

# Networks in ATLAS

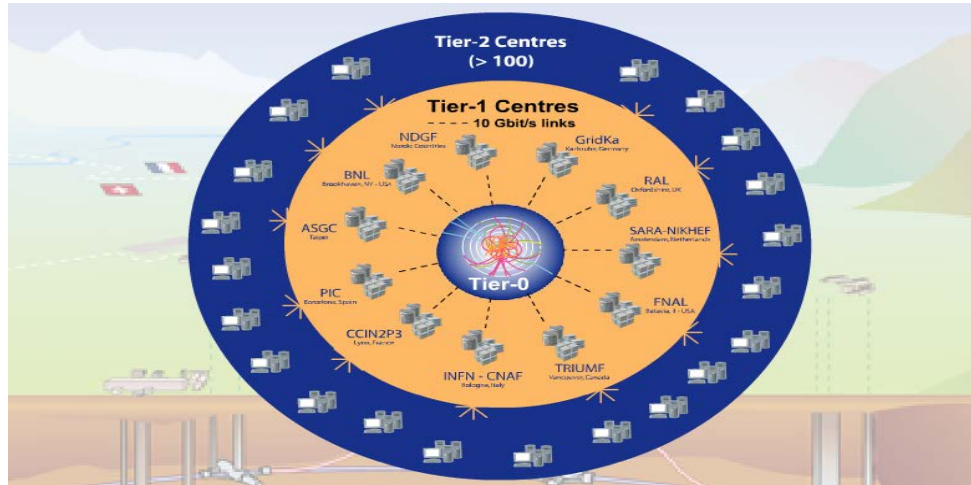
**Shawn McKee / University of Michigan**

***for the ATLAS Collaboration***

***GG C2 Track 3***

**CHEP 2016, San Francisco, October 11, 2016**

# Distributed Computing in ATLAS



ATLAS Computing Model :  
11 Clouds : 10 T1s + 1 T0 (CERN)

Cloud = T1 + T2s + T2Ds

T2D = multi-cloud T2 sites

2-16 T2s in each Cloud

Workload Management System



Task → Cloud : Task brokerage

Jobs → Sites : Job brokerage

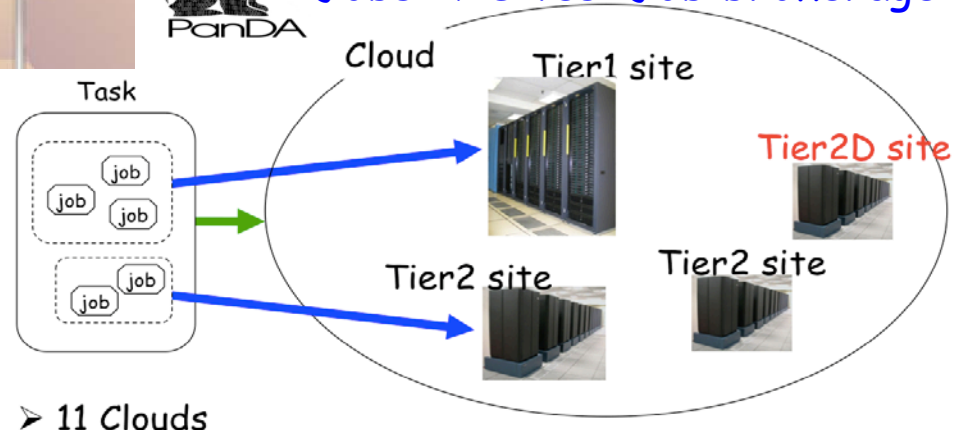
Basic unit of work is a job:

- Executed on a CPU resource/slot; may have inputs; produces outputs

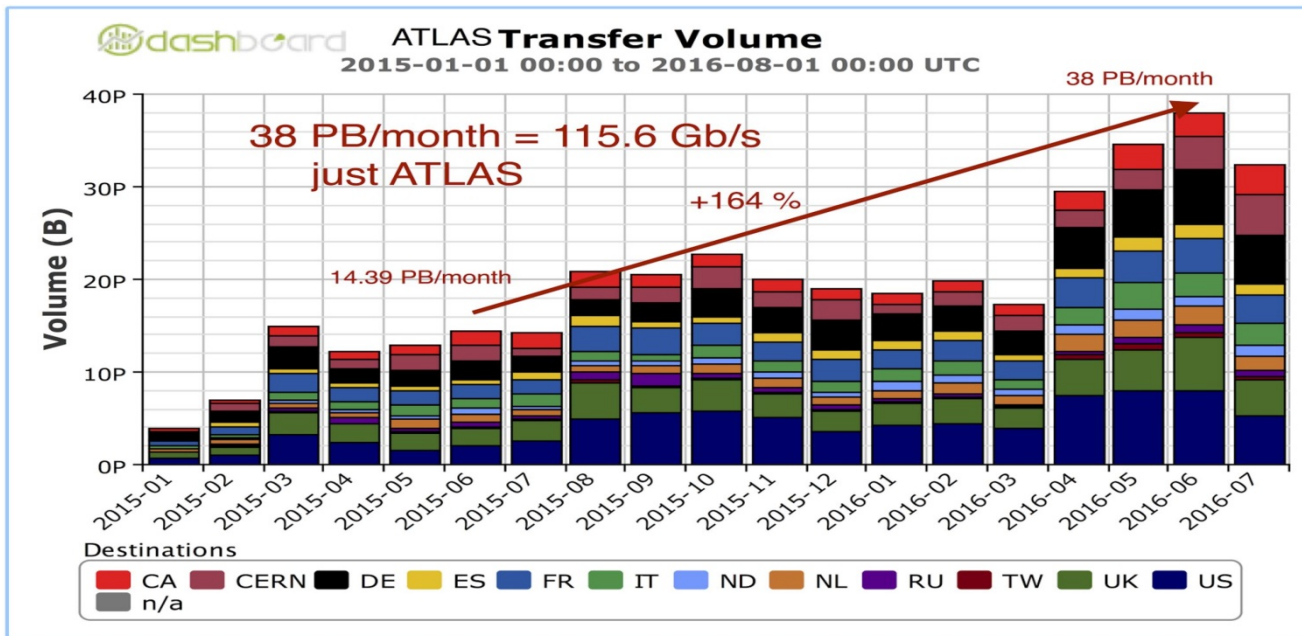
JEDI – layer above Panda creates jobs from ATLAS physics and analysis 'tasks'

Current scale – one million jobs per day

**The network ties this all together!**



# Network Use in ATLAS



ATLAS (and LHC in general) has been transferring an exponentially increasing amount of data since startup. This trend is likely to continue and is driven by increasing data volumes, more capable infrastructures and the excellent networks supporting our needs.

# Working on Networks for ATLAS



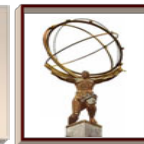
- There has been a small but long-term effort in the area of networks for LHC HEP (High Energy Physics)
  - Initial efforts started around the Internet2 HENP working group in 2001
  - The LHCOPN (and follow-on LHCONE) effort started in 2005 and focused on defining the LHC experiment networking needs and implementing them. It continues to meet twice per year.
  - USATLAS piloted perfSONAR in 2006, expanding to LHCONE in 2010 and WLCG wide in 2012
  - Open Science Grid(OSG) began a network focus area in 2012
    - OSG now provides a network service for WLCG/OSG, gathering perfSONAR metrics worldwide and making the available
  - WLCG has had a task-force and a working group in networks
    - perfSONAR deployment task-force which got ~250 perfSONAR toolkit innstances deployed globally in 2013-2014
    - WLCG Network and Transfer working group which organizes and maintains network and transfer data from perfSONAR and transfer data sources from 2015 to the present
  - A set of ATLAS collaborators working on network analytics

