# CERN Experience

D. Giordano / CERN-IT

WLCG Workshop

8-9 October 2016

# OpenStack Clouds at CERN



Total Cores in OpenStack Clouds at CERN

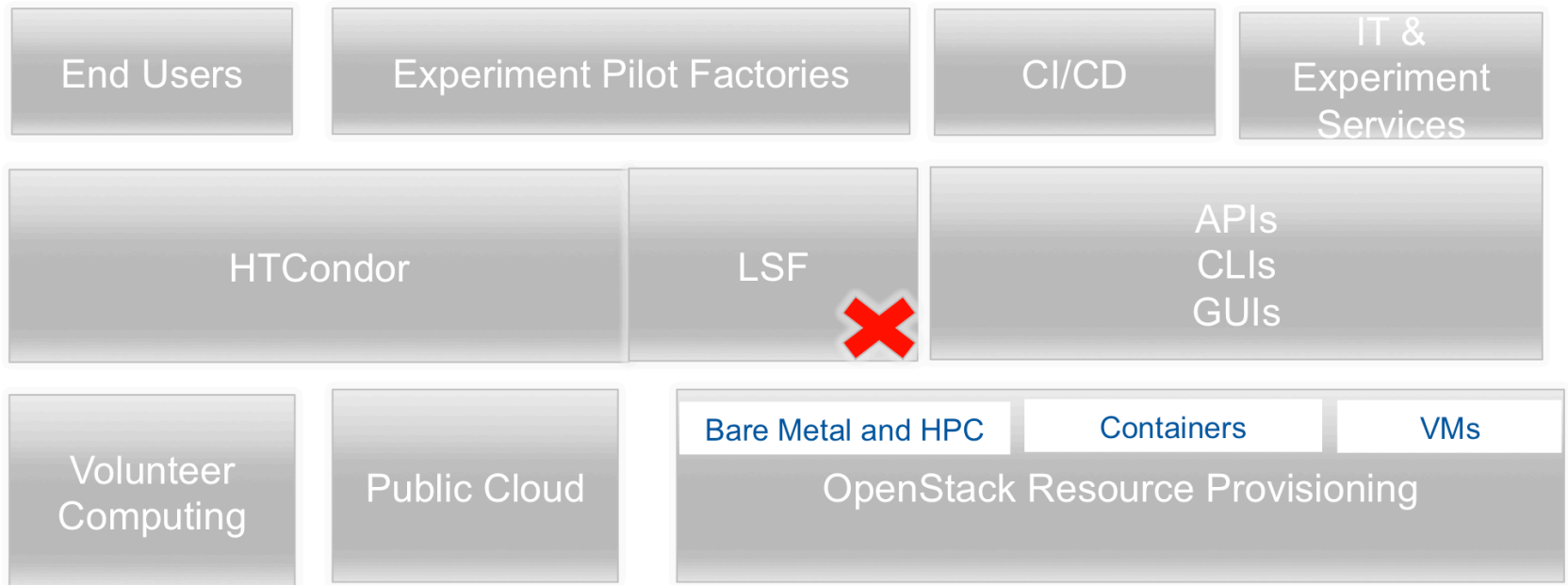Legend: IT, ATLAS HLT, CMS HLT, ALICE HLT

In production:

