades de las operadoras para solucionar problemas de gran escala de la red. Este movimiento nace en X. Uno de los objetivos que se pretenden desde el punto de vista de las operadoras, es reducir los costes y optimizar la red. Esto sucede ya que las redes actuales sufren de tres problemas: falta de escalabilidad en momentos precisos y específicos; excasa flexibilidad de la red, en momento específicos necesitaremos modificar funcionamientos de nodos en concreto, pero no se nos posibilita; falta de innovación en la gestión de los recursos, nuevos servicios, etc. Todo esto ocurre ya que las redes de las operadoras han evolucionado muy poco.

La solución que se propone es la virtualización de las redes. Esto nace en octubre de 2012, a partir del paper: https://portal.etsi.org/nfv/nfv_white_paper.pdf

Entendemos la virtualización ligera como un tipo de virtualización de un sistema. Dicha virtualización se hace a nivel de sistema operativo, aportando que existan diferentes espacios de usuarios aislados entre sí, y todo esto estando todos bajo el mismo kernel. Las diferentes implementaciones de

3 Interfaces de red virtuales en *Linux*

Linux dispone de una selección muy diferente de interfaces de red que nos permiten, de manera sencilla, el uso de máquinas virtuales o contenedores. En este apartado vamos a mencionar las interfaces más relevantes de cara a la virtualización ligera que proponemos para el despliegue de una red virtualizada. Para obtener una lista completa de las interfaces disponibles, podemos ejecutar el siguiente comando `ip link help`.

En este trabajo, vamos a comentar la siguientes interfaces:

- `eth0:{0,1,2...}`
- `eth0.{0,1,2...}` VLAN
- `eth0.{0,1,2...}` VLAN 802.1ad
- VETH pairs
- TUN/TAP

3.1 Nombres predecibles interfaces de red. **systemd networkd v197**

Para la explicación de las diferentes interfaces, vamos a suponer que estamos utilizando el nombrado de interfaces de red antiguo. Esto proviene de que en la últimas versiones del kernel, se ha cambiado la forma en la que las interfaces de red son nombradas por Linux. Es por esto por lo que antes podíamos tener interfaces tal que `eth0` y ahora nos encontramos con la siguiente nomenclatura `enps30`. Este cambio surge ya que anteriormente se nombraban las diferentes interfaces conforme el propio ordenador estaba en la etapa de `boot`, por lo que podría pasar que a lo que nosotros entendíamos como `eth1`, en el próximo arranque fuera `eth0`, dando lugar a incontables errores en el sistema. Es por esto por lo que se empezó a trabajar en soluciones alternativas. Por ejemplo, la que hemos utilizado de ejemplo, utiliza la información aportada por la BIOS del dispositivo para catalogarlo en diferentes categorías, con su formato de nombre para cada categorías. Dichas clasificación corresponde con las siguientes:

1. Nombres incorporando Firmware/BIOS que proporcionan un número asociado a dispositivos en la placa base. (ejemplo: `eno1`)
2. Nombres incorporando Firmware/BIOS proveniente de una conexión PCI Express hotplug, con número asociado al conector. (ejemplo: `ens1`)
3. Nombres que incorporan una localización física de un conector hardware. (ejemplo: `enp2s0`)
4. Nombres que incorporan una la MAC de una interfaz. (ejemplo: `enx78e7d1ea46da`)
5. Sistema clásico e impredecible, asignación de nombres nativa del kernel. (ejemplo: `eth0`)

3.2 enp2s0:{0,1,2...}

Todas las interfaces asociadas comparten la misma dirección MAC. Cada una de ellas, recibe el nombre de *alias*. La funcionalidad principal que tienen este tipo de interfaces es la de asignar varias direcciones de IP a una misma interfaz de red.

```
$ ip addr add 192.168.56.151/24 broadcast 192.168.56.255 dev enp2s0  
label enp2s0:1
```

Sin embargo, el comando `iproute2` admite esta misma funcionalidad sin tener que crear interfaces de red extra. Para ello, solo tenemos que asociar cada IP con la interfaz de red deseada.

```
$ ip addr add 192.168.56.151/24 dev enp2s0  
$ ip addr add 192.168.56.251/24 dev enp2s0
```

3.3 VLAN. enp2s0.{0,1,2...}

Mismo concepto que la interfaz anterior, sin embargo, utilizamos el estandar 802.1q, que permite etiquetar las tramas, para crear una red lógica independiente. Es necesario que la interfaz a la que estamos asignando, sea un puerto trunk, o bien sea tagged para una VLAN específica.

```
$ ip link add link enp2s0 name enp2s0.{num} type vlan id {num}  
$ ip addr add 192.168.100.1/24 brd 192.168.100.255 dev enp2s0.{num}  
$ ip link set dev enp2s0.{num} up
```

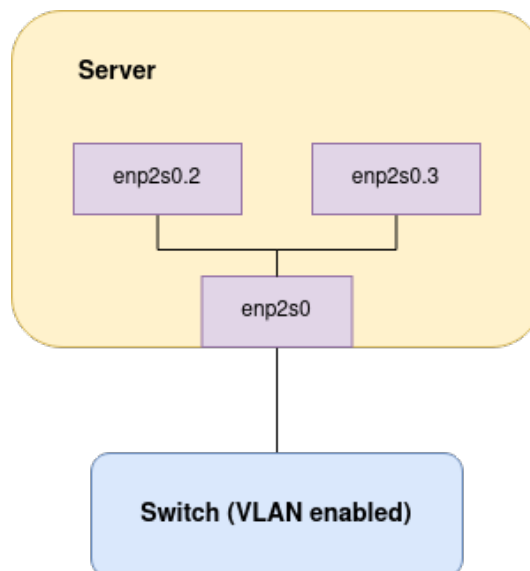


Figura 1: Diagrama conexión VLANs

3.4 VLAN 802.1ad. `enp2s0.{0,1,2...}`

3.5 Pares VETH.

Los VETH (Ethernet virtuales) son un dispositivo que forman un túnel local ethernet. El dispositivo se crea en parejas.

