

GWD-I (NMWG Internal Draft)
Network Measurements Working Group
<http://www.didc.lbl.gov/NMWG/>
<http://www.didc.lbl.gov/NMWG/measurements.pdf>

Bruce Lowekamp
College of William and Mary
Brian Tierney
Lawrence Berkeley National Lab
Les Cottrell
Stanford Linear Accelerator Center
Richard Hughes-Jones
University of Manchester
Thilo Kielmann
Vrije Universiteit
Martin Swany
UCSB
July 1, 2002

A Hierarchy of Network Measurements for Grid Applications and Services

DRAFT

Status of this memo

This memo is an internal draft for GGF5 and will be replaced prior to GGF6.

This memo provides information to the Grid community regarding current practices of network measurements used for network performance. It does not define any standards or make recommendations. Distribution is unlimited.

Copyright ©Global Grid Forum (2002). All Rights Reserved.

1 Introduction

This document describes a hierarchy of network measurements useful for Grid applications and services. The goal of this work is to categorize the various types of network measurements according to the network characteristic they are measuring. This hierarchy of network characteristics will facilitate the creation of common schemas for describing network monitoring data in Grid Monitoring and Discovery Services, and thus help to address portability issues between the wide variety of network measurements used between sites of a Grid.

The hierarchy presented in this document addresses the first step of this process by providing a common dictionary of terms and relationships between commonly used measurements. The hierarchy allows measurements to be grouped according to the network characteristic they are measuring. This document is the first product of the Network Measurements Working Group (NMWG). A future NMWG document will define mappings between available tools and the measurement methodologies they implement.

The NMWG focuses on existing and currently used network measurements. It does not attempt to define new standards or to define only the best measurement methodologies to use for grid applications. It does attempt to point out the advantages and disadvantages of different measurement methodologies. The NMWG is closely related to the IETF Internet Protocol Performance Metrics (IPPM) WG, however their focus is on defining best-practices metrics of use to network engineers, whereas the NMWG and this document focus on existing practices and requirements of grid applications and grid-related tools. Where possible, we adopt the terminology defined in the IPPM Framework [21], although due to the different goals of NMWG and IPPM, certain sections of that framework do not apply to this document.

The NMWG focuses on network measurements, however many measurement tools' results are influenced by bottlenecks at the hosts making the measurements. We work to identify these effects in this

document, but a more thorough approach to such bottlenecks is the focus of the Internet2 End-to-end Performance WG [15].

2 Sample Grid use of network measurements

As an example of how network measurements could be used in a Grid environment, we use the case of a Grid file transfer service. Assume that a Grid Scheduler [12] determines that a copy of a given file needs to be copied to site *A* before a job can be run. Several copies of this file are registered in a Data Grid Replica Catalogue [5], so there is a choice of where to copy the file from. The Grid Scheduler needs to determine the optimal method to create this new file copy, and to estimate how long this file creation will take. To make this selection the scheduler must have the ability to answer these questions:

- what is the best source (or sources) to copy the data from?
- should parallel streams be used, and if so, how many?
- what TCP window and buffer size should be used?

Selecting the best source to copy the data from requires a prediction of future end-to-end path characteristics between the destination and each possible source. Accurate predictions of the performance obtainable from each source requires measurement of *available bandwidth* (both end-to-end and hop-by-hop), *latency*, *loss*, and other characteristics important to file transfer performance.

Determining whether there would be an advantage in splitting up the copy, and, for example, copying the first half of the file from site *B* and in parallel copying the second half of the file from site *C*, requires hop-by-hop link availability information for each network path. If the bottleneck hop is a hop that is shared by both paths then there is no advantage to splitting up the file copy in this way.

Parallel data streams will usually increase the total throughput on uncongested paths, as shown by Hacker et al. [14]. However, on congested links, using parallel streams may just make the problem worse. Therefore, further measurements, such as delay and loss, are needed to determine how many parallel streams to use.