- ❑ 4 clouds
- ❑ >230K cores
- ❑ >8,000 hypervisors

90% of CERN's compute resources are now delivered on top of OpenStack

A further 42K cores to be installed in next few months

# Tier-0 Compute Services 2017

- Universal resource provisioning layer for bare metal, containers and VMs
- HTCondor as the single end user interface with LSF retirement by LS2
- Continue investing in automation and other communities for scaling with fixed staff
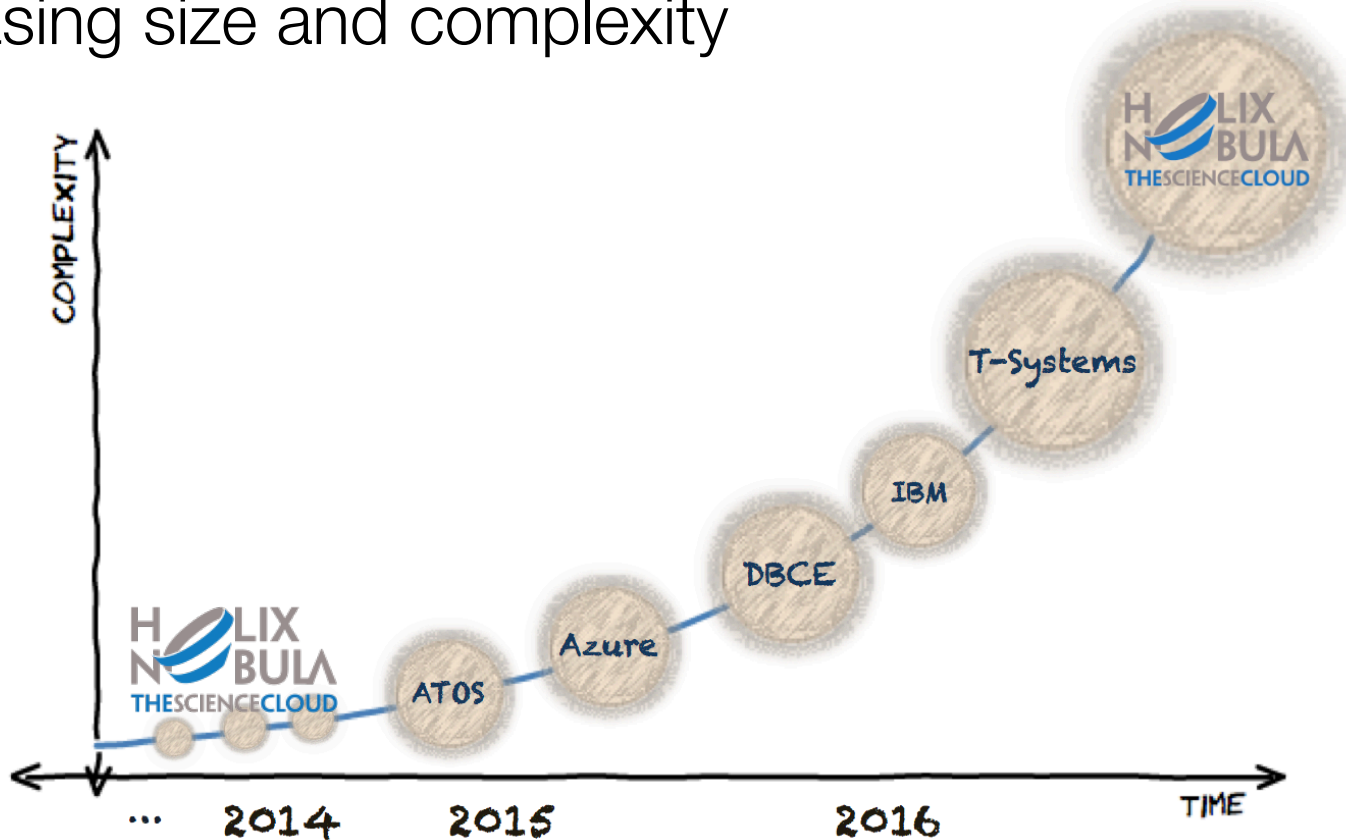- Self service for end users within the policies and allocations

| End Users | Experiment Pilot Factories | CI/CD | IT & Experiment Services |
|---|---|---|---|

| HTCondor | LSF ❌ | APIs CLIs GUIs |
|---|---|---|

| Volunteer Computing | Public Cloud | Bare Metal and HPC | Containers | VMs |
|---|---|---|---|---|
| | | OpenStack Resource Provisioning | | |

*courtesy of T. Bell*

# Evaluation of
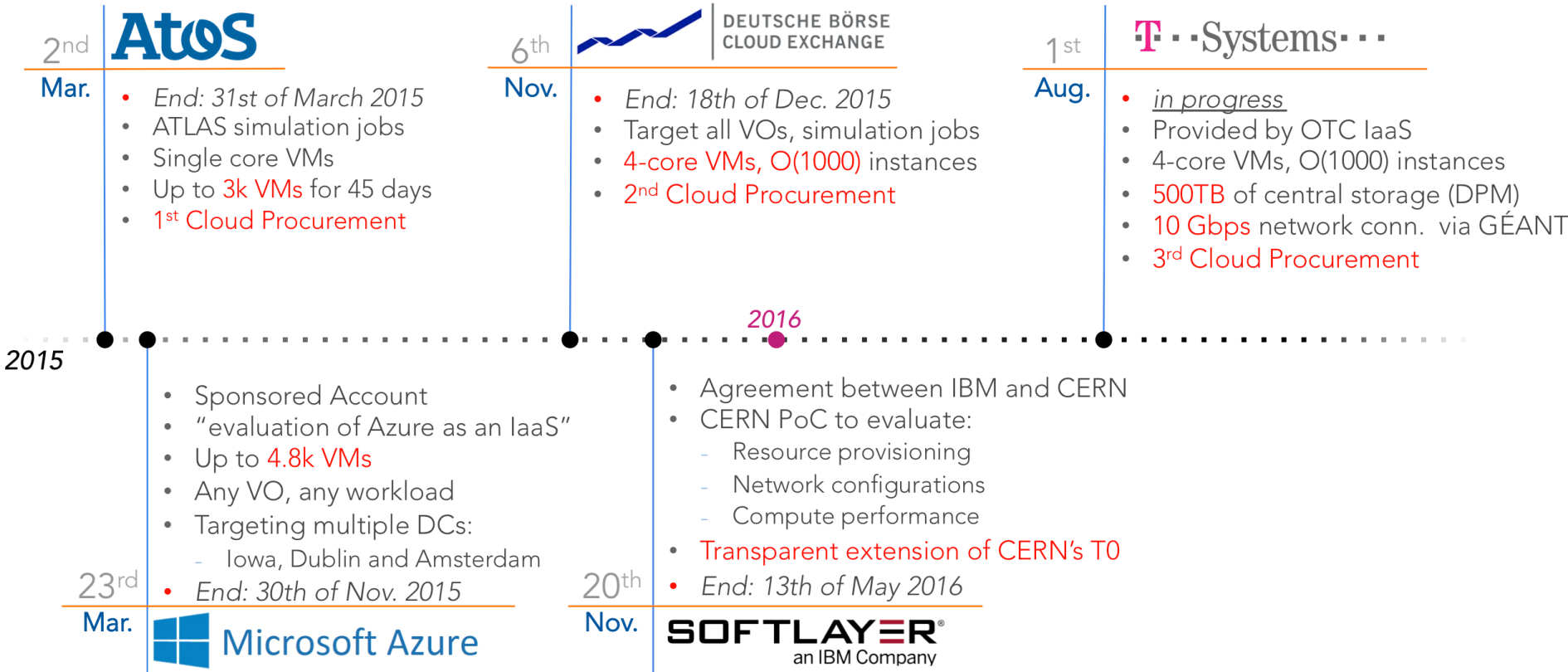# Public Cloud Service Providers (CSP)

# Roadmap

Started in 2011 with the EC funded project Helix-Nebula

Since 2015, series of short CERN <u>procurement projects</u> of increasing size and complexity

