



DB

Database Services

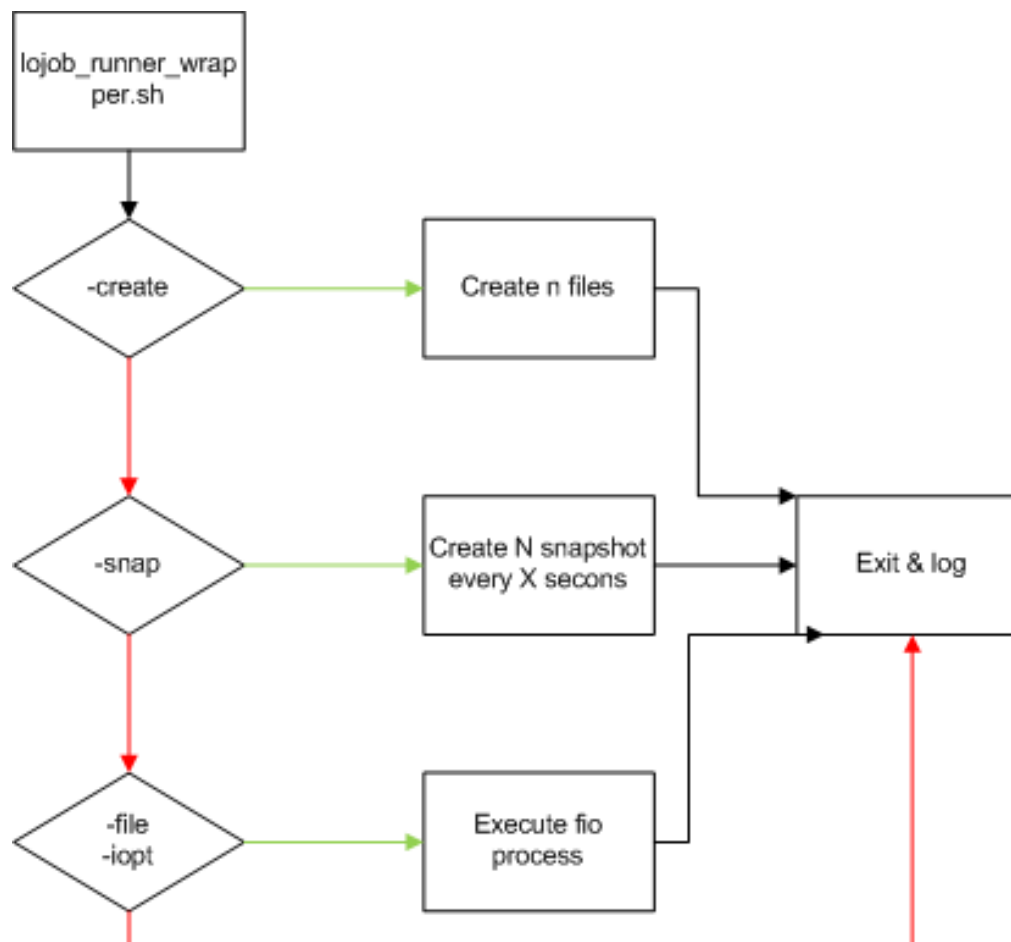
CERN IT
Department

Tests and evolution SSD as flash cache

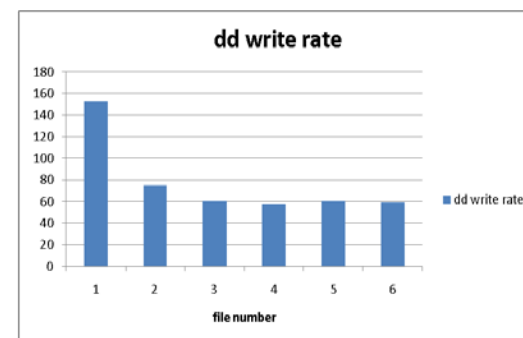
Ruben.Gaspar.Aparicio@cern.ch
IT/DB

- IO framework: IOTool
- IO Tests:
 - Intel X25-M 160GB SSD
 - Sun 7410 Storage: ZFS
 - Netapp: PAM II & FlexClone

