

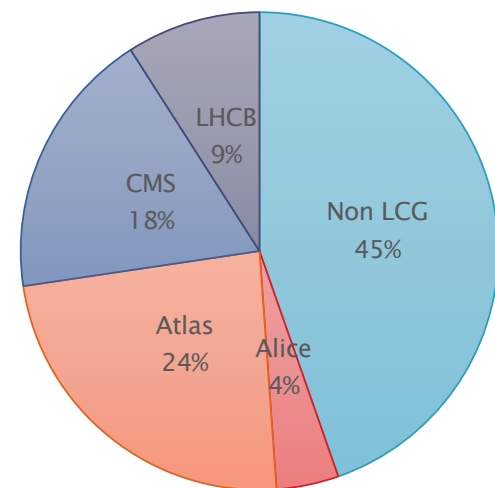
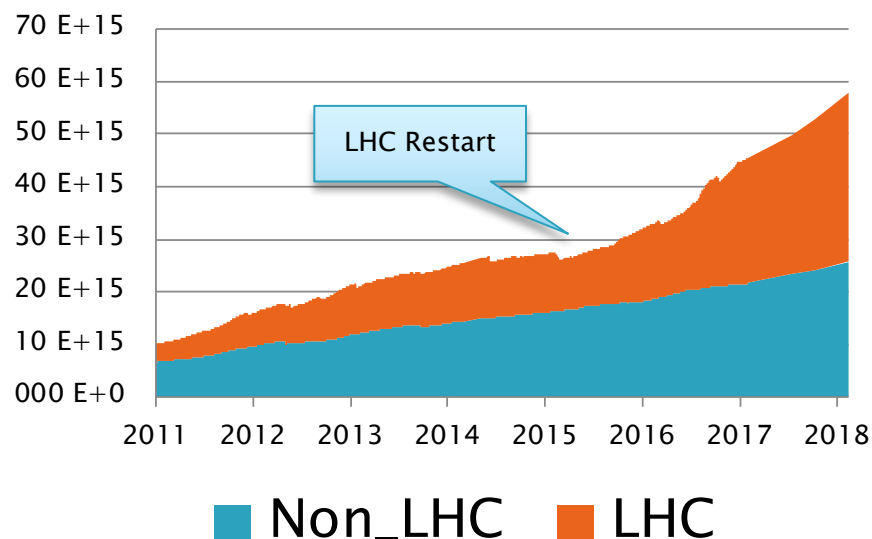
HPSS at IN2P3

Pierre-Emmanuel Brinette, 2018-04-16
2èmes rencontres HPSS France 2018

- ▶ HPSS and TREQS overview
- ▶ HPSS Monitoring
- ▶ Tape infrastructure and evolution
- ▶ HPSS 7.5 Migration
- ▶ Service Evolution