# Importance of Measuring Our Networks



- **End-to-end network issues are difficult to spot and localize**
  - Network problems are multi-domain, complicating the process
  - Standardizing on specific tools and methods allows groups to focus resources more effectively and better self-support
  - Performance issues involving the network are complicated by the number of components involved end-to-end.
- **Network problems can severely impact ATLAS's workflows** and have taken weeks, months and even years to get addressed!
- **perfSONAR provides a number of standard metrics we can use**
- **Latency measurements provide one-way delays and packet loss metrics**
  - Packet loss is almost always very bad for performance
- **Bandwidth tests measure achievable throughput and track TCP retries (using Iperf3)**
  - Provides a baseline to watch for changes; identify bottlenecks
- **Traceroute/Tracepath track network topology**
  - Measurements are only useful when we know the exact path they are taking through the network.
  - Tracepath additionally measures MTU but is frequently blocked

# Current perfSONAR Deployment



[http://grid-monitoring.cern.ch/perfsonar\\_report.txt](http://grid-monitoring.cern.ch/perfsonar_report.txt) for stats



**249 Active** perfSONAR instances

**199** Running latest version (3.5)

**95 sonars in latency mesh**

- 8930 links measured at 10Hz
- packet-loss, one-way latency, jitter, ttl, packet-reordering

**115 sonars in traceroutes mesh**

- 13110 links
- hourly traceroutes, path-mtu

**102 sonars in bandwidth mesh**

- 10920 links (iperf3)

[https://www.google.com/fusiontables/DataSource?docid=1QT4r17H](https://www.google.com/fusiontables/DataSource?docid=1QT4r17HEufkvnqhJu24nIptZ66XauYEIBWWWh5Kpa#map:id=3)

[EufkvnqhJu24nIptZ66XauYEIBWWWh5Kpa#map:id=3](https://www.google.com/fusiontables/DataSource?docid=1QT4r17HEufkvnqhJu24nIptZ66XauYEIBWWWh5Kpa#map:id=3)

- Initial deployment coordinated by WLCG perfSONAR TF
- Network commissioning by WLCG Network and Transfer Metrics WG

# Latency and packet loss matters



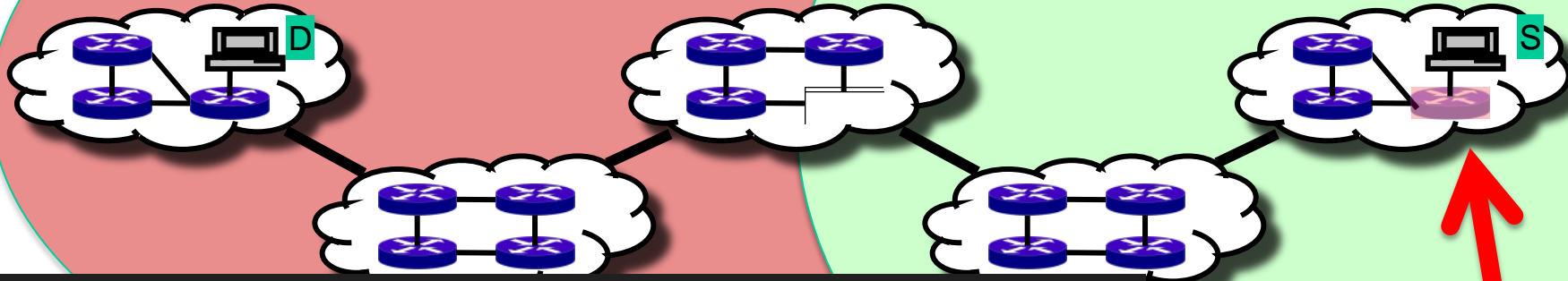
Performance is poor when  
RTT exceeds  $\sim 10$  ms

Performance is good when  
RTT is  $< \sim 10$  ms

Destination  
Campus

R&E  
Backbone

Source  
Campus



0.0046% loss (1 out of 22k packets) on 10G link

- with 1ms RTT: 7.3 Gbps
- with 51ms RTT: 122Mbps
- with 88ms RTT: 60 Mbps (factor 80)

