

## perfSONAR

Please provide short introduction, motivation and description of your use cases

Introduction:

perfSONAR description in detail can be found at <http://www.perfsonar.net/about/what-is-perfsonar/>

Motivation and use cases:

Ensure that in WLCG we're able to effectively use the network and quickly resolve network issues when and where they occur. In particular:

Identifying and localizing network problems

- Often this is very difficult and time-consuming for Wide-Area Network (WAN) problems

Scheduled perfSONAR bandwidth and latency metrics monitor WLCG network paths

- Significant packet-loss or consistent large deviation from baseline bandwidth indicate a potential network problem (see in GUI or via alarms).
- On-demand tests to perfSONAR instances can verify the problem exists. Different test points along the path can help pinpoint the location.
- Correlation with other paths sharing common segments can be used to localize the issue.
- The time things change is also very useful to find the root causes. Scheduled tests provide this.

The above use cases are based on the following considerations:

- Network problems can be hard to diagnose and slow to fix
- 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 (as well as benefiting from others work)
- Performance issues involving the network are complicated by the number of components involved end-to-end. We need the ability to better isolate performance bottlenecks

What network-based metrics are you providing, what are their characteristics (frequency, latency). Focus on raw metrics, but also add information on aggregations. Add additional details on coverage and granularity of measurements and data access (API if any)

Metrics:

Network path - We use perfSONAR's traceroute to track the network path between WLCG sites. Currently 1/hour between ALL WLCG perfSONAR latency instances.

Latency - We send 10Hz of one-way delay measurement packets between all WLCG sites. The packet statistics (avg,min,max delay) are summarized every minute and any packet losses (x/600 packets) are noted (a critical metric)

Bandwidth - We use perfSONAR's lperf tool to measure achievable bandwidth. Depending upon the grouping (mesh) we test every 6 hours both directions between src-dst. We additionally test all WLCG pairs 1/week)

Coverage and granularity:

214 hosts, 108 sites, coverage details at [http://grid-monitoring.cern.ch/perfsonar\\_coverage.txt](http://grid-monitoring.cern.ch/perfsonar_coverage.txt) (updated daily). Latency and bandwidth are measured sonar to sonar, in addition there is

Access/API:

Each perfSONAR has its own measurement archive that is accessible via REST API. We're also planning a centralized data store (more information at <https://twiki.opensciencegrid.org/bin/view/Production/OSGNetworkDatastorePlan>)

What is the operational status of your deployment, what are the outstanding tasks that need to be done to have production-ready system, what are the planned changes in the short-term (3-6 months) future.

Current operational report is available at [http://grid-monitoring.cern.ch/perfsonar\\_report.txt](http://grid-monitoring.cern.ch/perfsonar_report.txt)  
We're currently running a WLCG-wide campaign to update all sites to perfSONAR 3.4; there are regular perfSONAR operations meeting that follow up. Tracker(s) for commissioning and operations is at <https://its.cern.ch/jira/browse/METRICS-5> ; <https://its.cern.ch/jira/browse/METRICS-7> ; <https://its.cern.ch/jira/browse/METRICS-6>

Do you have any prior experience in correlating network and transfer metrics ?

Additional details to use cases:

What network or transfer metrics would you be interested in and what characteristics are expected (coverage, granularity, frequency, latency, etc.) ?

From the perspective of perfSONAR, it would be interesting to correlate latency metrics with link status provided by FAX and FTS as well as understand how to match/link the corresponding metrics wrt. coverage and granularity provided by FTS/FAX.

How do you plan to use them (optimize, what/how) ?

The results can be used to understand how to tune the current latency measurements, coverage and granularity of the data provided by the perfSONAR network.

## FTS

Please provide short introduction, motivation and description of your use cases

FTS

- Low-level data movement service, moves data sets from one site to another (SE to SE)
- Used for majority of LHC 3rd party transfers

Use cases (TBA)

What network-based metrics are you providing, what are their characteristics (frequency, latency). Focus on raw metrics, but also add information on aggregations. Add additional details on coverage and granularity of measurements and data access (API if any)

For each individual transfer publishes event messages to the message bus

- Raw data contain the following event types: start, complete, state
  - (Details at <https://fts3-service.web.cern.ch/content/install%26config>)

In addition, command line client reports aggregated snapshots

There is also an aggregated dashboard for FTS:

FTS dashboard service provides in-depth details (<http://dashb-fts-transfers.cern.ch/>)

- calculates transfer rates (throughput, volume, success/failure) per site, vo, vo's activities, host, country, token

- aggregates on FTS states

There is also FTS3 monitoring service (low-level, part of each FTS3 service deployed), e.g.

