

# Project Results

Bob Jones  
CERN

25 February 2019

# Helix Nebula Science Cloud Joint Pre-Commercial Procurement

Procurers: CERN, CNRS, DESY, EMBL-EBI, ESRF, IFAE, INFN, KIT, STFC, SURFSara

Experts: Trust-IT & EGI.eu

The group of procurers committed

- Procurement funds
- Manpower for testing/evaluation
- Use-cases with applications & data
- In-house IT resources

Resulting made available to end-users from many research communities

Co-funded via H2020 Grant Agreement 687614

Duration: 3 years - Jan 2016 to Dec 2018

**Total procurement budget >5.3M€**



# Challenges

Procure R&D for innovative IaaS level cloud services integrated with procurers in-house resources and public e-infrastructure to support a range of scientific workloads

## *☞ Compute and Storage*

- ☞ support a range of virtual machine and container configurations including HPC working with datasets in the petabyte range

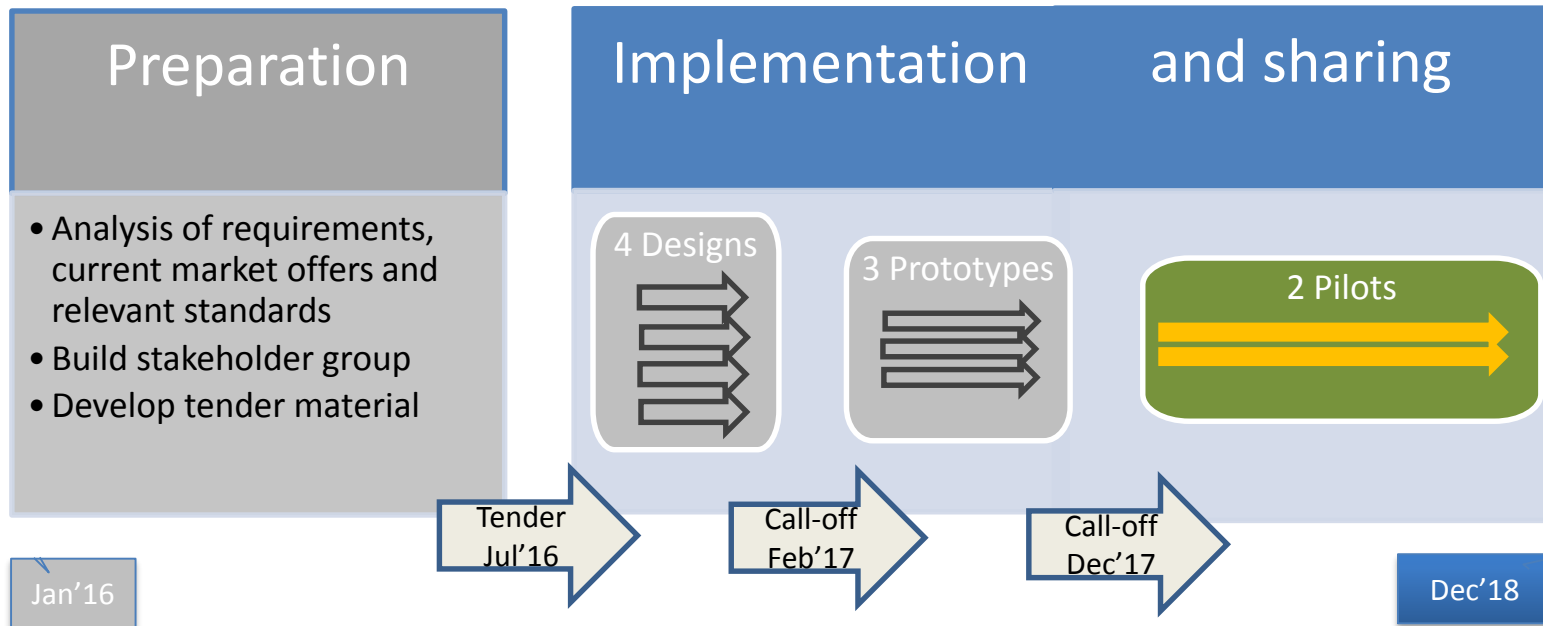
## *☞ Network Connectivity and Federated Identity Management*