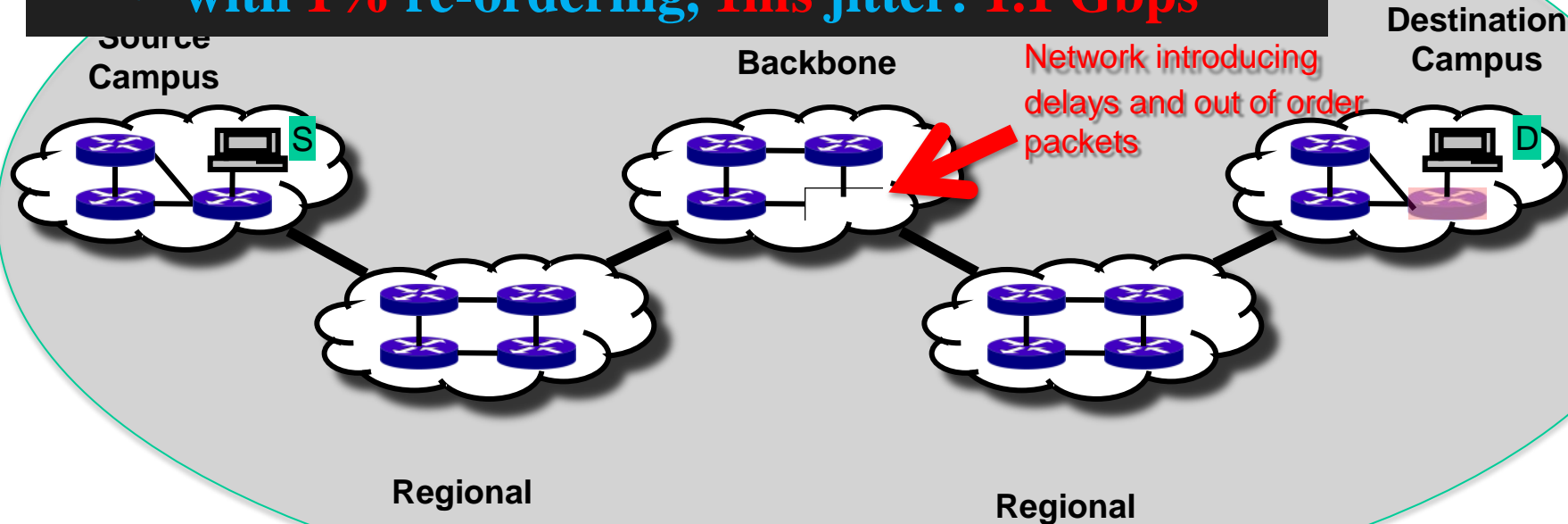


# Packet ordering and jitter



At 70ms RTT on 10G link, 60 seconds test

- with 1% re-ordering, 0.2 ms jitter: 8.45 Gbps
- with 1% re-ordering, 1ms jitter: 1.1 Gbps





# OSG and WLCG Network Efforts



- OSG is in its fifth year of supporting WLCG/OSG networking and is focused on:
  - Assisting its users and affiliates in **identifying** and **fixing** network bottlenecks
  - Supporting higher-level network services
  - Improving the ability to manage and use network topology and network metrics: Analytics Platform based upon ELK in use
  - Developing effective **Alarming and Alerting** for network problems
- The WLCG Network and Transfer Metrics working group has created a support unit to coordinate responses to potential network issues
  - Tickets opened in the support group can be triaged to the right destination
  - Many issues are potentially resolvable within the working group
  - Network issues can be identified and directed to the appropriate network support centers
  - Documented at [https://twiki.cern.ch/twiki/bin/view/LCG/NetworkTransferMetrics#Network\\_Performance\\_incidents](https://twiki.cern.ch/twiki/bin/view/LCG/NetworkTransferMetrics#Network_Performance_incidents)
  - Many issues resolved within hours mainly due to using perfSONAR information