<https://fts3.cern.ch:8449/fts3/ftsmon/#/>

<https://www-ftsmon.gridpp.rl.ac.uk:8449/fts3/ftsmon/#/>

<https://cmsfts3.fnal.gov:8449/fts3/ftsmon/#/>

Current granularity of raw FTS metrics are SE to SE (host-based)

Current coverage depends on the experiment (but all T1/T2 are included), links TBA

What is the operational status of your deployment, what are the outstanding tasks that need to be done to have production-ready system, what are the planned changes in the short-term (3-6 months) future.

TBA

Do you have any prior experience in correlating network and transfer metrics ?

TBA

Additional details to use cases:

What network or transfer metrics would you be interested in and what characteristics are expected (coverage, granularity, frequency, latency, etc.) ?

Understanding if perfSONAR link status (latency) and routing information could be used in the optimizer.

Ability to identify the closest available sonar to an SE would be needed to start with.

How do you plan to use them (optimize, what/how) ?

Potentially, fine tune the optimizer. The current optimizer has the following algorithm:

If the success rate is 100% and the most current throughput sample is higher than the previous sample than start one more transfer

If the success rate is 100% and the most current throughput sample is equal to the previous sample than number of active transfer remain steady

If the success rate is 100% and the most current throughput sample is less than the previous sample than number of active decreases by one

If success rate is  $< 100\%$  then decrease by 2

To smooth out the frequent bumps we use exponential moving average algorithm, and for getting the throughput of a link we run weighted average throughput.

## FAX

Please provide short introduction, motivation and description of your use cases

Federated XRootD services for ATLAS, introduction at

<https://twiki.cern.ch/twiki/bin/viewauth/AtlasComputing/AtlasXrootdSystems>

XRootD introduction at <http://xrootd.org/>

Mainly, how perfSonar information can be used to improve our understanding of the network links relevant to ATLAS. Our numbers are telling us all that we need to know concerning individual file transfers but nothing on the total bandwidth of the link or the total remaining bandwidth. Would be also nice to know if there are some ideas how to centrally map ATLAS/CMS sites to perfSonar instances.

What network-based metrics are you providing, what are their characteristics (frequency, latency). Focus on raw metrics, but also add information on aggregations. Add additional details on coverage and granularity of measurements and data access (API if any)

XRootD monitoring

- XRootD provides two types of monitoring: summary (xrd.monitor) and detail (xrd.detail)

([http://xrootd.org/doc/dev4/xrd\\_monitoring.pdf](http://xrootd.org/doc/dev4/xrd_monitoring.pdf))

- Monitoring information is sent over UDP (to provide near real-time with low impact, medium latency - 5 minutes or full buffer 1-8kB)

XrdMon

- Based on GLED, collects and aggregates information from multiple XRootD endpoints

- Processes the monitoring information and maintains in-memory representation of all ongoing user sessions and file transfers, highly configurable

- Supports multiple backends (ActiveMQ, Gratia, TTree, Http)

- Details at <http://www.gled.org/cgi-bin/twiki/view/Main/XrdMon>

FAX/AAA dashboards (<dashb-atlas-xrootd-transfers.cern.ch>, <dashb-cms-xrootd-transfers.cern.ch>)

- Raw data consumed from message bus (published by XrdMon),

- Events generated for every file close operation (many additional attributes are available, details at [https://twiki.cern.ch/twiki/bin/view/LCG/WLDataTransferMonitoring#RAW\\_Table](https://twiki.cern.ch/twiki/bin/view/LCG/WLDataTransferMonitoring#RAW_Table))

- Dashboard calculates statistics on transfer rates (throughput/volume) and active/finished transfer counts per site, country

There is also additional testing run in parallel.

