



A.Tsaregorodtsev, CPPM LCG-France, IRFU, 1 Nov 2014







- o Current operations
- Updated Computing Model for Run2
- o T2D status
- o Clouds, Vac, Vcycle
- Parallel calculations
- o Federated storage
- o Catalogs
- Data popularity and storage optimization



#### **Current operations**



# Grid usage in general

- Main metrics for this period stay unchanged
  - Dominated by Monte Carlo
  - 91 % of all payloads
     successfully executed
  - All Tier1 sites in the top of used sites (+ HLT & Yandex)
  - Main activities with very high CPU efficiency ~ 100 %





3 Nov '14

#### **Current operations**



# **Resource Usage**



- Low usage of resources in this period
  - Physicists waiting for new MC models to submit productions





# Normalized CPU by Country



Generated on 2014-10-27 20:38:59 UTC





### Lyon T1 center

- o Pledge/consumed Oct 2014
  - CPU: 21700 HEPSPEC'06 years, consumed 12500
    - Pledged ~18% of all the T1's
    - Consumed ~13% of all the T1's
  - Disk: 1.35 PB (allocated 1.4), consumed 1.15
  - Tape: 1.8 PB, consumed 1.35

- CC/IN2P3 preliminary pledges and percentage of the LHCb request
  - CPU: 23 000 HEPSPEC'06 (~19%)
  - Disk: 1.88 PB (~15%) (complemented by T2D storage)
  - Tape: 4.36 PB (~18%)
- The pledges 2015 are very satisfactory







- Currently running Stripping 21 campaign to produce "legacy Run I dataset"
  - As the number of stripping code lines doubled, fighting with memory requirements ( ~ 3GB/job ) – under control
- Using FTS3 for prestageing data in the BUFFER Storage Elements
  - Great simplification of operations and workflows
- Setting up Federated storage setup
  - HTTP/Webdav access





- RAW files will be produced in 2015 with a trigger rate of ~ twice
   2012 (10 kHz RAW + 2.5 kHz turbo)
  - Assume 25 TB of RAW / day
- Prompt reconstruction will be delayed for using final calibrations
  - Assume 2 weeks delay
- The size ratio FULL.DST/RAW will remain the same as 2012
  - 45 TB of FULL.DST / day
- o Stripping will hopefully be running almost in parallel with reco
  - Avoid having too large buffer used
- As soon as the RAW is in CERN-RAW, replicate to:
  - XXX-RAW (XXX being a Tier1, selected according to the RAW shares)
  - And in addition, to:
    - > YYY-BUFFER: if it is decided to reconstruct the run at YYY
- Main difference with now:
  - Processing input SE selected at replication time (by run)





#### • **Reconstruction**:

- Run from Tier1-BUFFER
- Upload FULL.DST to YYY-BUFFER
- Replicate to the same YYY-RDST (or another)
- o Stripping
  - Run from YYY-BUFFER replica of FULL.DST
  - Remove FULL.DST from BUFFER
    - Issue a removeReplica request in the job
  - Merging
    - As now, upload preferentially to selected YYY-DST but can be to any Tier1-DST (no need for FAILOVER)





- Jobs can run anywhere, preferentially where data is, but not necessarily if download policy:
  - Introduce new concept in JDL
    - Primary sites (SE is accessible by protocol for remote file opening)
    - Secondary sites (SE is allowed for downloading to the WN)
  - Matching:
    - If no job available in matching by a pilot using the site as as Primary Site, try and match as Secondary site
    - > If no specific explicit site matching available, match from ANY
  - User jobs can only use primary sites
- o This is a challenge for efficient matching





