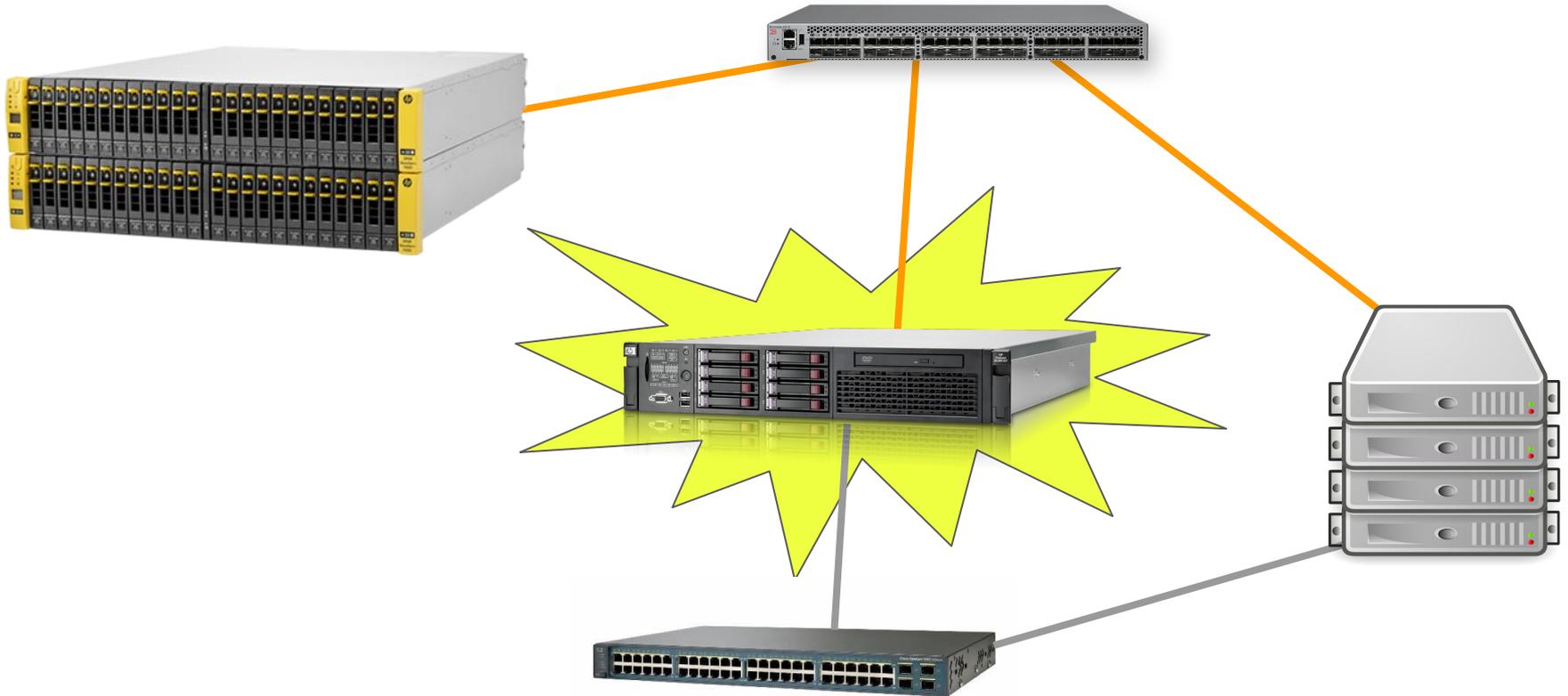


# Implantación de ZFS para servicio de ficheros corporativo en la UMU

Pablo López Mozas - Universidad de Murcia  
Jose Fco. Hidalgo Céspedes - Universidad de Murcia

- Introducción al entorno de almacenamiento
- Decisión del producto elegido: ZFS
- Descripción de la nueva infraestructura NAS
- Estructura de ZFS
  - Pools y recursos
  - Tareas habituales
  - Scrubs
  - Cuotas
  - Failover cluster
  - Snapshots y réplicas
- Monitorización
- Servicios alrededor de NFS
- Tareas futuras





