# The ATLAS Computing Model, Use Cases and Network Implications

**Shawn McKee - University of Michigan**

**Workshop on Transatlantic Networking for LHC Experiments**

**June 10th 2010**

ATLAS Detector Under construction
November 2005

# Talk Overview

❊ A brief outline of the ATLAS Computing Model

  ❑ How does data flow (in principle)

  ❑ How is ATLAS doing its scientific work?

❊ Use Cases, Recent Examples & Network Implications

❊ Summary and Conclusions

NOTE:  I will be presenting <u>my perspective</u> of ATLAS activities and use cases, which is US-centric!

# ATLAS Tiered Computing Model

* Within this model the each Tier has a set of responsibilities:

  - Tier-0 – First pass reconstruction, archive ALL data

  - Tier-1 – Data reprocessing and archiving, User/Group analysis

  - Tier-2 – Simulation and User analysis

* Implicit in this model and **central to its success** are:

  - High-performance, ubiquitous and robust **networks**
  - Grid middleware to securely **find**, **prioritize** and **manage** resources

* **Without either of these capabilities the model risks melting down or failing to deliver the required capabilities.**

* Efforts to date have (*necessarily*) focused on building the most basic capabilities and ensuring they can work.
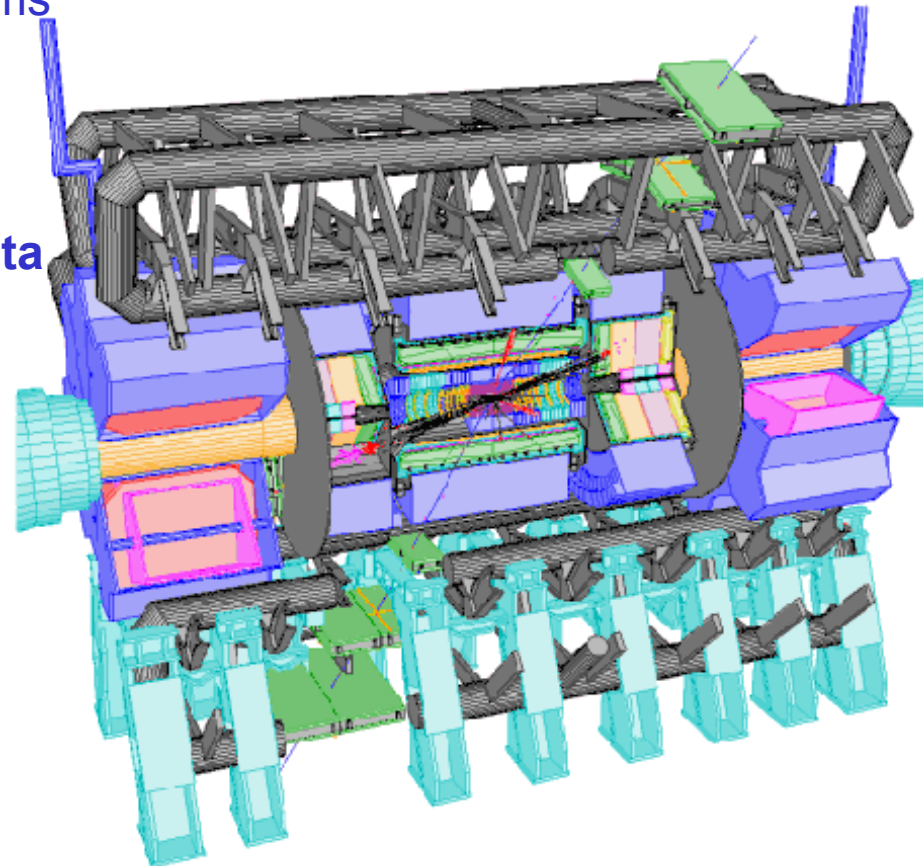
# ATLAS Physicist's Requirements