# In a nutshell

**ATOS**

2nd Mar.
- *End: 31st of March 2015*
- ATLAS simulation jobs
- Single core VMs
- Up to 3k VMs for 45 days
- 1st Cloud Procurement

**DEUTSCHE BÖRSE CLOUD EXCHANGE**

6th Nov.
- *End: 18th of Dec. 2015*
- Target all VOs, simulation jobs
- 4-core VMs, O(1000) instances
- 2nd Cloud Procurement

**T··Systems···**

1st Aug.
- *in progress*
- Provided by OTC IaaS
- 4-core VMs, O(1000) instances
- 500TB of central storage (DPM)
- 10 Gbps network conn.  via GÉANT
- 3rd Cloud Procurement

*2015*

*2016*

23rd Mar.
- Sponsored Account
- "evaluation of Azure as an IaaS"
- Up to 4.8k VMs
- Any VO, any workload
- Targeting multiple DCs:
  - Iowa, Dublin and Amsterdam
- *End: 30th of Nov. 2015*

**Microsoft Azure**

20th Nov.
- Agreement between IBM and CERN
- CERN PoC to evaluate:
  - Resource provisioning
  - Network configurations
  - Compute performance
- Transparent extension of CERN's T0
- *End: 13th of May 2016*

**SOFTLAYER** an IBM Company

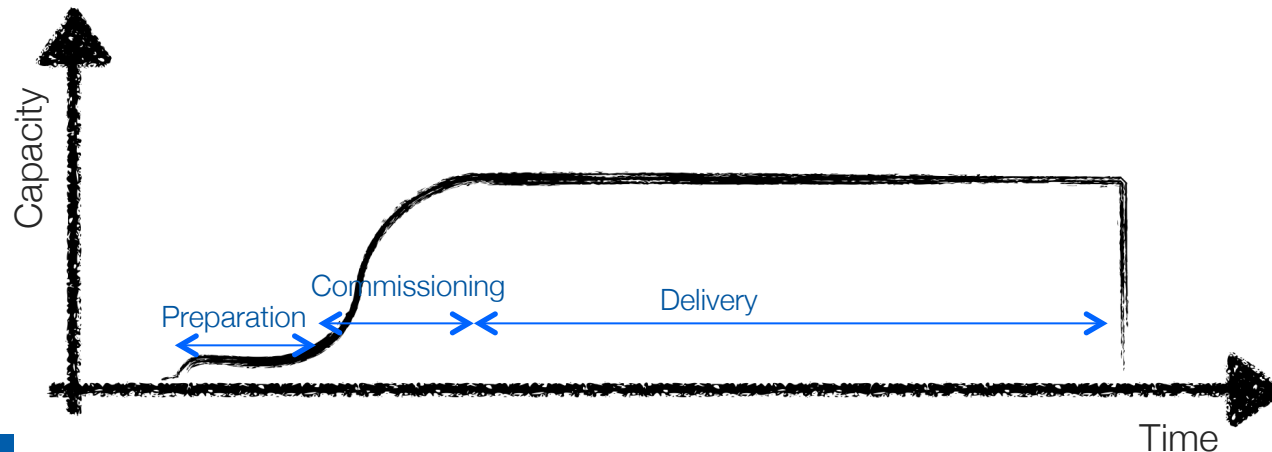*[see CHEP o-26]*

# Some Lessons learned

Cloud Brokerage not as effective as publicized

Overhead in managing several, small public CSP

Client-side accounting, monitoring and benchmarking are crucial

Cannot avoid preparation and commissioning phases
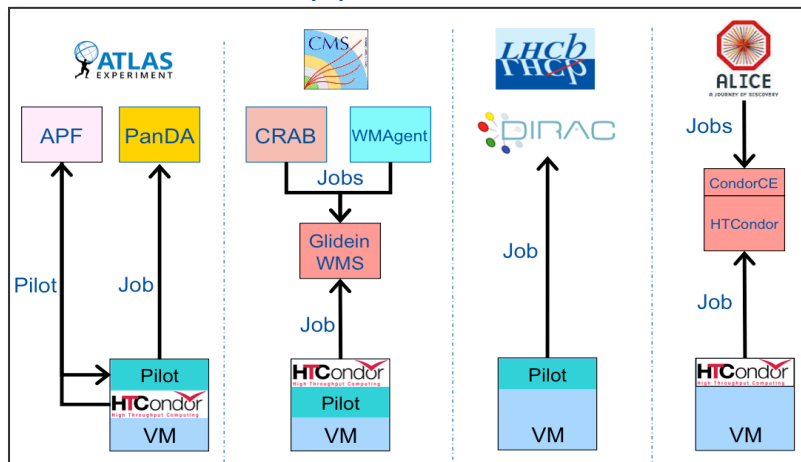- Ramp-up speed is often connected to the maturity and stability of the CSP