## Side Status Board

○ ○ ○ ⑤ Site Status Board × +									
💿 🕑 dashb-lhcb-ssb.cern.ch/dashboard/request.py/siteview#currentView=By+tier&sorting_SubTier=asc 🗸 C 🔞 Coogle							Q ☆ 自 🕹 合 🔲 🗄	=	
Help Login Register for site notifications						Site	Status for the LHCB sites, null 📌		
Show 200 - Print II Save View: By tier - Posearch_ entries									
Site Name 🗘	SubTier 🔷	SiteMask 🗘	Availability 🗘	Reliability \$	Storage Online	Storage Used 🗘	SE Status 🗘	Attached Tier-1 🗘	
LCG.RRCKI.ru	1: T1	ck	100	100	200.0	25.924	active		
LCG.SARA.nl	1.71	ck	100	100	1070.0	905.795	active		
LCG.CPPM.fr	2: T2-D	ck	100	100	109.201	0.637	active		
LCG.CSCS.ch	2: T2-D	ck	100	100	290.0	87.285	active		
LCG.IHEP.su	2: T2-D	ok	100	100	64.0	49.463	active	RRCKI	
LCG.LAL.fr	2: T2-D	ck	100	100	109.951	3.973	active		
LCG.LPNHE.fr	2: T2-D	ck	100	100	142.936	4,441	active		
LCG.Manchester.uk	2: T2-D	ck	100	100	300.0	88.444	active	CERN	
LCG.NCBJ.pl	2: T2-D	ck	100	100	300.0	68.217	active		
LCG.NIPNE-07.ro	2: T2-D	ck	100	100	219.902	1.734	active		
LCG.RAL-HEP.uk	2: T2-D	ok			455.0	129.058	active		
LCG.CNAF-T2.it	3: T2-R	ck	100	100				CNAF	
LCG.Dortmund.de	3: T2-R	ck	100	100				CERN	
LCG.JINR.ru	3: T2-R	ck	100	100				RRCKI	
LCG.KIAE.ru	3: T2-R	ck	100	100				RRCKI	
LCG.ARAGRID-CIENCIAS.es	4: T2-M	ck	100	100					
LCG.AUVER.fr	4: T2-M	ok	100	100					
LCG.BHAM-HEP.uk	4: T2-M	ck	100	100					
LCG.BIFI.es	4: T2-M	ok	100	100					
LCG.BMEGrid.hu	4: T2-M	ck	0	0					
LCG.Barcelona.es	4: T2-M	ck	100	100					
Showing 1 to 82 of 82 entries							First Previous 1	lext Last 🔼 📰 🖉 😒	







#### o 9 official T2D sites

- **dCache: CSCS.ch, IHEP.su, RAL-HEP.uk**
- DPM: CPPM.fr, LAL.fr, LPNHE.fr, Manchester.uk, NCBJ.pl, NIPNE-07.ro
- LAL.fr and LPNHE.fr share some infrastructure and offering 300TB together
- CPPM.fr offers ~100TB now
- T2D sites already provide 1991 TB disk storage for analysis jobs which meets the 2015 requirements
- DPM xrootd problem fixed and deployed at all DPM sites
- o Mainly used for user analysis of DST data







#### Storage used at T2-Ds





Generated on 2014-11-01 22:13:04 UTC



- o LHCb users Vac and Vcycle
- o Both VM Lifecycle Managers
- Vac is a standalone daemon you run on each worker node machine to create its VMs
- Vcycle manages VMs on IaaS Clouds like OpenStack
  - Can be run at the site, by the experiment, or by regional groups like GridPP
- Both developed at Manchester as part of GridPP Clouds/VMs effort
  - With help from Lancaster, Oxford, IC, CERN, LHCb and ATLAS
- o Both make very similar assumptions about how the VMs behave
  - The same LHCb VMs working in production on Vac and Vcycle









### Vac's Vacuum Model







Since we have the



#### LHCb jobs from Vac and Vcycle VMs





Generated on 2014-11-04 11:20:56 UTC



#### **Vac and Vcycle**

#### o Both Vac and Vcycle assume the VMs have a defined lifecycle

- The VM runs and its state is monitored
- VM executes shutdown -h when finished or if no more work available
- Can use machine/job features
  - Information on VM resources, e.g. number of cores
  - Information on remaining CPU, etc
- Machine/Job features are intended to be used in the DIRAC pilot running in the LHCb VMs
  - Knowledge of available resources
  - Graceful termination before the slot is exhausted





- o Both Vac and Vcycle can create multiple processor VMs
- LHCb VM architecture is a simple wrapper around dirac pilot script, so running one multiprocessor payload in a VM should be straightforward
  - Using DIRAC Pilot 2.0 architecture
    - > On the fly configurable pilots in terms of parameters and functionality
    - > Executing a customizable, e.g. VO dependent, list of commands/
- Managing multi-core slots is done by one or more JobAgents running inside the pilot
  - JobAgent is running a small local "batch system" for managing cores in its VM
  - Problems to solve
    - > How to reconcile multi-core job requirements and priorities
    - Masonry problem





## Multi-core node scheduling "Masonry" problem

## The Masonry Problem







GPU

- LHCb Gaudi based applications capable of using multiple cores
  - How to get such jobs scheduled on the grid
  - Mixture of jobs with different number of cores requirements
- The use of GPUs is explored in the context of the HLT software trigger
  - Gains and overheads





o Might be applied later for offline



Offloadir



- Currently setting up a prototype for http/webdav to all LHCb storages
  - HTTP/WEBDAV can become a main access protocol in the future
  - HTTP federation on top of these accesses
  - All but 5 sites available

000	And Whoth HCh/Collisio x	WMS history plots as hob	Data Storage plots as incb	WMS history plots as lhob.	. 🛃 4th LHCb Computing Wor	🄌 Miczo OLED Breakout - L 🕂 🕹		
					▽ C ☆ 白 (図 ~ Google Q) + - =			
Se 1 🗳	Apacha CouchOB (``) Read later	r 🔒 LocalTwiki 👫 Mydoba	★ ix 📋 CERN + 🚞 Prog +	📄 Privata = 🔹 v@Dirac	💱 pringr@Disc 🚫 BUM 🐴	atd 🚞 MachFinit + 🚺 CM 🚞 MasCPU + 30		