Los paquetes transmitidos por un extremo del VETH se reciben inmediatamente en el otro extremo. Si alguno de ellos se encuentra apagado, decimos que el link de la pareja esta también apagado.

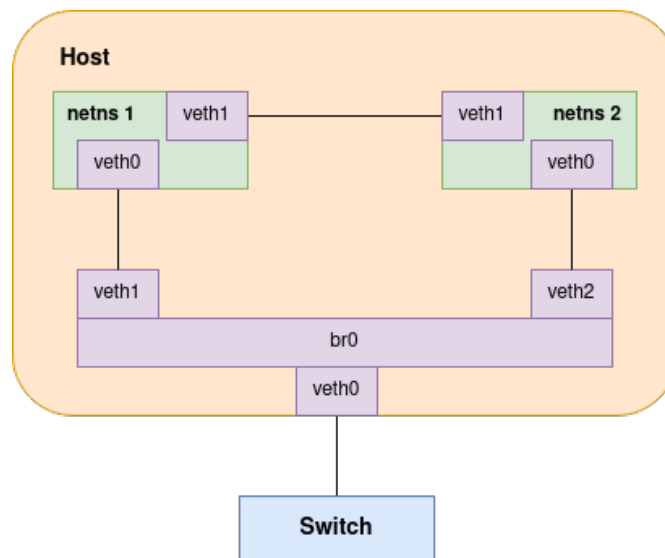


Figura 2: Diagrama VETH y netns

Es necesario utilizar VETH cuando un namespace necesita comunicarse con el host, o viceversa.

La configuración que tenemos que realizar para tener el diagrama de la Figura 2.

```
$ ip netns add net1
$ ip netns add net2
$ ip link add veth1 netns net1 type veth peer name veth1 netns net2
```

De esta manera, tendríamos creados los namespaces `net1` y `net2`, que estarían interconectados con las interfaces `veth1` y `veth2`. Nos quedaría asignar una IP a cada interfaz.

```
$ ip netns exec net1 ip addr add 10.1.1.1/24 dev veth1
$ ip netns exec net1 ip link set veth1 up
$ ip netns exec net2 ip addr add 10.1.1.2/24 dev veth1
$ ip netns exec net2 ip link set veth1 up
```

Ahora nos quedaría la parte de crear un bridge entre las diferentes interfaces virtuales que hemos creado. Para esto, podemos utilizar una interfaz tipo *bridge* de Linux, o bien podemos configurar un *bridge* usando *Open VSwitch*.

```
$ ip link add veth1_br type veth peer name veth0 netns net1
$ ip link add veth2_br type veth peer name veth0 netns net2
$ ovs-vsctl add-br ovsbr0
$ ovs-vsctl add-port ovsbr0 veth1_br
$ ovs-vsctl add-port ovsbr0 veth2_br
$ ovs-vsctl add-port ovsbr0 <outbound interface>
```

Añadimos las direcciones IP que nos faltan:

```
$ ip addr add 10.1.1.10/24 dev veth1_br
$ ip link set veth1_br up
$ ip addr add 10.1.1.20/24 dev veth2_br
$ ip link set veth2_br up
$ ip netns exec net1 ip addr add 10.1.1.15/24 dev veth0
$ ip netns exec net1 ip link set veth0 up
$ ip netns exec net2 ip addr add 10.1.1.25/24 dev veth0
$ ip netns exec net2 ip link set veth0 up
$ ip netns exec net1 ip link set lo up
$ ip netns exec net2 ip link set lo up
```

3.6 TUN/TAP

Permite conectar una aplicación en específico con el kernel, utilizando un dispositivo Ethernet. Una aplicación puede acceder al tuntap o bien crearlo ella misma.

TUN trabaja a nivel IP (capa 3), tiene por función principal definir túneles. No entiende nada de ARP (ni nada de nivel 2).

TAP trabaja solo a nivel de MAC (entiende ARP, etc.) y sirve directamente para un bridge virtual. Los túneles los puede crear directamente el comando ip, por lo que la interfaz TUN apenas se utiliza. Sin embargo, TAP es muy utilizado para crear un *bridge*.

Un *bridge* puede incorporar interfaces TAP. En este caso, la aplicación que lo monitoriza puede ser un hipervisor (Virtualbox, Qemu, libvirt), es decir, el *bridge* está dentro del kernel.

4 Espacio de nombres en *Linux*

4.1 ¿Qué es un *espacio de nombres*?

Los *espacios de nombres*, o también llamados, *namespaces*, son una característica del kernel de Linux que permite gestionar los recursos del kernel, pudiendo limitarlos a un proceso o grupo de procesos. Suponen una base de tecnología que aparece en las técnicas de virtualización más modernas (como puede ser Docker, Kubernetes, etc). A un nivel alto, permiten aislar procesos respecto al resto del kernel.