# Provisioning, Monitoring, Exploitation
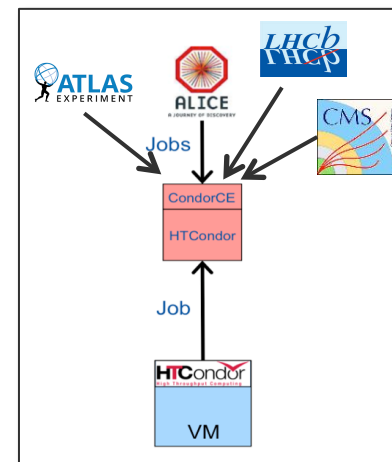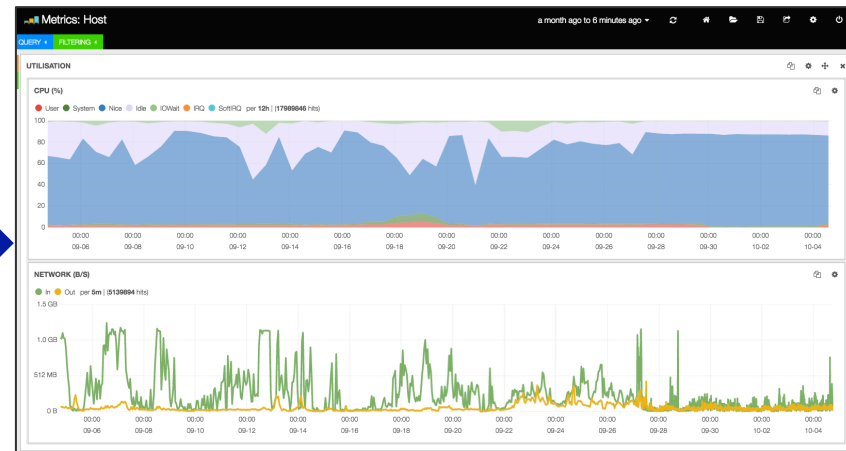
# Transparent Extension of CERN Resources

Consolidate the strategies adopted in the past cloud activities
– Manage and exploit external resources using same toolset and entry points as CERN on premises resources
  • *Puppet* configuration
  • **HTCondor** for scheduling and match-making
  • Infrastructure **monitoring**     *[see CHEP p-22]*
– Adopted *Terraform* for VM lifecycle management (N.B.: looking for long VM lifetime)
  • Open source toolkit, supports several cloud providers



2015 approaches

Now

Provisioning tools:
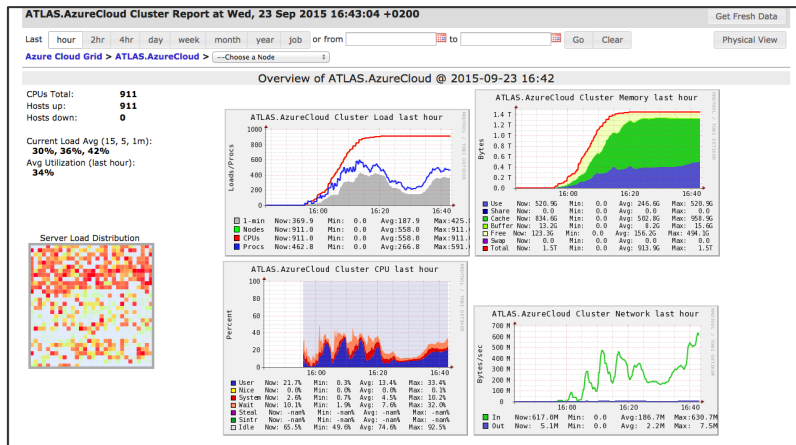*Vcycle; custom scripts for HNX-SlipStream, Azure RM*

*Terraform*

# Transparent Extension of CERN Resources

Consolidate the strategies adopted in the past cloud activities
- Manage and exploit external resources using same toolset and entry points as CERN on premises resources
  - *Puppet* configuration
  - **HTCondor** for scheduling and match-making
  - Infrastructure **monitoring**            *[see CHEP p-22]*
- Adopted *Terraform* for VM lifecycle management (N.B.: looking for long VM lifetime)
  - Open source toolkit, supports several cloud providers

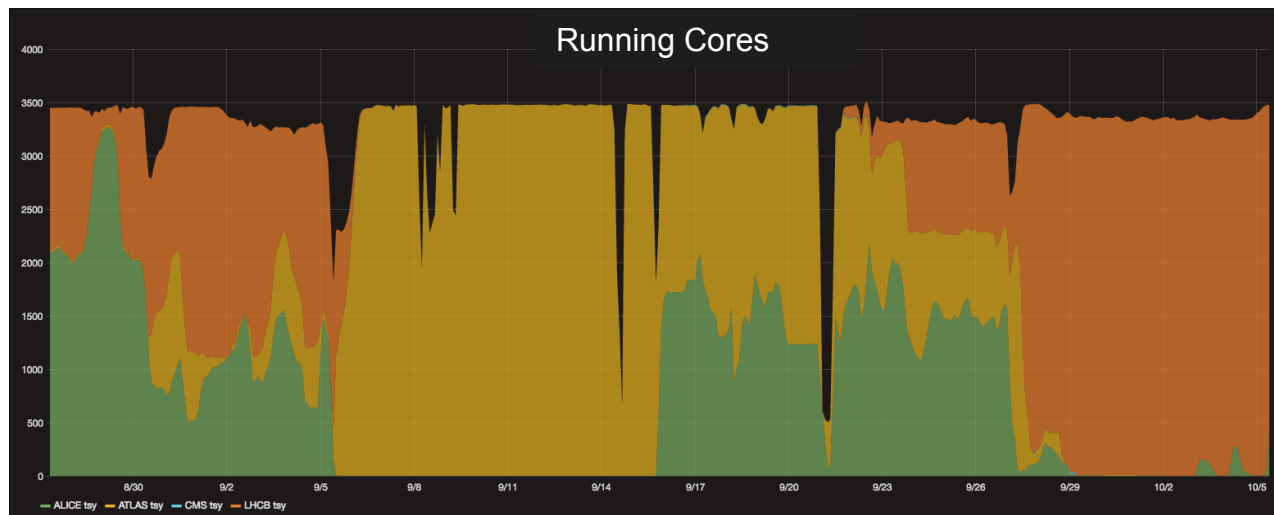## 2015 approaches                                          Now