/fed/lhcb/LHCb/Collision12/BHADRON.MDST/00034960/0000/

Mode	UID	GID	Size	Modified	Name
-rw-rw-r	0	0	13.2M	Tue, 04 Mar 2014 10:14:24 GMT	00034960 0000001 1.bhadron.mdst     00034960 0000001 1.bhadron.mdst
-rwxrwxrwx	0	0	41.4M	Tue, 04 Mar 2014 14:23:18 GMT	🎨 🛅 00034960_0000032_1.bhadron.mdst
-FW-FW-F	0	0	1.96	Tue, 04 Mar 2014 10:12:22 GMT	00034960_00000076_1.bhadron.mdst     00034960_00000076_1.bhadron.mdst
-IWXIWXIWX	0	0	58.4M	Wed, 26 Feb 2014 22:53:31 GMT	🎨 🛅 00034960 00000085 1.bhadron.mdst
-rwxrwxrwx	0	0	61.1M	Thu, 27 Feb 2014 00: 2 C GMT	00034960_0000095_1.bhadron.mdst
-IWXIWXIWX	0	0	937.7M	Tue, 04 Mar 2014 191:18	00034960_00000140_1.bhadron.mdst
-rwxrwxrwx	0	0	36.0M	Tue, 04 06 19 10 1135 GMT	00034960_00000150_1.bhadron.mdst
-rwxrwxrwx	0	0	2.0G	Tue, 4 Mar 2014 09:45:27 GMT	& 00034960 00000159 1.bhadron.mdst
-rw-rw-r	0	0	2.6G	Tue, D. Mar 2014 10:45:52 GMT	🎨 🖣 00034960 00000191 1.bhadron.mdst
	0	0	2.76	10 04 Nat 2014 11:13:47 GMT	8 00034960_00000288_1.bhadron.mdst
-metalink version="3.0" gener	rator="logdm-day"	pubdate="Tee	,04 Mar 2014 10:31:18	CMT'>	00034960_00000327_1.bhadron.mdst
- «files» - «file name="//hch'l.">					00034960 00000341 1.bhadron.mdst
<ul> <li>- size-983231858-(size)</li> <li>- sizesources&gt;</li> </ul>	•				00034960 00000351 1.bhadron.mdst
<ul> <li>- carl type="https"&gt; https://D1-060-124</li> </ul>	-e gridka de: 2880/p	nfi/gidla.de/I	00034960 00000383 1.bhadron.mdat		
- durbs			00034960 00000401 1.bhadron.mdst		
https://101-060-126	e gridka de 2880/p	pafs/gridlu.co/I	00034960 00000410 1.bhadron.mdst		
- duri type="https">-	2+3 <del>1-</del> 700/John	WTh/Callinian)	00034960_00000421_1.bbadron.mdst		
<urb></urb>	aport association.	execution in	CONTRACTOR OF CONTRACTOR	and a second sec	00034960_0000436_1_bbadron_mdst
office					C Starte Contraction
«/film>					00034960_00000450_1.bhadron.mdst
					D 00024960 00000469 1 bhadron mdat

Request by nobody (nobody) Powered by LCGDM-DAV 0.16.0





### LFC to DFC migration

- Even if Federated storage is a success, the File Catalog is still needed
  - Access control rules
  - Metadata standard and user defined
  - Smart data placement policies
  - Efficient storage usage reports
- LHCb has chosen to replace the LFC by the DIRAC File Catalog
  - The last among the LHC experiments
  - The migration will be done in January 2015, before the RUN 2
    - A hours pause in operations to copy over the contents of LFC to DFC
    - > Possibility to use LFC and DFC in parallel for a transition period
  - The tests, functionality and performance, are ongoing
    - Results are optimistic
  - DAV access would be needed to couple nicely with Federated storage





### Number of replicas vs Computing Model













### Data Popularity for storage usage optimization

# Storage space saving







### Data Popularity for storage usage optimization

# Storage saving strategies examples







- Very good efficiency of operations in 2014
- o 2015 pledges of French T1 resources are very satisfactory
- T2D sites are getting into production
  - Will be fully used with RUN2 start
- Using cloud virtualized resources with Vac and Vcycle in production on several sites
  - More to come
- o Using multi-core nodes is possible
  - More intelligent scheduling to increase efficiency is needed
- Using Federated storage is demonstrated integrating 13 out of 19 LHCb sites
- Replacing LFC by DFC is being prepared
- Analysis of data popularity can help to reduce the usage of disk storage