Even in the case of a single stream, accurate network measurements can greatly improve performance and host resource allocation. TCP has always used a “prediction” (or smoothed estimate) of the RTT to determine timeouts. Recent work in the Net100 [19] project aims to extend this by making more information available to the TCP stack and to external programs.

2.1 How the hierarchy helps

When a distributed application is designed, the designer makes decisions about what options the application has for adapting to the network, how to make a decision between them, and what measurement tools to use. Typically, an application-level interface, such as that provided by NWS [29], provides the necessary information.

Currently there exist a variety of APIs for collecting network information from network tools and Grid Information Services, none of which are compatible. Our goal is to define a hierarchy of characteristics that can be used to provide a common classification of measurement observations taken by various systems. We believe that our hierarchy would also be well suited to defining a set of schemata for disparate measurement data that allows for its discovery and presentation by a general-purpose Grid Information Systems interface such as that described in [26], however we believe that this hierarchy is as well-suited for annotating information contained in other schemas as it is for developing a new schema.

The natural strides of research support the development of multiple network measurement systems, the development of newer measurement tools, and the cooperation of various groups to share deployed infrastructure. Even with the cooperation of various groups building Grids, there will be different monitoring systems developed. The same monitoring system may be deployed using probes with different parameters. The development of new techniques as well as different needs requiring different tools or parameters will continue to guarantee the need for many different performance monitoring systems.

The network characteristics hierarchy presented here is not an attempt to unify all measurement systems under a specific set of measurements, nor is it an attempt to define standard measurement techniques for everyone to use. Instead, the hierarchy is aimed at allowing monitoring systems to classify the measurements they take. The independent development and deployment of the monitoring systems can continue, but the classification of the measurements they take will allow that information to be used in other ways.

With the proposed hierarchy, measurements can be classified and published. Other systems familiar with the specific measurement methodology will be able to use them as intended. Systems not familiar with that particular methodology can treat them as generic measurements of that particular characteristic. By maintaining both the original measurement and a generic classification, maximum information is available according to another system's ability to interpret it.

Full use of such a hierarchy requires continued work towards the NMWG's goals, as well as integration with the work of other GGF working groups. We submit this hierarchy to the Grid and Internet communities as an agreed-upon standard of the types of measurements in use, to allow both current and future measurement methodologies to classify their observations to maximize their portability.

3 Terminology

Before discussion and definition of network measurements can begin, we first define the relevant terms. The different research backgrounds of people participating in GGF necessarily bring slightly different terminology. We will define three terms, ranging from general to specific. *Network characteristics* are the intrinsic properties of a portion of the network that are related to the performance and reliability of the Internet. *Measurement methodologies* are the means and methods of measuring those characteristics. An *observation* is an instance of the information obtained by applying the measurement methodology.

It is important to note that most network characteristics are inherently hop-by-hop values, whereas most measurement methodologies are end-to-end. Therefore what is actually being reported by the measurements may be the result for the smallest (or "bottleneck") link. In this document we attempt to distinguish between links (IP-level hop-by-hop) and paths (end-to-end).

Aside: We originally adopted the terminology specified by the IETF IPPM RFC2330 [21]. However, in actual use we found that it wasn't ideally suited for our application. We defined the distinction between metric and measurement more strongly than in the IPPM framework, due to our desire to develop a hierarchy between the various characteristics and measurements, rather than simply establishing a flat dictionary of terms. IPPM has defined multiple "metrics" where our definitions would indicate only one characteristic, or subsets of a common characteristic, for instance "loss rate" versus "loss pattern". Furthermore, discussions with IPPM contributors have indicated some questions as to what the differences between metrics and measurement methodologies are. Our definitions also led to practical difficulties, where people who were unfamiliar with our vocabulary would misinterpret our use of the term "metric." In practice, we find that characteristic, measurement methodology, and observation are rarely misinterpreted, therefore we prefer that terminology and avoid the difficulty of conflicting

uses of the word “metric.” We preserve the use of the other terminology described by RFC2330 wherever possible.