El objetivo de cada *namespaces* es adquirir una característica global del sistema como una abstracción que haga parecer a los procesos de dentro del *namespace* que tienen su propia instancia aislada del recurso global.

4.2 ¿Cuántos *namespaces* hay?

El kernel ha estado en constante evolución desde que 1991, cuando Linus Torvalds comenzó el proyecto, actualmente sigue muy activo y se siguen añadiendo nuevas características. El origen de los namespaces se remonta a la versión del kernel 2.4.19, lanzada en 2002. Conforme fueron pasando los años, más tipos diferentes de namespaces se fueron añadiendo a Linux. El concepto de *User namespaces*, se consideró terminado con la versión 3.9.

Actualmente, tenemos 8 tipos diferentes de namespaces, siendo el último añadido en la versión 5.8 (lanzada el 2 de Agosto de 2020).

1. UTS (hostname)
2. Mount (mnt)
3. Process ID (pid)
4. Network (net)
5. User ID (user)
6. Interprocess Communication (ipc)
7. Control group (cgroup)
8. Time

4.2.1 UTS namespace

El tipo más sencillo de todos los namespaces. La funcionalidad consiste en controlar el hostname asociado del ordenador, en este caso, del proceso o procesos asignados al namespace. Existen tres diferentes rutinas que nos permiten obtener y modificar el hostname:

- *sethostname()*
- *setdomainname()*
- *uname()*

En una situación normal sin namespaces, se modificaría una String global, sin embargo, si estamos dentro de un namespace, los procesos asociados tienen su propia variable global asignada.

Un ejemplo muy básico de uso de este namespace podría ser el siguiente:

Listado 1: Example usage UTS namespace

```
$ sudo su                                # super user
$ hostname                               # current hostname
> arch-linux
$ unshare -u /bin/sh                     # shell with UTS namespace
$ hostname new-hostname                  # set hostname
$ hostname                               # check hostname of the shell
> new-hostname
$ exit                                   # exit shell and namespace
$ hostname                               # original hostname
> arch-linux
```

En el ejemplo planteado, vemos que utilizamos el comando `unshare`. Utilizando la documentación de dicho comando, `man unshare`. Podemos deducir lo siguiente:

- Ejecuta un programa con algunos namespaces diferentes del host.
- En los parámetros podemos especificar cual o cuales namespaces queremos desvincular.
- Tenemos que especificar la ruta del ejecutable que queremos aislar
- La sintaxis sería tal que: `unshare [options] <program> [<argument>...]`

4.2.2 Mount namespace

Un *mount namespace* (*mnt*) supone otro tipo de espacio de nombres, en este caso relacionado con los *mounts* de nuestro sistema. Lo primero es entender a que nos referimos cuando hablamos de *mount*. *Mount*, o montaje, hace referencia a conectar un sistema de archivos adicional que sea accesible para el sistema de archivos actual de un ordenador. Un *mount*, tiene asignado lo que se llama *mount point*, que corresponde con el directorio en el que está accesible el sistema de archivo que previamente hemos montado.

Por lo tanto, un namespace de tipo *mount* nos permite modificar un sistema de archivos en concreto, sin que el host pueda ver y/o acceder a dicho sistema de archivos. Un ejemplo básico de esta funcionalidad podría ser la siguiente:

Listado 2: Example of usage mount namespace

```
$ sudo su                                # run a shell in a new mount namespace
$ unshare -m /bin/sh
$ mount --bind /usr/bin/ /mnt/
$ ls /mnt/cp
> /mnt/cp
$ exit                                  # exit the shell and close namespace
$ ls /mnt/cp
> ls: cannot access '/mnt/cp': No such file or directory
```

Como vemos en el ejemplo, dentro del namespaces lo que hacemos es crear un *mount* de tipo *bind*, que tiene por función que un archivo de la máquina host se monte en un directorio en específico, en este caso, un directorio unicamente del programa que hemos asignado al namespace. Otro ejemplo de uso de estos namespaces es crear un sistema de archivos temporal que solo sea visible para ese proceso.

4.2.3 Process ID namespace

Para entender en que consiste este namespace, primero tenemos que conocer la definición de *process id* dentro del Kernel. En este caso, *process id* hace referencia a un número entero que utiliza el Kernel para identificar los procesos de manera unívoca.

Concretando, aísla el namespace de la ID del proceso asignado, dando lugar a que, por ejemplo, otros namespaces puedan tener el mismo PID. Esto nos lleva a la situación de que un proceso dentro de un *PID namespace* piense que tiene asignado el ID "1", mientras que en la realidad (en la máquina host) tiene otro ID asignado.

Listado 3: Example of usage process id namespace

```
$ echo $$                # PID de la shell
$ ls -l /proc/$$/ns      # ID espacios de nombres
$ sudo unshare -f --mount-proc -p /bin/sh
$ echo $$                # PID de la shell dentro del ns
$ ls -l /proc/$$/ns      # nuevos ID espacio de nombres
$ ps

$ ps -ef                 # ejecutar en una shell fuera del ns. Comparar PID
$ exit
```

Si ejecutamos el ejemplo, lo que podemos comprobar es que el ID del proceso que está dentro del namespaces (`echo $$`), no coincide con el proceso que podemos ver de la máquina host (`ps -ef | grep /bin/sh`). Más concretamente, el primer proceso creado en un PID namespace recibirá el pid número 1, y además de un tratamiento especial ya que supone un `init process` dentro de ese namespace.

