

BCP - Campus and RON Guidelines

Performance WG BCP Task Force conference call notes

This document is essentially the LHC Tier 2 Best Common Practices draft, along with notes summarizing the discussion from the 9/4/2008 conference call of the "BCP Task Force." The notes from the call are represented as text formatted the way this section appears (red text, red box, shady background). Participants on the call: Jeff Boote, Carla Hunt, Joe Metzger, Siegrid Rickenbach, Mike Van Norman, ken lindahl.

Section heading numbers have been added, to make it easier to refer to specific parts of the document in subsequent discussion.

Some section headings are marked with a boxed number like **3**. These are an indication of the degree to which the section needs to be rewritten:

3 = major rewriting will be required

2 = medium rewrite

1 = minor rewrite

0 = perfect, no rewriting

(Please note that these indications were not determined during the conference call; they are solely ken's determinations and should be interpreted with appropriately sized grains of salt.)

Introduction **3**

This documentation is a recommendation from the US PerfSONAR community to the US LHC community that describes how the US LHC Tier-1, Tier-2 and Tier-3 centers can effectively utilize the network measurement infrastructure tools developed by the perfSONAR collaboration to debug, monitor and manage the network paths critical to their center.

The LHC computing model is a complex distributed workflow system that relies on a large number of compute, storage, and network services provided and supported by many different organizations around the globe. Many of the components of this system are new, and have never operated as production infrastructure on this scale before.

The way the LHC community is going to use the global networks is significantly different than most prior large science experiments. It is expected that this new fashion of using and stressing the research and education network infrastructure will probably bump into previously unknown problems or limitations in some parts of the global infrastructure. Simple faults where a single system fails completely until it is repaired are usually easy to diagnose and repair. However, transient faults and subtle partial failures in a system as large and complex as the LHC computing model can be very difficult to track down. Deploying a perfSONAR network measurement infrastructure in the LHC Tier 2 centers will make the network components of the workflow system more predictable and deterministic. It will make it trivial to determine if the network services are up and functioning correctly, or suffering some impairment. It should reduce the effort required to diagnose complex workflow problems, and it should allow the LHC scientists to focus their time on other parts of the computing model, or the LHC science.

Goals **2**

Allow LHC Scientists to easily:

1. Characterize and track network connectivity between their center, and the centers they serve or rely on.
2. Characterize and quantify network performance problems to accelerate diagnosing and fixing them.
3. Differentiate between application and network performance problems.
4. Differentiate between local and remote network problems.
5. Identify, understand and respond effectively to changes in the underlying network.

Use Cases **1**

This section seems pretty long; shorten? Move to a different section, e.g. appendix?

There are many use-cases for robust network measurement capabilities. We expect common use cases in the US LHC community to include:

1.
 - a.
 - i. Monitoring the set of paths important to a particular LHC Center.
 - ii. Debugging lack of performance such a slow data transfers or poor video conference capability between LHC centers
 - iii. Tracking raw network availability statistics to ensure they remain within acceptable bounds
 - iv. Monitoring the impact of LHC traffic flows on the network for capacity planning & community relations

How these use cases apply to the LHC Center management process can be illustrated in the following scenarios. A Tier 2 center manager has heard that there is a network performance problem between his center and 4 different sites. He might use the measurement infrastructure to confirm the problem, isolate it to a particular domain, and gather the diagnostic information ad described below. Note: general information about the diagnostic processes , flow charts, etc. will be in the Usage Guide section of this document.

1. He looks at the regularly scheduled latency tests and sees the following:

- 20% packet loss to Site A.
- Increased jitter to Site B.
- Step changes in the latency to Site C
- No change to Site D.

2. He investigate the performance to Site C.

1. Checks the regularly scheduled bandwidth tests.

They show no change in performance

2. Checks the interface utilization data on the path to Site C

All interfaces are running below capacity

3. So, he concludes that the network to Site C has failed over to an alternate path, but it is working correctly.

He determines that this is not a critical problem and makes a note to ask about status on the next regular conf call.