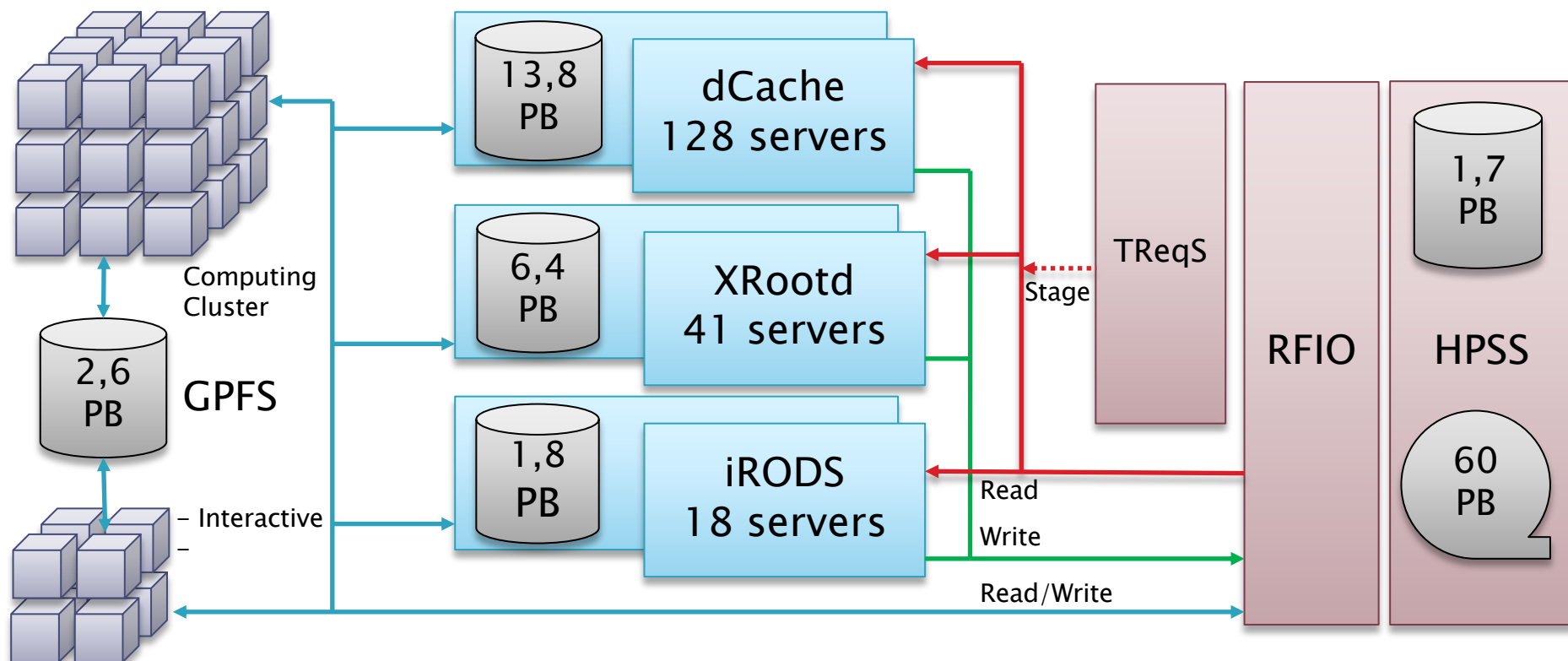
- ▶ HPSS is the main repository for scientific data
 - 80 different VO (groups) store data in HPSS
 - 55 % used for LHC data (Alice, Atlas, CMS, LHCb)
- ▶ Usage (Apr 2018)
 - 60 PB stored
 - 75 M of files
- ▶ Evolution over last year +11,7 PB (+26 %)
 - LCG : +8 PB (+34 %)
 - Non LCG : +3,7 PB (+ 17%)
- ▶ Forecast for 2018 : + 16 PB (~ 2000 tapes)

