Computing trends for WLCG

Ian Bird

ROC_LA workshop

CERN; 6th October 2010



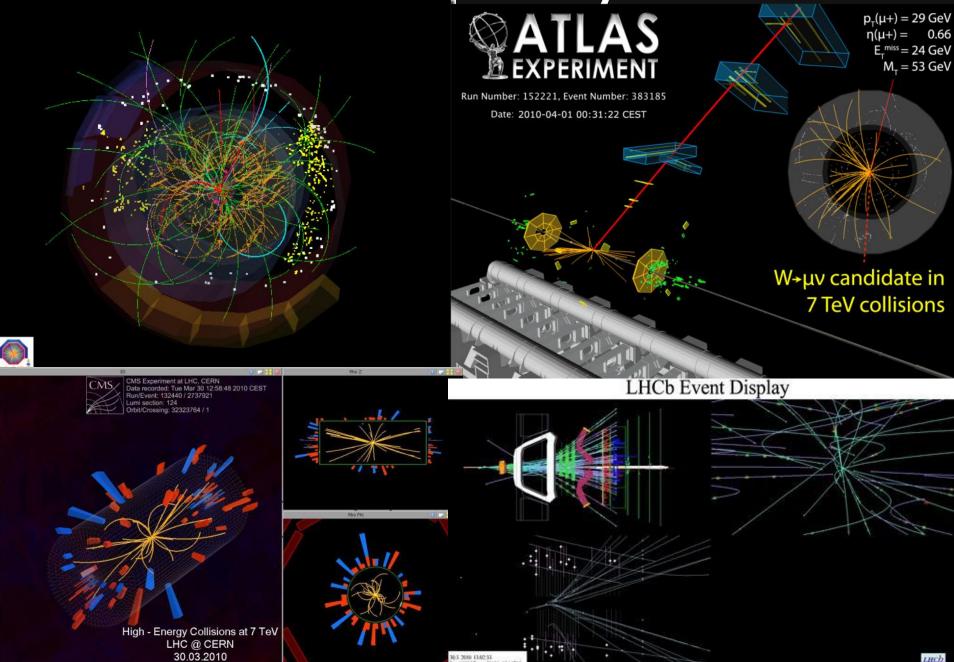






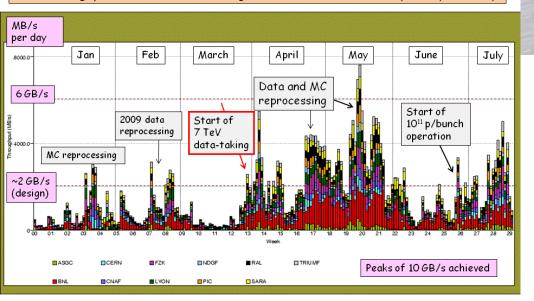


... And now at 7 TeV



Worldwide data distribution and analysis

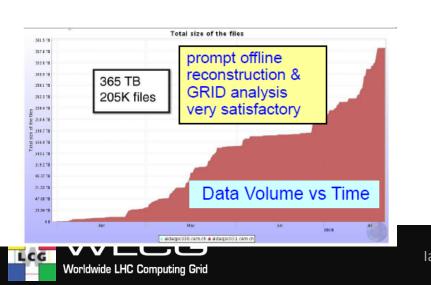
Total throughput of ATLAS data through the Grid: from 1st January until yesterday



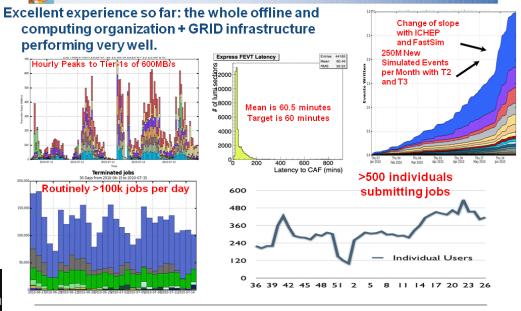
Progress

- Overall is clear physics output in very short time
- Huge effort: Combination of experiment sw & computing and grid infrastructures
- And a lot of testing!

