

2494-8

**Workshop on High Performance Computing (HPC) Architecture
and Applications in the ICTP**

14 – 25 October 2013

HPC Storage Management

C. Onime & E. Corso
ICTP, Trieste

HPC Storage Management

Overview

- Global plan
- RAID & data protection
 - Hardware, BIOS & software
 - Provisioning (expansion)
- File-systems
 - Local & remote
- Network File System
 - Configuration, quota, mounting on demand
 - Performance tuning
- Troubleshooting storage problems

sustainability

GLOBAL PLANNING

Comprehensive plan

- Multiple plans for
 - Short term
 - Now to 3 months
 - Medium term
 - 6 months to 1 year
 - Long term
 - 1 year to 4 years
- Phased deployment

Lifecycle & protection

- Data storage lifecycle
 - Online -> near line, off-line & discard
- Data protection (backup)
 - What
 - Can it be recreated or downloaded?
 - When
 - Daily, weekly, monthly
 - How
 - Remote mirroring
 - Tape device

User/group needs

- Space
 - Growth plan
- Performance
 - Improving performance
- Partitioning
 - Security
 - Capacity

Important items

- Data center (environment)
 - Power (consumption)
 - Sleeping disks
 - Physical volume/size of device
 - Cooling needs
- Maintenance
 - Warranty & post-sales support/service
 - Spares (in-house, on-demand)
 - MTBF (disks)

Data protection

RAID

All about RAID

- Redundant Array of Independent Disks (RAID)

Level	Useable capacity	Data protection
RAID0	$\text{Size}_{\min} * n$	None
RAID1	Size_{\min}	Failure of one single disk
RAID5	$\text{Size}_{\min} * (n - 1)$	Concurrent failure of one single disk
RAID6	$\text{Size}_{\min} * (n - 2)$	Concurrent failure of two disks
RAID1+0	$\text{Size}_{\min} * (n/2)$	Concurrent failure of more than two disks

RAID types

Characteristics	Hardware RAID	BIOS RAID	Software RAID
Cache RAM	dedicated	shared	shared
Battery backup unit	Yes (48 hours)	No	No
Raw data disk Portability	Not recommended <i>(Works for same controller family)</i>	Not sure	Yes <i>(works for same O.S)</i>
Configuration tool	Dedicated firmware based	Firmware+Host O.S	Host O.S
Hot disk replacement	yes	No recommended	Not recommended
Performance enhancement	Yes (faster)	none	none

RAID Volumes

- Typical unit presented to O.S
 - Provisioning (mostly ability to expand)
 - Reduction may require destroying and make a new one
- States
 - NORMAL
 - DIRTY
 - DEGRADED

Choosing

FILE SYSTEMS

Characteristics of file-systems

- Journaling
 - Data protection mechanism for faster consistency check
- Snapshots
 - Frozen image of data typically used for backup
- IO scheduler
 - Delayed write & Read ahead
- Others
 - Overall capacity, metadata handling, quota management, etc.

Examples

- Local file systems
 - FAT, NTFS
 - (de-fragmentation issues)
 - ext2, ext3, ext4, xfs, jfs...
- Network file systems (*single servers*)
 - NFS & CIFS
- Distributed file systems (multiple servers)
 - AFS (encrypted) & pNFS
- Clustered file systems (*parallel access*)
 - GFS2, GPFS & Lustre