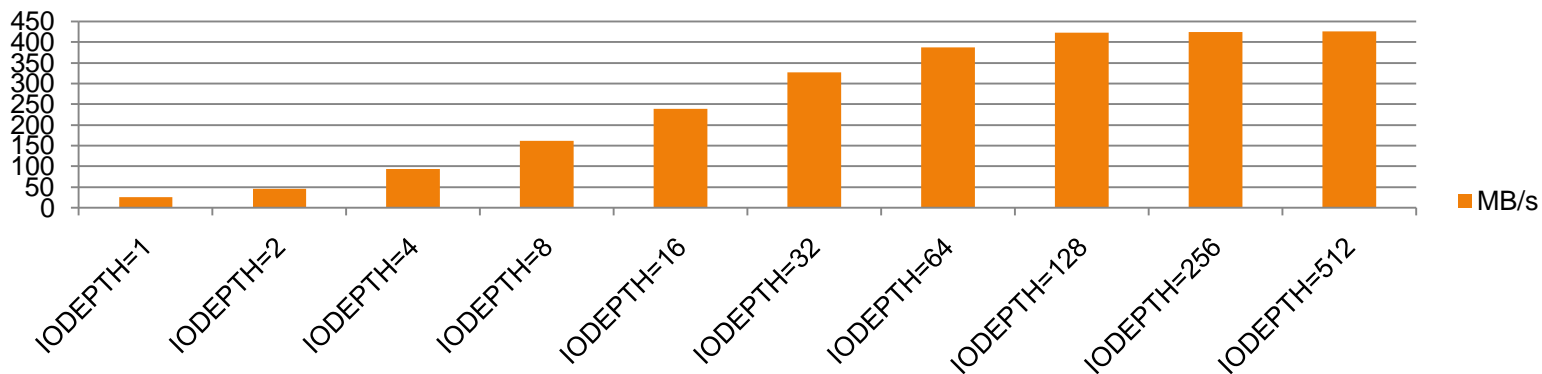
- Based on Open source project fio:
<http://freshmeat.net/projects/fio/>
- It can be used to stress NAS or SAN storage
- Run from a central server it monitors several IO servers
 - To stress the storage behind
 - To avoid limitation from a single server
 - Central point to gather logs
 - Dashboard from a central machine
- It can be run standalone on a IO server
- Perl & Bash



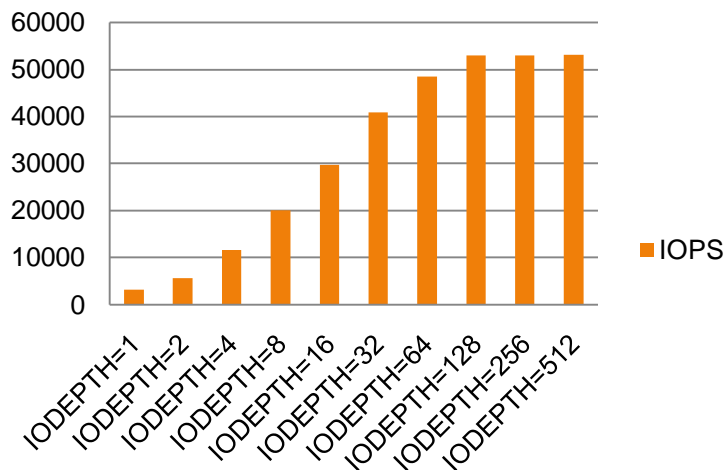
- [IOIntelSSDBenchmark.pdf](#) (external/CERN account required)
- Five Intel X-25 M in RAID0. Ext2 file system. Locally attached.
- Server: 64 cores, 125Gb RAM, RHE5
- 100GB files
- Different workload size
 - Write performance decreases as soon as SSD gets full
- IO scheduler: noop vs cfq (Completely Fair Queuing)