–HammerCloud continuously submits jobs to ~40 ATLAS analysis queues

(<http://hammercloud.cern.ch/hc/>)

–These jobs copy special files from each of the FAX endpoints (every 15min, including traces)

–Results reported to SSB and FSB

[<http://dashb-atlas-ssb.cern.ch/dashboard/request.py/siteview#currentView=FAX+cost+matrix>,

<http://waniotest.appspot.com/>] and via AGIS to JEDI for brokering

What is the operational status of your deployment, what are the outstanding tasks that need to be done to have production-ready system, what are the planned changes in the short-term (3-6 months) future.

- All services mentioned are deployed and operated in production
- There are two GLED/XrdMon operated in US, one for ATLAS, one for CMS
- Additional GLEDs for EOS, ActiveMQ brokers and FAX/AAA dashboards are operated at CERN

Do you have any prior experience in correlating network and transfer metrics ?

TBA

Additional details to use cases:

What network or transfer metrics would you be interested in and what characteristics are expected (coverage, granularity, frequency, latency, etc.) ?

TBA

How do you plan to use the results (optimize, what/how) ?

TBA

## PhEDEx

Please provide short introduction, motivation and description of your use cases

PhEDEx stands for Physics Experiment Data Export (more details at <https://github.com/dmwm/PHEDEX>)

Use cases TBA

What network-based metrics are you providing, what are their characteristics (frequency, latency). Focus on raw metrics, but also add information on aggregations. Add additional details on coverage and granularity of measurements and data access (API if any)

PhEDEx currently provide throughput graphs (e.g. <https://cmsweb.cern.ch/phedex/prod/Activity::RatePlots>) - counts the rate at which files arrive at a destination. Plots are coarsely binned, for any transfer that ends in a given one-hour window the total transfer is ascribed to that window. That can be inaccurate for low-rate transfers, but for high-rate, it's close enough to tell us if things are working or not.

PhEDEx also counts the error-rate for file transfers, and this translates into the 'link quality' (<https://cmsweb.cern.ch/phedex/prod/Activity::QualityPlots>).

PhEDEx maintains the history of these metrics internally, and uses them to choose a source-site for a given transfer (users only specify that they want dataset XYZ transferred to a given destination, PhEDEx chooses the source for them).

PhEDEx has only a very high-level view of the network topology, so instead of monitoring performance along a sequence of internet hosts a la traceroute, it only knows about 'nodes' such as 'T1\_US\_FNAL\_Buffer' or 'T0\_CH\_CERN\_Export' or 'T2\_US\_Florida\_Disk' as network endpoints. It uses hostnames when building transfer-URLs, of course, but they're opaque data as far as PhEDEx is concerned, it has no understanding of what a hostname means.

N.B. a 'link' to PhEDEx is simply a logical connection between two nodes, it has nothing to do with physical network links.. We 'enable' links (set a bit in our DB) to allow traffic between two sites, and we disable links that we don't want to permit (so, for example, we don't have outbound links from MSS systems to every T2/T3 in CMS).

All this data is stored in the PhEDEx Oracle DB at CERN, and can be accessed using APIs (see <http://cmsweb.cern.ch/phedex/datasvc/doc> for the documentation). The data itself is produced by PhEDEx file-download agents at the sites when they report the success or failure of a given file transfer, and aggregated centrally into metrics per-hour.

What is the operational status of your deployment, what are the outstanding tasks that need to be done to have production-ready system, what are the planned changes in the short-term (3-6 months) future.

We have three separate instances of PhEDEx; the Production instance (real + MC data), the Dev instance (for testing new code) and the Debug instance, for debugging and commissioning sites and links.

Do you have any prior experience in correlating network and transfer metrics ?

Additional details to use cases:

What network or transfer metrics would you be interested in and what characteristics are expected (coverage, granularity, frequency, latency, etc.) ?

At least the following is needed to progress on this:

- 1) Ensure that we have full coverage of the same mesh
- 2) Map the real network topology to the 'PhEDEx topology'
- 3) Understand what perfSONAR measurements mean in terms of our concepts of throughput

How do you plan to use the results (optimize, what/how) ?

As described above



## RUCIO

Please provide short introduction, motivation and description of your use cases

What network-based metrics are you providing, what are their characteristics (frequency, latency). Focus on raw metrics, but also add information on aggregations. Add additional details on coverage and granularity of measurements and data access (API if any)

What is the operational status of your deployment, what are the outstanding tasks that need to be done to have production-ready system, what are the planned changes in the short-term (3-6 months) future.