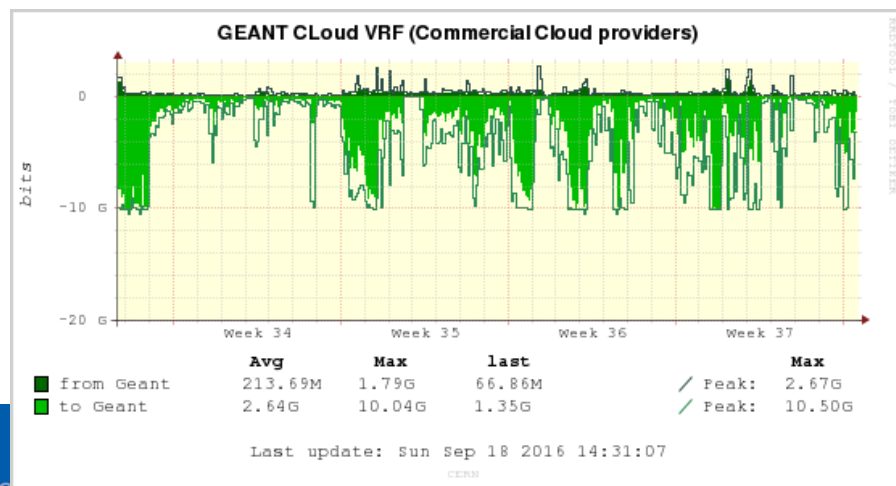Monitoring:
*Ganglia*

*AI Monitoring*

# Recent activity: T-Systems

- Batch resources
  fully loaded
  - shared among VOs



Running Cores

- Mixture of "CPU-intensive"
  and "network-intensive" tasks
  - MC workloads tend to
    dominate: easier to manage?

| | Max | Avg ▾ |
|---|---|---|
| LHCB tsy | 99.05 | 85.04 |
| ALICE tsy | 93.83 | 75.98 |
| ATLAS tsy | 100.00 | 64.13 |

- WAN largely used
  - Sometimes even saturated



GEANT CLoud VRF (Commercial Cloud providers)

| | Avg | Max | last | | Max |
|---|---|---|---|---|---|
| from Geant | 213.69M | 1.79G | 66.86M | Peak: | 2.67G |
| to Geant | 2.64G | 10.04G | 1.35G | Peak: | 10.50G |

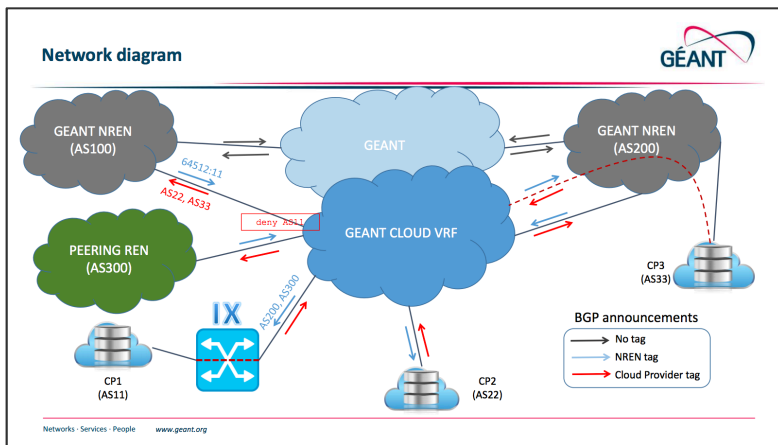Last update: Sun Sep 18 2016 14:31:07

# Network Connectivity
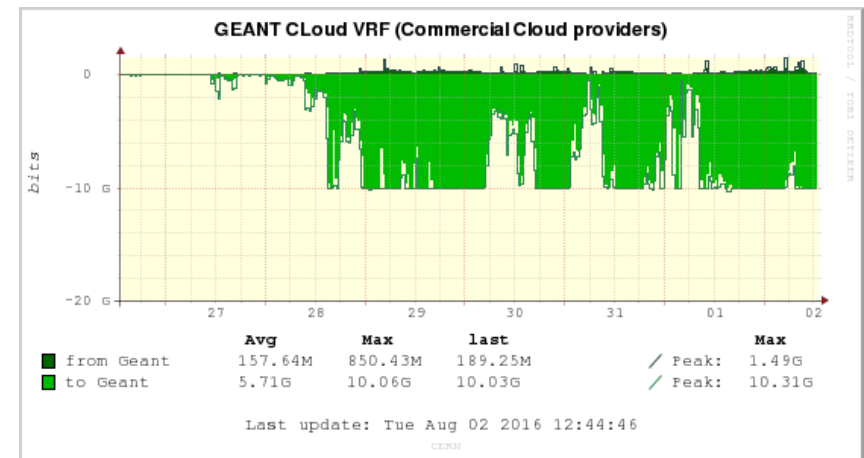
# WAN connectivity over GÉANT network

Requirement for CSP since the first procurement (early '15)