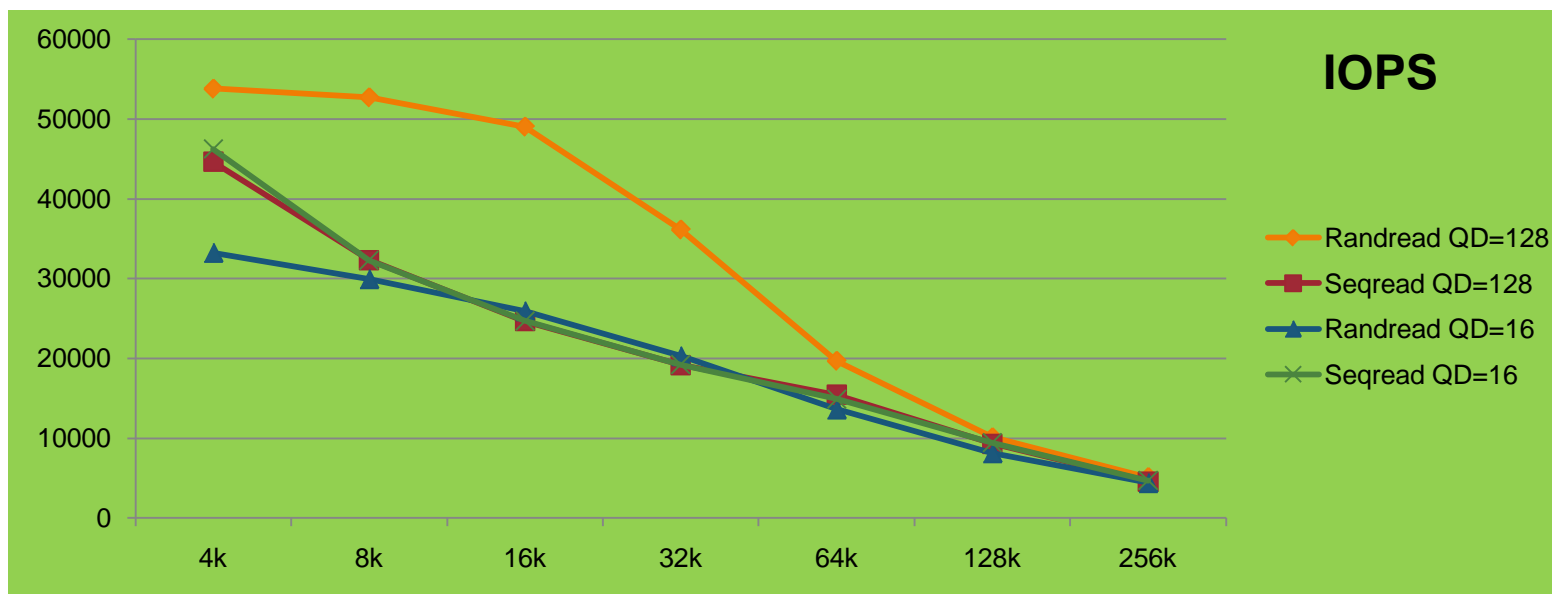
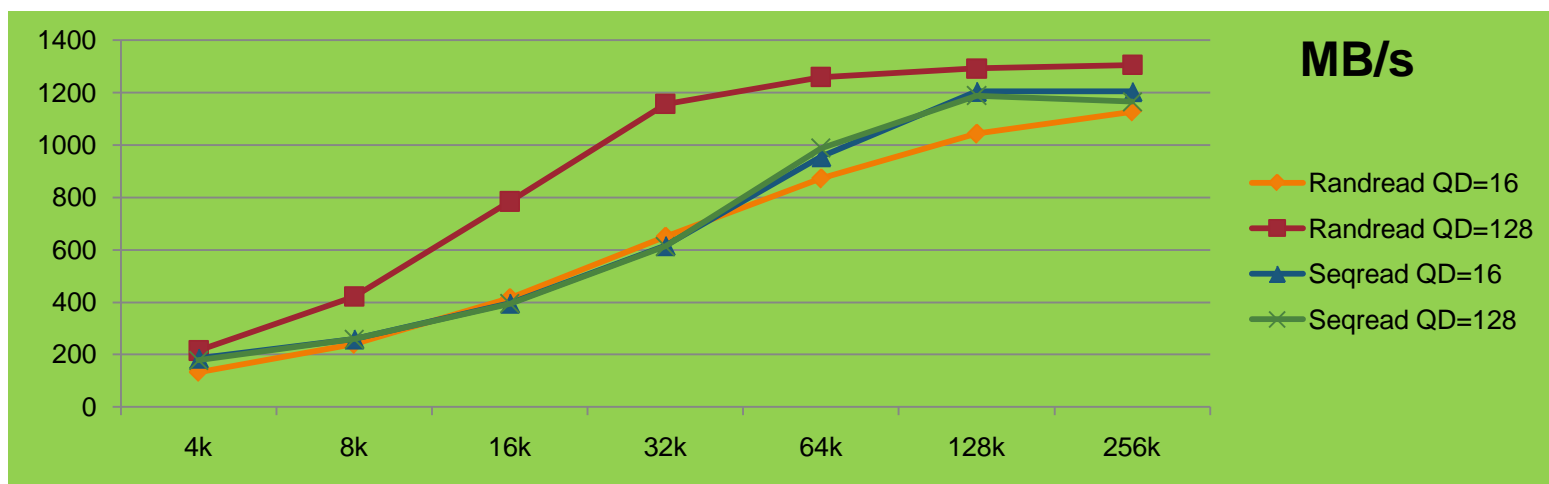
Sustained Rate Random Reads