* ATLAS physicists need the software and physical infrastructure required to:
  - Calibrate and align detector subsystems to produce well understood data
  - Realistically simulate the ATLAS detector and its underlying physics
  - **Provide timely access to ATLAS data globally**
  - Accurately reconstruct data to allow new physics discoveries
  - Define, manage, search and analyze data-sets of interest

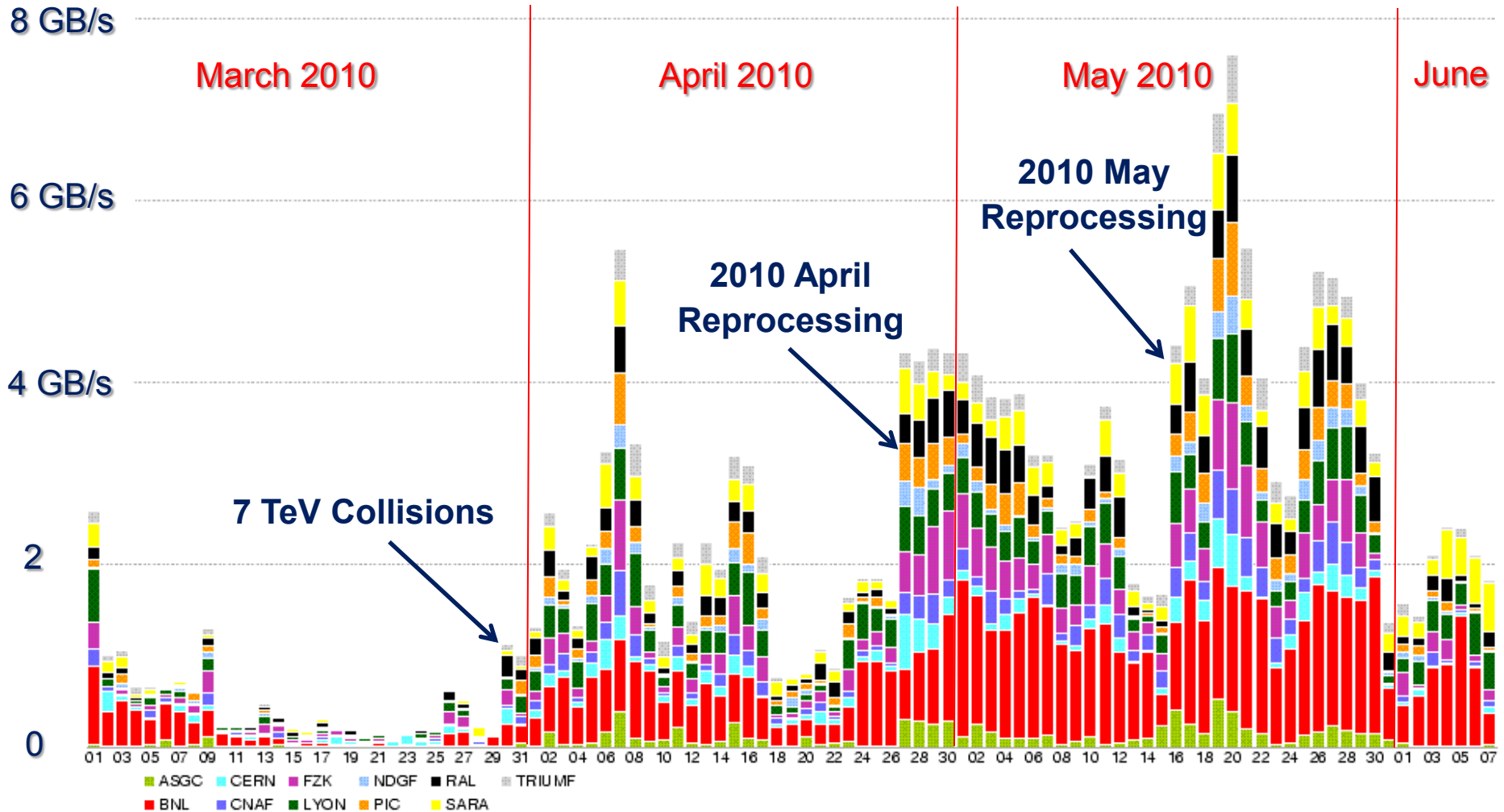* Networking plays a fundamental role for all of these activities

ATLAS

# Recent ATLAS Data Movement Results

* Since the LHC restarted in March we have gotten a quick look at how the DDM infrastructure has worked (next slide)

* Summary is "Very well"…but there are some details which are useful to cover and some issues we need to address.