GÉANT Cloud VRF is currently connecting CERN and T-Systems (via DFN)

- 10 Gbps of total reserved peak bandwidth available
- The VRF is configured to only allow traffic between CSPs and NRENs; no CSP-CSP traffic is allowed



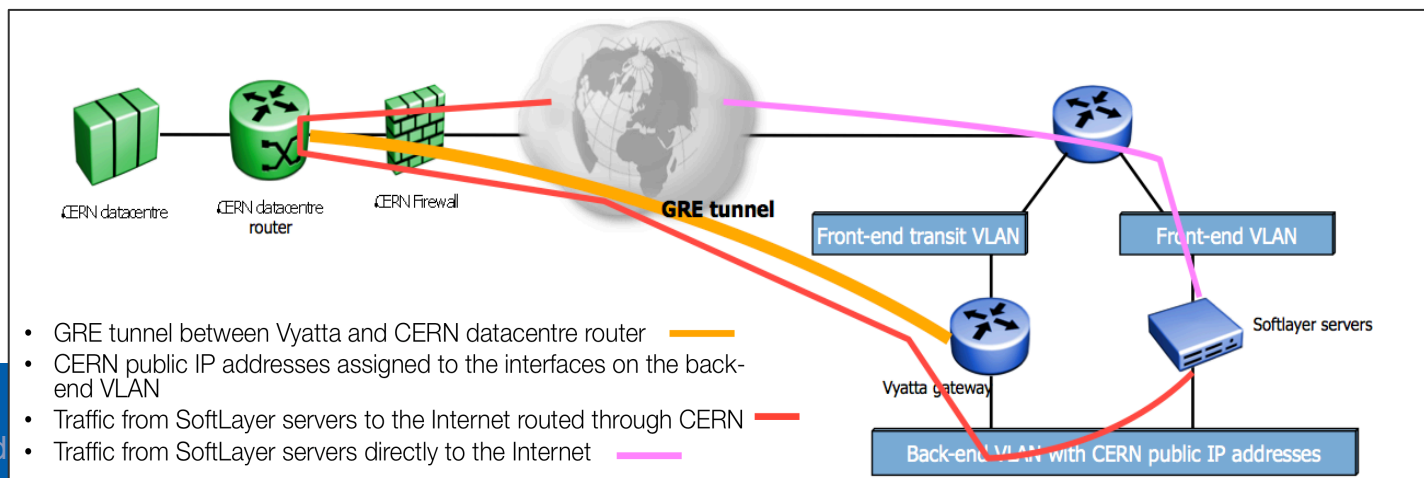*source V. Capone: GÉANT approach to Cloud R&E traffic*

# VPN Evaluation

Performed during the IBM-SoftLayer activity over GÉANT

- Established connectivity @ Amsterdam in co-located PoP
- Performance measured using **perfSONAR** (throughput, latency, loss, routing)

VPN evaluation outcome

- <u>Performance</u>: not necessary the best option for maximum throughput
- <u>Management</u>: overhead due to additional configuration
- <u>Security</u>: VPN implies full access to CERN from outside
  - On the contrary w/o VPN going through the CERN firewall. No access to full CERN from outside. Only ports open and to be secured are for Puppet CA.
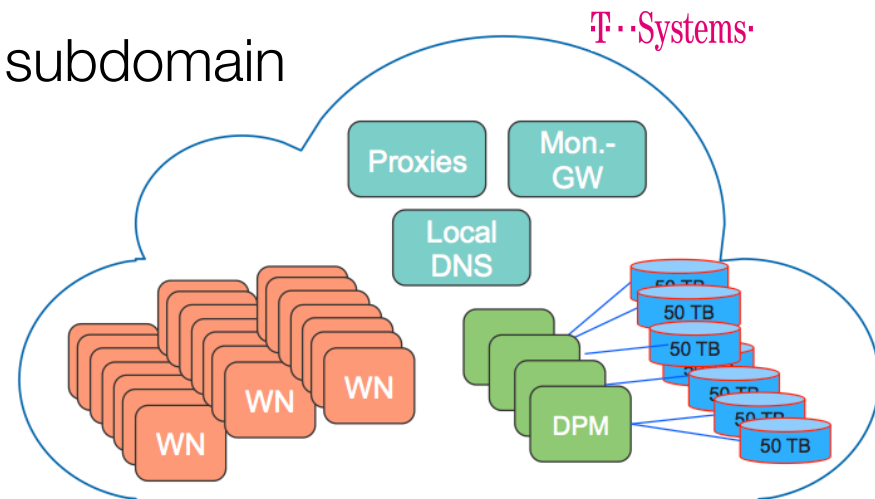