HPSS growth over last 7 years



■ Non LCG ■ Alice ■ Atlas ■ CMS ■ LHCb

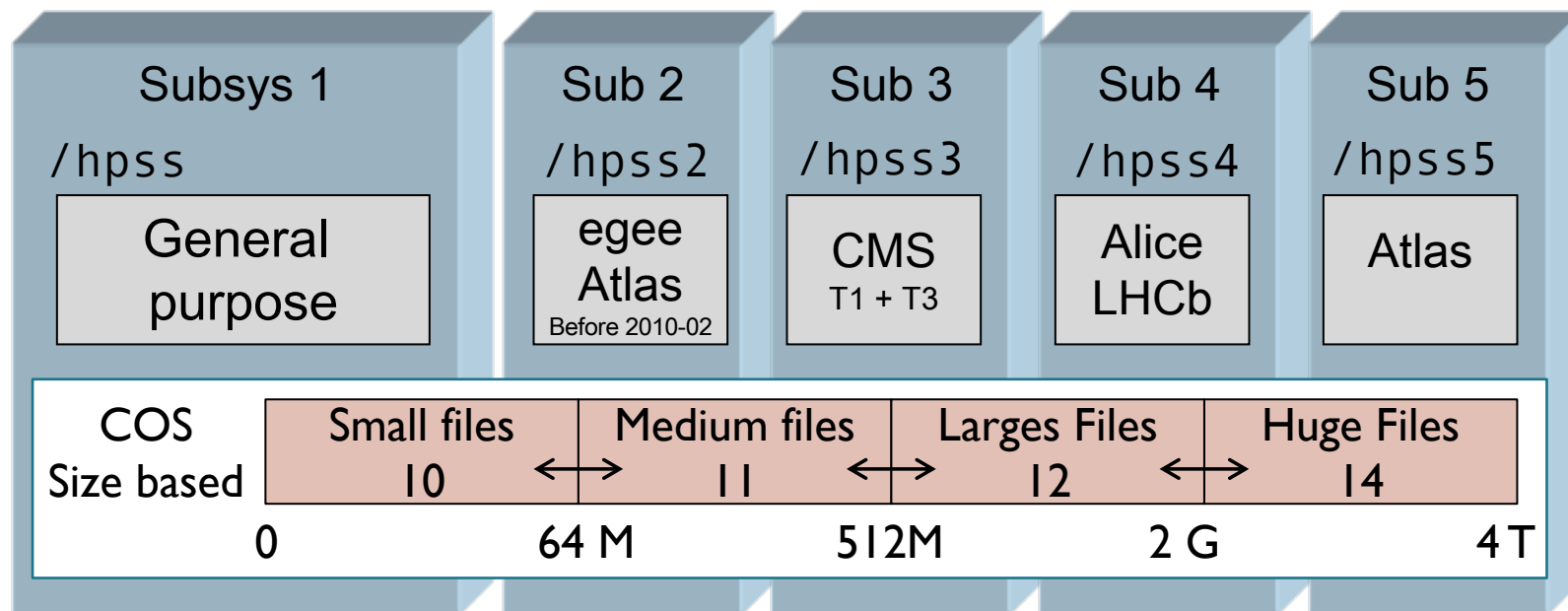
HPSS Overview



- ▶ HPSS v7.4.3p2
- ▶ HPSS Interface : RFIO with HPSS extensions
- ▶ 85 % of HPSS access are performed through storage middleware
 - **dCache** (LCG/egee),
 - **Xrootd** and **iRods**
- ▶ Still some direct access to HPSS but decreasing
- ▶ Disk cache renewed in 2017
 - + 8 new movers (DELL R730xd)
 - Total 12 movers (1,7 PB) @ 10Gbits
- ▶ Read operations from storage middleware are handled by TREQS 2

HPSS Storage policy

- ▶ 5 subsystems, 4 COS Only (selected by size), 10th file families
- ▶ Different tape resources per COS (ie. Small files on “Sport” tapes)