ON LINUX



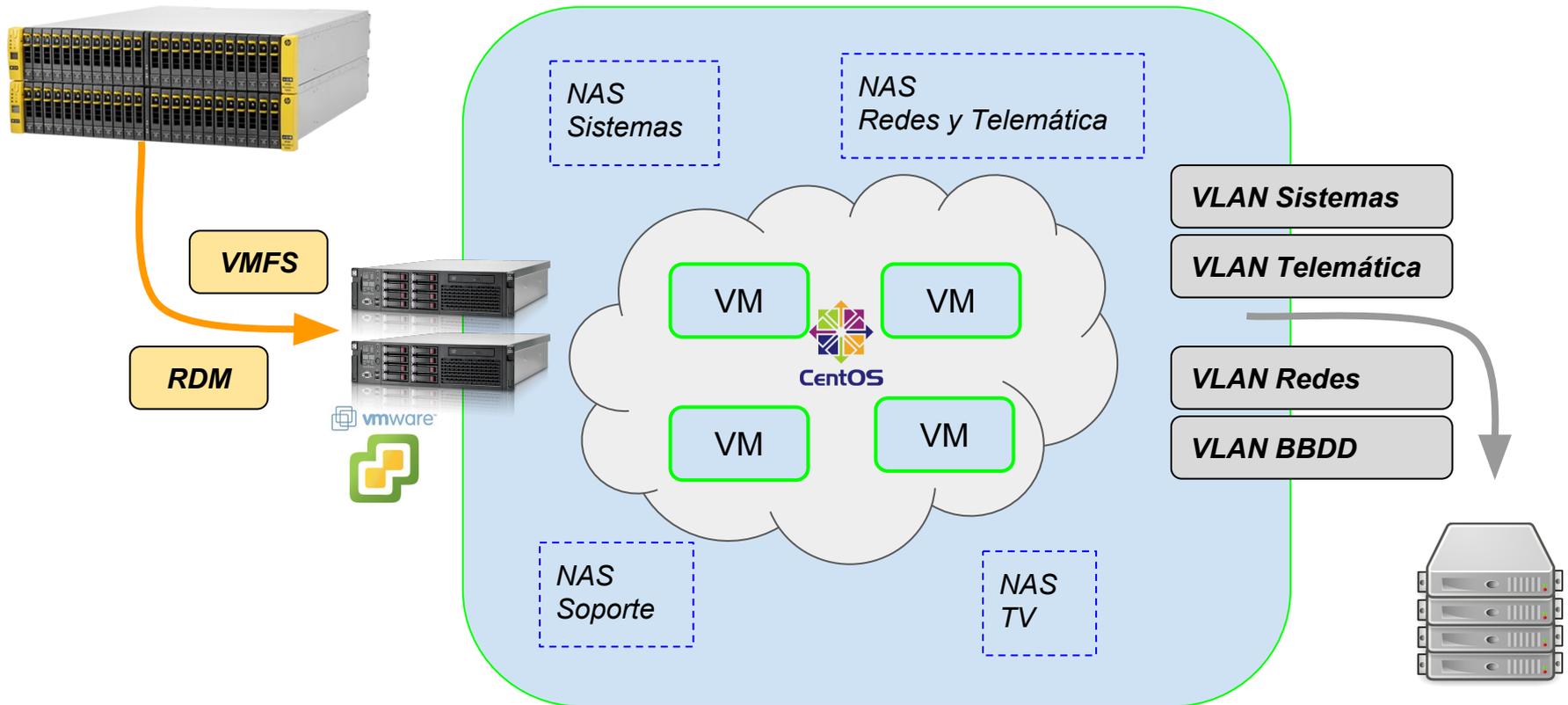
## A favor de ZFS

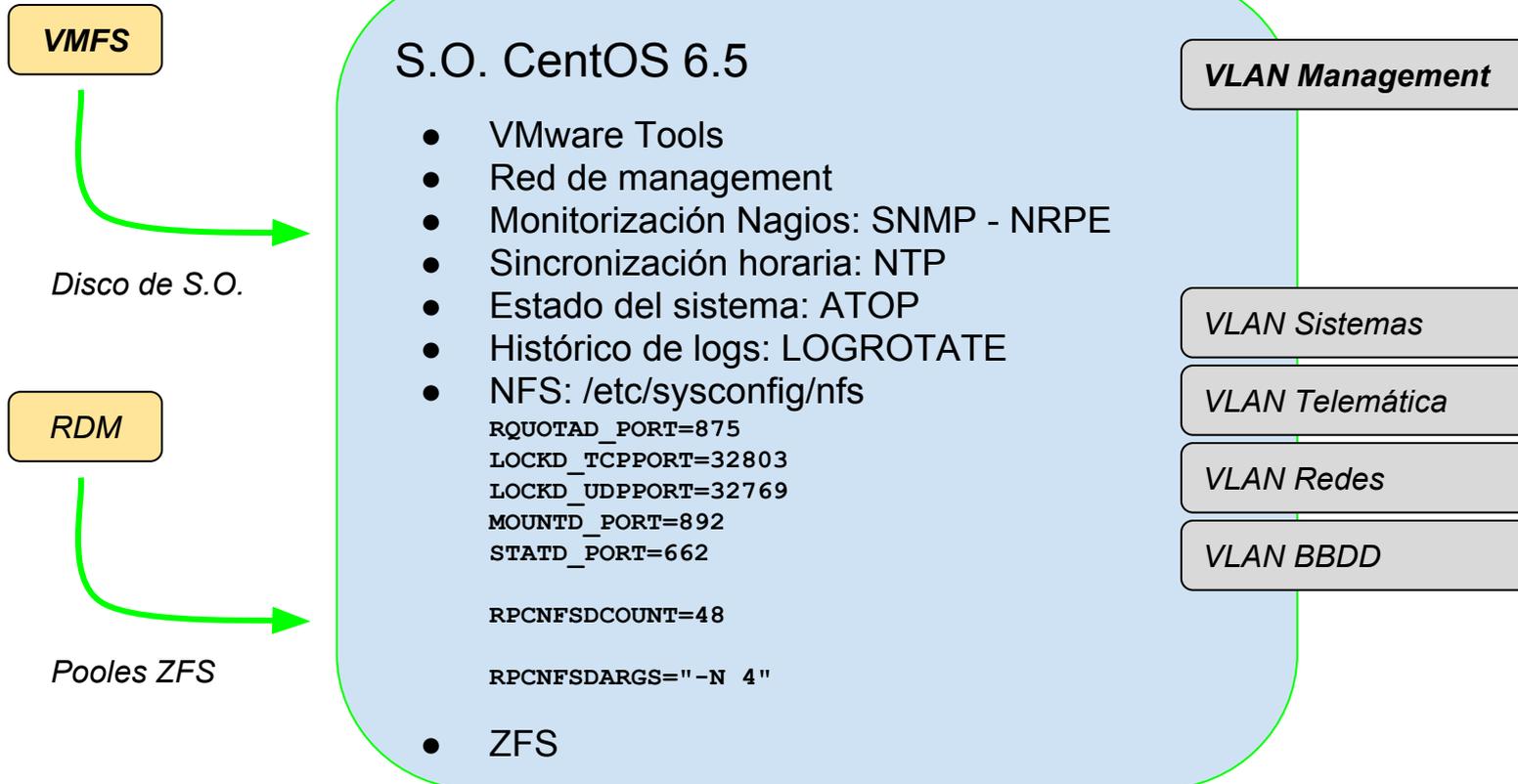
- Software libre
- Conocimiento previo
- Características esperables
- Buenas críticas y recomendaciones