IOPS



	Completion Time (us)	Standard Deviation (us)
IODEPTH=1	288.783124	126.295424
IODEPTH=2	320.510819	135.099729
IODEPTH=4	318.405563	205.260997
IODEPTH=8	372.784486	149.152536
IODEPTH=16	517.758054	184.776984
IODEPTH=32	769.693771	265.919942
IODEPTH=64	1307.207774	382.034242
IODEPTH=128	2405.593282	592.249235
IODEPTH=256	4809.014412	850.041343
IODEPTH=512	9610.620898	1344.863408



- [SunStorage7410Benchmark.pdf](#)
- ZFS: dynamic stripping, compression, deduplication, snapshot, dynamic file system sizing, 128-bits file system ...
- Supports NFS v3 and v4
- Web interface (BUI) & CLI
- Analytics based on Dtrace
- Several possibilities to configure storage pool(s)



Show Details

dbsunos1



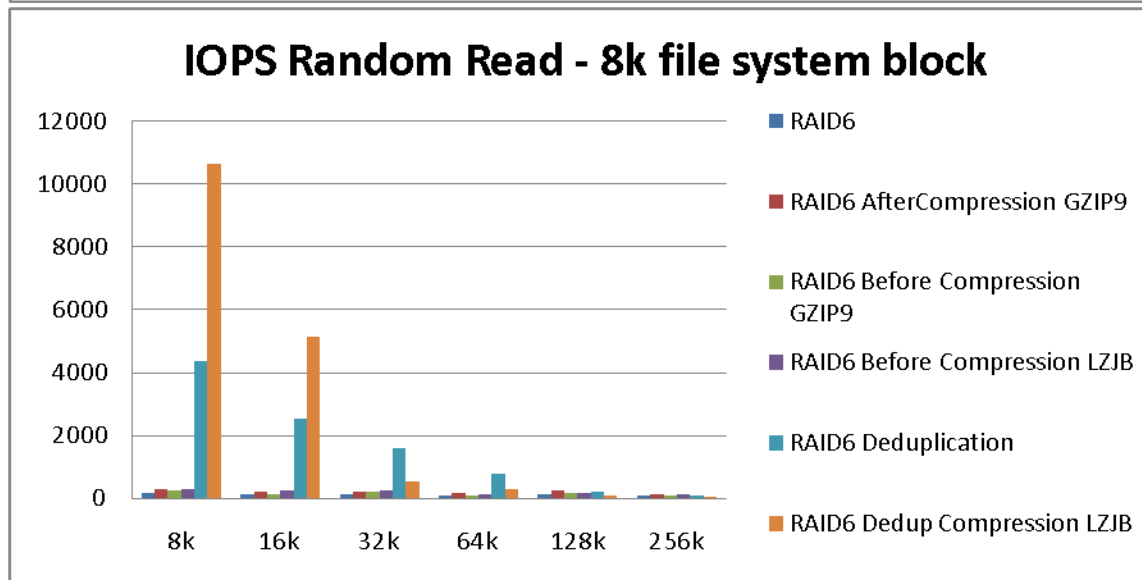
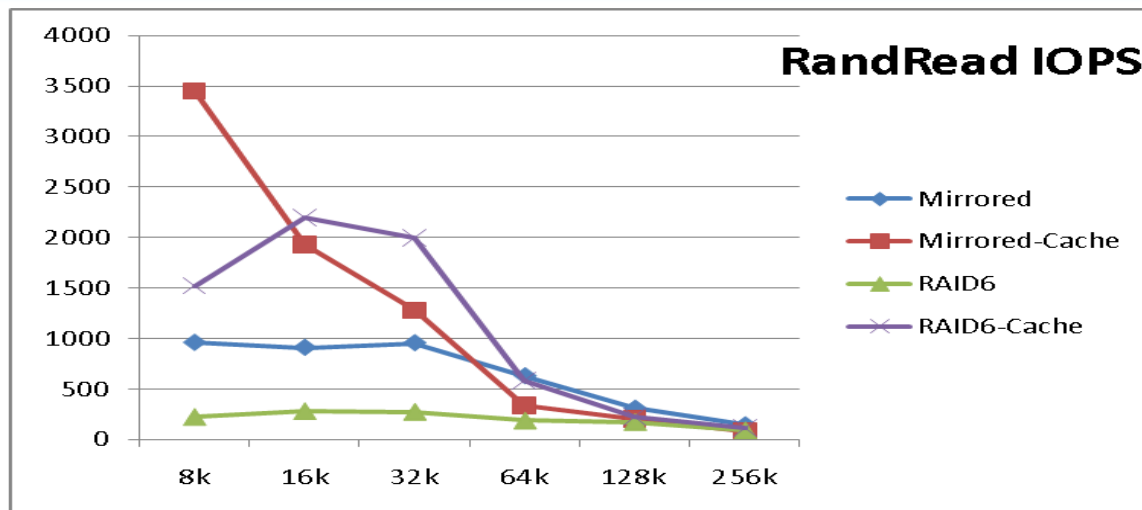
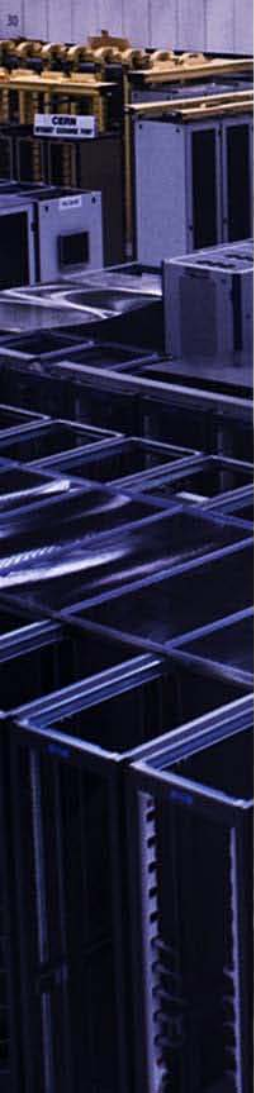
Manufacturer Sun Microsystems, Inc.
 Model Sun Storage 7410
 Serial 0852QAF028
 Processors 2x2.30GHz Quad-Core AMD
 Opteron(tm) Processor 2356
 Memory 64GB