* Generally there is a large amount of data to move and not all of it is equally interesting for physicists.

* After data reprocessing, many physicists want immediate access to certain datasets and submit a **large** number of grid jobs targeted at this new data.

* Remember that ATLAS generally tries to send jobs to the data…so, as long as the data exists where the CPU slots exist, we are good…if not…( I will cover this in a few slides)

# Recent ATLAS Data Distribution

# Use Case: ATLAS Datasets to Physicists

* In ATLAS we have a goal of getting datasets to Tier-2 centers quickly (~ 4 hours).

* This is especially important for **"interesting"** or **"hot"** datasets that generate a large number of user/group jobs requiring these datasets as inputs (will discuss later)

* The ATLAS "cloud" based distribution model previously described makes timely access to 'hot' new datasets challenging:

  ❑ Get to destination cloud: transfer from source to local Tier-1

  ❑ Next:  transfer from local Tier-1 to destination site in cloud

  ⇒ **Increases both I/O and latency vs direct src-dst move**

  ⇒ **However it is better controlled and easier to debug problems**
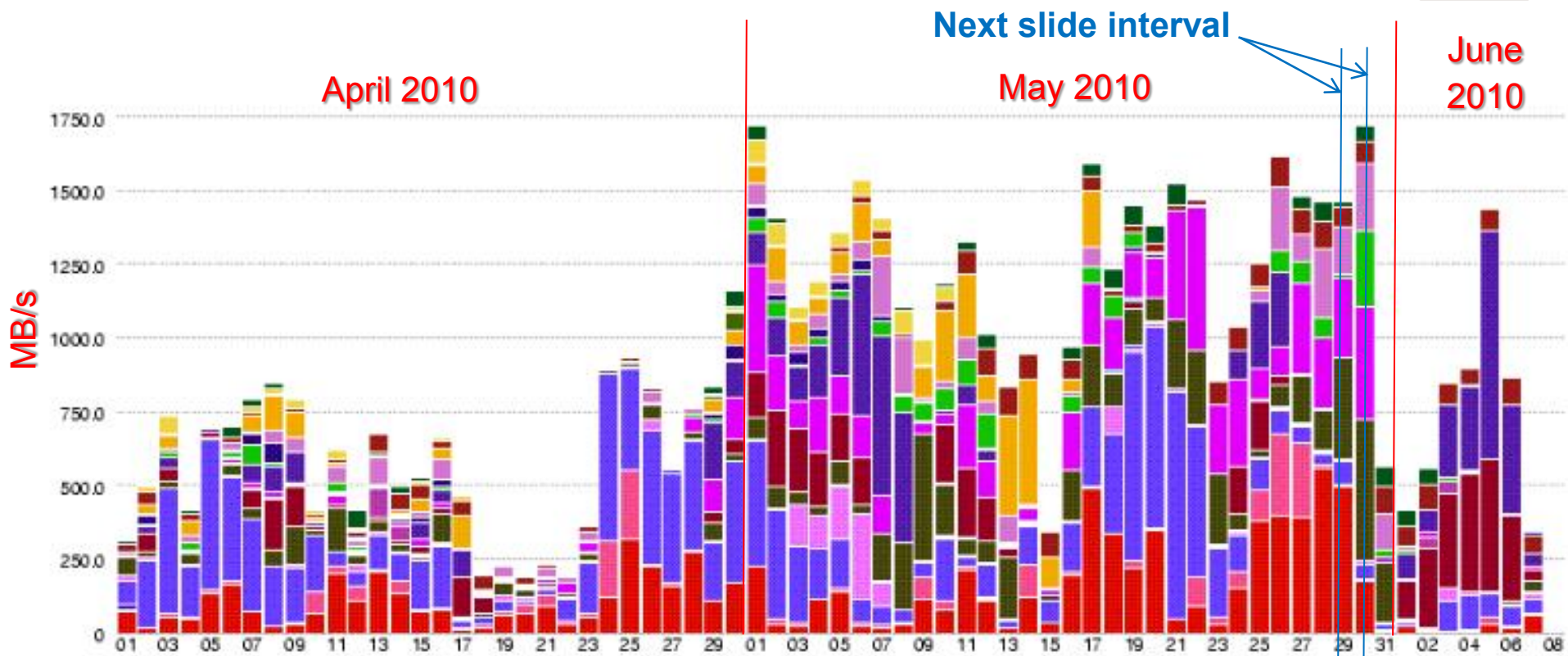
# US Tier-2 Transfer Capabilities