GRID-based analysis in June-July 2010: > 1000 different users, ~ 11 million analysis jobs processed



Data Processing, Transfer and Analysis Activities



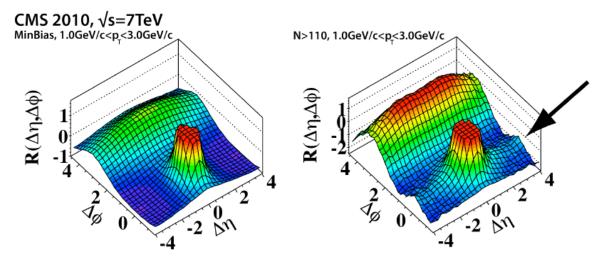
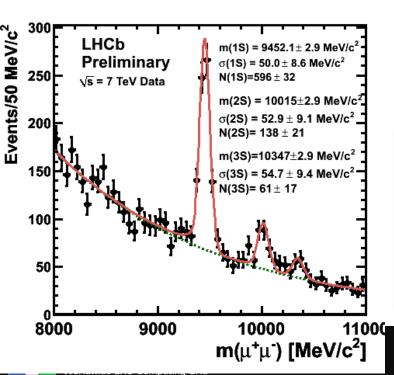
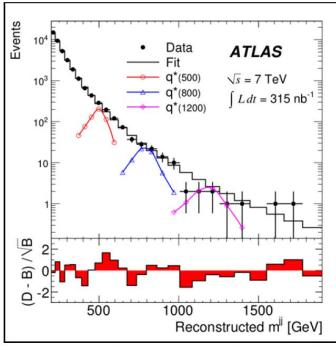


Figure 2: The variation of R with $\Delta \eta$ and $\Delta \phi$, for proton-proton collisions in CMS. Left: for minimum bias collisions; Right: for collisions that produced at least 110 charged particles.



ATLAS sets world's best limits on q*

7 September 2010



Reconstructed dijet mass (filled points) fitted to a smooth background distribution, with predicted q* signal indicated at three different masses.

- Progress == real physics output in a very short time ...
- "rediscovery" of Standard Model
- Starting to see improvements on state of the art, and hints ...

Worldwide LHC Computing Gric

From testing to data:

Independent Experiment Data Challenges

2004

Service Challenges proposed in 2004

To demonstrate service aspects:

- -Data transfers for weeks on end
- -Data management
- -Scaling of job workloads
- -Security incidents ("fire drills")
- -Interoperability
- -Support processes

2005

2006

2007

2008

2009

e.g. DC04 (ALICE, CMS, LHCb)/DC2 (ATLAS) in 2004 saw first full chain of computing models on grids

SC1 Basic transfer rates

SC2 Basic transfer rates

SC3 Sustained rates, data management, service reliability

SC4 Nominal LHC rates, disk→ tape tests, all Tier 1s, some Tier 2s

CCRC'08 Readiness challenge, all experiments, ~full computing models

STEP'09 Scale challenge, all experiments, full computing models, tape recall + analysis

• Focus on real and continuous production use of the service over several years (simulations since 2003, cosmic ray data, etc.)

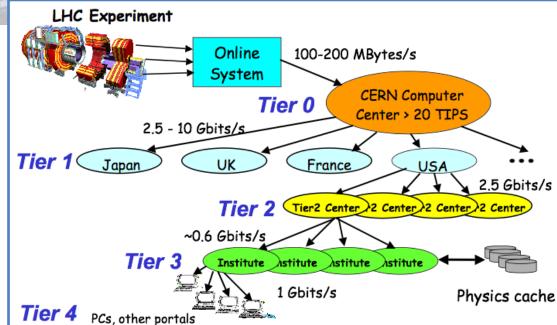
 Data and Service challenges to exercise all aspects of the service – not just for data transfers, but workloads, support structures etc.

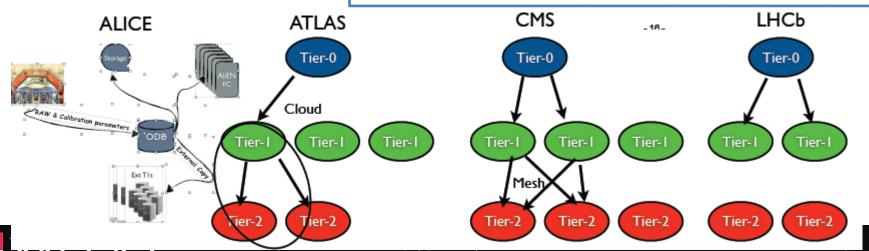
LCG

2010 Ian Bird, CE

Experiment models have evolved

- Models all ~based on the MONARC tiered model of 10 years ago
- Several significant variations, however