System 932GB (2 disks)
 Data -
 Cache 186GB (2 disks)
 Log -
 Total 1.09TB (4 disks)

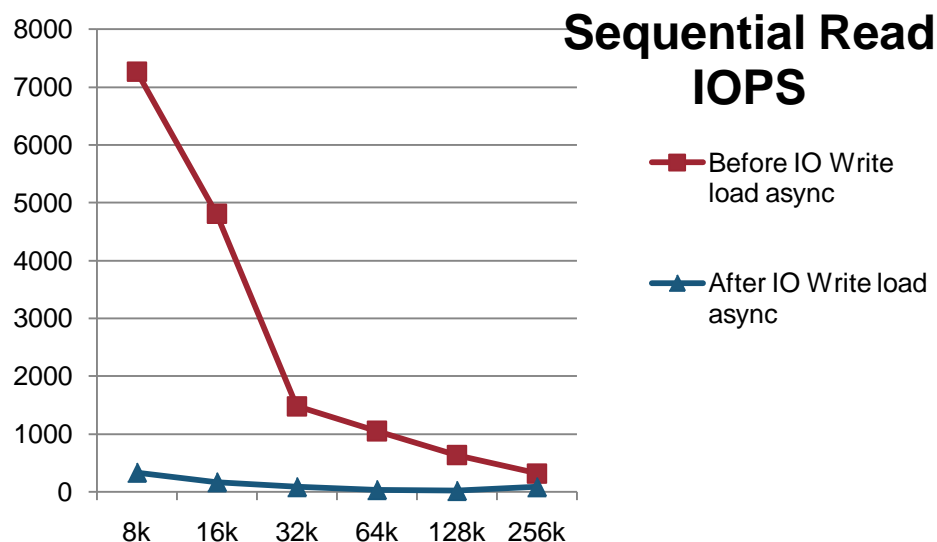
Disk Shelves

	NAME	MANUFACTURER	MODEL	DATA	CACHE	LOG	PATHS	
➔ ⓘ	0903QBK011	Sun Microsystems, Inc.	J4400	9.10TB	-	34GB	2	⊙