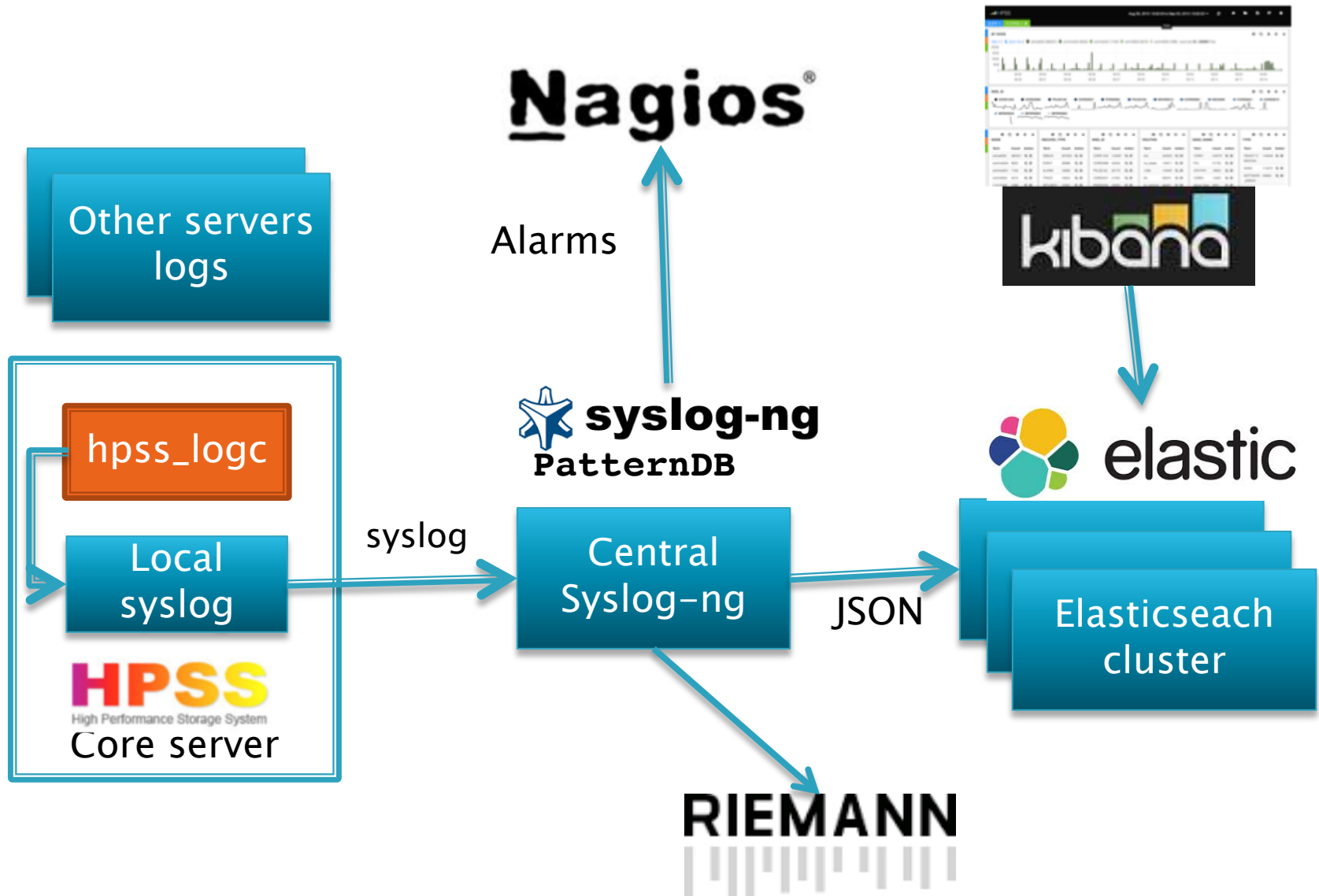


- ▶ **Historical**
 - 26 PB
 - ~2000 UID
 - 50 M files

- **Newly created**
 - 34 PB
 - 20 M files
 - Mainly used for LHC Data

- **Dedicated subsystem**
 - Allow to dedicate DISK resources for specific set of users when using **automatic COS selection**
 - Specific database for a set users → faster query
 - Subsys 1 : 40 GB
 - Subsys [2-5] : 1.5 to 6 GB

- ▶ TREQS 2 is the IN2P3 tape scheduler for HPSS
 - Optimize **read** operations by sorting files by tapes and positions
 - Reduce the number of mounts / dismounts of the same tape.
 - Limit the number of drives used for staging
- ▶ Fully in production since June 2017
 - 8 M files / 14,5 PB proceed
 - 2 M files on cache
- ▶ Features detailed at HUF 2017 [1]
- ▶ Product stable, no new development since the HUF.
- ▶ Code available for the HPSS community
 - <https://gitlab.in2p3.fr/cc-in2p3-dev/treqs2>
 - License : GPLv3
 - Account opened on request



Treqs Kibana based dashboard :