3. He investigates the performance to Site B.

1. He notes that the base latency hasn't changed significantly

So, the path probably hasn't changed either.

2. He checks the regularly scheduled bandwidth tests

They show a gradual drop over the last month

3. He checks the Interface Utilizations along the path.

They show a bottle neck link in domain X is now running at capacity.

He asks his network administrator to contact Domain X and ask about upgrade plans.

4. He investigates the performance to site A.

1. He checks the scheduled bandwidth tests

Yesterdays results are good, but todays results are showing 100 Kbps.

2. He checks the scheduled end to end latency tests.

They are showing 20% packet loss starting yesterday at 10:37 PM.

3. He checks the interface data

It is showing almost no traffic after 10:37 PM yesterday.

4. He calls his network administrator to open a ticket with their network provider.

He sends links to the latency, bandwidth and interface utilization data.

5. The upstream network provider correlates the performance change with a maintenance event on an adjacent piece of equipment and believes a fiber might have been bumped.

They send out a tech to re-seat the fiber connection which fixes the problem.

6. He investigates the performance to Site D.

1. The latency hasn't changed

2. The regularly scheduled bandwidth test results haven't changed.

3. The Utilization along the path is low.

4. The PhEDEx graphs show the many queued and failed connections He instructs the grid support team to look at application changes.

Measurements

The network measurement infrastructure will support the following network measurements. These can be used to characterize the network between points of interest, and in some cases all along the path.

General Diagnostics 1-2

Continuously measure end to end delay

What

- Manage a local star **as opposed to a mesh, but neither term is defined in the document** configuration of continuous latency measurements from a local measurement point to remote collaborators and store results in an Measurement Archive (MA)
- Publish results via a standardized web service interface

Why

- Measure & document actual availability
- Provide time references for when problems occurred and when they were fixed
- Detect & assist in diagnosing common causes of performance degradation
 - Packet Loss
 - Congestion related
 - Non-Congestion related
 - Queuing & Jitter caused by congestion
 - Routing Issues: changes, asymmetry, flapping, etc

Make regular scheduled bandwidth measurements across paths of interest

What

- Manage a local star configuration of regularly scheduled tests from a local measurement point to remote collaborators and store results in an MA
- Publish results via a standardized web service interface

Why

- Detect performance problems
- Identify when problems appeared
- Document performance delivered
- Demonstrate the performance achievable with tuned systems.

Monitor up down status of cross domain circuits

Move this section to later in the document, discuss monitoring layer 3 links first.

What

- Determine the status of a circuit
- Publish status via a web services interface

Why

- Determine when circuits are available
- Simplify debugging of end to end circuit problems

Monitor Link Circuit Capacity Utilization and Errors

Need to specify inclusion of RON-to-campus links (layer 3 and below).

What

- Publish interface utilization & error statistics via a web services interface

Why

- Allow determining available capacity to tune target data transfer rates.
- Simplify throughput problem diagnosis
- Capacity Planning

Diagnostics to look for specific known performance problems 2

What

- Make diagnostic servers available (a la NDT/NPAD)

Both NDT and NPAD use TCP. NPAD is designed for looking at short paths. Jeff: work is underway for a common scheduler, so that NDT and NPAD cannot step on each other.

Why

- T-2 sites will be accessed by end-users. These tools make it possible to diagnose paths not related to LHC T-X traffic patterns, but all the way to the scientist.
- Helps inform the scientist what performance they should expect.

Tools 1-2

There are several tools that can be used to provide the measurements listed above. The first priority in this recommendation is to ensure that diagnostic measurements can take place, but those diagnostics are most useful given historical data for comparison. Therefore, the specific tool recommendations in this list were chosen based on the ability of the tool to work in both an on-demand mode for diagnostics as well as a scheduled mode for on-going monitoring and historical analysis. Additionally, because network performance diagnostics is inherently a multi-domain effort, the perfSONAR framework is used as a medium to share the results of the measurements.