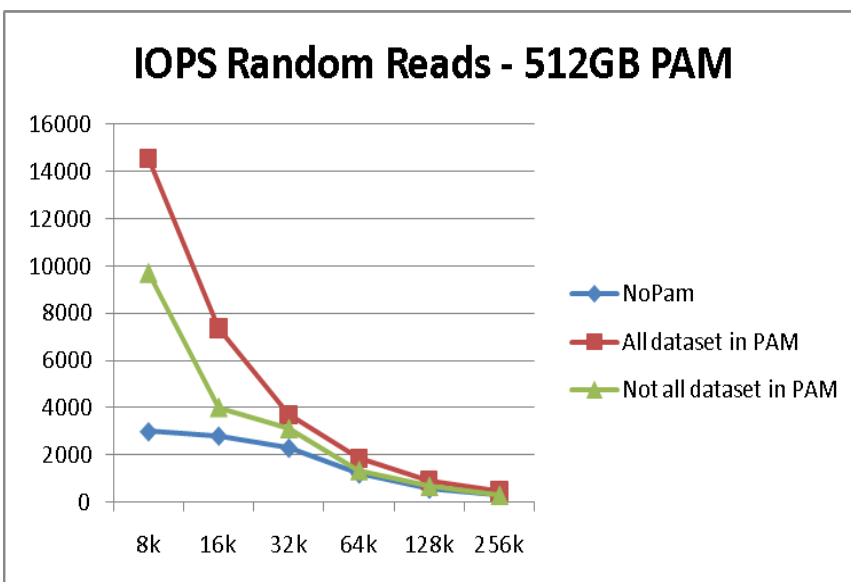


- Copy on write file system
- After 30 hours of random writes

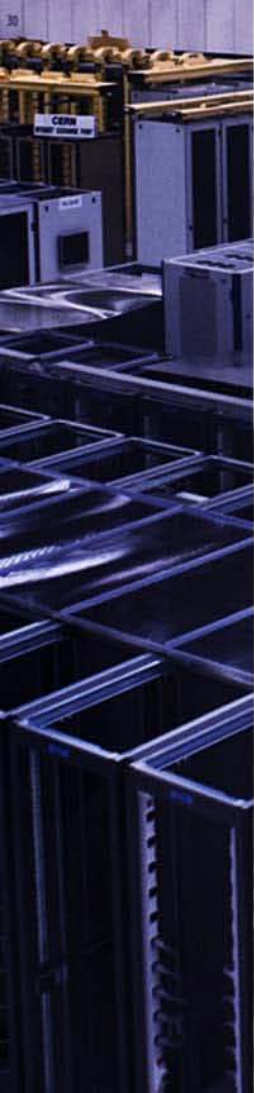


- Defrag functionality will be included in next release.

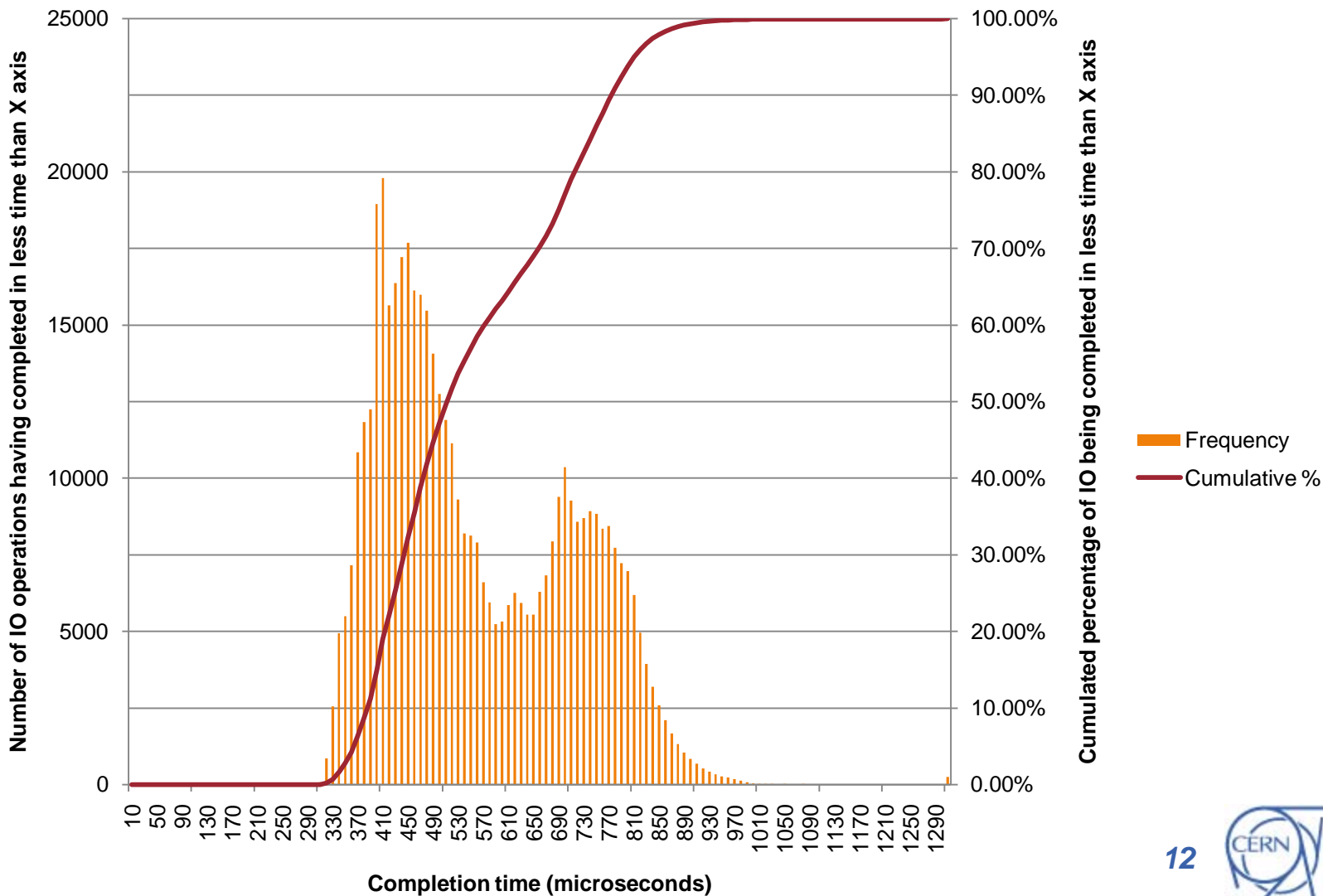
- FAS3170 running Ontap 7.3.4 with 27 FC data disks 15 rpm, raid_dp.
- Predictive cache stadistics:
 - flexscale.lopri_blocks
 - flexscale.normal_data_blocks