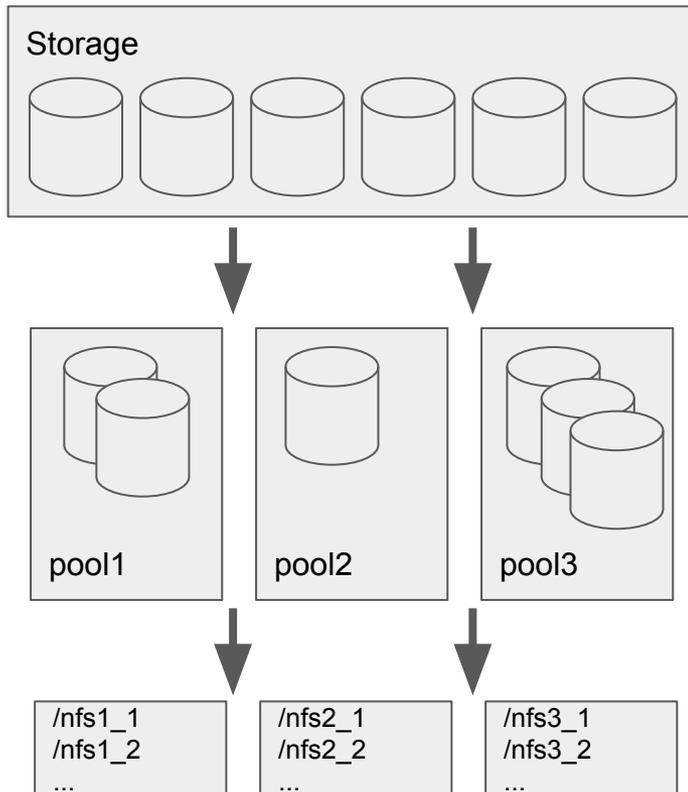
## En contra de otros

- EXT4: antiguo y superado
- XFS: pocas características
- LVM: sólo particionado
- Red Hat Cluster: complejo, pocos FS
- Appliances: cerrados, funciones de pago









- Un “filer” diferente para cada servicio: Sistemas, Telemática (RyT), Soporte-Novell, TV
- Pool como agrupación de uno o más volúmenes, mismo tipo de disco, y con un determinado fin

```
zpool create -f pool1 sdb sdc
zpool add -f pool1 sdj
```

- Los recursos se definen sobre los pools

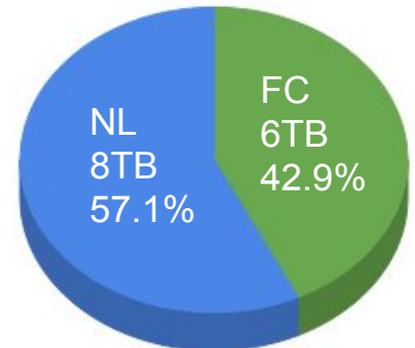
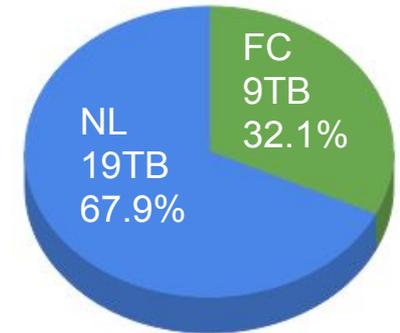
```
zfs create -o quota=200G -o mountpoint=/nfs1_1 pool1/nfs1_1
zfs set sharenfs=rw=155.54.208.17:155.54.208.25 pool1/nfs1_1
```

- Mountpoint no tiene que ser igual que el nombre
- Tamaño flexible (quota)

```
zfs set quota=300G pool1/nfs1_1
```