- ☞ provide high-end network capacity via GEANT for the whole platform with common identity and access management

## *☞ Service Payment Models*

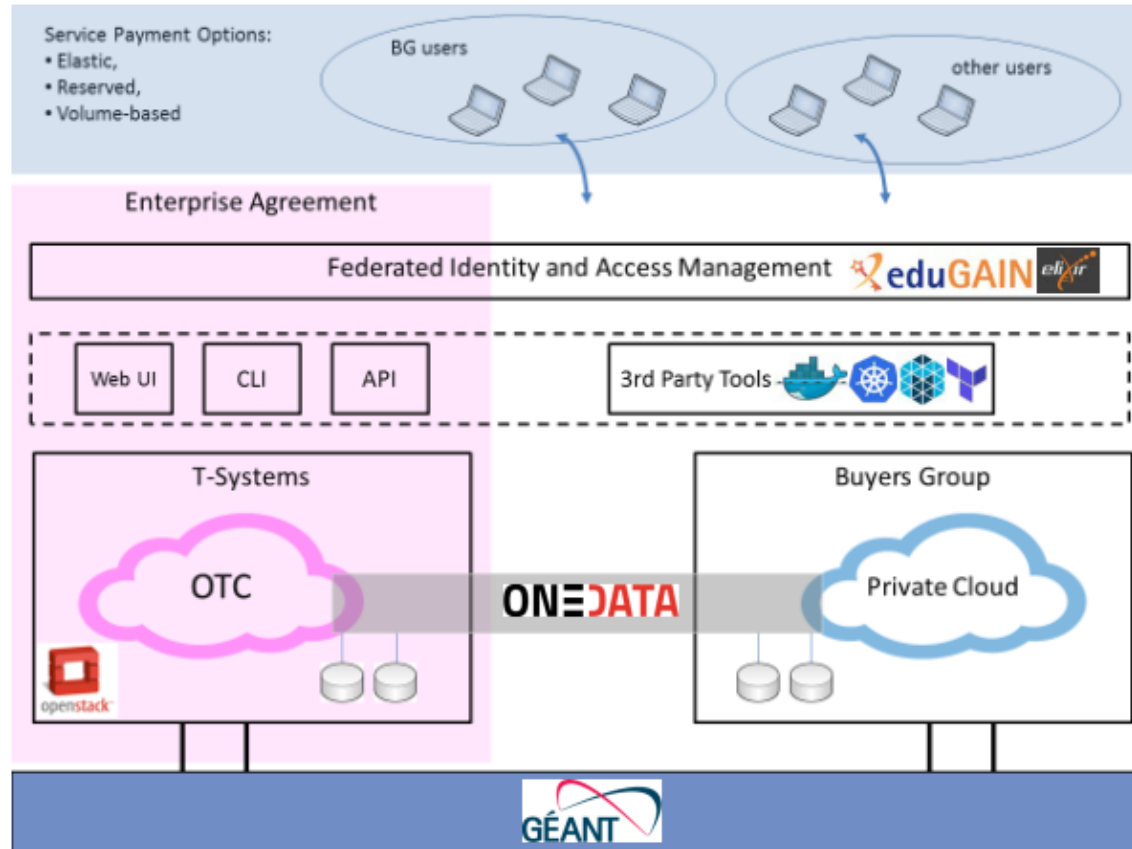
- ☞ explore a range of purchasing options to determine those most appropriate for the scientific application workloads to be deployed

# HNSciCloud project phases



Each step is **competitive** - only contractors that successfully complete the previous step can bid in the next