81,035

Requests

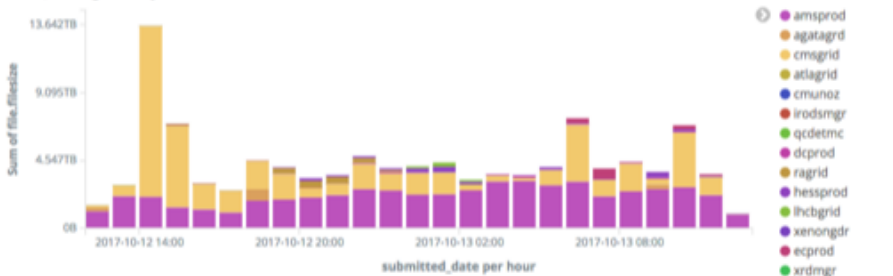
109.947TB

Total Size

TREQS2: Requetes par utilisateurs



TREQS2: Stage rate by users



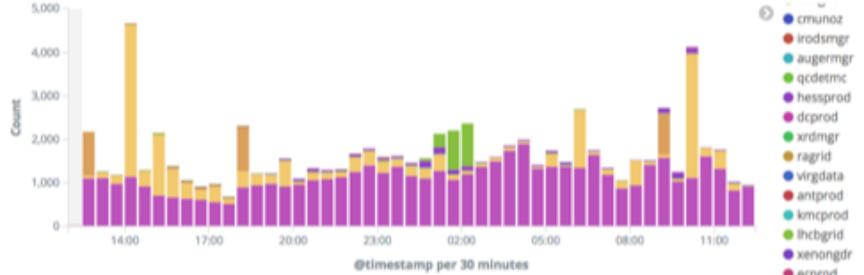
TREQS2: Cache Hints per user



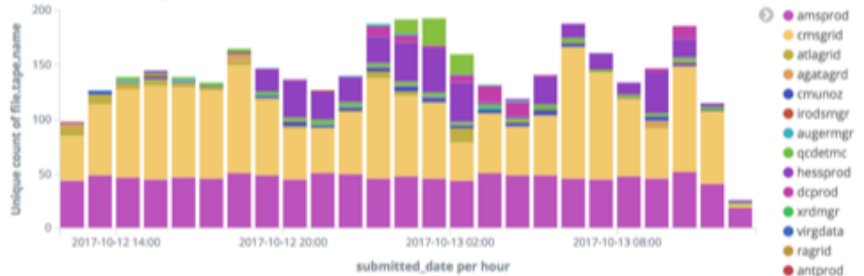
Status	Count	File size
STAGED	64,885	91.115TB
ALREADYONDISK	16,108	18.79TB
FAILED	42	43.245GB

Export: [Raw](#) [Formatted](#)