4.2.4 Network namespace

Este namespaces nos permite aislar la parte red de una aplicación o proceso que nosotros elijamos. Con esto conseguimos que el *stack* de red de la máquina host sea diferente al que tenemos en nuestro namespace. Debido a esto, el namespace crea una interfaz virtual, conjunto con el resto de necesidades para conformar un stack de red completo (tabla de enrutamiento, tabla ARP, etc...).

Para crear un *namespace* de tipo *network*, y que este sea persistente, utilizamos la *tool* ip (del *package* iproute2).

Listado 4: Creation persistent network namespace

```
$ ip netns add ns1
```

Este comando creará un network namespace llamado ns1. Cuando se crea dicho namespace, el comando ip realiza un montaje tipo bind en la ruta /var/run/netns, permitiendo que el namespace sea persistente aún sin tener un proceso asociado.

Listado 5: Comprobar network namespaces existentes

```
$ ls /var/run/netns
or
$ ip netns
```

Como ejemplo, podemos proceder a añadir una interfaz de *loopback* al namespace que previamente hemos creado:

Listado 6: Asignar interfaz loopback a un namespace

```
$ ip netns exec ns1 ip link dev lo up
$ ip netns exec ns1 ping 127.0.0.1
> PING 127.0.0.1 (127.0.0.1) 56(84) bytes of data.
> 64 bytes from 127.0.0.1: icmp_seq=1 ttl=64 time=0.115 ms
```

La primera línea de este ejemplo, corresponde con la directiva que le dice al namespace que "levante" la interfaz de loopback. La segunda línea, vemos como el namespace ns1 ejecuta el ping a la interfaz de loopback (el loopback de ese namespace).

Es importante mencionar, que aunque existen más comandos para gestionar las redes dentro de linux (como pueden ser ifconfig, route, etc), el comando ip es el considerado sucesor de todos estos, y los anteriores mencionados, dejarán de formar parte de Linux en versiones posteriores. Un detalle a tener en cuenta con el comando ip, es que es necesario tener privilegios de administrador para poder usarlo, por lo que deberemos ser root o utilizar sudo.

Por lo tanto, utilizando el comando ip, podemos recapitular que si utilizamos la siguiente directiva, podemos ejecutar el comando que nosotros indiquemos, pero dentro del network namespace que previamente hemos creado.

Listado 7: Ejecutar cualquier programa con un network namespace

```
$ ip netns exec <network-namespace> <command>
```

4.2.4.1 Ejemplo práctico

Una de las problemáticas que supone el uso de los network namespaces, es que solo podemos asignar **una interfaz** a **un namespace**. Suponiendo el caso en el que el usuario root tenga asignada la interfaz eth0 (identificador de una interfaz de red física), significaría que solo los programas en el namespace de root podrán acceder a dicha interfaz. En el caso de que eth0 sea la salida a Internet de nuestro sistema, pues eso conllevaría que no podríamos tener conexión a Internet en nuestros namespaces. La solución para esto reside en los **veth-pair**.

Un veth-pair funciona como si fuera un cable físico, es decir, interconecta dos dispositivos, en este caso, interfaces virtuales. Consiste en dos interfaces virtuales, una de ellas asignada al root namespace, y la otra asignada a otro network namespace diferente. Si a esta arquitectura le añadimos una configuración de IP válida y activamos la opción de hacer NAT en el eth0 del host, podemos dar conectividad de Internet al network namespace que hayamos conectado.

Listado 8: Ejemplo configuración de NAT entre eth0 y veth

```
# Remove namespace if exists
$ ip netns del ns1 &>/dev/null

# Create namespace
$ ip netns add ns1

# Create veth link
$ ip link add v-eth1 type veth peer name v-peer1

# Add peer-1 to namespace.
$ ip link set v-peer1 netns ns1

# Setup IP address of v-eth1
$ ip addr add 10.200.1.1/24 dev v-eth1
$ ip link set v-eth1 up

# Setup IP address of v-peer1
$ ip netns exec ns1 ip addr add 10.200.1.2/24 dev v-peer1
$ ip netns exec ns1 ip link set v-peer1 up
# Enabling loopback inside ns1
$ ip netns exec ns1 ip link set lo up

# All traffic leaving ns1 go through v-eth1
$ ip netns exec ns1 ip route add default via 10.200.1.1
```

Siguiendo el ejemplo propuesto, llegamos hasta el punto en el que el tráfico saliente del namespace ns1, será redirigido a v-eth1. Sin embargo, esto no es suficiente para tener conexión a Internet. Tenemos que configurar el NAT en el eth0.

Listado 9: Configuración de NAT para dar Internet a un network namespace

```
# Share internet access between host and NS

# Enable IP-forwarding
$ echo 1 > /proc/sys/net/ipv4/ip_forward

# Flush forward rules, policy DROP by default
$ iptables -P FORWARD DROP
$ iptables -F FORWARD

# Flush nat rules.
$ iptables -t nat -F

# Enable masquerading of 10.200.1.0 (ip of namespaces)
$ iptables -t nat -A POSTROUTING -s 10.200.1.0/255.255.255.0 -o eth0
    -j MASQUERADE

# Allow forwarding between eth0 and v-eth1
$ iptables -A FORWARD -i eth0 -o v-eth1 -j ACCEPT
$ iptables -A FORWARD -o eth0 -i v-eth1 -j ACCEPT
```

Si todo lo hemos configurado correctamente, ahora podríamos realizar un ping hacia Internet, y este nos debería resultar satisfactorio.

```
$ ip netns exec ns1 ping google.es
> PING 8.8.8.8 (8.8.8.8) 56(84) bytes of data.
> 64 bytes from 8.8.8.8: icmp_seq=1 ttl=50 time=48.5ms
> 64 bytes from 8.8.8.8: icmp_seq=2 ttl=50 time=58.5ms
```

Aún así, no resulta muy cómodo el utilizar `ip netns exec` seguido de la aplicación a utilizar. Es por esto por lo que es común ejecutar dicho comando para asignar el network namespace a una shell. Esto sería tal que así:

```
$ ip netns exec ns1 /bin/bash
```

Utilizaremos `exit` para salir de la shell y abandonar el network namespace.