# HNSCICLOUD - T-SYSTEMS SOLUTION

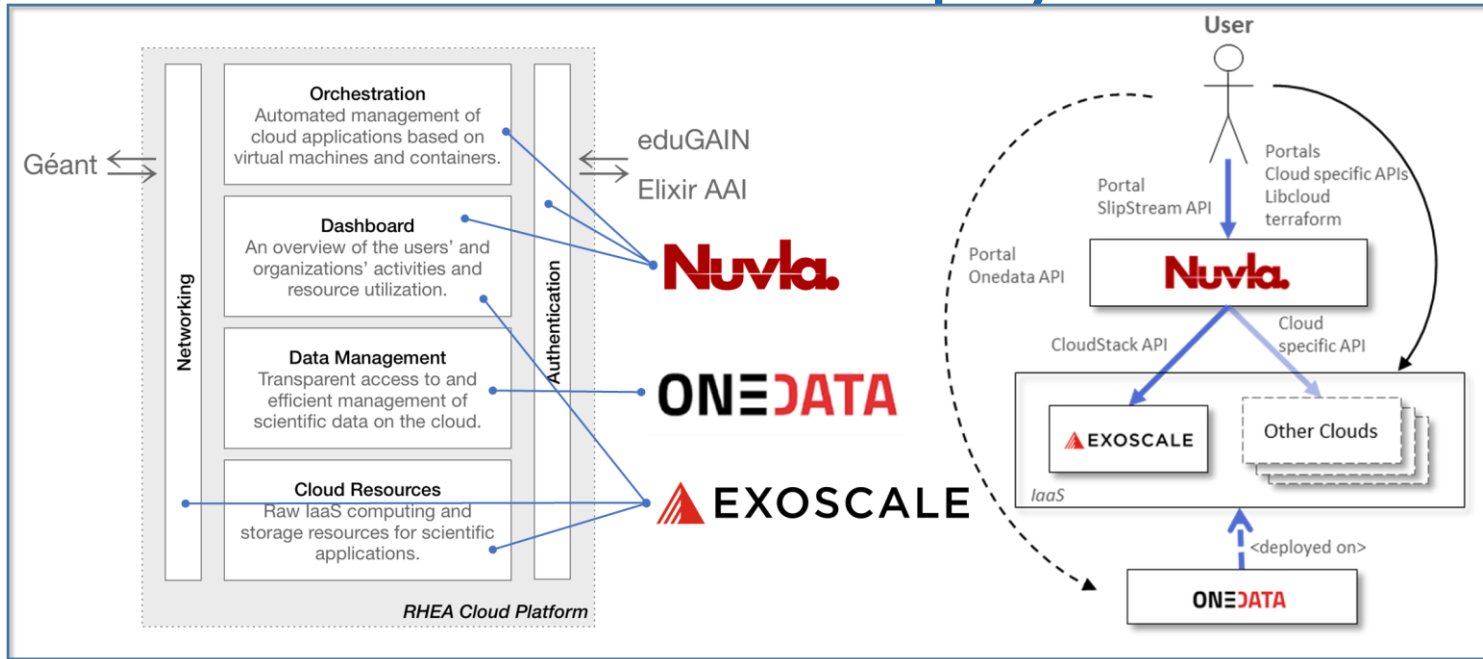


## HIGHLIGHTS

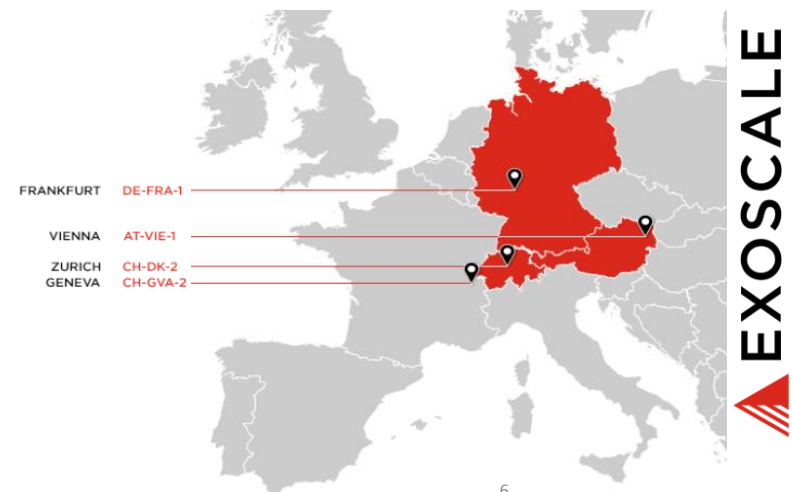
- HYBRID CLOUD**
- FEDERATED IDENTITY**
- CLOUD-NATIVE**
- LARGE-SCALE DATA**
- HPC SERVICE**
- GÉANT ACCESS**
- DASHBOARDS**
- ENTERPRISE SERVICE**
- OPENSTACK**