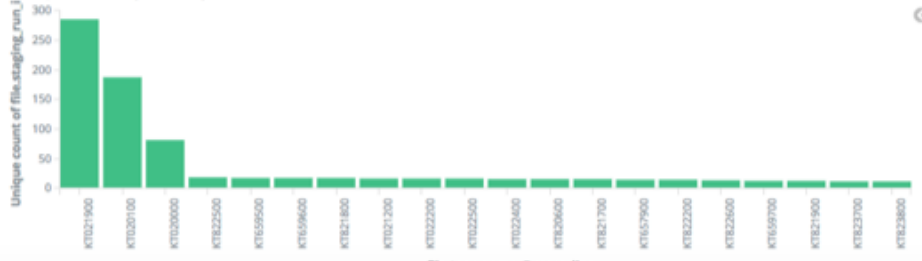
TREQS2: File requests by hour

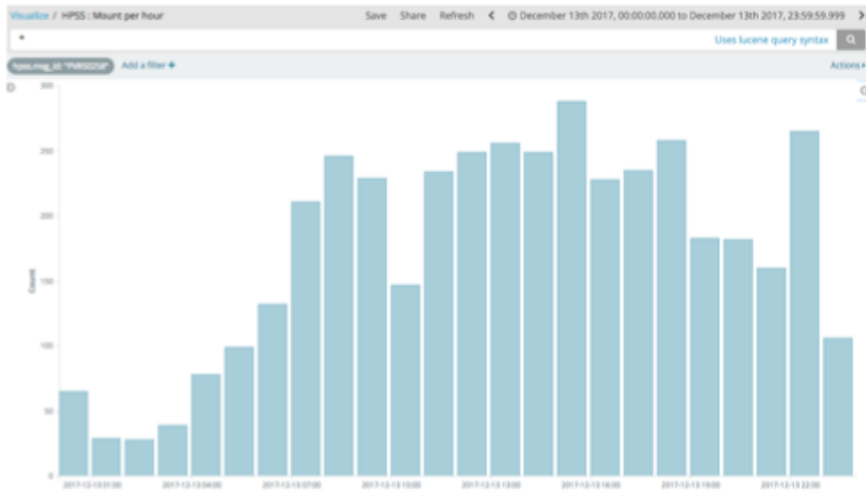


TREQS2: Tape count by users



TREQS2: Most requested tapes





- ▶ HPSS Dashboard :
 - Mount per hour
 - Migration/purge stats

Visualize / HPSS : Migrations / purges stats

Save Share Refresh ◀ December 13th 2017, 00:00:00.000 to December 13th 2017, 23:59:59.999 ▶

Uses lucene query syntax

hps.class_id:"61" hps.class_id:"11" Add a filter +

Actions ▶

Subsys	SC	Migrated bytes	Migrated Files	Purged bytes	Purged Files
1	10	1.727GB	256	0B	0
1	11	125.176GB	509	0B	0
1	12	367.804GB	347	0B	0
1	14	15.897TB	1,224	64.634TB	5,260
2	10	21.608GB	2,755	153.75GB	10,285
2	11	102.348GB	870	250.344GB	1,766
2	12	302.106GB	267	0B	0
2	14	2.621TB	810	0B	0
3	10	0B	0	0B	0
3	11	0B	0	0B	0
3	12	0B	0	0B	0
3	14	786.308GB	198	17.675TB	4,762
4	10	9.294GB	386	10.403GB	227
4	11	18.218GB	101	82.109GB	464
4	12	673.765GB	905	0B	0
4	14	4.369TB	1,152	25.136TB	8,266
5	10	12.325GB	1,280	0B	0
5	11	100.044GB	1,045	0B	0
5	12	1,007.754GB	1,257	0B	0
5	14	43.466TB	5,833	56.903TB	22,356
-		69.8TB	19,195	168.833TB	53,386

- ▶ Lack of overall performance monitoring
 - What is the actual bandwidth delivered by the system
 - What is the maximum capability of the system ?
 - How much data can be ingested by second / hour / day ?
 - How much data can be delivered by second / hour / day ?

- ▶ Count tape remounts
 - How many time the same tape has been remounted within 24h ?
 - Need to compare TREQS and ACSLS logs.

Tape infrastructure

- ▶ **Tape Libraries**
 - 4 Oracle SL8500 Libraries
 - Interconnected (with PTP)
 - Collocated with TSM (backup)
- ▶ **130 Tapes drives**
 - T10K-B/C out of warranty used on tests system
 - LTO 4/6 used for TSM
- ▶ **56 Tapes drives in production for HPSS**
 - 56 T10K-D (8,5 TB on T10K-T2)
 - +6 T10K-D (installed Q1-2018)
- ▶ **1 LTO 8 tape drive for test**
- ▶ **22 000 Tapes**
 - 11500 T10000T2 (8,5 TB)
 - 5 000 LTO 4
 - 2 000 LTO 6
 - 3 500 T10000T1 (to destroy)
- ▶ **Daily tape mounts:**
 - 2 000 average
 - > 6 000 peak
- ▶ **HPSS Repacks**
 - 23,000 T1 → T2 proceed in 2 years
 - 2,000 T10K-C → T10K-D in 2017