3.1 Characteristic

A characteristic is an intrinsic property that is related to the performance and reliability of a network path or link. More specifically, a characteristic is a primary property of the network, or of the traffic on it. A characteristic is the property itself, not an observation of that characteristic. An example characteristic is link capacity.

Note that a characteristic is not necessarily associated with a single number. For instance, packet loss is an important characteristic of paths and links. However, as discussed in Section 8, loss can be expressed generally as a fraction of all traffic sent, or specifically with detailed statistical properties.

3.2 Measurement methodology

A measurement methodology is a technique for recording or estimating a characteristic. Generally, there will be multiple ways to measure a given characteristic. Measurement methodologies may be either “raw” or “derived.” Raw measurement methodologies use a technique which directly produces a measurement of the characteristic, while derived measurements might be an aggregation or estimation based on a set of low-level measurements, such as using a statistical analysis of bursts of packets to estimate bandwidth capacity (e.g.: pchar and pathrate).

Consider roundtrip delay as a characteristic to be measured. Roundtrip delay may be measured directly using ping, calculated using the transmission time of a TCP packet and receipt of a corresponding ACK, projected from separate one-way delay measurements, or estimated from link propagation data and queue lengths. Each of these techniques is a separate measurement methodology for calculating the roundtrip delay characteristic.

3.3 Observations

An instance of output from a measurement methodology is an observation. An observation can be a *singleton*, which is an atomic observation, a *sample*, which is a number of singletons together, or a *statistical observation*, which is derived from a sample observation by computing a statistic on the sample.

Because network characteristics are highly dynamic, each reported observation needs to be attributed with timing information, indicating when a certain observation has been made. For singleton observations, a simple timestamp may be sufficient. For statistical observations, both the beginning and end of the observation time interval need to be reported.

3.4 Relationships between the terms

The intuitive use for these definitions is to establish a hierarchy, where the internal nodes are characteristics, and the leaf nodes are measurement methodologies. All measurement methodologies under a particular characteristic should be measuring “the same thing,” and could be used mostly interchangeably, while measurements methodologies under separate characteristics are not directly interchangeable (although may be used along with other characteristics to derive those values).

Two common characteristics are capacity and utilization. Note the relationship between these two characteristics and available bandwidth. Available bandwidth can be defined in terms of capacity and utilization (using utilization in a general sense to represent traffic utilizing the path). Using the above definitions of characteristic and measurement methodologies, available bandwidth is a characteristic because it has a

well-specified place in the hierarchy and because its measurements are not equivalent to those of any other characteristic. This definition is consistent with our intuitive sense of the hierarchy, which provides for characteristics being derived from a set of related characteristics. Note that available bandwidth is not unique in having measurement methodologies that measure it directly as well as methodologies used to estimate it from other characteristics.

To determine if a particular concept is a characteristic or a measurement methodology, the most important factor is whether the technique used to make the observation has any influence on the value itself. In particular, if there are different ways to observe identical or similar concepts, resulting in different values, then the concept may be a characteristic, but the techniques are measurement methodologies.

For example, consider the question of whether TCP-capacity (as opposed to link-capacity) is a characteristic or a measurement methodology. Although the maximum bandwidth that a TCP connection can achieve over a particular link is important to many applications, it is not truly a characteristic. In particular it is a function of the path capacity and delay, the TCP implementations used by both sending and receiving machines, and the power of the machines at each endpoint. Therefore, link-level capacity is the true characteristic, while the other factors determine how an application-level observation relates to the link capacity characteristic.

4 Statistical Representations

There are several issues in determining the statistical representations for network characteristics. Most observations are made as either packet traces or periodic samples of a particular characteristic. In either case, most applications do not use the actual sequence of measurements, but instead rely on some sort of statistical representation of those observations. This section will detail different representations that are in use.