**T-Systems**

# HNSciCloud Pilot Deployments



- ✓ Four regions: **Geneva**, **Frankfurt**, Zurich, and Vienna
- ✓ Based on **Apache CloudStack**
  - ✓ Supports CloudStack API, **Libcloud**, **Terraform** and **Docker Machine**
- ✓ **Simple portal** and API to manage workloads
- ✓ Powerful **object storage** for cloud native applications
- ✓ Géant: aggregate **40 Gb/s bandwidth**



Test/Deployment completion	74.07%	Last Update:							
Test/Deployment on-going	26.86%	26/10/19							
TEST/Deployment NAME	PROCURER	STATUS	Expected Results	Summary of results obtained	Running Provider	Score if applicable	USE OF PROCURER'S DATA?	Additional Comments	
Batch Service Deployments	OERN	Completed	Use of 80% of the quota on both providers to deploy WLOG workloads	Functional deployment OK. Now scaling out.	Both	-	No		
TOTEM Deployment test		On-going	Fully functional data analysis software stack. Scaling out the installation and testing a complete analysis example.	Security challenge executed in both providers	Both	-	No	Test to be completed before December 19th 2019	
Security test challenge		Completed	Ag Latency (in/out): Reference value for OERs located in Europe 30ms Latency will be averaged across the 10 Buyers institutes: 0- Latency > 60 ms (Failed) 1- 60 < Latency < 50 (Failed) 2- 50 < Latency < 45 (Failed) 3- 45 < Latency < 40 (Passed) 4- 40 < Latency < 30 (Passed) 5- Latency < 30ms (Passed) Ag. Bandwidth (in/out): Reference value observed was ~ 5Gbps max. end to end. Bandwidth will be averaged across the 10 Buyer Group institutes: 0- Bandwidth < 100 Mbps 1- 100 Mbps < Bandwidth < 500 Mbps 2- 500 Mbps < Bandwidth < 1 Gbps 3- 1 Gbps < Bandwidth < 2 Gbps 4- 2 Gbps < Bandwidth < 5 Gbps 5- Bandwidth > 5Gbps	Network performance was satisfactory for both providers. T8systems ag latency across all buyers group was 25.9 ms (max for POC at 46.1), ag throughput across all buyers was 4.23 Gbps (min was for POC at just 100Mbps). Biscache ag latency was 24.99ms (max for POC at 31ms), ag throughput across all buyers was 1.38 Gbps (min was for POC at just 103 Mbps).	Both	Passed	No		
PERSONAR tests	Completed	Parallel training of Generative Adversarial Networks. Scaling experiments and comparison to HPC clusters. Biject (linear) scaling	T8systems: initial run on 2 VMs (NVIDIA P100). Scaling benchmark on 8 NVIDIA V100 (on a BMS). Exoscale: Scaling benchmark using container Kubernetes workload on 18 GPUs.	Both	-	No	Dedicated PERSONAR on public network was recently established by T8systems and results from it will be taken into account if available		
Deep Learning with GPUs	On-going	Deployment of a fully functional and containerized WLOG site (OERs/CE, Triage Batch System, Triage Worker nodes and later, HTC/Condor as well) using SIMPLE Grid Framework.			RHEA/Exoscale	-	No	Extend to larger number of VMs as soon as they become available. Test scheduling with SLURM.	
Deployment of Lightweight WLOG sites	CNR8	Completed		In Atlanta we had issue with the operabi. network configuration. We could not assign EP on our VM. The issue solved after intervention of the provider. Our model works well for medium number of VM (0-100) (with low IO network demand) and we can deploy with efficient manner on OERs (OTC and Biscache). The availability of Public IP facilitates a lot the deployments and the operation mode.	Both	-	No	Awaiting release of SIMPLE Framework. (tentative date: 1st week of December 2019)	
IaaS access via BluGain and local accounts		Completed		marked as completed but was not able to reach 500 IBS. Should be re-run as if problems with network performance are solved	Both	-	No	The Test consists of a set of subtests (AMI subtests and subtests for IaaS capabilities as are described in the initial document of the test cases) 0 Subtest: Failed 1 Subtest: Partial failed 2 Subtest: Partial Succeeded 3 Subtest: Almost Succeeded 4 Subtest: Full Succeeded 5 Subtest: Full Succeeded and recovers status by end of the expected requirements The final score will be the average of the average of subtests for AMI subtests and average of subtests for IaaS	
Onedata Wave 1: HDFS_IO		Completed				Both	-	No	
Onedata Wave 2: Scaling HDFS_IO	DEBY	Completed	100TB shared data (500 IBS average file size, min 70KB, max 4GB), more than 100 (parallel) instances of HDFS_IO on 20 nodes minimum. Single scientific app performance of 400MB/sec (10 mins peak, 100MB/sec total average). In case DEBY needs to base the performance metrics of the local (at OTC) caching infrastructure. It's expected to be less than 8% overhead in terms of IO through OneClient when compared to the direct access to the storage system.	Currently trying to run tests with multiple files. Experience problems with importing data to DEBY to OneData - some files are only partially imported. Need to investigate further. Update 06/09/2019 - Release 10 - problem is still there. Access to OneProvider on our local instance is provided to Cylontel (although RHEA experts seeing similar problem on their environment as well). T8system experts do not perform scaling tests with a significant amount of files.	Both	-	No		
		Completed	In case of the final configuration (reference scenario II) DEBY expects the following for data not on cache (only one local DEBY copy exists, application runs at OTC): max file (deleting delay) (file could be opened in OneData - known to space - after local creation): 100 seconds same max delay is expected for the way back (to DEBY storage system) Metadata scan rate (from DEBY storage into OneData Space and the reverse direction): 50 MB File (initial) access rate: not less than 85% of the slowest network segment (subtest)	On Exoscale - only B3 is really scales, but it does not allow to reach desired bandwidth (just several IBS per application) Works, still need to check updated LUMI A version and replication scripts Update 26/11/2019 - works as expected, score 5 for OTC, 3 for Biscache	Both	-	No		
SLURM integration	Completed				Both	-	No	0- Cannot execute IO with Docker in the cloud	