Some observations

- ✓ Experiments have truly distributed models
- ✓ Network traffic is close to what was planned and the network is extremely reliable
- ✓ Significant numbers of people (hundreds) successfully doing analysis at Tier 2s
- ✓ Physics output in a very short time unprecedented
- ✓ Today resources are plentiful, and not yet full; This will surely change ...
- ✓ Needs a lot of support and interactions with sites
 - heavy but supportable



Other observations

- Availability of grid sites is hard to maintain...
 - 1 power "event"/year at each of 12 Tier0/1 sites is 1/month
 - DB issues are common use is out of the norm?
 - Still takes considerable effort to manage these problems
- Problems are (surprisingly?) generally not middleware related ...
- Actual use cases today are far simpler than the grid middleware attempted to provide for
- Advent of "pilot jobs" changes the need for brokering
- Hardware is not reliable, no matter if it is commodity or not; RAID controllers are a spof
 - We have 100 PB disk worldwide something is always failing
- We must learn how to make a reliable system from unreliable hardware
- Applications must (really!) realise that:
 - the network is reliable,
 - resources can appear and disappear,
 - data may not be where you thought it was
 - even at 0.1% this is a problem if you rely on it at these scales!

Providing reliable data management is still an outstanding problem

The EEE Production Infrastructure



LCG

Test-beds & Services

Production Service

Pre-production service

Certification test-beds (SA3)

Training infrastructure (NA4)

Support Structures & Processes

Operations Coordination Centre

Regional Operations Centres

Global Grid User Support

EGEE Network Operations Centre (SA2)

Operational Security Coordination Team

Training activities (NA3)

Security & Policy Groups

Joint Security Policy Group

Grid Security Vulnerability Group

EuGridPMA (& IGTF)

Operations Advisory Group (+NA4)

What is WLCG today?

An infrastructure demonstrated to be able to support LHC data processing and analysis: consisting of (in order of importance!)

- Above all a collaboration:
 - Single point of coordination, communication, requirements synthesis, ...
 - Vehicle to coordinate resources and funding
 - Extremely useful in organising with technology and service providers
 - In EC-speak WLCG is a highly structured community via the MoU
- A service:
 - WLCG provides a common operations coordination and management on top of EGEE/EGI, OSG, and others
 - Security coordination operational and policy development
 - World-wide trust federation of CA's and VO's
 - The Policy framework was indispensable in actually deploying WLCG across the world
- An implementation of a distributed computing infrastructure
 - Today with grid technology and higher level (WLCG and experiment-specific) middleware
 - Tomorrow ...





Evolution and sustainability

- Need to adapt to changing technologies; e.g.:
 - Major re-think of storage and data access
 - Use of many-core CPUs (and other processor types?)
 - Filesystems, etc.
 - Virtualisation as a solution for job management
 - Brings us in line with industrial technology
 - Integration with public and commercial clouds
- Network infrastructure
 - This is the most reliable service we have
 - Invest in networks and make full use of the distributed system
- Grid Middleware
 - Complexity of today's middleware compared to the actual use cases
 - Evolve by using more "standard" technologies: e.g. Message
 Brokers, Monitoring systems are first steps





Areas for evolution – 1

- End to end usable and transparent networks and data movement:
 - for the 100 Gigabit era;
- Data issues:
 - Data management and access
 - How to make reliable systems from commodity (or expensive!) hardware
 - Fault tolerance
 - Data preservation and open access
- Global AAI:
 - **–** SSO
 - Evolution/replacement/hiding of today's X509
 - Use existing ID federations?
 - Integrate with commercial/opensource software?





Evolution – 2

- Support and evolution of the current grid middleware services for maintainability and effectiveness (e-infrastructure):
 - Evolve towards sustainable and externally provided/supported services for fabric layer
 - Improvement in efficiencies for use, support and maintenance
- Technology evolution from research to end-to-end production deployment in research/scientific codes:
 - Multi-core
 - Virtualized environments.
 - Commercial clouds.
 - GPUs
- Green computing
 - Innovations in end-to-end application, middleware and fabric design, technology and process that increase energy efficiency
 - Use of remote data centres







- 1st workshop held in June
 - Recognition that network as a very reliable resource can optimize the use of the storage and CPU resources
 - The strict hierarchical MONARC model is no longer necessary
 - Simplification of use of tape and the interfaces
 - Use disk resources more as a cache
 - Recognize that not all data has to be local at a site for a job to run allow remote access (or fetch to a local cache)
 - Often faster to fetch a file from a remote site than from local tape
- Data management software will evolve
 - A number of short term prototypes have been proposed
 - Simplify the interfaces where possible; hide details from end-users
- Experiment models will evolve
 - To accept that information in a distributed system cannot be fully up-todate; use remote access to data and caching mechanisms to improve overall robustness
- Timescale: 2013 LHC run