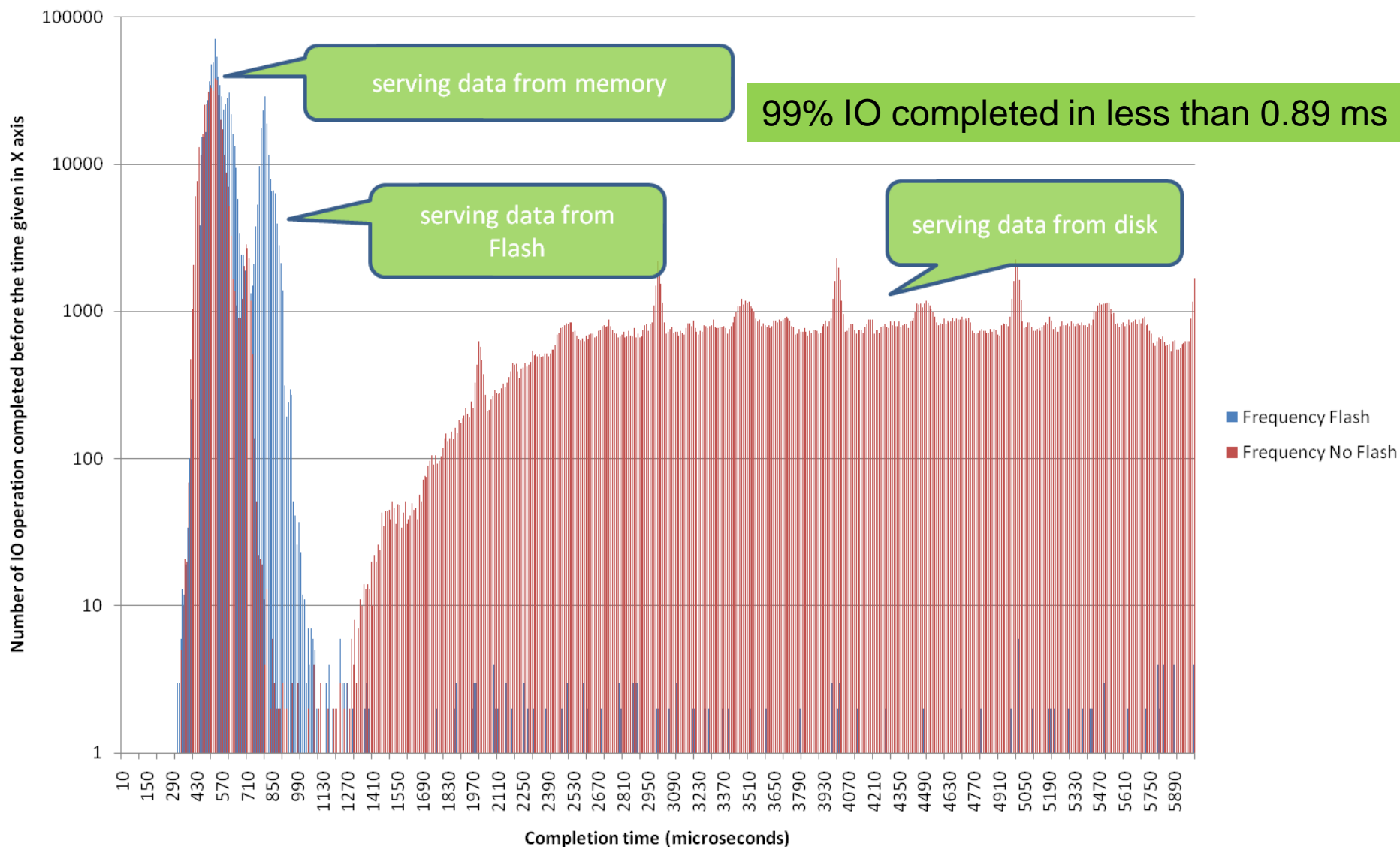
MB/s	C.Time(us)	MB/s	C.Time(us)	MB/s	C.Time(us)
23.99121	5311.89312	116.4121	1080.90882	77.72949	1628.3538
45.06348	5665.95807	117.6846	2161.13857	64.24219	3971.2753
73.64453	6936.01354	118.6426	4299.63039	99.80176	5110.8363
76.7373	13331.3421	118.6475	8617.18462	86.19629	11866.698
71.56543	28602.0238	118.2695	17300.5632	86.34668	23703.034
73.63184	55609.3087	118.542	34532.2281	74.13574	55230.67