Buyers Group assembled a suite of approx. 30 tests and executed them against commercial cloud provider solutions

# Use cases deployed



Buyers



	Tested features			Fields of Research					
	AAI Federated Identity	Storage blocks/objects	Computing capacity (VMs)	Photon / Neutron Science	High Energy Physic	Astronomy	Life sciences		
				FDMNES	CrystFEL	DODAS	LOFAR	PanCancer	WeNMR/HADDOCK
CERN	●●●	●●	●●●●						
CFP	●●●	●●●	●●			●●●			
DESY	●●●	●●	●		●				
EMBL	●●●	●●	●●					●	
ESRF	●●●	●	●	●					
IFAE	●●●	●	●						
INFN	●	●	●●●			●●●			
KIT	●●	●	●						
SURF SARA	●●●	●●	●				●		●

● useful, they can ease the way I do my job   ●● relevant, they can really improve the way I do my job   ●●● strategic, they are fundamental to execute my job





# HEP Deployments



- WLCG
- ALICE, ATLAS, CMS and LHCb
- Daniele Spiga (INFN) talk about DODAS, Track 7:
  - <https://indico.cern.ch/event/587955/contributions/2937198/>
- Matthias Schnepf (KIT) talk about Dynamic Integration of resources, Track 8:
  - <https://indico.cern.ch/event/587955/contributions/2937900/>