There are several domains within which these tools will be used: RONs, campuses, application communities. A "lookup service" will be needed to find perfSONAR deployments. Development is already underway; should be ready sometime soon.

Delay Measurements

diagnostics:

Two main tools are recommended in this area:

1. ping (ICMP)

Advantages:

No cooperation needed from remote site

Disadvantages:

round-trip metrics only

ICMP echos may be blocked

2. owamp

ken: how useful/important is owamp? some campus folks from CENIC community have expressed disinterest in supporting owamp, in part because of the need for an external clock. Jeff: owamp is the most useful tool for watching delay variation – indicating queuing. NDT is sufficient; an external clock is advantageous but not required.

Advantages:

one-way metrics such as jitter, hops available

direction of performance problems is isolated

Disadvantages:

Remote site must be running an OWAMP daemon (owampd)

Sites need to have reasonably synchronized clocks (NTP)

on-going monitoring:

Two tools can be used for on-going monitoring with the decision upon which depending largely on the amount of cooperation between the sites. (It is expected that sites that are serious about monitoring will implement both.)

1. pingER along with the perfSONAR-PS pingER-MA

pingER is used to set up a regular set of ping measurements and creates an archive of them. The pingER-MA is used to share that archive using perfSONAR compliant interfaces. The advantages/disadvantages from the ping diagnostics are relevant.

Many campuses (and perhaps some RONS) have already have smokeping installations, some of them fairly extensive. Guidelines should discuss integration of smokeping into the perfSONAR environment. Jeff: this is already on the list of things to be developed.

2. perfSONARBUOY

perfSONARBUOY is used to setup a regular set of owamp measurements and creates an archive of them. perfSONARBUOY exposes a perfSONAR compliant interface to share that archive. The advantages/disadvantages from the owamp diagnostics are relevant.

Bandwidth Measurements

Here, and perhaps elsewhere in the document, we need to address the IPv4 / IPv6 issues. One approach has been to use v4- and v6-specific names so that the person running the test can be assured of testing over the desired version of IP. (Also, there is an issue with the current bwctl software: the client connects to both endpoints (bwctl servers) to set up the connection. If an endpoint is v4- and v6-capable, v6 is selected, even if the other endpoint is v4-only. But that needs to be fixed independent of these guidelines.)

diagnostics: bwctl

Advantages:

- - Scheduling and authorization allow on-demand and scheduled measurements to coexist
 - well known tools are supported (iperf/nuttcp/thruly)
 - Synchronous administrator coordination not required

Disadvantages:

up down status of cross domain circuits

TBD

link capacity utilization errors

TBD

Best Practices

Introduction

This section describes the network measurement infrastructure that should be deployed at all major US LHC centers, and the set of regularly scheduled network measurements that should be made to all other US LHC centers of interest.

Measurement Infrastructure

General Guidelines

The measurement infrastructure contains multiple components that may influence each other making results analysis more difficult. The primary example is that bandwidth tests run on the same computer as continuous latency measurements will affect the latency measurement results. In order to simplify the analysis process, bandwidth measurement points and latency measurement points SHOULD NOT be deployed on the same physical machine.

Measurement points SHOULD be deployed as close to the network administrative boundaries as possible. *inside or outside of firewalls? or both?* The reason for this is to facilitate diagnosing problems using path decomposition techniques and to make the resulting data as actionable as possible.

If the measurement point is outside the firewall, some sites will not take responsibility for the security of the device, they will want someone else to manage it. This will be a problem: local ownership is really needed, otherwise the system will eventually die – e.g. AMP project.

Bandwidth Measurement Infrastructure

The bandwidth measurement infrastructure will measure achievable TCP bandwidth using memory to memory transfers over tuned TCP sessions between bandwidth measurement points.

Bandwidth measurements are useful to detect various network problems that may not affect delay measurements. Since bandwidth measurements are intrusive, they should be used with restraint as described in the schedule section below.