4.2.5 User ID (user)

Cada sistema dispone de una manera de monitorizar que usuario es el dueño de cada archivo. Esto permite al sistema restringir el acceso a aquellos archivos que consideramos sensibles. Además, bloquea el acceso entre diferentes usuarios dentro del mismo sistema. Para el usuario, este identificador de usuarios se muestra como el usuario que en ese momento está conectado, sin embargo, para nuestro sistema, el identificador de usuario esta compuesto por una combinación arbitraria de caracteres alfanuméricos. Con el fin de mantener el monitoreo correctamente, hay un proceso encargado de transformar esos caracteres a un número específico de identificación (UID), como por ejemplo sería 1000. Es este valor el que se asocia con los archivos creados por este usuario. Esto nos aporta la ventaja de que, si un usuario cambia su nombre, no es necesario reconstruir el sistema de archivos, ya que su UID sigue siendo 1000.

Si por ejemplo queremos ver el UID del usuario que estamos usando en este momento, podemos ejecutar: `echo $UID`, el cual nos devolverá el número asociado a nuestro usuario, en mi caso es el 1000.

Además de diferenciar entre los IDs de usuarios (UID), también se nos permite separar entre IDs de grupos (GID). En linux, un grupo sirve para agrupar usuarios de modo que un grupo puede tener asociado un privilegio que le permite usar un recurso o programas.

Por lo tanto, el namespace de UID, lo que nos permite es tener un UID y GID diferente al del host.

Listado 10: Ejemplo de uso UID namespace

```
$ ls -l /proc/$$/ns                                # espacios de nombres originales
$ id
> uid=1000(user) gid=1000(user) groups=1000(user), ...
$ unshare -r -u bash                                # Crea un namespace de tipo usuario, programa bash
$ id
> uid=0(root) gid=0(root) groups=0(root),65534(nobody)
$ cat /proc/$$/uid_map
>          0          1000          1
$ cat /etc/shadow                                # No nos deja acceder
> cat: /etc/shadow: Permission denied
$ exit
```

Como vemos en el ejemplo, el UID de usuario difiere de la máquina host. Dentro del namespace, tenemos UID 0, sin embargo, eso no significa que podamos acceder a los archivos con UID 0 de la máquina host, ya que en verdad lo que hace el namespace es *mapear* el UID 1000 al 0.

4.2.6 Interprocess Communication namespace (IPC)

Este namespace supone uno de los más técnicos, complicados de entender y explicar. IPC (Inter-process communication) controla la comunicación entre procesos, utilizando zonas de la memoria que están compartidas, colas de mensajes, y semáforos. La aplicación más común para este tipo de gestión es el uso en bases de datos.

4.2.7 Control group (cgroup)

Los grupos de control, cgroups, de Linux suponen un mecanismo para controlar los diferentes recursos de nuestro sistema. Cuando CGroups están activos, pueden controlar la cantidad de CPU, RAM, acceso I/O, o cualquier faceta que un proceso puede consumir. Además, permiten definir jerarquías en las que se agrupan, de manera en la que el administrador del sistema puede definir como se asignan los recursos o llevar la contabilidad de los mismos.

Por defecto, los CGroups se crean en el sistema de archivos virtual `/sys/fs/cgroup`. Si creamos un namespace de tipo CGroups, lo que estamos haciendo es mover el espacio de archivos virtual de dicho CGroup. Un ejemplo de esto sería, creamos un CGroup namespace en el directorio `/sys/fs/cgroup/mycgroup`. El host verá lo siguiente `/sys/fs/cgroup/mycgroup/{group1, group2, group3}`, sin embargo, el namespace solo verá `{group1, group2, group3}`. Esto es así ya que aporta seguridad a un namespace ya que los procesos del namespace solo pueden acceder a su sistema de archivos.