4.1 Sampling Techniques

Following standard terminology, a single observation of a measurement is referred to as a *singleton*. Because there is little interest in single measurements, typically measurements are collected over a period of time at particular intervals. Such a series of measurements is referred to as a *sample*. A statistical measurement is a representation of a characteristic derived from a sample of measurements. As an example, consider the case of making 12 pings starting at 10:00 and recording the average. Each ping is a singleton observation, all 12 pings form the sample and the average that is recorded is a statistical measurement.

The first issue in building a representation is the sampling pattern used to collect individual observations of the characteristics. Even for packet traces, while they capture all details during the trace, for most networks it is impossible to gather those traces continuously. Therefore, for either technique the interval at which the observations are made must be specified.

A number of techniques are in common use:

- Periodic intervals, beginning each observation at a consistent interval
- Aperiodic intervals, typically distributed according to a Poisson or geometric distribution.

There are positives and negatives associated with each sampling technique, the details of which are beyond the scope of this document.

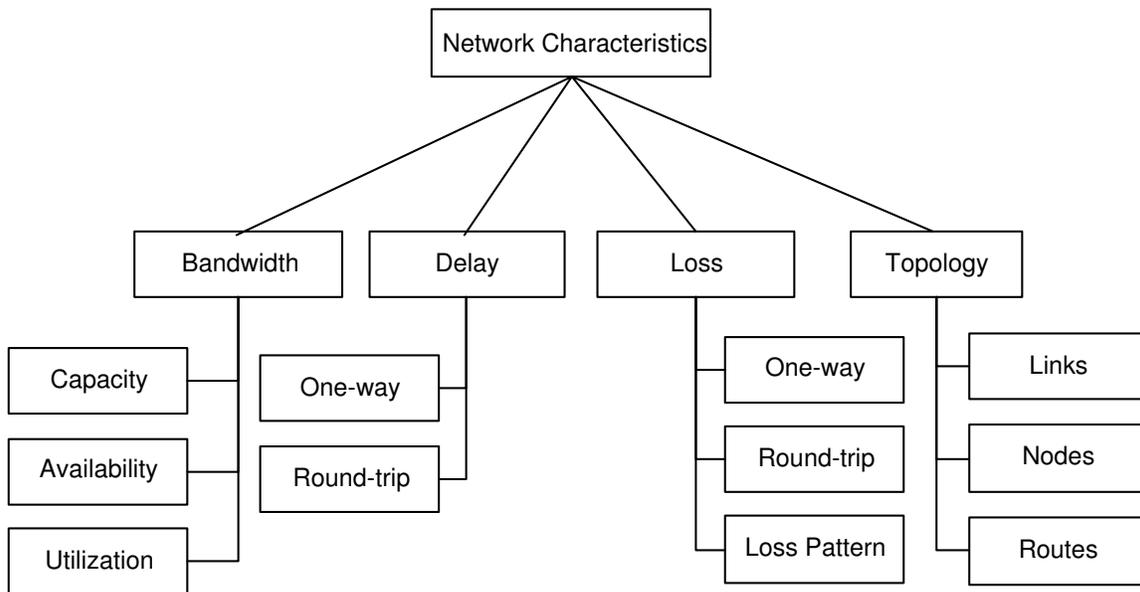


Figure 1: Hierarchy of network characteristics

4.2 Representing variability

Once the sampling technique has been determined, the next issue is how to represent the varying values observed for the characteristic at each interval. Although the simplest technique is to provide the sample in raw time series form, most customers of this information desire higher-level summary information.

Statistical representations span a wide range. Simple techniques, such as mean and variance, may capture important information, but are not necessarily appropriate for all data. More detailed representations might use quartiles or percentiles, or Box Plots [23] or histograms to represent the entire distribution of measurements. More advanced representations such as wavelets are more complex than classical statistics, but capture the temporal component associated with many network behaviors.

5 Overview of Characteristics Hierarchy

We have organized the network characteristics described in this document into a hierarchy shown in Figure 1. The following sections will discuss these characteristics, and some aspects related to their measurement, in the order presented in this figure.