One of the issues that must be addressed when deploying bandwidth measurement tools is should the servers be capable of saturating the network? This issue is an area of active debate. There is a general consensus that test machines should be at least 1 Gbps capable to detect the most common problems. Some domains are actively deploying 10G capable bandwidth test systems so they can easily identify and debug problems that only appear in networks faster than 1 Gbps.

Hardware

- Bandwidth measurement infrastructure SHOULD consist of a dedicated server supporting BWCTL.
- The server SHOULD have at least a 1 Gbps interface into the local network.
- The server MAY be 10GE connected if there is a desire to diagnose and debug problems with streams faster than 1 Gbps.

Protocols

- Bandwidth measurements SHOULD be coordinated with BWCTL. <http://e2epi.internet2.edu/bwctl/>
- Bandwidth measurement tests currently SHOULD use the iperf transfer protocol. Iperf is chosen initially since more network administrators have experience with this tool. This may be extended later to include support for other tools such as thrulay, nuttcp, pathload, etc.
- Bandwidth tests SHOULD use TCP

Operations

- Scheduled measurements SHOULD NOT be started until the test system is configured and tuned so it can consistently sustain 950 Mbps or faster achieved throughput on single stream TCP tests to other test systems.
- The bandwidth measurement system SHOULD allow the following tests without requiring authentication from members of your target community:
 - Sourcing or Sinking TCP tests up to 60 seconds in duration with at least a 1 Gbps maximum bandwidth threshold.
 - Sinking UDP tests up to 60 seconds in duration at up to 50% of network uplink speed or 1Gbps, which ever is lower.

Delay Measurement Infrastructure

Delay measurements can provide very sensitive light-weight indications of many different network changes or pathologies.

- ICMP echo request and reply SHOULD be enabled to support the round trip delay measurement infrastructure.
- The delay measurement infrastructure SHOULD consist of a dedicated server synchronized to a stratum one time source.

Clock Synchronization

One way delay measurements protocols rely on the servers having both stable and accurate clocks. The protocols also require the ability to estimate the accuracy of their time synchronization. Therefore the delay measurement system must have a stable and accurate clock. The problem is that configurations tuned for stability alone are not very accurate and vice versa. The engineering compromise that MUST be maintained is as follows:

- The one way delay measurement system MUST have 4-5 stratum 1 NTP peers with as divergent of network paths as possible.
- Clock synchronization Accuracy MUST be maintained within 500 microsecond of true time and error bounds within 500 microseconds.
- Clock synchronization Accuracy SHOULD be within 5 microseconds and error estimates within 250 microseconds.

Obtaining this level of clock accuracy is not that difficult but it does require some planning. The Accuracy requirement can be achieved by synchronizing to a Stratum 1 time source such as GPS or CDMA synchronized hardware clocks or NTP synchronizing with a Stratum 1 time source over a low jitter network path. Maintaining the error bounds within the recommended range requires NTP synchronization with 4 or 5 other stratum 1 time sources over low jitter paths. This should be straight forward to achieve if most Delay Servers have their own hardware clocks, and they NTP peer with the Delay Servers that they are making regularly scheduled tests against, or a set of public clocks maintained by the community.

NDT and NPAD Measurement Infrastructure

Passive Measurement Point

Scheduled Measurements

Delay

One way delay measurements are more valuable because they essentially perform a first level path decomposition by measuring each direction unidirectionally. Therefore, one-way latency measurements SHOULD be used when ever possible, and round-trip measurements SHOULD only be used when one-way measurement infrastructure is not available.

One Way Delay Measurements

Protocols

One way delay measurements MUST be made using the OWAMP protocol. 4656

Schedule

Delay test MUST be run on a regular basis with a probe sent at least once every 5 minutes.

Delay tests SHOULD average at least 1 packet per second over a 1 minute interval.

Delay test packets between a single set of test points MUST NOT exceed 100 Kbps over a 1 minute interval

Delay test probes should consist of XXX packets 1000 Bytes in length in a ZZZ distribution.