* Most USATLAS Tier-2 centers have 10GE connectivity (SWT2-UTA (last 1GE site) upgraded by September 2010)

* Our baseline assumes Tier-2's should be able to ingest data at >400MB/sec.   Bursts up to line capacity (1.25 GB/s).

* Recently data reprocessing distributions have shown Tier-2s capable of >800 MB/sec continuously (some ~1.25 GB/s)

* Assuming 800 MB/sec, we can move 1 TB datasets in about **21 minutes** or 10 TB datasets in **3 ½  hours**

* Note a dataset larger than **11.52 TB** has a total transfer time **> 4 hours…in other words**, there is a lower limit on dataset latency determined by size & achievable bandwidth

* For reference **10 Gbps** data transfer => **4.5 TB/hour**
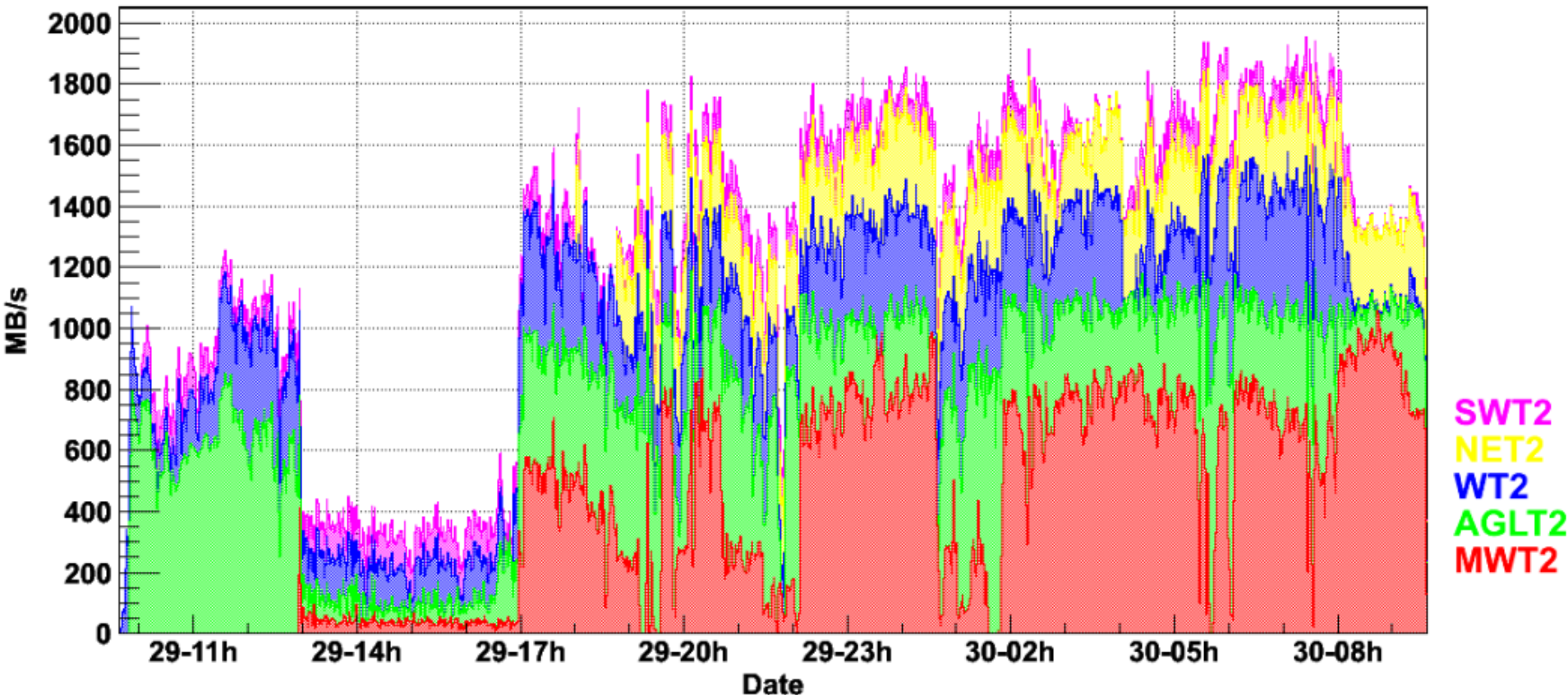
# ATLAS Data Transfers and Latency