Do you have any prior experience in correlating network and transfer metrics ?

Additional details to use cases:

What network or transfer metrics would you be interested in and what characteristics are expected (coverage, granularity, frequency, latency, etc.) ?

How do you plan to use the results (optimize, what/how) ?

## PANDA

Please provide short introduction, motivation and description of your use cases

What network-based metrics are you providing, what are their characteristics (frequency, latency). Focus on raw metrics, but also add information on aggregations. Add additional details on coverage and granularity of measurements and data access (API if any)

What is the operational status of your deployment, what are the outstanding tasks that need to be done to have production-ready system, what are the planned changes in the short-term (3-6 months) future.

Do you have any prior experience in correlating network and transfer metrics ?

Additional details to use cases:

What network or transfer metrics would you be interested in and what characteristics are expected (coverage, granularity, frequency, latency, etc.) ?

How do you plan to use the results (optimize, what/how) ?

## Experiments/ATLAS

### terminology:

**link** = the transfer route from SRC to DST (both, either reserved or non-reserved)

**mesh** = a map of the ATLAS sites where all the sites are linked to each other

**full mesh** = a map of all possible ATLAS sites and links between them

**The Sonar** = networking tool developed by ADC. Each site stores at its datadisk 10 sonar files of 1GB which are sent once per week to every other site. Transfers are equally spread over whole grid. Rucio is used as transfer manager and fts logs are fetched to find out the information about given transfer.

<https://twiki.cern.ch/twiki/bin/view/AtlasComputing/AtlasSonarTests>

Please provide short motivation and description of use cases (expectations on the outcome of this WG)

- to investigate weak links within full mesh and to make a diagnostics of that based on several networking tools (perfSonar, FAX, Sonar + FTS3 data including all transfers, ...)
- to provide an analysis on top of the first bullet. Understanding the slow transfers. To understand if there are weak sites or any other sources of problems.
- to define a procedures to assign the ownership and responsibility for the issue. We have to be sure that all the relevant parties are involved.
- to define the storage topology in a common way is mandatory if we want to use the same tools for all experiments (i.e closeness of the SEs between different experiments, for example can well be that SE assigned to CMS take place in the same site as different SE assigned to ATLAS). Common place where all urls to the SEs has to be stored is needed to be setup.
- cost matrix: optimize the capability of sending jobs reading over WAN and not LAN data

What systems providing network and transfer metrics already described are used in the experiment's workflow (add references if available) ? What systems are missing ?

please follow Site Status Board, network measurements view:

<http://dashb-atlas-ssb.cern.ch/dashboard/request.py/siteview#currentView=Network+measurements&highlight=false>

there are three metrics describing the Sonar+FTS3 data for every channel. The data for every channel are taken as 1 week average throughput. Sonar transfers are running exactly once per week using RUCIO to be sure that we have at least some transfer for given channel. There is also FAX metric at very right column.

- perfSonar is needed to be included and combined with existing measurements. Note that combination of perfSonar and existing metrics can provide **deeper look into potential issue** coming from SE.
- analysis tools for investigating **weak sites** from networking point of view, something similar what have been done in the past with T2D candidates.

We suggest two layers of metrics:

- metrics dedicated to each single link for all networking tools
- metrics dedicated to each site (after analysis)

→ definition of “site closeness” to the rest of the sites from networking point of view  
(simply how well it is connected to the rest of the grid)

What are the experiment’s operational use cases for correlating network and transfer metrics.  
What is currently missing, what is critical and what would be good to have.

TBA

- correlation between Sonar and perfSonar to see whether and how often similar performance is observed.
- correlation between Sonar and FAX.

Note: (JUST NOTE, TO REMOVE)

All the measurements are slightly different. The full mesh includes approx. 17 000 links. If we test each of the link we have sufficient statistic to see whether networking tools give correlated results (this can be visualised as 2d histogram with for example perfSonar on y axis and Sonar on x axis). This is particularly interesting for future decisions which metrics we need for monitoring.

Any other comments

## Experiments/CMS

Please provide short motivation and description of use cases (expectations on the outcome of this WG)

What systems providing network and transfer metrics already described are used in the experiment's workflow (add references if available) ? What systems are missing ?

What are the experiment's operational use cases for correlating network and transfer metrics. What is currently missing, what is critical and what would be good to have.