CERN datacentre    CERN datacentre router    CERN Firewall    **GRE tunnel**

Front-end transit VLAN    Front-end VLAN

Softlayer servers

- GRE tunnel between Vyatta and CERN datacentre router
- CERN public IP addresses assigned to the interfaces on the back-end VLAN
- Traffic from SoftLayer servers to the Internet routed through CERN
- Traffic from SoftLayer servers directly to the Internet

Vyatta gateway

Back-end VLAN with CERN public IP addresses

# Cloud Storage

# Cloud Storage in addition to Compute

Opportunity to study requirements and performance of "data-intensive" workloads, i.e. more than MC simulation

**500 TB** of Central Storage included in the last contract, awarded to T-Systems

- Block storage, managed via **DPM**
  (Tier-2 like model)
- Nodes registered in tsy.cern.ch subdomain
- Local DNS service needed for performance reasons with 1:1 NAT

# Storage Vs Network
## "Block" vs "Object" Storage

Two open questions

1. Is storage (mainly a cache) needed in public cloud?

   - Or is it better to buy more WAN bandwidth?

   - Different requirements from the VOs:

     – Disk-less resources (LHCb)

     – Cache for Minimum Bias (ATLAS and initially CMS the before pre-mixed PU approach)

     – Tier-2 like storage for Analysis (ALICE)

2. Status of adoption of Object Storage in WLCG workloads

   - CSPs prefer to leverage Object Storage

     – More scalable, reliable and cost effective than block storage
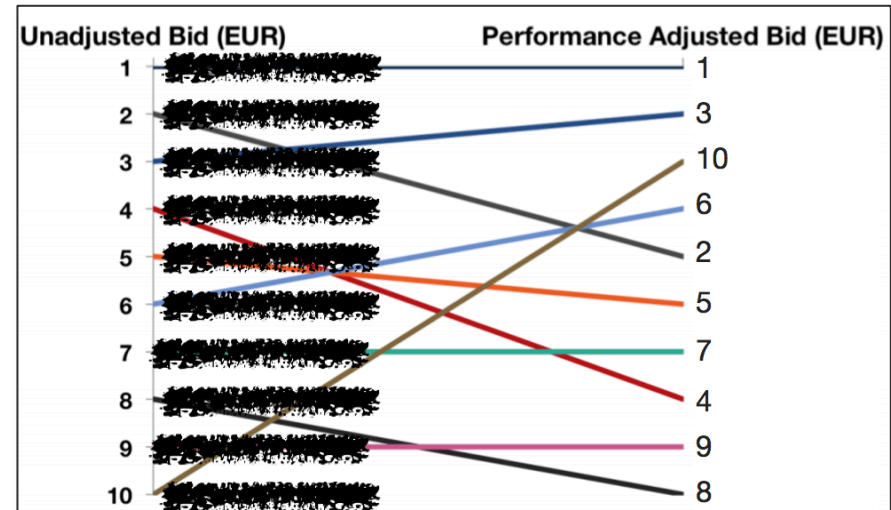
# Procurement Actions

# CERN Procurements

## Procurement rules apply

- No credit card to buy cloud services
- Request for Tender with upfront detailed Technical Specifications
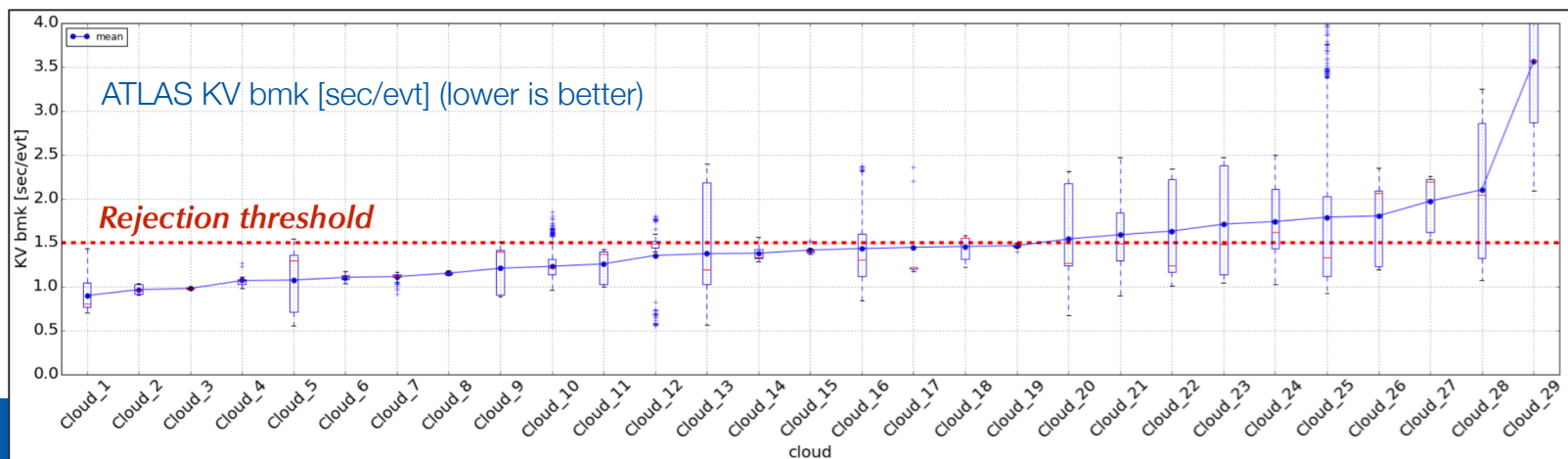- Award to the <u>lowest compliant bid</u>