Data management - 1

- Can probably improve use of existing resources with some reasonable steps
- Use the network to access data not found locally – don't insist that all data be present
- Storage:
 - Separate archive from cache and allow only organised access to archive
 - Random user read of archive is not supportable
 - Can simplify archive and cache interfaces
 - Can then consider industrial solutions for the archives and uses robots/drives in more normal manner





Data management - 2

Data access:

- Job should not assume all files will be local get missing files (either remote access or cache)
- Catalogues may not be fully up to date (i.e. Need to be able to correct by remote access)
- Need effective caching mechanisms and organised data placement tools
 - Policies/Algorithms to be flexible
- User access model is that of filesystem
 - No complexities like SRM should be visible
- Reliability is critical. "Reliable" hardware, software adaptation





Data management – 3

Data transfer

- Support organised data placement
- Support data caching
- Support for remote access to data not local, or bring data locally as needed
- Asynchronous, reliable data movement (e.g. From MC to archive)
 - Ensure data is delivered and catalogued

– FTS:

- Well defined recovery from failure
- Partial transfers, or partial files. Recovery from failure
- Manipulate datasets as a single entity





Data management - 4

- Namespaces, catalogues, auth², etc.
 - Need catalogues to reflect storage contents and be synchronized (asynchronously)
 - Applications must recognise that information may not be 100% correct (or up to date)
 - Synchronization could use e.g. Message bus, or by more dynamic (DHT, etc.)
 - Need for global namespace with both LFNs and **GUIDs**
 - ACLs must be implemented (but only once!) avoid back-doors





Data management - 5

- Global home directory
 - Is required (e.g. Drop boxes, Amazon S3, etc.)
- Demonstrator projects:

D1	ATLAS Dynamic data placement	Graeme Stewart	
D2	LHCb Dynamic data placement	Philippe Charpentier	
D3	ARC Caching	David Cameron	
D4	Proxy caches	Dirk Duellmann	
D5	CoralCDN	Jeff Templon	
D6	MSG/catalogue synchronisation	J-P Baud	
D7	MSG/ACL propagation	J-P Baud	
D8	NSF4.1 as access protocol	Patrick Fuhrmann	
D9	Xrootd-global: CMS	Brian Bockelman	
D10	Xrootd-global: ATLAS/IT large-scale tests	Dirk Duellmann	
D11	Cassandra/Fuse as LFC/SRM alternative	Oscar Koeroo	
D12	CHIRP Ian.Bird@cern.ch	Rod Walker	

Networking

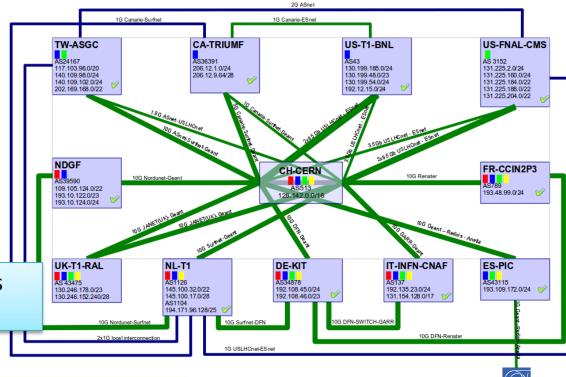
Data transfer:

•SW: gridftp, FTS (interacts with endpoints, recovery), experiment layer

•HW: light paths, routing, coupling to storage

Operational: monitoring

+ the academic/research networks for Tier1/2!



LHC PN

- Requirements wg to discuss with network communities
- How to control traffic?
 - Today in FTS, but what if a VO does not use FTS?

T1-T1 traffic only Not deployed yet

(thick) >= 10Gbps

At the network layer? Over provision?