Any other comments

## Experiments/LHCb

Please provide short motivation and description of use cases (expectations on the outcome of this WG)

Debugging of data transfer issues

Checking network link quality between sites for optimizing data processing workflows and data transfer activities.

What systems providing network and transfer metrics already described are used in the experiment's workflow (add references if available) ? What systems are missing ?

LHCbDIRAC provides information about Data Transfers on metrics such as

- number of files transferred
- amount of data transferred
- "quality", i.e. amount of successful/failed transfers
- throughput

What are the experiment's operational use cases for correlating network and transfer metrics. What is currently missing, what is critical and what would be good to have.

The use case would be for debugging data transfer issues to see whether the link between two endpoints could have an impact on a seen issue. Also when commissioning storage areas to see what are the network capabilities between the new site and others would be good to have.

Any other comments

## Experiments/ALICE

Please provide short motivation and description of use cases (expectations on the outcome of this WG)

The end goal is to provide sufficient amount of monitoring data to optimize data access. Due to the data-intensive aspect of user and organised analysis the storage performance and partially the network can be limiting factors and the monitoring information should help identifying the weakest links in the chain and assist the site administrators in addressing these.

By only having handles on the end machines we usually have to guess the reason and/or iterate with all the parties involved to sort out an inefficiency visible, for example, on application level. A significant positive step would be to provide a uniform access to monitoring data of the network equipment along the data paths, including computing centres internal and WAN network. We expect this to be addressed in the WG.

What systems providing network and transfer metrics already described are used in the experiment's workflow (add references if available) ? What systems are missing ?

There are several data sources for network and transfer metrics:

- Monitoring data from `xrdcp` operations (an [ApMon](#)-instrumented version of it). The data is sent to the site-local VoBox [MonALISA](#) instance and provides the data volume and rate for the transfer along with the source and destination;
- Monitoring data reported by the Xrootd data servers themselves (*xrd.report*). Also sent to the site-local MonALISA instance where data is aggregated on several levels: network class, site name, LAN and WAN and total rates;
- Host monitoring of the Xrootd data servers used to correlate machine status with transfer patterns to identify the overloaded/underperforming data servers;
- Host monitoring of the WNs running ALICE jobs, including network interfaces. The information is not stored persistently but used on demand for debugging;
- Monitoring data from the FTD (File Transfer Daemon) service. For centrally performed transfers (typically raw data replication, storage elements migration, calibration data distribution) the 3rd party xrootd transfers are monitored in detail to assure that critical data is safely replicated. The data rates of the individual transfers are also collected.

To optimize the data transfers between computing centres, we perform regular network tests. These are coordinated from the [central](#) MonALISA repository and done between pairs of VoBoxes. They execute 1 TCP stream throughput tests (using FDT: <http://fdt.cern.ch/>) and also record the traceroute/tracepath between the machines, the RTT and the relevant kernel configuration parameters (buffer sizes, congestion algorithm). This information is used for debugging purposes during site commissioning, after upgrades or when issues are spotted. We only use 1 TCP stream as this reflects what a job reading data from remote storage element uses. Archived results are available here: <http://alimonitor.cern.ch/speed/> .

The VoBox-to-VoBox network tests also allow for [automatic](#) discovery of the network topology. Combining this with the storage monitoring (functional tests, occupancy, utilization) we build a distance metric between any client and each storage element. This function is used to dynamically point the client to the closest replica to read from, or the best several storages where to upload the results to. Since jobs are dispatched to sites that have a copy of the input data to begin with, most data reading is done locally, but

the system automatically compensates for broken file replicas and temporarily unavailable storage or portion of it.

The most data-intensive activity - the analysis jobs - are instrumented to record all file access parameters and job characteristics: logical file name, physical copy that was used, data access rate, CPU efficiency and data access rates. This information is aggregated per analysis cycle and helps identifying network and storage IO bottlenecks in correlation with the network and machine monitoring details.

What are the experiment's operational use cases for correlating network and transfer metrics.  
What is currently missing, what is critical and what would be good to have.

We currently have enough monitoring data to analyse and optimize job performances and to spot majority of the the problems on SE and WN level. What is lacking is uniform access to the infrastructure monitoring, at least at the site level where most of the problems usually occur.

Any other comments