- Permisos de acceso directamente en ZFS

	<i>Tipo disco</i>	<i>Tamaño</i>	<i>Usado</i>	<i>Nº vols</i>	<i>Nº recursos</i>	<i>Compresión</i>
<b>SISTEMAS</b>		<b>28TB</b>	<b>12.2TB</b>		<b>171</b>	
pool_prod_fc	Fast class	4TB	1.8TB	2x2TB	47	1.07%
pool_prod_nl	Near line	4TB	540GB	1x4TB	32	1.09%
pool_desa_fc	Fast class	2TB	186GB	1x2TB	34	1.22%
pool_desa_nl	Near line	2TB	503GB	1x2TB	51	1.26%
pool_sakai_fc	Fast class	3TB	2.37TB	2 (1+2TB)	3	1.11%
pool_bak_nl	Near line	12TB	6.5TB	3x4TB	3	off
pool_inter_nl	Near line	1TB	252GB	1x1TB	1	off
<b>TELEMATICA</b>		<b>14TB</b>	<b>6.12TB</b>		<b>35</b>	
pool_ryt_fc	Fast class	2TB	173MB	1x2TB	1	1.06%
pool_ryt_nl	Near line	8TB	4.15TB	2x4TB	27	1.00%
pool_box_fc	Fast class	4TB	1.97TB	2x2TB	7	1.14%
<b>TV</b>		<b>20TB</b>	<b>14.5TB</b>		<b>15</b>	
pool_tv_nl	Near line	20TB	14.5TB	5x4TB	15	off
<b>SOPORTE</b>		<b>2TB</b>	<b>985GB</b>		<b>2</b>	
pool_nov_nl	Near line	2TB	985GB	1x2TB	2	off



```
# zfs get all pool_box_fc/umubox_dg_a1
NAME                                PROPERTY                                VALUE                                SOURCE
pool_box_fc/umubox_dg_a1           type                                   filesystem                           -
pool_box_fc/umubox_dg_a1           creation                               mar nov 11 12:23 2014                -

pool_box_fc/umubox_dg_a1           mounted                                yes                                   -
pool_box_fc/umubox_dg_a1           mountpoint                             /umubox_dg_a1                       local
pool_box_fc/umubox_dg_a1           sharenfs                                rw=10.14.2.6:10.14.1.240/28         local
pool_box_fc/umubox_dg_a1           sharesmb                                off                                   default

pool_box_fc/umubox_dg_a1           quota                                   400G                                  local
pool_box_fc/umubox_dg_a1           used                                    306G                                  -
pool_box_fc/umubox_dg_a1           available                               94,3G                                  -

pool_box_fc/umubox_dg_a1           compression                             lz4                                    inherited from pool_box_fc
pool_box_fc/umubox_dg_a1           compressratio                             1.11x                                  -

pool_box_fc/umubox_dg_a1           atime                                    off                                    inherited from pool_box_fc

pool_box_fc/umubox_dg_a1           snapdir                                 hidden                                 default

pool_box_fc/umubox_dg_a1           readonly                                off                                    default

pool_box_fc/umubox_dg_a1           dedup                                    off                                    default
```

## Operaciones sobre pooles

- Crear un pool
- Ampliar pool existente
- Consultar estado de pooles y tareas en curso
- Exportar o importar (¡cuidado!)
- Eliminar un pool (¡cuidado!)

## Operaciones sobre recursos

- Crear un recurso: nombre, tipo de disco (pool), tamaño, mountpoint, shares
- Modificar propiedades: shares, tamaño
- Modificar propiedades: mountpoint, nombre
- Eliminar un recurso (¡cuidado!)

- Chequeo de integridad del sistema de ficheros, recorre el contenido del pool en busca de errores
- Escanea los datos, no el total del pool
- Tarea en segundo plano, sólo una por nodo

```
# zpool status pool_name
pool: pool_name
state: ONLINE
  scan: scrub in progress since Fri Nov 11 12:52:35 2016
        1,22G scanned out of 40,9G at 47,9M/s, 0h14m to go
        0 repaired, 2,97% done
```

- Script genérico para realizar scrub de todos los pooles del filer
  - Planificado en CRON para un sábado del mes a las 0:00, según filer

- ZFS soporta cuotas de usuario
  - Establecimiento de cuotas y consulta mediante comandos locales
- No se integra con el comando “quota” para consultas remotas (via NFS)
- Creamos script para consulta y modificación de cuotas
  - Definir usuarios locales, certificados para autenticación SSH
  - Script en /home/usuario
  - Configurar /etc/sudoers para que puedan ejecutar comando /sbin/zfs

- ZFS no ofrece mecanismo para clusterización
- Pooles montados exclusivamente en un nodo, pero visibles por todos
  - Usar “zpool import” y “zpool export” para montar los pooles en el nodo deseado
  - OJO: Cuidado con montar un pool en dos nodos a la vez, posible corrupción de datos
- Configuración de red equivalente en todos los nodos
- Scripts sencillos para failover
- El propio cliente NFS se recupera ante desconexiones puntuales

## nas\_telematica start

```
ifconfig eth2 10.59.0.4 netmask 255.255.252.0
ifconfig eth3 10.14.0.7 netmask 255.255.252.0 mtu 9000

zpool import pool_ryt_fc
zpool import pool_ryt_nl
zpool import pool_box_fc
```