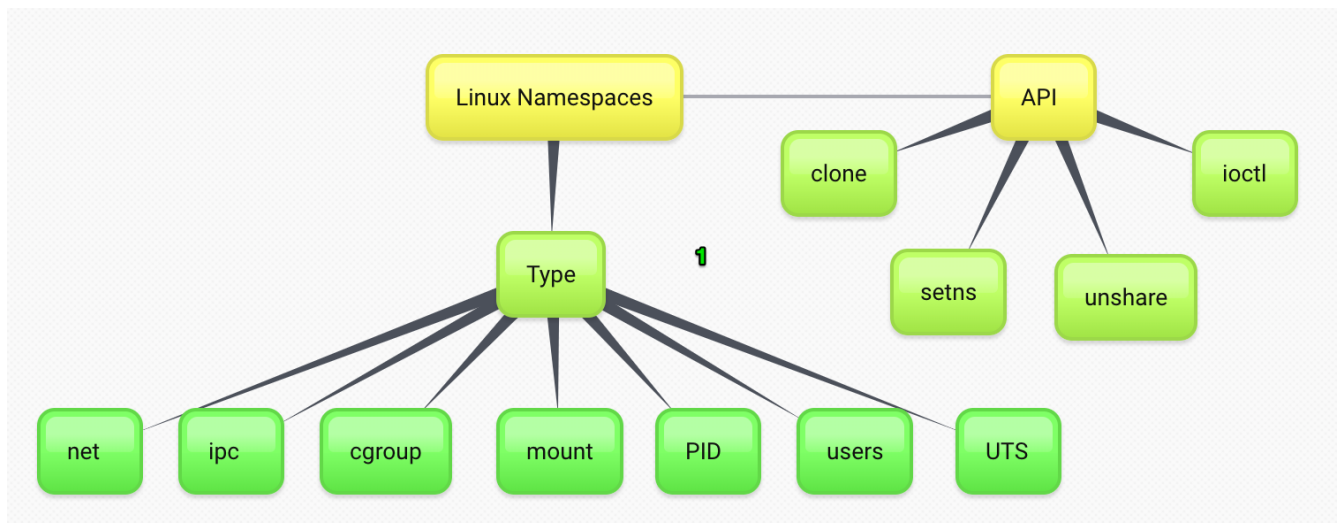


Figura 3: Diferentes namespaces en Linux y su API de acceso. (1)

4.2.8 Time

Por último, nos queda el namespaces asociado al tiempo. Este namespace fue propuesto para que se incorporara al kernel de Linux en 2018 y en enero de 2020 fue añadido a la versión mainline de Linux. Apareció en la release 5.6 del kernel de Linux.

El namespace time, permite que por cada namespace que tengamos, podamos crear desfases entre los relojes monotónicos (CLOCK_MONOTONIC) y de boot (CLOCK_BOOTTIME), de la máquina host. Esto permite que dentro de los contenedores se nos permita cambiar la fecha y la hora, sin tener que modificar la hora del sistema host. Además, supone una capa más de seguridad, ya que no estamos vinculando directamente la hora a los relojes físicos de nuestro sistema.

Un namespace de tipo time, es muy similar al namespace de tipo PID en la manera de como lo creamos. Utilizamos el comando `unshare -T`, y mediante una `systemcall` se nos creará un nuevo time namespace, pero no lo asocia directamente con el proceso. Tenemos que utilizar `setns` para asociar un proceso a un namespace, además todos los procesos dependientes también tendrán asignado dicho namespace.

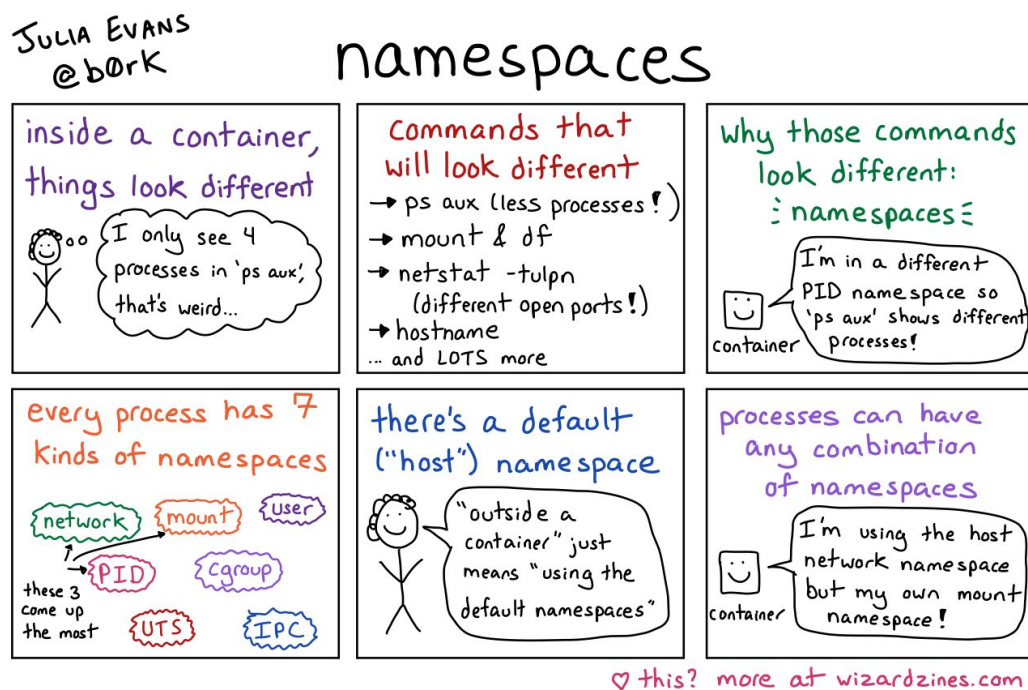


Figura 4: Como funcionan los contenedores. (2)

4.3 Ejemplo de uso de 'netns' usando comando **ip**

En este apartado, vamos a detallar un ejemplo de como funciona el comando IP manejando los 'network namespaces'.

Creamos los *network namespaces*, en este caso, con nombre h1 y h2. El sistema no crea directamente el namespace, lo que en realidad hace es definirlos en el sistema. El *network namespace* se crea cuando una aplicación se asocia a el.

```
$ ip netns add h1
$ ip netns add h2
```

Si utilizamos el comando `ip netns`, nos mostrará los netns existentes. Como es un sub-comando de `ip`, muestra los ns que el comando `lsns` no muestra.