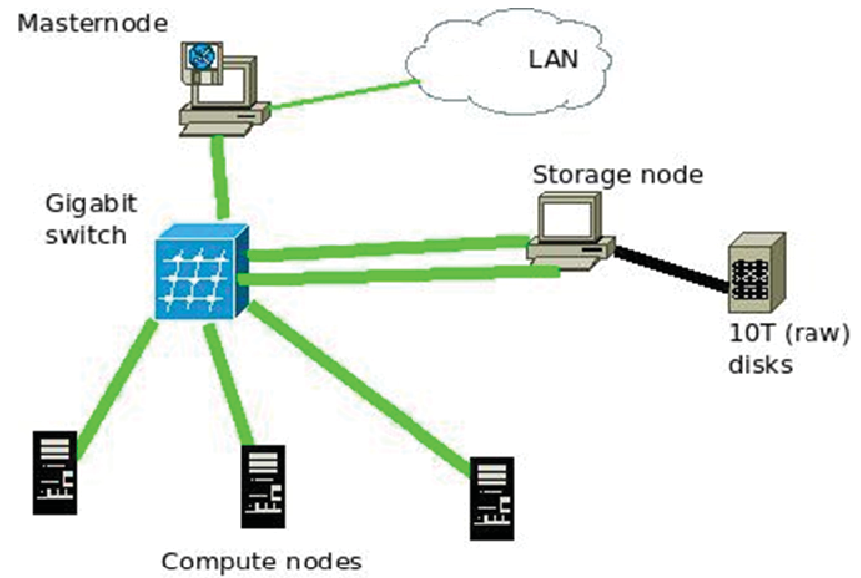
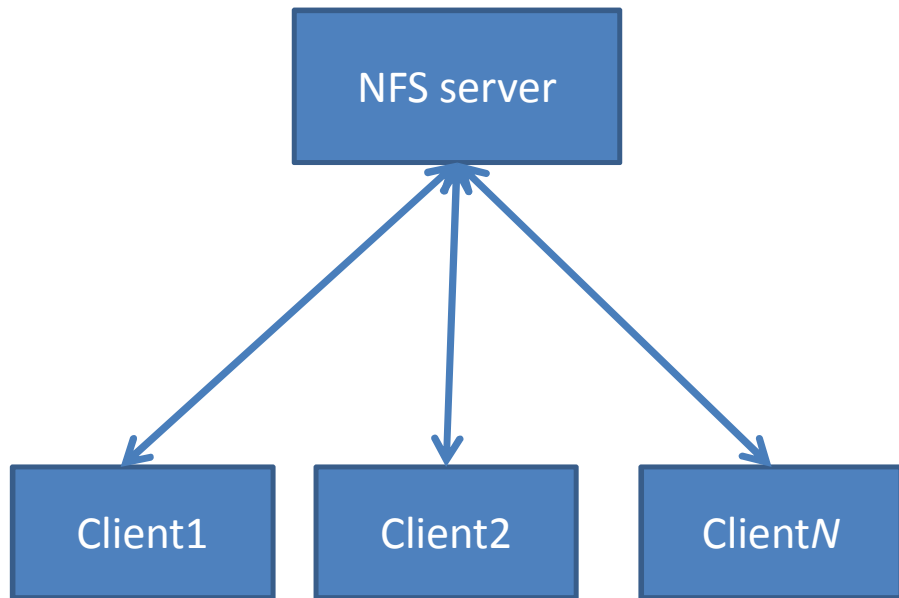
(Choice should be governed by intended usage)

NETWORK FILE SYSTEM

Network File System

- Version 3 (NFS or NFSv3)
 - Most widely deployed implementation
 - Simple security system
 - IP address based
 - UID user authentication with POSIX/Unix permissions and ability to exclude UID=0 (root user).
- Version 4 (NFSv4)
 - Improved security thanks to kerberos 5 user authentication
- Version 4.1 (pNFS)
 - Improved performance: Separating metadata from data

HPC storage (NFS) architecture



Server side configuration

- Access control
 - /etc/exports
 - None, read-only, read-write (with wildcards)
- Quota control
 - Mount local file-system with quota options
 - in /etc/fstab
 - Use local tools to set and manage quota

Client side configuration

- Mount
 - /etc/fstab
- On-demand mount
 - Automounter
 - auto.master
 - auto.home
- Performance tuning
 - mount options
 - rsize and wsize

Server performance tuning

- Vertical scaling *(bigger single server)*
 - More or faster RAM (and/or CPU)
 - More network connections
 - More daemons
- Horizontal scaling *(more physical servers)*
 - Requires partitioning of data
 - Works best with automounter based client mounting

Troubleshooting

STORAGE PROBLEMS

Possible problems

- faulty connectivity (network)
- Bad/faulty disk
- Failing disk
- Bad power supply unit
- Server crash (needed hard power cycle)
- Slowdown in performance (scalability)
 - High load average on server
 - other processes/services on server
 - Overloaded, too many clients

Visual/physical inspection

- *Photo of CED storage lights in dark room*
- LED /Lights
 - On disks
 - Network ports (both computer & network device)
 - Power supply
- Damaged/broken cables
 - Broken heads, old cables
- High temperatures can also degrade the MTBF

TOP

```
top - 09:11:59 up 3 days, 22:06, 1 user, load average: 0.80, 0.45, 0.51
Tasks: 177 total, 1 running, 175 sleeping, 0 stopped, 1 zombie
pu0 : 1.0%us, 1.3%sy, 0.0%ni, 97.3%id, 0.0%wa, 0.0%hi, 0.3%si, 0.0%st
pu1 : 0.7%us, 1.3%sy, 0.0%ni, 97.7%id, 0.0%wa, 0.0%hi, 0.3%si, 0.0%st
pu2 : 0.6%us, 0.6%sy, 0.0%ni, 98.4%id, 0.0%wa, 0.0%hi, 0.3%si, 0.0%st
pu3 : 0.7%us, 0.7%sy, 0.0%ni, 98.4%id, 0.0%wa, 0.0%hi, 0.3%si, 0.0%st
pu4 : 0.3%us, 1.0%sy, 0.0%ni, 98.7%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
pu5 : 1.0%us, 0.7%sy, 0.0%ni, 98.3%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
pu6 : 0.3%us, 0.6%sy, 0.0%ni, 71.2%id, 27.9%wa, 0.0%hi, 0.0%si, 0.0%st
pu7 : 0.0%us, 0.0%sy, 0.0%ni, 100.0%id, 0.0%wa, 0.0%hi, 0.0%si, 0.0%st
Mem: 41283036k total, 27924908k used, 13358128k free, 543676k buffers
Swap: 12586892k total, 204k used, 12586688k free, 21806484k cached
```