Histogram of IO completion (Oracle trace, 10046 event)



Distribution of IO operations completion time



- Based on snapshot technology
- Fast and easy to duplicate production volumes. No extra space involved (0.5% of original parent volume size).

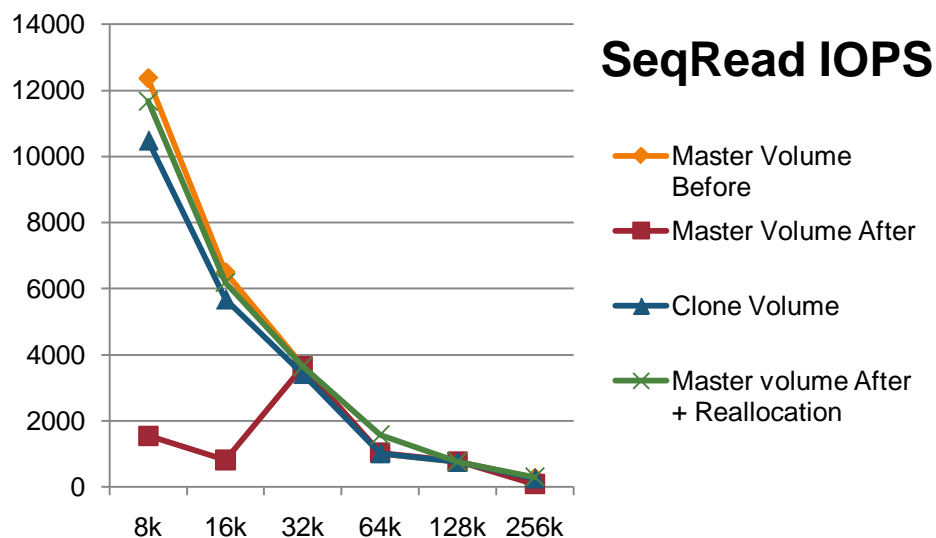
```
dbnasb301> vol clone create clonetest_clone -b clonetest
```

```
dbnasb301> snap list  
Volume clonetest  
working...
```

<u>%/used</u>	<u>%/total</u>	<u>date</u>	<u>name</u>
0% (0%)	0% (0%)	Sep 10 15:58	clone_clonetest_clone.1 (busy,vclone)

- Possibility to split the relationship:
 - vol clone split

- No difference among parent and clone volumes
- Defragmentation should be taken into account on both clone and parent



Questions?