## Compare quality & cost

- Requiring firms to execute benchmarks during the tender procedure
- Adopt CPU benchmarks representative of the experiments' workloads *[see CHEP p-28]*



Ranking changes significantly when price is adjusted for expected performance (1 is cheapest)
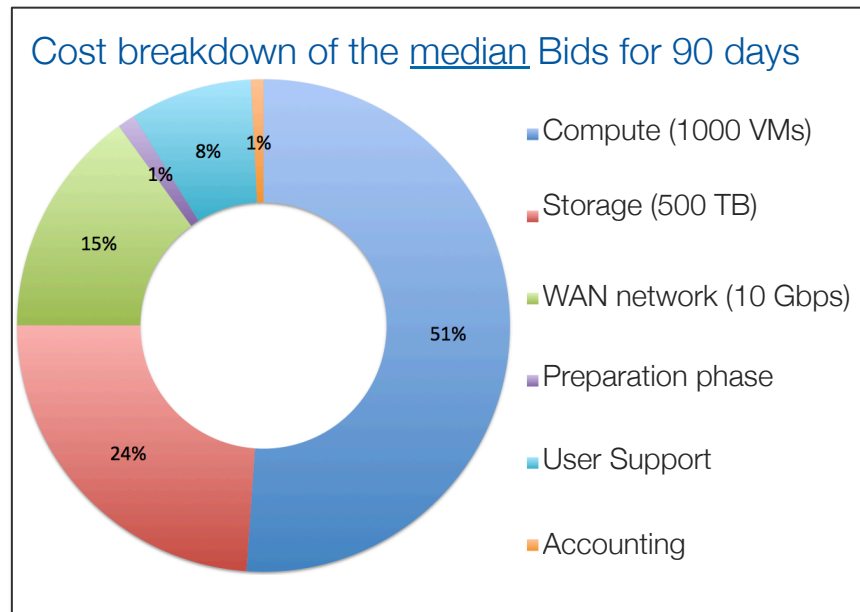


ATLAS KV bmk [sec/evt] (lower is better)

Rejection threshold

# Cost Evaluation

Main focus on the procurement process and features of cloud deployments: functionality, stability, maintainability, …
– Prefer duration (several weeks) to high peak capacity

Side effect: economy of scale not reached
– Shown by the large cost spread per CSP and per service (compute, storage, network, etc)

Cost breakdown of the <u>median</u> Bids for 90 days



- Compute (1000 VMs)
- Storage (500 TB)
- WAN network (10 Gbps)
- Preparation phase
- User Support
- Accounting

51%
24%
15%
1%
8%
1%

Median Price for compute
– 0.02 CHF/core/hour
– Including RAM (2 GB/core) and Disk (25 GB/core)

Storage (500 TB) represents 24% of the cost
– If not needed can invest more in network

# HNSciCloud Pre-Commercial Procurement

- 5.3 MEur procurement of R&D services   *[see CHEP o-397,o-399, p-401]*
- Build a Hybrid Cloud platform for the European research community
  - Combining services at IaaS to be integrated with in-house resources
- HNSciCloud challenges, triggering some of the next questions:

**Compute and Storage**
- support a range of virtual machine and container configurations working with datasets in the petabyte range

**Network Connectivity and Federated Identity Management**
- provide high-end network capacity for the whole platform with common identity and access management

**Service Payment Models**
- explore a range of purchasing options to determine the most appropriate ones for the scientific application workloads that will be deployed
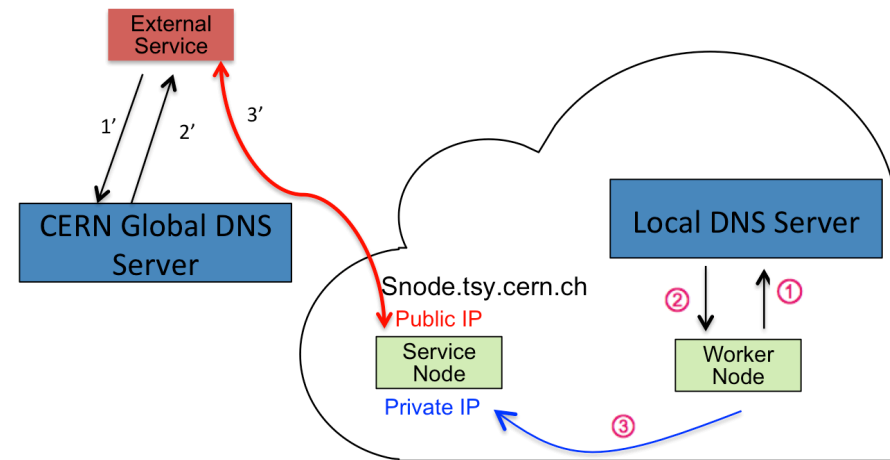
**www.cern.ch**

# DNS setup

## VMs are generated with private IP addresses
– Public Elastic IP bound to the VM via 1:1 NAT
– VMs assigned to CERN subdomain tsy.cern.ch

## NAT
– Different behavior for internal/external clients
– Requires cloud-local DNS



## DNS reverse lookups
– Grid clients need this for both public and private IPs
– T-Systems, as authoritative entity for reverse DNS, has implemented the Global DNS reverse function for the CERN-OTC service nodes