## nas\_telematica stop

```
zpool export pool_ryt_fc
zpool export pool_ryt_nl
zpool export pool_box_fc

ifconfig eth2 down
ifconfig eth3 down
```

- Gestión de snapshots integrada en el sistema de ficheros ZFS
  - Propiedad “snapdir” (hidden / visible) para acceder a los snapshots desde el recurso raíz
  - Snapdir: /fs\_test/.zfs/snapshot/snap1
- Nombre de snapshot es igual al del recurso más un sufijo
- Propiedades del recurso snapshot reducidas
  - Propiedad “type” cambia de “filesystem” a “snapshot”
- Snapshots sucesivos siempre incrementales
- Eliminación de snapshots consolida para mantener la cadena
- Etiquetado de snapshots (hold, release)

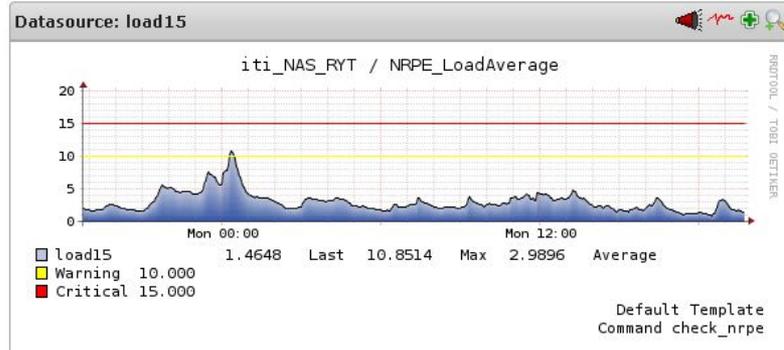
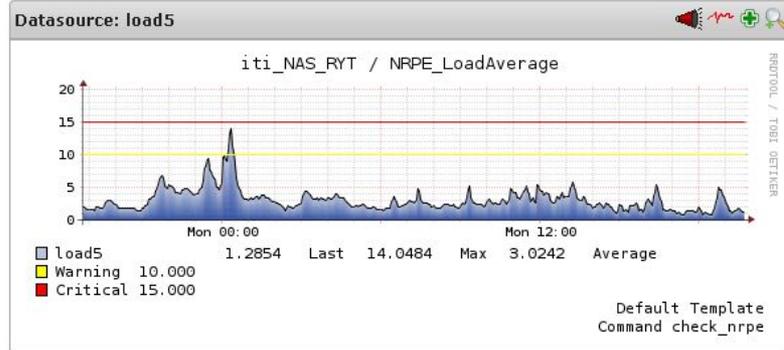
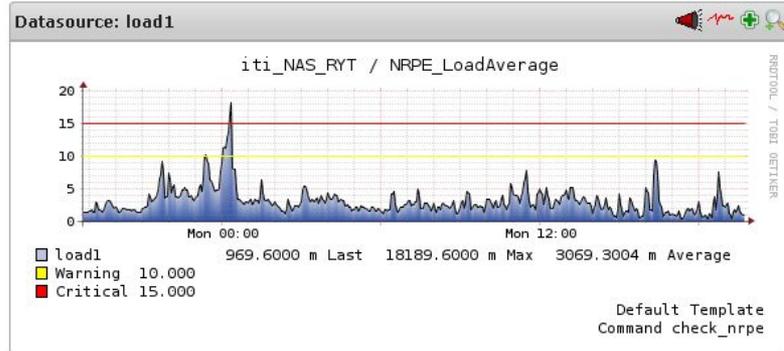
- Mecanismo de réplicas basado en dos comandos
  - zfs send: convierte un recurso en un flujo de bytes
  - zfs receive: convierte un flujo de bytes en un recurso
- Los comandos se concatenan para realizar una réplica
- Cuando hay snapshots, se pueden hacer réplicas recursivas e incrementales
- Se puede encadenar con un SSH para enviar a un filer remoto
- Script para mantener estructura de snapshots y réplicas
  - Crear y eliminar snapshots, etiquetas de snapshot diario, semanal, etc
  - Réplicas con filer remoto (opcional)
  - Diferente retención en local y en remoto

- No hay check SNMP para ZFS, creamos scripts NRPE
- Monitorización de pools, recursos, tráfico de red, procesos, memoria...