= internet backup available

p2p prefix: 192.16.166.0/24 edoardo.martelli@cern.ch 20100916



Multi-core jobs Grid-wide

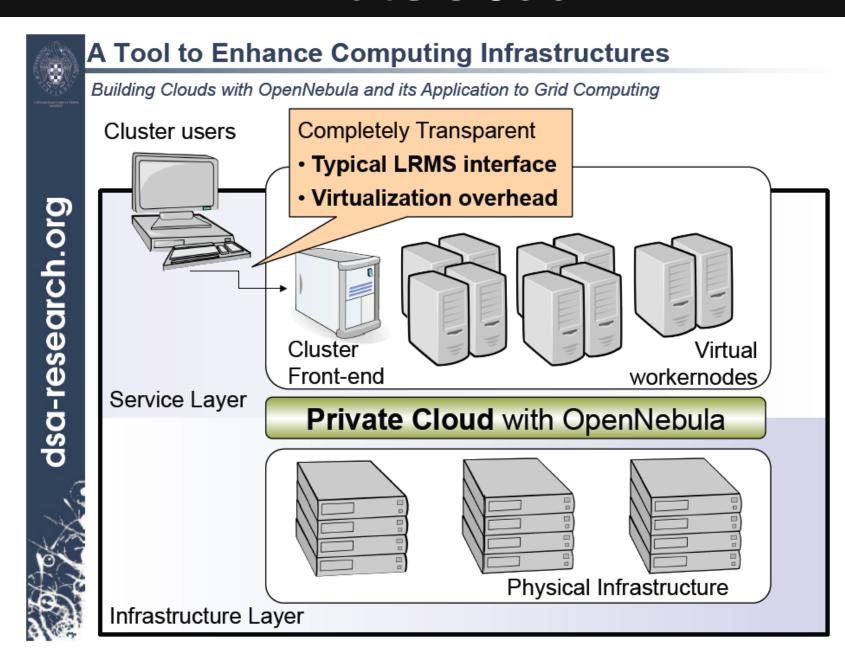
- Experiments may want to request entire nodes
 - To better manage the cores by using pilot frameworks to optimize the total workload on a box
- Needs some additional brokering
 - To request entire nodes
 - To request a number of cores
 - Very similar use case to MPI jobs
 - Should be supported by CREAM



Virtualization and "Clouds"

- Virtualization potentially provides:
 - Better optimization of resource usage at a site
 - And hence power usage
 - Dynamic provisioning of services on demand
 - Breaking of dependencies OS<->m/w<->application
 - Etc.
- "Clouds"
 - Can mean anything
 - Private cloud → A site using virtualisation to better manage its resources; provides a cloud interface (which one?)
 - Public cloud → commercial resource e.g. Amazon EC2
 - Both are being deployed and experimented with

Private cloud









Virtualisation Activities at **CERN**



- In many areas not all directly relevant for LHC computing
- 3 main areas:
 - Service consolidation
 - "VOBoxes", VMs on demand, LCG certification testbed
 - "I X" services
 - Worker nodes, pilots, etc (issue of "bare" WN tba)
 - Cloud interfaces

Rationale:

- Better use of resources optimise cost, power, efficiency
- Reduce dependencies esp between OS and applications (e.g. SL4 → SL5 migration), and between grid software
- Long term sustainability/maintainability → can we move to something which is more "industry-standard" ?
- don't forget WLCG issues of how to expand to other sites
 - which may have many other constraints (e.g. may require virtualised WN)
 - Must address trust issue from the outset







Service consolidation



- VO Boxes (in the general sense of all user-managed services)
 - IT runs the OS and hypervisor; user runs the service and application
 - Clarifies distinction in responsibilities
 - Simplifies management for VOC no need to understand system configuration tools
 - Allows to optimise between heavily used and lightly used services
 - (eventually) transparent migration between hardware: improve service availability
- VMs "on demand" (like requesting a web server today)
 - Request through a web interface
 - General service for relatively long-lived needs
 - The user can request a VM from among a set of standard images
 - E.g.:
 - ETICS multi-platform build and automated testing
 - LCG certification test bed. Today uses a different technology, but will migrate once live checkpointing of images is provided.



CVI: CERN Virtualisation Infrastructure

- Department
- Based on Microsoft's Virtual Machine Manager
 - Multiple interfaces available
 - 'Self-Service' web interface at http://cern.ch/cvi
 - SOAP interface
 - Virtual Machine Manager console for Windows clients
- Integrated with LANdb network database
- 100 hosts running Hyper-V hypervisor
 - 10 distinct host groups, with delegated administration privileges
 - 'Quick' migration of VMs between hosts
 - ~1 minute, session survives migration
- Images for all supported Windows and Linux versions
 - Plus PXE boot images







CVI: some usage statistics



Today, CVI provides 340 VMs...