❊ Of course the previous slide implicitly assumed we were ONLY moving the "important" dataset and had the "input" dataset already at the Tier-1 in our cloud, neither of which is typically true when new interesting datasets are first ready.

❊ In practice, Tier-2s usually are busy either receiving or sending data (MC results, calibration/conditions data, user jobs, etc)…**there is usually competing activity**.

❊ Data from outside the cloud requires the Tier-1 to transfer the dataset first  (but it can "overlap" the transfer).  **Note I/O is doubled: write-then-read, which impacts throughput**

❊ **Therefore meeting a 4 hour latency can be difficult…**

# BNL Tier-1 Cloud Data Transfers



Shown are the BNL Tier-1 cloud data transfers for April till early June.  Many days average above 10 Gb/s (1250 MB/s). (Different colors denote destination sites)

# USATLAS Tier-2 Data Transfers



The US Tier-2 sites are currently capable of using ~1.5x10 Gb/s of network bandwidth for data distribution.

Once the SWT2 is upgraded to a 10GE path we anticipate regularly filling 2x10GE for Tier-2 data distributions

# Challenges with Dataset Availability

❄ During/after the April reprocessing in early March we had a large number of users submit jobs needing these datasets.

❄ The BNL Tier-1 relatively quickly acquired the datasets and begin redistributing them to the US Tier-2 sites

❄ However a significant amount of other Monte-Carlo data was also being transferred and reprocessed datasets were arriving at Tier-2's much too slowly

❄ Because of job-to-data matchmaking, BNL quickly had ALL users jobs queued up (~100,000)

❄ **Tier-2's had empty analysis slots waiting for jobs**

❄ The problem was identified and a series of manual "fixes" were applied to allow ONLY reprocessed data to transfer to the Tier-2s to resolve the backlog. Not a long-term solution

# Implications for the Future

❋ The transfer capabilities of the Tier-1 and Tier-2s are very good.  Some Tier-2 sites can fill a 10GE link on their own.

❋ When large amounts of data are being distributed, the ATLAS DDM system *performs well in ensuring all data is transferred…<u>eventually</u>*.

❋ However, in most cases datasets are NOT equally important and have different urgencies (in terms of being ready for users to access).  This importance changes in time.

❋ We need the capability of expressing relative "importance" by dataset AND have an infrastructure that can allocate available resources accordingly.