- CERN Batch Service
- Container Federation
- Openstack Summit, reference talk:
  - <https://www.openstack.org/summit/vancouver-2018/summit-schedule/events/20768/cern-experiences-with-multi-cloud-federated-kubernetes>



- Belle II
- Silvio Pardi (INFN) poster about experience of Belle II with commercial clouds, Track 7:
  - <https://indico.cern.ch/event/587955/contributions/2937060/>



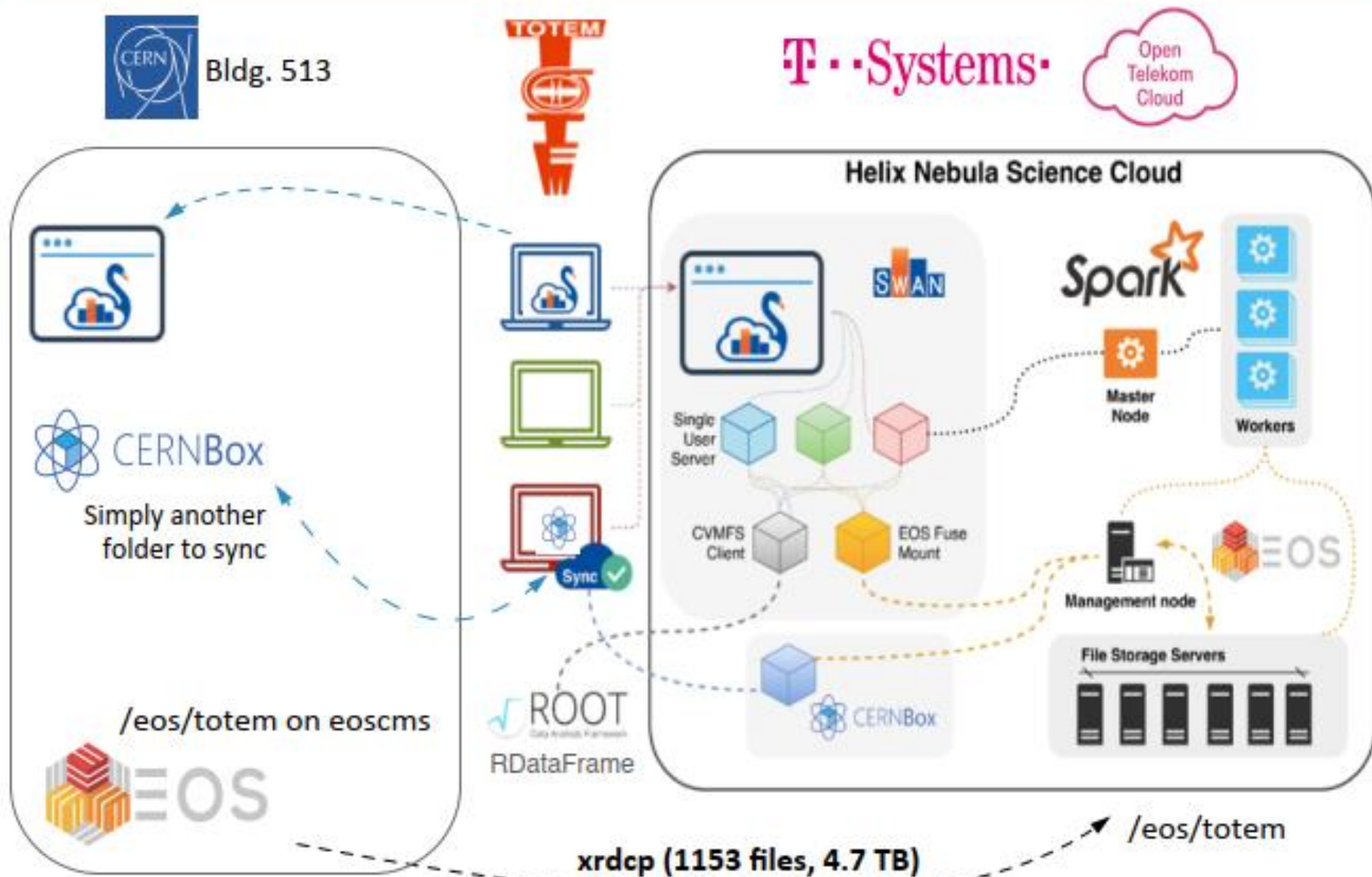
- Interactive Analysis for End Users for TOTEM
  - [https://indico.cern.ch/event/727193/contributions/3039091/attachments/1667076/2674030/TotemTest\\_HNSciCloud.pdf](https://indico.cern.ch/event/727193/contributions/3039091/attachments/1667076/2674030/TotemTest_HNSciCloud.pdf)