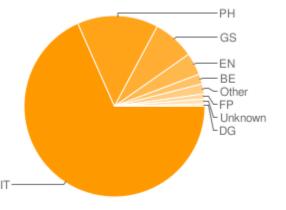
70% Windows, 30% Linux

... for different communities

- Self-service portal
 - 230 VMs from 99 distinct users
 - Mixture of development and production machines



- Engineering Services: 85 VMs
- Media streaming: 12 VMs
- 6 Print servers, 8 Exchange servers
- ... etc





"LX" Services



LXCloud

- See presentation of Sebastian Goasguen + Ulrich Schwickerath at workshop
- Management of a virtual infrastructure with cloud interfaces
- Includes the capability to run CernVM images
- Scalability test are ongoing
- Tests ongoing with both Open Nebula and Platform ISF as potential solutions

LXBatch

- Standard OS worker node: as VM addresses dependency problem
- WN with a full experiment software stack
 - user could choose from among a standard/certified set. These images could e.g. Be built using the experiment build servers.
- As the previous case but with the pilot framework embedded
- CERNVM images
- Eventually: LXBatch → LXCloud







Evolution



	Today	(SLC = Scientific Linux CERI
	Batch	
Physical SLC4 WN	Physical SLC5 WN	

Near future:

Batch							
SLC4 WN SLC5 WN	Physical	Physical SLC5 WN					
hypervisor cluster	SLC4 WN						

(far) future?

Batch	TO	development	other/cloud applications					
Internal cloud								







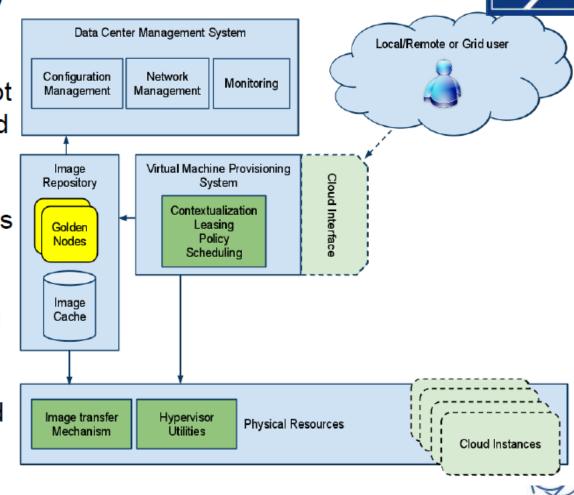
LXCloud



CERN's Ixcloud architecture

- Image repository with Golden nodes.
- VM instances not quattor managed have finite lifetime
- Specific IP/MACs are pinned to hypervisors
- Currently testing two provisioning system:
 Opennebula and

Opennebula and Platform ISF.



CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it



Ongoing work



- Integration with existing management infrastructure and tools
 - Including monitoring and alarm systems
- Evaluating VM provisioning systems
 - Opensource and commercial
- Image distribution mechanisms
 - With P2P tools
- Scalability tests
 - Batch system (how many VMs can be managed)
 - Infrastructure (network database, etc.)
 - Image distributions
 - VM performance
- To be understood:
 - I/O performance, particularly for analysis jobs
 - How to do accounting etc.
 - Virtualised infrastructure vs allocate "whole node" to application