Round Trip Delay Measurements

Protocols

Round trip delay measurements SHOULD be made using ICMP Echo Request & Echo Reply

Schedule

Delay tests MUST be run on a regular basis with a set of probe packets sent at least once every 5 minutes

Delay tests MUST NOT exceed 10 Kbps or 100 packets per second over a 10 second interval.

The time that test probes are launched within a 5 minute interval should be varied to prevent synchronization effects.

Bandwidth Measurements

Schedule

These specifications are probably appropriate for RON-to-RON testing, but may need adjustment for RON-to-campus and campus-to-campus testing, depending on what applications, and what kind of research, the campus (or pair of campuses) is supporting.

- Bandwidth tests MUST be run at least once a day.
- Bandwidth tests SHOULD be run once every 6 hours.
- Bandwidth tests SHOULD transfer data for 60 seconds.
- Bandwidth tests SHOULD NOT not transfer data for more than 120 seconds.
- Bandwidth tests between 2 endpoints SHOULD NOT run more often than once an hour.

Server Configuration

- Bandwidth tests servers MUST be connected to the local infrastructure with at least Gigabit interfaces, unless this exceeds the sites nominal uplink speed.
- Bandwidth test servers MAY be connected at 10 Gigabit speeds.
- Gigabit Ethernet attached bandwidth test servers MUST have tuned TCP stacks that can achieve at least 900 Mbps TCP throughput over 80 millisecond long round trip paths or better.
- 10 Gigabit Ethernet attached test servers SHOULD have tuned TCP stacks that can achieve at least 7 Gbps TCP throughput over 80 millisecond long round trip paths or better.
- In order to support interactive tests, the BWCTL service should be configured to allow unlimited bandwidth ad-hoc TCP tests.
 - Note that the TCP protocol ensures that this will not exceed the available bandwidth of the remote endpoint.

Passive Measurements

Interface Statistics

Utilization

Errors and Discards

Adhoc Measurements

NDT

NPAD

Bandwidth

Latency

Legal Issues

Most countries have privacy laws regarding the publication of information about people. They range from the relaxed US laws to the UK requirement that information should be accurate to the Norwegian law that says that you can't publish individually identifiable information unless you get specific permission from the individual. Every maintainer of network performance information should publish data according to the national law of the country in which the local database which holds the information resides.

In general, individually identifiable information is not required for network performance monitoring, analysis and debugging. It is recommended that organizations do not publish network performance information about interfaces, flow records, or network attributes that can be identified with a single individual.

Organizations should also consider any other legal restrictions on their network performance data before publication. For example, some commercial network provider contracts explicitly prohibit publication of network performance data. It is recommended that organizations attempt to negotiate any such terms to allow as broad of network performance data publication as possible.

Security Considerations

Considerations for Deploying Measurement Systems

Considerations for Sourcing and Sinking Active Measurements

Considerations for Publishing Measurement Results

Implementation Guide

The implementation guide section of this document describes how to deploy a network measurement infrastructure.

This section probably will need to address differing guidelines for RONS vs campuses. Again, this is likely to be highly influenced by what applications are being supported, what kinds of research are being supported.

Based on what we expect to be available at Tier-1 sites, these tools and configurations would be useful at Tier-2 sites. This should remain fairly high-level, with the expectation that we will create very detailed instructions for the LHC community accepted portions.

Setup local infrastructure so others can perform robust measurements to your site

The hardware for a typical perfSONAR installation should contain at least 2 systems, one for a bandwidth measurement point and one for the latency measurement point, so the different measurements do not affect each other. These measurement points should be placed as close to the administrative borders of the network as possible.

We anticipate 2 main deployment options. One option is to use a bootable CD with all of the tools already installed. Another option is to use a set of Red Hat Enterprise Linux 5 RPMs.

Bootable CD Installation

Insert URL to knoppix install here.

RHEL5 RPM Installation

Basic Configs

Install and configure NTP on both the latency and bandwidth measurement points first.

Get clock synchronization working within on the latency measurement point to the tolerances identified in the Best Practices section.

Get clock synchronization working on the bandwidth measurement point to within 1/2 second or better.