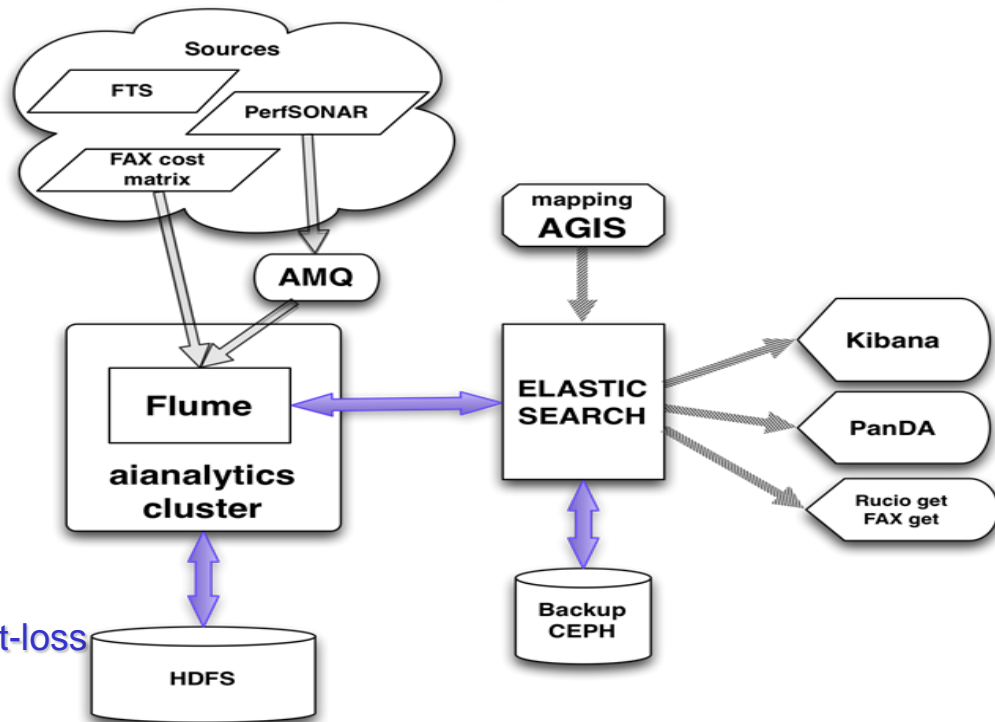
# ATLAS Network Analytics



- Ilija Vukotic/U Chicago has led the effort to get network metrics into an analytics platform.
- This analytics service indexes historical network related data while providing predictive capabilities for network throughput. (See Ilija's talk Thurs Track 5)

## Primary functions:

- Aggregate, and index, network related data associated with WLCG “links”
- Serve derived network analytics to ATLAS production, DDM & analysis clients
- Provide a generalized network analytics platform for other communities in the OSG
- Initial “Alarm” query prototyped and tested for Source-Destination paths with high packet-loss
- More details at: <http://tinyurl.com/gt92zwb>



# Making the Most of our Networks



- Much of our ATLAS infrastructure is NOT tuned to take the best advantage of the networks we currently have
  - There are a wide range of mis-configurations, non-optimal tunings and incorrect application, firmware and hardware settings that lead to inefficient use of our networks
  - As you have seen we have a wealth of data now available and analyzable to help identify bottlenecks and poor performance.
- With this infrastructure we need to take the next step and work to improve ATLAS resources ability to effectively utilize the network
  - Doesn't require SDN, new hardware or new networks but can make a huge difference
  - Should we organize a near-term workshop to share best practices, tools and tuning information?

# PanDA and Networking



- PanDA is ATLAS's workload manager
  - PanDA automatically chooses job execution site
    - Multi-level decision tree – task brokerage, job brokerage, dispatcher
    - Also predictive workflows – like PD2P (PanDA Dynamic Data Placement)
  - Site selection is based on processing and storage requirements
    - **Why not use network information in this decision?**
    - Can we go even further – network provisioning?
  - Network knowledge useful for all phases of job cycle
- Network as resource
  - Optimal site selection should take network capability into account
    - We do this already – but indirectly using job completion metrics
  - Network as a resource should be managed (i.e. provisioning)
    - We also do this crudely – mostly through timeouts, self throttling
- Longer-term goal for PanDA
  - Direct integration of networking with PanDA workflow – never attempted before for large scale automated WMS systems