- ▶ **Nombreuses erreurs de relectures**
 - Messages CORE3136 : Read error flag set in tape KT553700
 - Quelques dizaines de fichiers de certaines bandes sont illisibles
 - Fichiers écrits entre mi 2016 et début 2017 (?)
 - Corruption silencieuse à l'écriture due à 1 ou plusieurs lecteurs
 - Lecteur(s) impacté(s) non identifié
- ▶ **Les erreurs apparaissent lors de la relecture des données**
 - FSC 37F6 / FSC 3F21
 - Impossible de connaître a priori les fichiers corrompus
 - L'état des dégâts sera connu lorsque toutes les bandes T10K-D seront repackées (2020 ?)
- ▶ **Parc de lecteur assaini suite à une mise à jour de FW en mai 2017**
- ▶ **Etat actuel :**
 - ~ 40 bandes identifiées
 - Plusieurs centaines de fichiers, toute VO confondue.
- ▶ **Bandes envoyées chez OpenStorage pour analyse/restauration**
 - Faible probabilité de récupérer les fichiers.



- ▶ Oracle stopped developing “Enterprise drives” (T10000)
 - T10000-E drives won't be marketed
 - Need to move to a new technology

- ▶ 2 scenarios :
 - Move to IBM Enterprise class tapes drives (Jaguar)
 - Keep our libraries and use LTO drives.

- ▶ IBM Enterprise tapes (Jaguar) :
 - Native capacity :
 - 15 TB on a JD cartridge (TS1155)
 - 20 TB on a JE cartridge (TS1160)
 - Short media (“Sport” Tape) for storing small files.
 - Drive support latest's advanced features
 - 64 landing zone allowing fast positioning
 - Tape Ordered Recall and End To End Data integrity
 - Drive is NOT supported on Oracle libraries → Need to purchase new libraries

- ▶ LTO 8
 - Native Capacity : 12 TB
 - Media cost 25% lower than Enterprise tape and may decrease quickly.
 - Use the same R/W head than Jaguar (TMR) head and BeFe media.
 - But Only 2 landing zones → Performance lower on random recall.
 - Advanced features not supported (TOR)

- ▶ LTO 9
 - Native capacity : 20 TB (?)
 - Not available before 2019/2020
 - Could we wait until these days ?

▶ Choice not evident

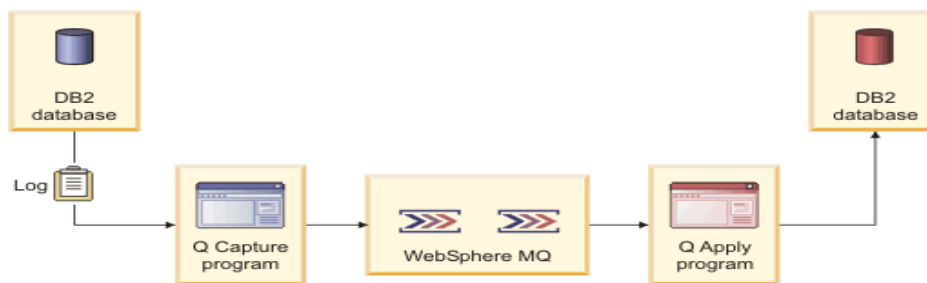
- Reliability/performances of the LTO drives / media ?
 - LTO tapes can support our workload (6000 mount/day) ?
 - Today, we “break” about 8 to 10 drives T10K-D per month.
- Service and support ?
 - Today, T10K-D drives are monitored by Oracle SDP2
 - Service Request opened automatically when a drive fail.
- Our libraries getting old (10 years)
 - Maintenance cost will increase by 50 %
- How long Oracle will continue in the tape business ?