6 Bandwidth Characteristics

Bandwidth is defined most generally as data per unit time. However, the “bandwidth of a link (or path)” is not a precisely defined term and specific characteristics must be clearly defined prior to discussing how to measure bandwidth.

There are three characteristics that describe bandwidth:

Capacity: The maximum amount of data that a link or path can carry

Utilization: The aggregate traffic currently on that link or path.

Available Bandwidth: The maximum throughput that the path can provide to an application, given the link or path's current load (utilization). Can be measured directly or estimated from capacity and utilization.

Each of these characteristics can be used to describe characteristics of an entire path as well as a path's hop-by-hop behavior.

6.1 Capacity

Capacity is the theoretical maximum link-layer throughput that the link or path has available, when there is no competing traffic. The capacity of the bottleneck link is the upper bound on the capacity of a path that includes that link. This information is needed to determine how to set optimal TCP buffer sizes.

6.1.1 Capacity Measurements

Link capacity can sometimes be obtained directly via SNMP queries to the network switches and routers along the path. However in general it is not possible get access to this information in a commercial ISP's network.

There are a variety of techniques for inferring capacity, most of which are based on "packet trains." Bursts of carefully-spaced datagrams can be injected into the network, and the difference between the separation of packets at the source and the destination (or other point), referred to as the dispersion, is observed. Packet dispersion techniques evaluate the dispersion of packets based on analytical models of network behavior, and estimate path characteristics based on this information. For example, the bottleneck capacity link separates the packets by the time it takes each packet to be cross that hop. If this separation is maintained through the remainder of the path, the bottleneck capacity can be derived by observing the closest spacing of the packets. Actual implementations based on these techniques apply a variety of statistical techniques to filter out the noise and random behavior of real networks.

Although they are quite useful, there are a number of challenges associated with implementing and applying these measurement techniques. They include, but aren't limited to:

- Host timing issues—As link speeds increase, intra-host latencies make a larger difference in measurements. Interrupt coalescing and driver and kernel implementations can make appreciable differences.
- Differential Queuing—There are many techniques for so-called "traffic shaping," and it is difficult to be certain that UDP or ICMP are treated the same way as TCP in the network infrastructure.

6.2 Utilization

Utilization is the aggregate capacity currently being consumed on a link or path. As a singleton observation, utilization is merely a boolean value. More useful is a statistical observation, which summarizes utilization as a fraction between 0 and 1 over a particular time interval.

More complex representations of utilization are possible. For example, a profile of traffic on a link during a particular time interval is also a way of observing utilization, however at significantly more detail than the simple proportion of link capacity in use. **It is an open question for this document how to organize traffic profile vs proportion within the characteristic hierarchy. We propose that traffic profile be represented underneath utilization, because the simpler proportion can always be derived from a more detailed representation of the traffic utilizing the link.**

6.2.1 Utilization Measurements

Utilization measurements are generally collected passively. Like capacity, utilization can sometimes be obtained directly via SNMP queries to the network switches and routers along the path. However in general it is not possible to get access to this information for a commercial ISP's network.

6.3 Available Bandwidth

Available bandwidth is the maximum IP-layer throughput that a link or path can provide to a flow or a set of flows given the current traffic load. In the context of a path consisting of several links, the link with the minimum transmission rate determines the capacity of the path, while the link with the minimum unused capacity limits the available bandwidth.

The IETF IPPM defines a TCP measurement of the available bandwidth characteristic: “Bulk Transfer Capacity” (BTC) [18]. The specific definition of the bulk transfer capacity is:

$$\text{BTC} = \text{datasent} / \text{elapsedtime}$$

This effectively tries to capture the “steady state” of a long-lived flow—amortizing out the constant of overhead.