Procedemos a asociar una aplicación a cada netns. Utilizamos `bash`.

```
$ ip netns exec h1 bash
$ ip netns exec h2 bash
```

Ahora, si que podemos utilizar el comando `lsns`. Comprobamos que si nos aparecen los ns que hemos creado, cosa que antes de asociar una aplicación al ns, no pasaba.

El comando `ip` crea automaticamente un *nsfs* para poder colocar los archivos de configuración del *netns*. Para ello, debe crearse el directorio `/etc/netns/h1` y poner en el los archivos de configuración de la red.

```
$ mkdir /etc/netns/h1
$ echo "nameserver 8.8.8.8" > /etc/netns/h1/resolv.conf
```

En este momento, tenemos el ns configurado con DNS. Nos quedaría realizar la conexión entre el ethernet físico de nuestro *host* y las interfaces de nuestros namespaces. Para ello, vamos a utilizar un conmutador virtual, en este caso Open vSwitch.

```
$ systemctl enable --now openvswitch.service
```

Creamos un *brige* utilizando el OpenvSwitch.

```
$ $ ovs-vsctl add-br s1
```

Utilizando el comando `ip`, creamos las interfaces virtuales de ethernet y las asignamos a sus namespaces.

```
$ ip link add h1-eth0 type veth peer name s1-eth1
$ ip link add h2-eth0 type veth peer name s1-eth2
$ ip link set h1-eth0 netns h1
$ ip link set h2-eth0 netns h2
```

Utilizando el comando `ovs-vsctl`, asignamos al *bridge* el otro par ethernet que hemos creado para cada namespace.

```
$ ovs-vsctl add-port s1 s1-eth1
$ ovs-vsctl add-port s1 s1-eth2
```

Verificamos que el controlador sea *standalone*, así el switch se comportará como un *learning-switch*.

```
$ ovs-vsctl set-fail-mode br0 standalone
```

Como la conexión es desde localhost al exterior, entendemos que es una conexión fuera de banda.

```
$ ovs-vsctl set controller br0 connection-mode=out-of-band
```

En este momento, tenemos todos configurado a falta de habilitar las diferentes interfaces de nuestra topología.

```
$ ip netns exec h1 ip link set h1-eth0 up
$ ip netns exec h1 ip link set lo up
$ ip netns exec h1 ip add add 10.0.0.1/24 dev h1-eth0
$ ip netns exec h2 ip link set h2-eth0 up
$ ip netns exec h2 ip link set lo up
$ ip netns exec h2 ip add add 10.0.0.2/24 dev h2-eth0
$ ip link set s1-eth1 up
$ ip link set s1-eth2 up
```

Ahora tenemos todas las interfaces configuradas, el switch activado y el sistema interconectado, por lo que podemos ejecutar un ping en una de las terminales de los namespaces para verificar la topología.

```
$ ip netns exec h1 ping -c4 10.0.0.2
```

Si queremos revertir todas las configuraciones que hemos hecho, lo que tenemos que hacer es ejecutar los siguientes comandos:

```
$ ovs-vsctl del-br s1
$ ip link delete s1-eth1
$ ip link delete s1-eth2
$ ip netns del h1
$ ip netns del h2
```


4.4 Ejemplo de uso de 'netns' usando comando **unshare**

En el ejemplo anterior, utilizabamos el comando `ip` para manejar los netns, sin embargo, eso nos limitaba los tipos namespaces que queriamos asignar a nuestro namespace. En contra partida a esto, el comando `unshare` nos da más libertad a la hora de crear los namespaces.

Utilizando `unshare`, no podemos ponerle un nombre, pero sí que nos permite asociarlo a un archivo, que montará con tipo `bind`. Esto nos permitirá utilizar en namespace aunque no haya ningún proceso corriendo en el, para ello podemos utilizar el comando `nsenter`.

```
$ touch /var/net-h1
$ touch /var/uts-h1
$ unshare --net=/var/net-h1 --uts=/var/uts-h1 /bin/bash
```

Utilizando el comando `nsenter` podemos ejecutar comandos dentro del namespace.

```
$ nsenter --net=/var/net-h1 --uts=/var/uts-h1 hostname h1
$ nsenter --net=/var/net-h1 --uts=/var/uts-h1 ip address
```

Para destruir el namespace, lo que tendremos que hacer es desmontar los archivos asignados a dicho namespace.

```
$ umount /var/net-h1
$ umount /var/uts-h1
```

Como será común necesitar más de un namespace, de ahora en adelante tendremos que utilizar los comandos `unshare` y `nsenter`.

5 Virtualización

5.1 Contenedores usando **unshare**

5.2 Contenedores LXC

5.3 Contenedores Docker

6 Caso práctico: Virtualización para simulación de redes

6.1 Evaluación de prestaciones

6.2 Interconexión física de diferentes red virtuales

6.3 Openflow OKO

7 Conclusiones

Anexos

Anexo 1. Instalación de **ansible** para automatizar una VM

En el caso de que queramos utilizar la herramienta ansible para configurar una VM, tendremos que seguir los siguientes pasos.