▶ Preliminary tests started on LTO-7

- Tape filled with 2GB files
- Good performances on LTO-7 at migration (writing)
 - Close to 300 MB/s
- Read operations slower on LTO-7
 - Positioning slower on LTO-7 vs T10K-D (-10% to -30%)
 - But performance similar using Treqs (!)
- Tests has to be made with small / medium files size (10 to 100 MB) and aggregates

▶ LTO 8 Tests planned in Q2-2018

- ▶ HPSS 7.5.1 is the new major HPSS version
 - Features presented by J. Procknow at HUF 2017 [2]
 - Database partitioning
 - End To End Data Integrity
 - Tape Ordered Recall + 'Quaid'
 - Many changes in the metadata schema
 - Redesigned for improving NS performances (files creation / deletion)
 - SOID reduced from 32 bytes to 19 bytes
- ▶ Migration based on QREP
 - Designed to reduce downtime
 - Metadata converted while HPSS running



- ▶ Two scenarios :
 - In place metadata conversion (on the same machine)
 - Server to server conversion (data replicated and converted on a target server)

- ▶ **Migration of the test environment completed**
 - Source : HPSS 7.4.3p2 on Openstack VM (RHEL 6.9)
 - 3 subsystems and about 1.1 millions of files
 - Scenario 2 : Migration on a new machine (RHEL 7.4)
 - Target HPSS 7.5.1.2
 - Documentation and tools provided by HPSS support
 - QREP and a set of python scripts
 - IBM Websphere + DB2 licence

- ▶ **My feedbacks :**
 - Some mistakes in the documentation
 - It's not clear which commands has to be run on the source or target server
 - Files and directories permissions has to be tuned
 - Kernel value for hpss 7.5.1.2 (kernel.sem = 4096 2048000 32 4096)
 - Many component need to be deployed on servers
 - Python 2.7.5 must be compiled for RHEL 6.9 servers
 - DB2 python module > 2.0.4 doesn't works
 - Websphere MQ use 10 GB is on the root filesystem
 - Need to create a dedicated partition
 - All the DB must be catalogued on both nodes
 - Both servers are able to access to source an target DB
 - But databases must be catalogued in different way depending the host
 - DB2 Instance need to be restarted anyway
 - To upgrade DB2 v10.5 fp8
 - To set Federated mode
 - Hard to troubleshoot : Sometime no errors messages, but nothing happens

▶ My feedbacks (cont)

- Bug detected at “Verify” step
 - Problem due to default collating sequence of the DB that change the “ORDER BY” results
 - On source DB, default values is “SYSTEM_819” and on target DB, default value is “INDENTITY”
 - Problem quickly identified the HPSS support and a fix was delivered
- Some operations take lot time :
 - Ie : Initial load of the DB (“activate” step)
 - Almost 11 hours for 72 M files on the production database.
 - Qverify step : ~ 12 hours
- Some commands are confusing :
 - ie : stop capture
`./manage_qrep.py -c -s 1 -s 2 -s 3 --stop_capture`
 - ie : restart replication after a reboot :
`./manage_qrep.py -c -s 1 -s 2 -s 3 --stop_capture --start_capture`

▶ Final status of the conversion in qverify output file :

`/hpss_src/QRep/templates/7.4.3.2a_to_7.5.1/replication_logs/qverify_run_cfg_s1_s2_s3_s4_s5.log`

▶ Current status :

- Target databases synced with the sources databases
- Each changes on the source (while hpss running) is applied within ms on the target
- Next step : Stop the replication and switch HPSS to the target server

▶ Schedule for the production :

- March 2018 : Setup QREP and start replication
- June 2018 : Migrate to HPSS v7.5.1p2

- ▶ API change :
 - div_cl64m removed from the API
- ▶ Isvol :
 - Now display the full path of file
 - Fileset is not displayed anymore

- ▶ Propose an alternative to RFIO
 - RFIO use deprecated readlist/writelist API
 - May be remove from HPSS API in HPSS 8

- ▶ Alternate tools
 - HSI / HTAR ?
 - Xrootd ?

- ▶ Expose HPSS filesystem ?
 - GHI ?
 - Expensive !
 - VFS over NFS
 - Problem with COS selection when creating file
 - Lustre HSM ?
 - Conflict with GPFS client ?

Thank you

- [1] <https://conference-indico.kek.jp/indico/event/28/session/10/contribution/25>
- [2] <https://conference-indico.kek.jp/indico/event/28/session/6/contribution/9>