❋ **Interaction of network services and the DDM system will be required to deliver this capability**

# Additional Implications

※ The USATLAS Tier-2 centers are large and planned to grow to meet there MOU requirements.

※ Network needs scale with processing power and local storage.  Currently a typical US site has ~1500 processors and 1 Petabyte of storage and this will grow.

※ The current ATLAS cloud model restricting transfers between clouds to the Tier-1's needs re-evaluation.

❑ Original intent was to provide well defined and managed inter-cloud links to facilitate debugging and manage "load"

❑ As Tier-2s become more powerful we need to look at the cost in latency and additional I/O impact for the store and forward model.

❑ Data transfer decisions should be based on resources capabilities

※ Changes would have implications for transatlantic networks

# Need for Pervasive Monitoring

�֍ Many of you are probably aware that all problems of unknown origin are **"network"** problems

✷ It is easy to attribute such problems to the "network" because of its black-box nature and its potentially large set of administrative domains for a typical end-to-end path.

✷ In practice problems in the "network" or more likely to be local problems at the source or destination…but how can we know?

✷ Having "standardized" monitoring that can **identify current and past performance** as well as the capability of **isolating the location of performance or connectivity issues** is critical for managing wide-area science.

# Network Monitoring: perfSONAR

❋ There is a significant, coordinated effort underway to instrument the network in a standardized way. This effort, call perfSONAR, is jointly supported by DANTE, Esnet, GEANT2, Internet2 and numerous University groups.

❋ Since the network is so fundamental to our work on the ATLAS, we targeted implementation of a perfSONAR instance at all our primary facilities.

❋ **perfSONAR's primary purpose is to aid in network diagnosis** by quickly allowing users to isolate the location of problems. **In addition it can provide a standard measurement of various network performance related metrics over time.**

❋ Has already proven very useful in USATLAS!

# Example: AGLT2's perfSONAR



## pS-Performance Node For AGLT2-MSU In East Lansing, MI, USA

### User Tools
- Services On This Node
- Global Set Of Services
- Java OWAMP Client
- Reverse Traceroute
- Reverse Ping
- PingER Web GUI

### Service Graphs
- Throughput
- One-Way Latency
- Ping Latency
- SNMP Utilization
- Cacti Graphs

### Toolkit Administration
- Administrative Information
- External BWCTL Limits
- External OWAMP Limits
- Enabled Services
- NTP

### Services Offered

**Bandwidth Test Controller (BWCTL)[1]** — Running
- tcp://psmsu02.aglt2.org:4823

**Lookup Service[1]** — Running
- http://psmsu02.aglt2.org:8095/perfSONAR_PS/services/hLS

**Network Diagnostic Tester (NDT)[1]** — Running
- tcp://psmsu02.aglt2.org:3001
- http://psmsu02.aglt2.org:7123

**Network Path and Application Diagnosis (NPAD)[1]** — Running
- tcp://psmsu02.aglt2.org:8100
- http://psmsu02.aglt2.org:8200

**One-Way Ping Service (OWAMP)[1]** — Disabled
- tcp://psmsu02.aglt2.org:861

**perfSONAR-BUOY Regular Testing (Throughput)[1]** — Running

**perfSONAR-BUOY Measurement Archive[1]** — Running
- http://psmsu02.aglt2.org:8085/perfSONAR_PS/services/pSB

**perfSONAR-BUOY Regular Testing (One-Way Latency)[1]** Disabled

**PingER Measurement Archive and Regular Tester[1]** — Disabled
- http://psmsu02.aglt2.org:8075/perfSONAR_PS/services/pinger/ma