0

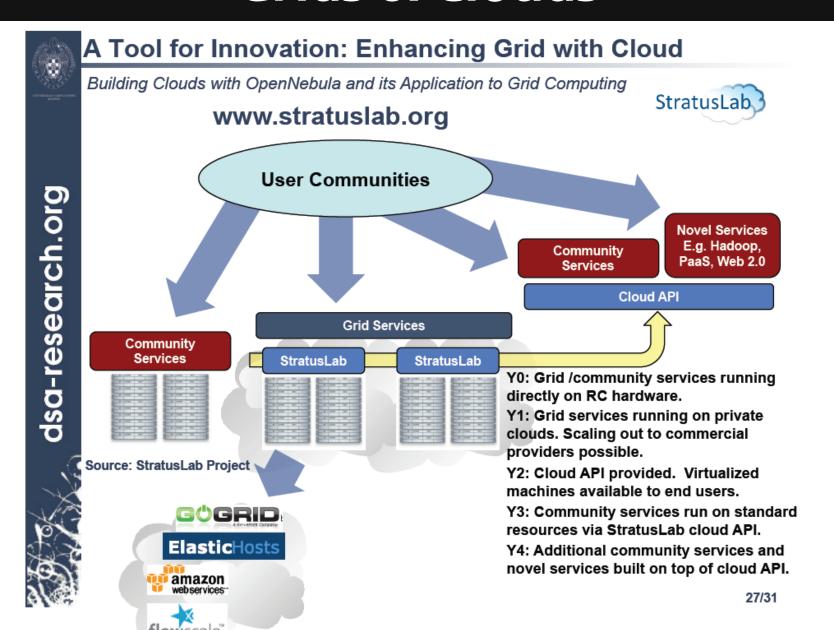








Grids & Clouds



Hybrid clouds

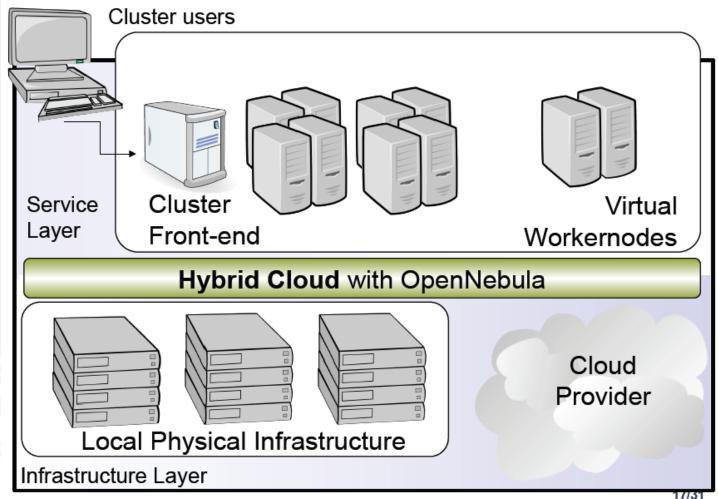








Building Clouds with OpenNebula and its Application to Grid Computing





Hybrid clouds

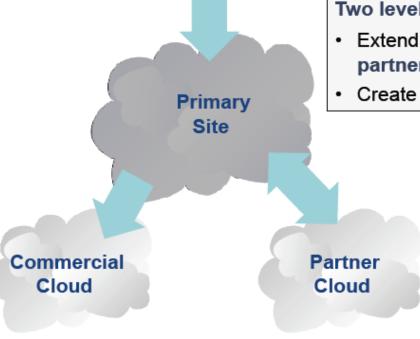


Innovations: The Business Perspective

Innovation in Cloud Computing Architectures

Hybrid Cloud Computing and Federation

- Cloudbursting at infrastructure layer, fully transparent to users
- Scale-out decisions are taken by infrastructure administrators according to business policies



Two levels of Collaboration

- Extend the private cloud using both partner and commercial clouds
- Create a federation of clouds

Sounds familiar?

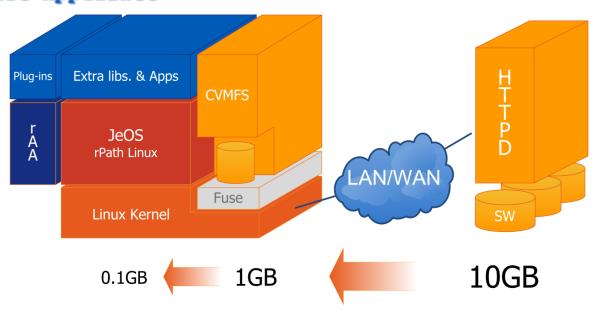


WP9 - Virtualization R&D

- Aims to provide a complete, portable and easy to configure user environment for developing and running LHC data analysis locally and on the Grid independent of physical software and hardware platform (Linux, Windows, MacOS)
 - Code check-out, edition, compilation, local small test, debugging, ...
 - Grid submission, data access...
 - Event displays, interactive data analysis, ...
 - Suspend, resume...
- Decouple application lifecycle from evolution of system infrastructure
- Reduce effort to install, maintain and keep up to date the experiment software (CernVM-FS)
- CernVM 1.x (SLC4) and CernVM 2.x (SLC5) released
 - Small (200-350 MB) image
 - Available for all popular hypervisors and on Amazon Cloud (EC2)