PID	USER	PR	NI	VIRT	RES	SHR	S	%CPU	%MEM	TIME+	COMMAND
3527	root	20	0	204m	2136	1156	S	3.3	0.0	196:37.32	rsyslogd
3256	named	20	0	266m	5928	1772	S	2.7	0.0	169:54.61	named
3776	ganglia	20	0	477m	158m	2216	S	2.3	0.4	215:36.90	gmond

Total DISK READ: 0.00 B/s				Total DISK WRITE: 0.00 B/s			
TID	PRI	USER	DISK READ	DISK WRITE	SWAPIN	IO>	COMMAND
1	be/4	root	0.00 B/s	0.00 B/s	0.00 %	0.00 %	init
2	be/4	root	0.00 B/s	0.00 B/s	0.00 %	0.00 %	[kthreadd]
3	be/4	root	0.00 B/s	0.00 B/s	0.00 %	0.00 %	[ksoftirqd/0]
5	be/0	root	0.00 B/s	0.00 B/s	0.00 %	0.00 %	[kworker/0:0H]

Atop - argo 2013/10/14 09:16:13												-----		10s elapsed			
PRC	sys	0.60s	user	0.39s		#proc	173		#zombie	0	clones	15		#exit	17		
CPU	sys	5%	user	3%	irq	1%		idle	785%	wait	7%		steal	0%	guest	0%	
cpu	sys	0%	user	0%	irq	0%		idle	93%	cpu007	w 7%		steal	0%	guest	0%	
cpu	sys	1%	user	0%	irq	0%		idle	99%	cpu002	w 0%		steal	0%	guest	0%	
cpu	sys	1%	user	1%	irq	0%		idle	97%	cpu000	w 0%		steal	0%	guest	0%	
cpu	sys	0%	user	0%	irq	0%		idle	100%	cpu003	w 0%		steal	0%	guest	0%	
cpu	sys	1%	user	1%	irq	0%		idle	98%	cpu001	w 0%		steal	0%	guest	0%	
cpu	sys	1%	user	1%	irq	0%		idle	98%	cpu004	w 0%		steal	0%	guest	0%	
cpu	sys	1%	user	1%	irq	0%		idle	99%	cpu005	w 0%		steal	0%	guest	0%	
cpu	sys	0%	user	0%	irq	0%		idle	100%	cpu006	w 0%		steal	0%	guest	0%	
CPL	avg1	0.29	avg5	0.36	avg15	0.46			csw	32541		intr	39447		numcpu	8	
MEM	tot	39.4G	free	12.7G	cache	20.8G	dirty	159.5M	buff	531.1M	slab	2.9G					
SWP	tot	12.0G	free	12.0G													
DSK			busy	10%	read	3	write	763	KiB/r	5	KiB/w	8	MBr/s	0.00	vmlim	31.7G	
NET	transport		tcpi	52	tcpo	65	udpi	7212	udpo	2881	tcpao	6	tcppo	2	tcprs	0	
NET	network		ipi	7287	ipo	2959	ipfrw	0	deliv	7269			icmpi	0	icpmo	16	
NET	eth0	0%	pcki	5290	pcko	218	si	460 Kbps	so	12 Kbps	erri	0	erro	0	drpi	0	
NET	eth1	0%	pcki	23	pcko	842	si	1 Kbps	so	139 Kbps	erri	0	erro	0	drpi	0	
NET	eth2	0%	pcki	22	pcko	20	si	1 Kbps	so	1 Kbps	erri	0	erro	0	drpi	0	
NET	lo	----	pcki	1958	pcko	1958	si	153 Kbps	so	153 Kbps	erri	0	erro	0	drpi	0	
PID	RUID	EUID	THR	SYS CPU	USR CPU	VGRW	RGROW	RDSK	WRDSK	ST	EXC	S	CPUNR	CPU	CMD	1/2	
3527	root	root	6	0.20s	0.08s	OK	OK	OK	124K	--	--	S	5	3%	rsyslogd		
3256	named	named	11	0.15s	0.10s	OK	OK	OK	OK	--	--	S	0	3%	named		
3776	ganglia	ganglia	2	0.12s	0.07s	OK	OK	OK	OK	--	--	S	7	2%	gmond		

– Identify bottlenecks & other processes

- top (CPU, RAM or IO)
- iotop (which processes are generating IO?)
- atop (CPU, RAM, network IO or disk IO)

Periodic monitoring

- Monitoring
 - Use smartd for SMART monitoring of disks
 - periodic self testing of disks, predict disk faults and sends e-mail notifications
 - Use NAGIOS or CACTI
 - Monitor hardware, status and occupancy/capacity
- Periodic benchmarking (iozone/bonnie++)
 - File-system on both server and client side.
 - Can show trends

Questions??

Thank you
&
Now the hands-on