- Machine Learning/Deep Learning for Fast Detector Simulation using GPUs
- Sofia Vallecorsa (CERN openlab), Track 2 and Jean-Roch Vlimant (CMS), Track 6:
  - <https://indico.cern.ch/event/587955/contributions/2937595/>
  - <https://indico.cern.ch/event/587955/contributions/2937513/>



# Deployment (June-December 2018)



# Total Cost of Ownership study

- Understand the costs of using commercial cloud services as part of a hybrid cloud model
- 2 use-cases selected with different requirements
  - ALICE**: single core jobs, up to 50.000 at any time (monte-carlo, reconstruction, analysis)
  - PANCANCER**: burst pattern, with minimal resources constantly used (few VMs) and periods of ramp-up (up to 400 VMs)

# TCO results: T-Systems

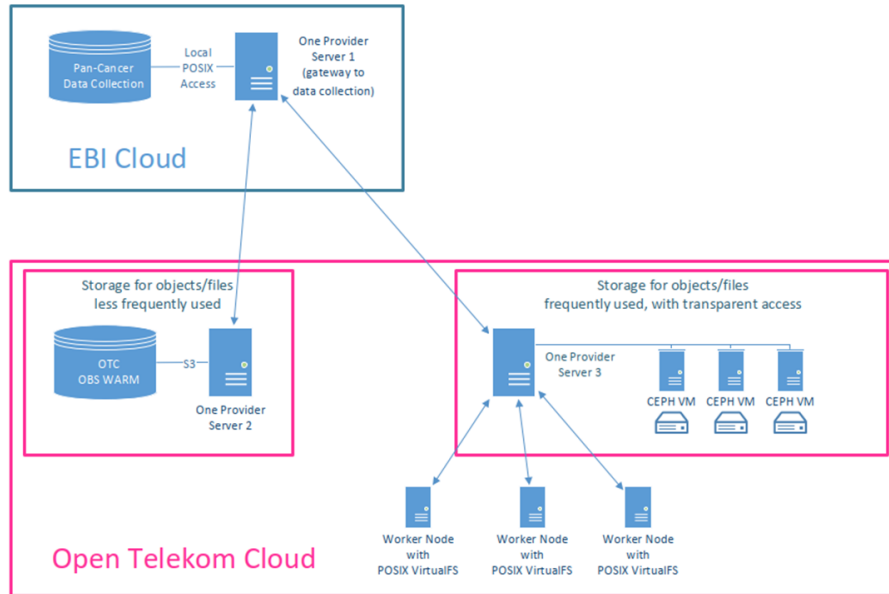


Figure 1: Pan-Cancer Data Management scenario

“The impact of flavour and VM size on TCO can be significant.”