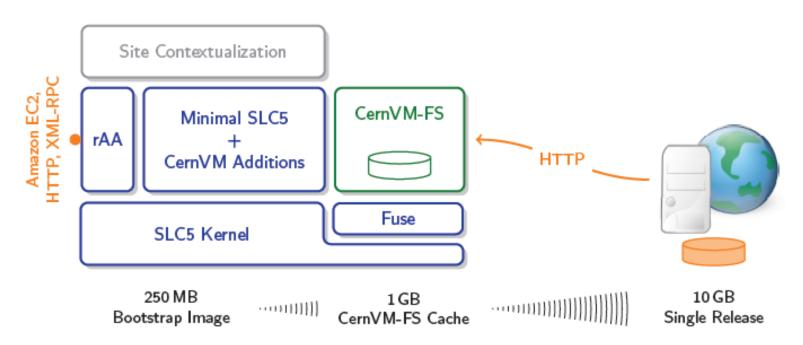
CernVM Model



- Minimal Operating System (common platform) sufficient to satisfy the most basic use cases of LHC experiments
- 2. File system based on HTTP protocol and optimized for software distribution using aggressive caching and capable of off-line operations
- Appliance Agent providing a simple Web UI for configuration and maintinance

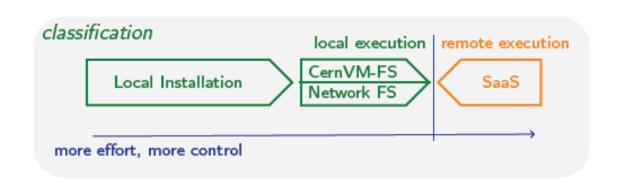
Software Distribution for the CernVM

Principle: Virtual software installation by means of an HTTP File System



Essential Properties:

- Defined platform
- 2 Read-only files
- 3 Public files





Summary



- Ongoing work in several areas
 - Will see benefits immediately (e.g. VOBoxes)
 - Will gain some experience in (e.g. LXCloud) test environment
 - Should be able to satisfy a wide range of resource requests – no longer limited to the model of a single job/cpu
- Broader scope of WLCG
 - Address issues of trust, VM management; integration with AA framework etc
 - Interoperability through cloud interfaces (in both directions) with other "grid" sites as well as public/commercial providers
 - Can be implemented in parallel with existing grid interfaces





Outlook

- Cloud technology can implement a "grid"
 - Industry-standard remote interfaces
 - Standard technology not (necessarily) supported by academia
 - Integration with commercial resource provisioning
 - Virtualisation can enable more efficient resource utilisation and deployment
- Features of a grid (that we still need):
 - Multiple administrative domains (no cloud solution)
 - Single sign-on/world-wide AAI (no cloud solution)
 - Concept of virtual collaborations (V0) (no cloud equiv)
 - Distributed resource providers (by definition)
 - Etc.



Middleware: Baseline Services

The *Basic* Baseline Services – from the TDR (2005)

- Storage Element
 - Castor, dCa
 SRM is too complex
 - Storm added in 2007
 - SRM 2.2 deployed in production Dec 2007
- Basic transfer t OK but why not HTTP?
 OK for some use cases
- File Tran OK, but must sync with storage
- LCG File No need for distributed catalogue
- LCG data mgt tools lcg-utils
- "Posix" I/O -
 - Grid File Access Library (GFAL)
- Synchroniced databases Tn←→T1s
 - 3D pro Frontier/Squid for many use cases

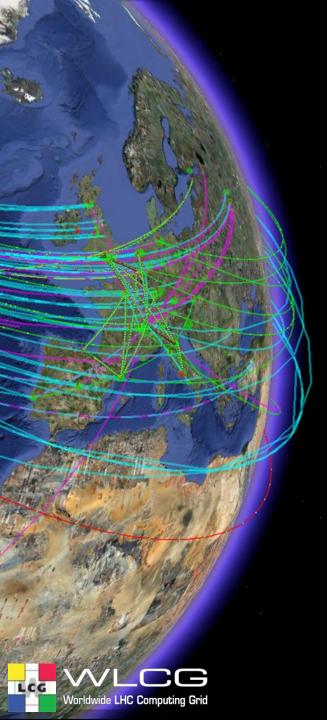
- Information
 LDAP → messaging?
 Static vs dynamic info
- Compute E Still have LCG-CE
 - Globus/ CREAM v slow coming
 - web ser MUPJs!
 - Support for multi-user pilot jobs

Actual LHC use cases much simpler Pilot frameworks will supercede it

- WIVIS, LB
- VO Management System (VOMS), MyProxy
- VO Boxes -> Virtual machine
 - Application > CVMFS or Squid
 - → MSG, Nagios, etc

APEL etc.





Conclusions

- Distributed computing for LHC is a reality and enables physics output in a very short time
- Experience with real data and real users suggests areas for improvement –
 - The infrastructure of WLCG can support evolution of the technology