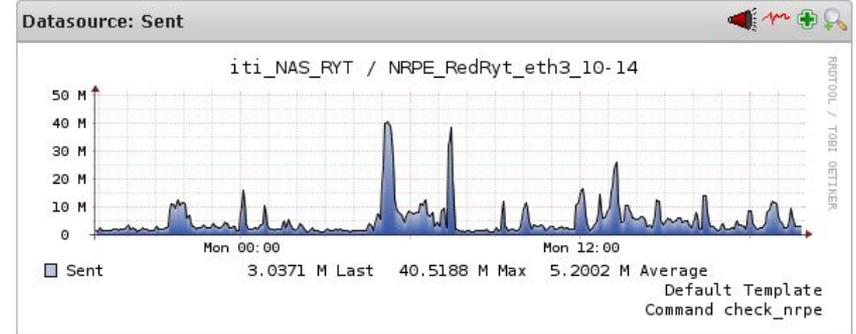
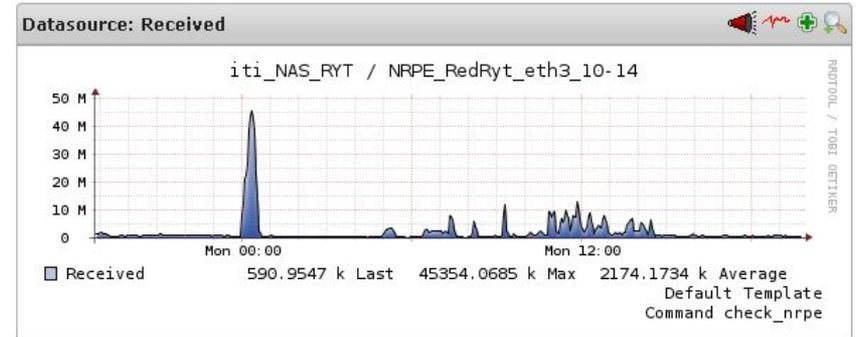
Host: iti\_NAS\_RYT Service: NRPE\_LoadAverage

25 Hours 13.11.16 18:42 - 14.11.16 19:42



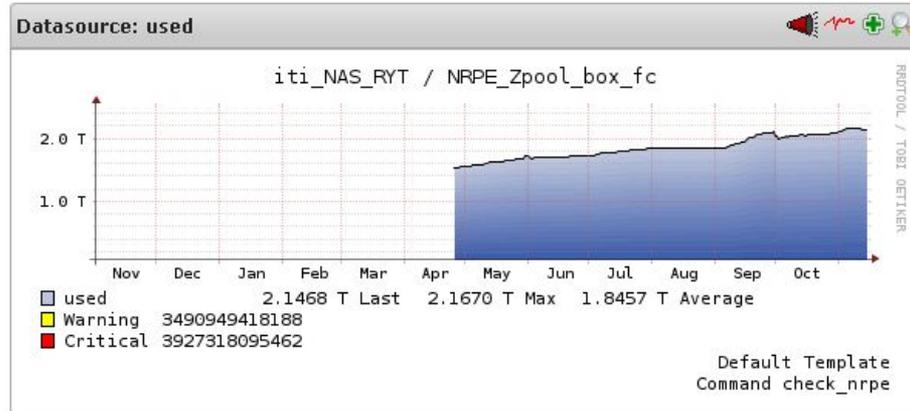
Host: iti\_NAS\_RYT Service: NRPE\_RedRyt\_eth3\_10-14