OWAMP

Install OWAMP on the latency measurement point system

Detailed instructions on deploying OWAMP can be found at

<http://e2epi.internet2.edu/npw/binder-docs/owamp-cookbook.pdf>

BWCTL

Install BWCTL on the bandwidth measurement point system.

Detailed instructions on deploying BWCTL can be found at

<http://e2epi.internet2.edu/npw/binder-docs/bwctl-cookbook.pdf>

PerfSONARBUOY

Install perfSONARBUOY on both measurement points to manage scheduled measurements.

Detailed instructions for deploying PerfSONARBUOY are at

<https://wiki.internet2.edu/confluence/display/PSPS/Deploying+perfSONAR-BUOY>

Pinger

Install the Pinger tools on the latency measurement system.

Detailed instructions for deploying the Pinger are at

<https://wiki.internet2.edu/confluence/display/PSPS/Deploying+perfSONAR-PS+PingER>

PerfSONAR-PS Utilization MA

1. Install the utilization MA on the latency test system (or other web services platform.)
2. Detailed instructions for deploying the utilization MA will be developed

NDT

NPAD

Etc.

Identify important collaborators

1. Organizations that provide important services to you

1. Tier 1 sites that serve important data
2. Tier 2 sites that you collaborate with
3. Cern?

2. Organizations you provide services to

1. Tier 3 sites that you service
2. Tier 2 sites that you collaborate with.

Configure Local Measurements

Once you have identified your collaborators, you need to identify which collaborators are participating by deploying local measurement infrastructure and which are not participating at this time. It is expected that all Tier 1's will be participating before the LHC goes online.

1. Participating Collaborators

1. Setup continuous latency measurements to the peers OWAMP service. (HOW DO WE DO THIS USING PERFSNAR SERVICES?)

2. Setup 1 minute bandwidth tests 4 to 6 times a day with the peers BWCTL service. (HOW DO WE DO THIS USING PERFSNAR SERVICES?)

2. Non-Participating Collaborators

1. Send a note to the remote site administrator asking if they could recommend two reliable servers that you could ping to monitor site availability.

1. Externally accessible Grid service nodes or storage server frontdoors may be good candidates.

2. Routers are typically not a good idea

3. NEED TO ADD DETAILED PING TARGET LOAD EXPECTATIONS HERE. IE How many packets per day will the default config send? How does this compare to normal background junk levels from Internet?

2. Identify 2 hosts per remote location that you are going to measure.

3. Configure a local perfSONAR Pinger system to track performance to the hosts identified

Example Configuration Files

Conclusion

It is possible to participate in the network measurement infrastructure at different levels. Your organization will get different levels of benefits depending on the level of participation.

EXPAND STRAWMAN BELOW....

Non Participant

Need a better descriptor for this class of participant.

Non-participants do not expend any effort **"do not expend any effort" is probably not the right description**, and have no control of network measurements made from remote sites into the local infrastructure.

Participating in the measurement infrastructure at this level provides you information about, and some level of control over the measurements made to your local site from remote locations.

You need to do the following to participate at this level:

use the <INSERT NAME HERE> Knoppix CD.

Or, install the following (a, b, c) from (src)

Normal Participant

(NEED A BETTER TITLE HERE...)

Participating at this level will allow a site to measure, document and understand the network characteristics between the local site and the important customers and providers, simplifying problem identification and resolution, and capacity planning.

You need to do the following to participate at this level:

...

Measurement Champion

Participants at this level are expected to assist others in their community with deploying and maintaining the measurement infrastructure.

They will host the Web visualization tools allowing inspection and analysis of the measurement data collected at their local site, as well as at other sites...

You need to do the following to participate at this level:

Install basic tools ...

Configure tests to collaborators

Install visualization tools

...

Usage Guide

The Usage Guide section of this document will describe how to use an operating network measurement infrastructure to detect, diagnose and confirm resolution of network performance problems.

Authors

Joe Metzger (for the LHC Tier2 BCP document on which this document is based.)