**SNMP Measurement Archive[1]** — Not Running
- http://psmsu02.aglt2.org:8065/perfSONAR_PS/services/snmpMA

Deployed at Tier-1 and all Tier-2s in the USATLAS
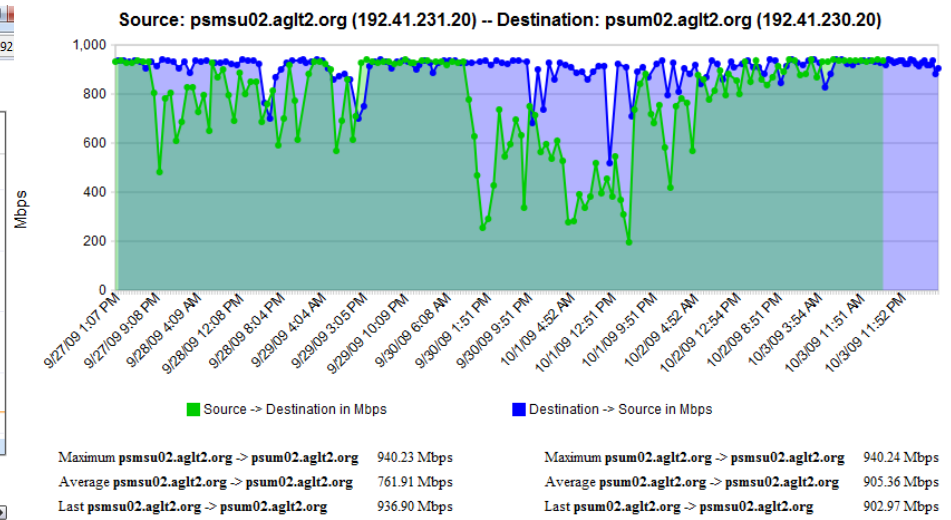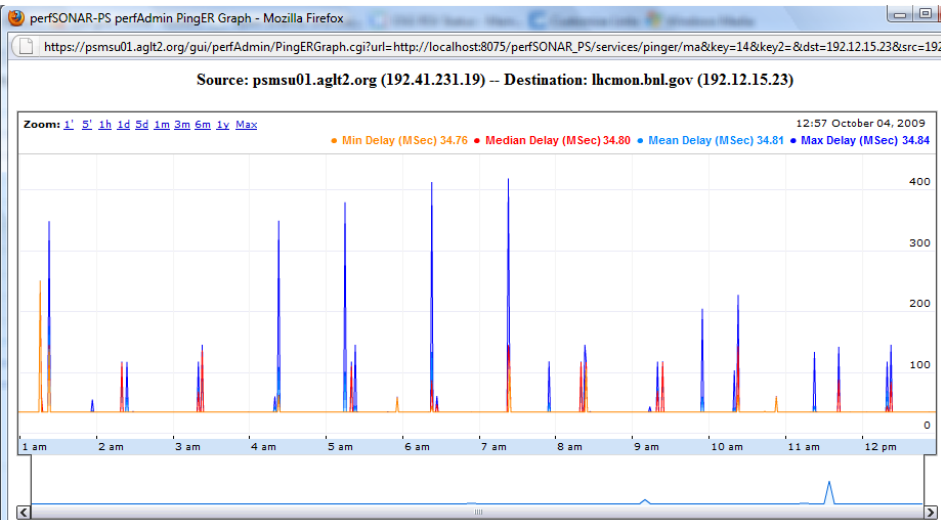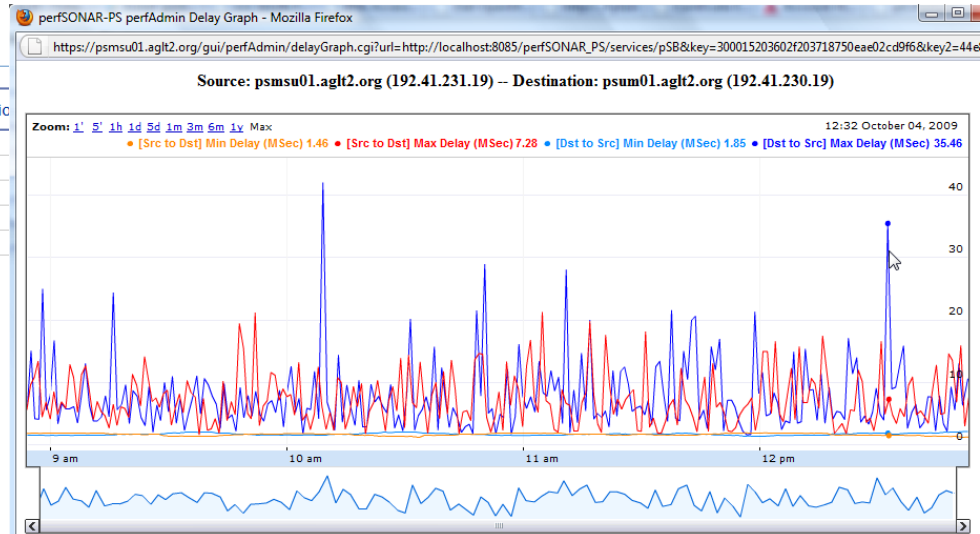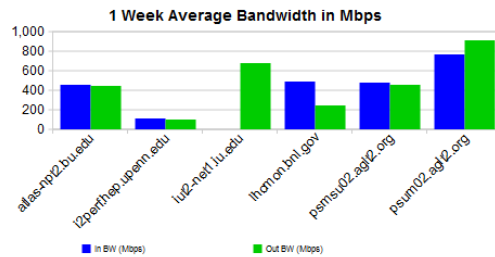
Provides throughput and latency measurements