25 Hours 13.11.16 18:47 - 14.11.16 19:47



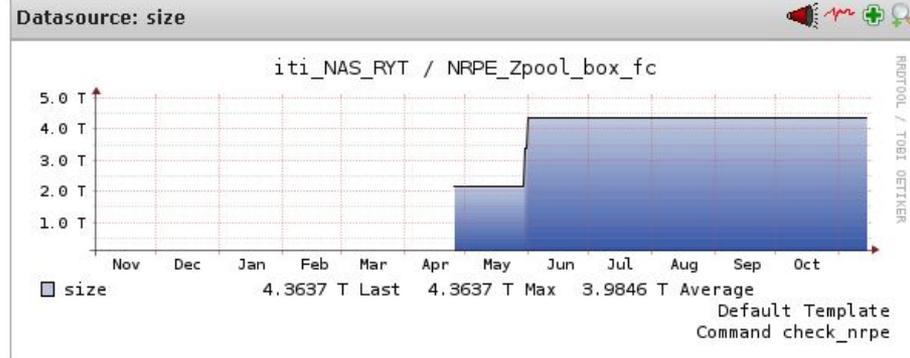
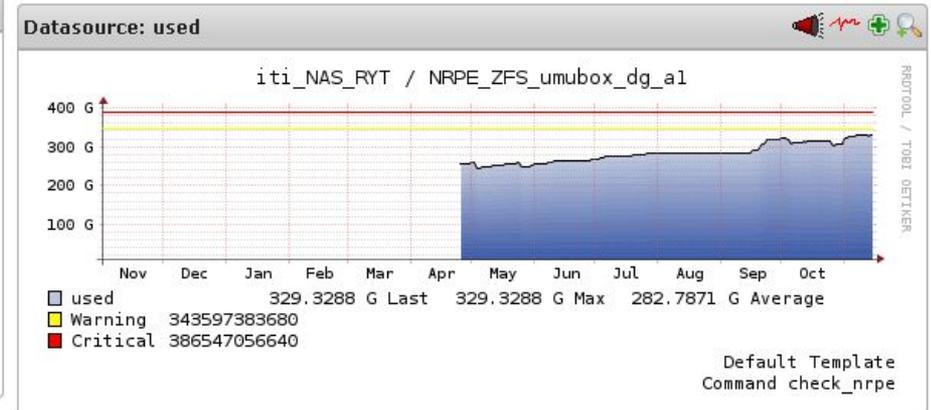
Host: iti\_NAS\_RYT Service: NRPE\_Zpool\_box\_fc

One Year 31.10.15 20:01 - 14.11.16 20:01



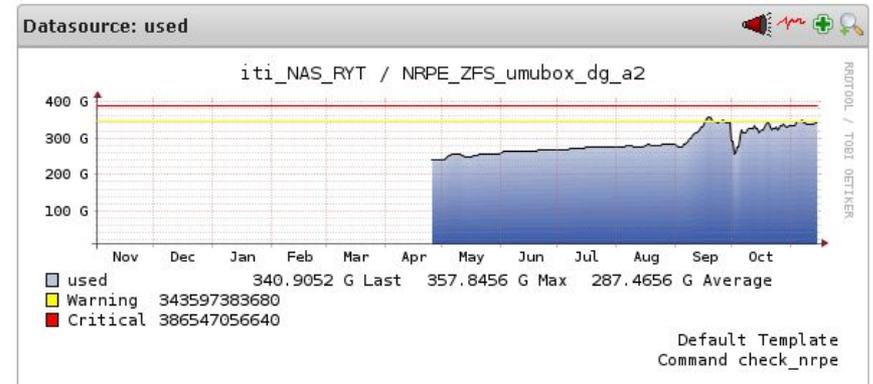
Host: iti\_NAS\_RYT Service: NRPE\_ZFS\_umubox\_dg\_a1

One Year 31.10.15 20:09 - 14.11.16 20:09



Host: iti\_NAS\_RYT Service: NRPE\_ZFS\_umubox\_dg\_a2

One Year 31.10.15 20:11 - 14.11.16 20:11



- Araneo (web corporativa [www.um.es](http://www.um.es))
- Configuraciones tomcat, OAS...
- Aula Virtual (SAKAI)
- DALI (impresión centralizada)
- Consigna (envío de archivos grandes por email)
- MySQL Telemática
- Socrates (disco compartido)
- Umubox (disco en la nube)
- Webs particulares departamentos y usuarios
- Concentrador de logs
- Diferentes entornos (producción y desarrollo)
- Archivo de videos y streaming en TV.UM

- Utilizar SSD para caché
- Snapshots y réplicas con CPD remoto
- Deduplicación
- Despliegue de nodos y configuración con Puppet
- Actualización de versiones de ZFS
- Automatización del cluster
- Posible entorno web para administrarlo