# Playing with SDN in ATLAS



- Future networks won't just have larger capacity
- A group of people in the US from AGLT2, MWT2, SWT2 and NET2 are planning to explore SDN in ATLAS
  - Working with the LHCONe point-to-point effort as well
- The plan is to deploy Open vSwitch on ATLAS production systems at these sites (<http://openvswitch.org/>)
  - IP addresses will be move to virtual interfaces
  - No other changes; verify no performance impact
  - Traffic can be shaped accurately with little CPU cost
- The **advantage** is the our data sources/sinks become **visible** and **controllable** by OpenFlow controllers like OpenDaylight
  - **BENEFIT:** Traffic shaping can result in significantly improved use of the WAN for some paths
- Follow tests can be initiated to provide experience with controlling networks in the context of ATLAS operations.
- Interest from UVic, KIT and SurfSARA in participating
- Possible partnership with ESnet/CORSA in ~Dec 2016 timeframe
- *For more details talk to Rob Gardner or Shawn McKee*

# Future Directions



- The WLCG efforts at CERN are being reorganized and this is an opportunity to chart future directions for the our networking efforts.
- We have a number of project areas we are considering and we need to understand where these efforts should be housed (Stay in WG, move to GDB, to LHCONE, elsewhere?)
  - It is important to note there is currently very little manpower for networking (much, much less than computing and storage)
  - To undertake all our plans will require identifying new effort
- We are planning a Pre-GDB meeting on January 10<sup>th</sup> focusing on networking. Save the date!

# Summary



- We have a working infrastructure in place to monitor and measure our networks in use for ATLAS
- perfSONAR provides lots of capabilities to understand and debug our networks
- Work on new applications is underway
  - Notifications/alerting
  - Predictive capabilities
  - Current utilization and capacity planning
- Analytics on network and transfer metrics now possible along with the chance to fix non-optimal infrastructure once we identify it
- Network capabilities will evolve based upon commercial goals...ATLAS should be ready to take advantage of what becomes available if it make sense.

**Questions or Comments?**



# References



- Network Documentation <https://www.opensciencegrid.org/bin/view/Documentation/NetworkingInOSG>
- Deployment documentation for OSG and WLCG hosted in OSG  
<https://twiki.opensciencegrid.org/bin/view/Documentation/DeployperfSONAR>
- Measurement Archive (MA) guide [http://software.es.net/esmond/perfsonar\\_client\\_rest.html](http://software.es.net/esmond/perfsonar_client_rest.html)
- Modular Dashboard and OMD *Prototypes*
  - <http://maddash.aglt2.org/maddash-webui> [https://maddash.aglt2.org/WLCGperfSONAR/check\\_mk](https://maddash.aglt2.org/WLCGperfSONAR/check_mk)
- **OSG Production instances for OMD, MaDDash and Datastore**
  - <http://psmad.grid.iu.edu/maddash-webui/>
  - [https://psomd.grid.iu.edu/WLCGperfSONAR/check\\_mk/](https://psomd.grid.iu.edu/WLCGperfSONAR/check_mk/)
  - <http://psds.grid.iu.edu/esmond/perfsonar/archive/?format=json>
- Mesh-config in OSG <https://oim.grid.iu.edu/oim/meshconfig>
  - Being updated to a new standalone mesh-config application (ready for v3.6?)
- Use-cases document for experiments and middleware  
<https://docs.google.com/document/d/1ceiNITUJCwSuOuvbEHZnZp0XkWkwkPQTQic0VbH1mc/edit>
- Big Data Analytics Tools (CHEP 2016) <https://indico.cern.ch/event/505613/contributions/2228332/>

# Back up slides



# Possible Future Project Areas



- **Title:** LHCONE Traffic engineering
- **Areas:** LHCONE, routing, debugging, network orchestration
- **Title:** LHCONE L3VPN Looking Glass
- **Areas:** LHCONE, monitoring, debugging
- **Title:** Integration of network and transfer metrics to optimize experiments workflows
- **Areas:** FAX/Phedex, Rucio, perfSONAR, DIRAC
- **Title:** Advanced notifications/alerting for network incidents
- **Areas:** WAN, Advanced Notifications/Alerting, perfSONAR, Hadoop/Spark
- **Title:** Network performance of the commercial clouds
- **Areas:** Clouds, WAN connectivity, WAN performance (perfSONAR), establishing and testing network equipment at the cloud provider (VPN)
- **Title:** Software Defined Network Production Testbed
- **Areas:** WAN, SDN, LHCONE/LHCOPN, Storage/Data nodes