Table 5: ALICE job costs compared between CERN cloud and OTC public cloud

Job Type	TCO Benefit OTC Public Cloud
Monte Carlo	
Reconstruction	
Analysis	

Table 4: Pan-Cancer qualitative factors

Factor	Importance (H, M, L)	Effectiveness	Notes	Additional Information
Agility	H	+1	The solution provides easy and quick deployment and changes to the processing and storage functionality	
Contract Review and Negotiation	M	+1	Can be implemented with minimum review	
Elasticity and Scalability	H	+1	Solution facilitates easy and quick expansion of available processing and/or storage capacity	
Regulatory and Policy Requirements	H	+1	Solution adequately enables compliance with external regulations	User data can be stored encrypted in the public cloud. Through Onedata, data privacy can be maintained between public, restricted and confidential data. Confidential data would only be hosted on-premise.
Security	H	+1	Solution provides effective mechanisms by which constantly escalating security threats are prevented and security events or breaches are constantly monitored	Cloud service is certified e.g. for relevant ISO standards, EU GDPR and German national regulations.

Factor	Importance (H, M, L)	Effectiveness	Notes	Additional Information
Service Levels	M	+1	Service level targets were not provided. The solution would be able to match or improve the current SLA based on on-premise services. Especially with regards to performance, the cloud offers more granularity and diversity to tune performance to requirements when compared to on-premise resources.	Service Availability of 99,95% is achievable when using 2 availability zones and load balancing. Support is provided 24/7.

# TCO results: RHEA System

Table 15: TCO Cost Summary for Each Use Case

Use-Case	Exoscale	AWS	CloudFerro
PANCANCER (no storage)	123,170	147,575	155,247
PANCANCER (with storage)	239,562	-	-
ALICE (all three use-cases) – 50000 jobs*	6,216,252	7,907,449	9,282,265
ALICE Monte Carlo – 1 job**	0.082	0.106	0.125
ALICE Raw Data Reconstruction – 1 job**	0.096	0.138	0.159
ALICE Analysis Trains – 1 job**	0.038	0.034	0.042

“The PANCANCER and ALICE use-cases each have their own specific requirements for cloud deployment. Both use-cases have more than one job type requiring different VM flavours and/or pricing options. Also, use-cases can be supported without having to storage large volumes of data in the cloud, minimising storage costs. Therefore, we have not included the cost data management solutions, since our assumptions for use-cases are that the data can be ‘streamed’ to the VMs. “

# Lessons Learned

- Framework agreements provide a convenient structure for service procurements in the scientific community*
- Volume and requirement aggregation across a group of scientific organisations brings advantages*
- Use a small number of cloud providers in parallel to avoid lock-in and ensure service continuity*
- Need to repatriate data at the end of contracts*
- Commercial clouds offer opportunities to rapidly scale cutting edge technology for R&D deployments*
- Vouchers/Credits are a practical means to provide limited-scale access to commercial cloud services for end-users*
- Commercial clouds providers offer services that are certified against international standards (e.g. ISO 27000 family) and consistent with legislation (e.g. GDPR)*

- ***Procurement of digital services for the European Open Science Cloud***
- ***Procurement Budget: 9.5M€***
- ***Starting Date: 1<sup>st</sup> January 2019***
- ***Duration: 36 Months***
- ***Coordinating Partner: GÉANT***



**OCRE will become the procurement vehicle of the European Open Science Cloud (EOSC)**

**Procuring commodity services from multiple cloud providers across IaaS/PaaS/SaaS**

**Suppliers readiness will be technically validated with a test-suite**

**Procurement Framework: Oct 2019, Call-offs: Q1 2020 and Q1 2021**

**OCRE will provide opportunities for research communities to adopt commercial cloud services as part of their IT computing strategy**



- **Pre-Commercial Procurement**
- **Focus: Archiving and Data Preservation Services in commercial clouds**
- **Procurement Budget: 3.4M€**
- **Starting Date: 1st January 2019**
- **Duration: 36 Months**
- **Coordinating Partner: CERN**



EMBL-EBI



addestino  
innovation delivered.



**ARCHIVER Pre-Commercial Procurement will run an open tender process for R&D in the next generation of digital archiving and Long-Term Data Preservation services**

**Experience gathered from HNSciCloud about the PCP process will be applied in ARCHIVER**

**Open Market Consultation Kick-off event: 8th of April 2019:**  
<https://goo.gl/rhWTmH>