1. Asegurarnos de que tenemos instalado el programa en el host. La principal dependencia de ansible es Python.

Para instalar ansible en diferentes distribuciones sería tal que:

- Ubuntu 21.04:

```
$ sudo apt install -y software-properties-common
$ sudo add-apt-repository --yes --update ppa:ansible/ansible
$ sudo apt update
$ sudo apt install -y ansible
```

- Fedora 32:

```
$ sudo dnf update
$ sudo dnf install ansible
```

- OpenSUSE:

```
$ zypper in ansible
```

- Arch Linux:

```
$ sudo pacman -S ansible
```

2. Creamos una clave SSH para utilizarla como método de autenticación con la maquina virtual.

```
(host)$ ssh-keygen -t ed25519 -C "Host Ansible ssh key"
```

3. Utilizando nuestro gestor de máquinas virtuales, encendemos la VM. Es importante que nos aseguremos que tiene internet, por lo que configuramos que la interfaz de red sea de tipo *NAT*. Dentro de la máquina, buscaremos que ip tiene asignada con el comando `ip address`. Desde este momento, dejaremos la máquina virtual encendida.

4. Procedemos a copiar la clave pública SSH que hemos generado en el paso dos. Utilizamos el siguiente comando, sustituyendo con el usuario e IP específico de la máquina. Tendremos que configurar una red tipo *Host-only network* para poder comunicarnos correctamente con nuestra máquina virtual, para ello podemos seguir los pasos detallados en [15]

```
(host)$ ssh-copy-id -i $HOME/.ssh/id_ed25519.pub <user>@<IP VM>
```

5. En este momento ya podríamos conectar con la máquina virtual utilizando `ansible`. Sin embargo, por comodidad, lo que vamos a hacer es crear un archivo que funcione con inventario de servidores. Nos servirá para guardar las direcciones IP de los servidores en los que queremos ejecutar comandos remotos con `ansible`. Para ello, creamos un archivo `inventory`. Como ejemplo, dicho archivo puede ser tal que:

```
## VMs locals
[virtualbox]
10.0.100.1
10.0.100.2

## Server SSH
[server]
192.168.1.200
```

6. Para probar la conectividad. Nos aseguramos de que tenemos la máquina virtual encendida y que además, hemos copiado la clave `ssh` y tenemos su IP añadida a nuestro inventario. Después, procedemos a ejecutar el siguiente comando en la máquina `host`:

```
(host)$ ansible -i inventory -m ping
```

En la consola nos aparecerá la información de cada uno de las IP que habrá probado para ejecutar el comando remoto `ping`.

Anexo 2. Configuración de *Guest Network* para comunicar con la VM

Para comunicar nuestra máquina host [15], con una *virtual machine* es necesario que estén en una misma red, es por esto por lo que vamos a crear una red en específico para ello. En el programa *Virtualbox*, esta funcionalidad recibe el nombre de *Host-only*.

Es importante que cuando queramos configurar una máquina virtual con esta funcionalidad que siempre pongamos el adaptador de red 1 como el que tendrá la comunicación *guest*, mientras que el adaptador 2 será el que tendrá la conexión a internet via *NAT*.

Procedemos a enumerar los pasos a seguir para configurar una máquina virtual con conectividad con el host, utilizando *Virtualbox*.

1. Abrimos el programa *Virtualbox*. Navegamos a la sección: **File** → **Preferences**.
2. Seleccionamos la tabla de ***Host-only Networks***. Procedemos a pulsar el botón + para añadir una nueva red.
3. Asignamos las diferentes IP que consideremos. Es importante que dejemos activado el servidor DHCP, así la máquina virtual tendrá una IP válida dentro de la red.
4. Ahora, lo que tenemos que hacer es añadir una interfaz de red a nuestra máquina virtual. Para ello, vamos a las **Settings** de la máquina virtual. En el apartado de Network, cambiamos el **Adapter 1** a **Host-only Adapter**, y asignamos una segunda interfaz de red con NAT (Adapter 2 attached to NAT).
5. Para comprobar que lo hemos hecho correctamente. Arrancamos la máquina virtual, ejecutamos el comando `ip address`. Nos deberían aparacer las diferentes interfaces que previamente hemos configurado, es decir, una interfaz de loopback, otra que corresponde al Host-only network, y por último, una que corresponda con la interfaz de NAT.

Anexo 3. Playbooks **ansible**

Glosario de términos

- NFV. *Network Function Virtualization*, o bien, Virtualización de funciones de red.
- Namespace. *Espacio de nombres*
- Linux. *Sistema operativo tipo UNIX, de código abierto, multiplataforma, multiusuario y multitarea*
- Kernel de linux. *Núcleo del sistema operativo Linux*
- PID. *Process Identifier*
- root. *Cuenta superusuario del sistema operativo Linux*
- veth-pair. *Virtual Ethernet Pair*
- UID. *Identificador de usuario*
- GID. *Identificador de grupo*
- ns. *Network namespace*

Bibliografía

Enlaces y referencias

1. *Namespaces*
2. Tutorial: Espacio de nombres en Linux
3. *Time namespaces coming to linux*
4. *Container is a lie. Namespaces*
5. *Namespaces. Uso de cgroups.*
6. *Introduction to Network Namespaces*
7. *Build a container by hand: the mount namespace*
8. Identificador de procesos (*process id*)
9. *Linux PID namespaces work with containers*
10. *Network Namespaces*
11. *Introduction to Linux interfaces for virtual networking*
12. *Introducción a los grupos de control (cgroups) de Linux*
13. *Time namespaces*
14. *Network namespaces. Assign and configure*
15. How to configure Network Between Guest VM and Host in Virtualbox
16. How to install ansible on fedora for it and server automation
17. *Herramientas de virtualización libres para sistemas GNU/Linux*
18. Fundamentos de Docker

Imágenes

1. Namespaces y API de acceso.
2. How containers work.