Used for problem isolation and performance monitoring

# perfSONAR Example Information

# Status of perfSONAR for USATLAS

❋ Fully deployed at all Tier-2 (and BNL Tier-1) "sites" (most Tier-2s are comprised of more than one physical site)

❋ Original hardware specified in 2008.  Inexpensive system (1U) from KOI Computing. Two boxes deployed: latency and bandwidth measurement roles.  Has been problematic:

- ❑ Boxes have had some driver issues and exposed perfSONAR bugs
- ❑ Systems seem underpowered at the scale of use for USATLAS
  - ⌘ Difficult to look at results (slow or timeouts)
  - ⌘ Some measurements hang (size of DB related?)

❋ **Primary missing component**: Automated monitoring of results with ALERTING.   Ongoing project for USATLAS

❋ New hardware purchased: Dell R4101U, Intel E5620, 12GB, 10GE Myricom, 2x1GE.  Possible 2x1U replacement

# ATLAS Networking Needs

❄ There isn't anything unique about ATLAS networking needs compared with LHC networking needs. ATLAS requires:

❑ **Robust networks**, end-to-end.   Extended loss of connectivity can be extremely disruptive.

❑ **Sufficient bandwidth** to support our physics needs.  This varies with time and source/destination but currently is:

⌘ 20-30 Gb/s for the Tier-1

⌘ 10 Gb/s for each Tier-2 (Tier-1=>Tier-2s at 20Gb/s)

⌘ These values support the **current** peak usage…this will grow as processors and storage at sites ramp-up (factor of ~2 by 2013?)

❑ Ability to **prioritize** traffic to match our needs.  High-demand datasets need higher priority to meet user needs/expectations.

❑ **Monitoring** to identify and isolate problems and verify normal operation (baseline setting)

# Status and Conclusions

※ ATLAS transatlantic networking has worked well as the LHC has started physics operations.

※ Current ATLAS cloud model certainly needs re-examination. A change to a more grid-like data access model may be facilitated by better, more pervasive monitoring, e.g., perfSONAR.

※ Having prioritization mechanism's for data distribution is needed. This may involve network services to support this capability.

※ Depending upon how ATLAS DDM evolves there may be more transatlantic traffic (burst-wise) because of Tier-2 related data transfers.   The Tier-2s in the US are already large and are planned to grow significantly in both storage and processors.

※ Robust, well monitored transatlantic networks are required for US Physicists to be able to effectively participate in ATLAS

# ?Questions?