This definition underscores a problem with our generalizations. Users tend to think of a measure *Availability* as identical to “Bulk Transfer Capacity.” The problem is how this characteristic is derived. There are many tools that measure what can be thought of as available bandwidth, although their results may vary significantly with their methodology. This is the strength of using a hierarchical arrangement in that measures of availability are grouped together by their common goal, but are differentiated closer to the leaves of the tree.

Tools to measure available bandwidth fall into two general categories:

- Flow based (or Connection-Oriented)
- Datagram based (typically packet trains)

Issues related to each of these measurement techniques are discussed below.

6.3.1 Flow Measurements

One of the most commonly used flow-oriented measurement is available TCP bandwidth. Available TCP bandwidth is both one of the most useful and yet most meaningless of all network measurements. A common method of measuring available bandwidth is to open up a TCP stream and send some data, thus simulating an application data stream. A number of tools have been developed over the years that do this, including *ttcp*, *iperf*, and so on. There are a number of problems with this technique:

- Results are greatly affected by TCP implementation. As is described in the BTC RFC, the TCP implementation on both the sending and receiving host operating systems can have a large influence on the achieved bandwidth. Any methodology that relies on a system's TCP implementation is therefore subject to its influence on its results. Furthermore, tuning the sending and receiving hosts, such as selecting the appropriate sized socket buffers, can have a profound influence on performance.
- Results are affected by the TCP slow-start algorithm. For a fairly typical high-speed link (e.g.: capacity = OC-12, RTT = 50 ms), slow start takes about 1 second. Therefore short tests are dominated by slow start, and longer tests are more intrusive. Some tools try to factor out the effect of slow start.

- TCP-based tools can be quite intrusive, putting a high load on the network. Experiments at SLAC have shown that to get a reasonable estimation of available bandwidth on a WAN using iperf requires about a 10 second test, which places a lot of unnecessary traffic on the network [7].
- Real applications are more bursty than most of these tools, and will therefore be more subject to router buffer overflows and queuing delays than tools like iperf.
- TCP-based tools only measure end-to-end bandwidth, where often one wants to have hop-by-hop information as well.

Other types of stream-oriented measurements use datagrams to simulate TCP flows or to quickly saturate the network. We make the distinction that if the aggregate behavior of the (connectionless) flow is considered, then this measurement technique is stream-oriented. This most correctly models UDP-based bulk transfer utilities and “streaming” applications as well.

So, a stream of suitably spaced UDP packets can be transmitted and the amount of data received and the time taken to receive that data measured at the destination. The spacing could be regular or follow a Poisson distribution, provided that the average generated packet rate does not exceed the transmission capability of the NIC used. If this were the case, packets could be queued or even lost in the sending Operating system, distorting the sequence presented to the network. Recording the time to transmit the data may only observe the time to move data from user space to the kernel, not the time to send it over the network.

6.3.2 Datagram Measurements

There is general consensus that packet train and packet dispersion Methodologies, described above, are well-suited for measuring path capacity. Packet dispersion methodologies are also used for measuring available bandwidth. There is, however, some question about their ability to measure availability in all situations [10], and even in the ideal situations, achieving an accurate result may require many measurements.

7 Delay Characteristics

This section draws heavily on RFC 2679 “One-way Delay Metric for IPPM,” G. Almes, S. Kalidindi, and M. Zekauskas and RFC 2681 “A Round-trip Delay Metric for IPPM,” G. Almes, S. Kalidindi and M. Zekauskas [1, 3]. For more details these references should be consulted.

As described in RFC 2679, delay is important because:

- Some applications do not perform well (or at all) if end-to-end delay between hosts is large relative to some threshold value.
- Erratic variation in delay makes it difficult (or impossible) to support many real-time applications.
- The larger the value of delay, the more difficult it is for transport-layer protocols to sustain high bandwidths.
- The minimum value of this metric provides an indication of the delay due only to propagation and transmission delay.
- The minimum value of this metric provides an indication of the delay that will likely be experienced when the path traversed is lightly loaded.
- Values of this metric above the minimum provide an indication of the congestion present in the path.

A general definition of delay, following RFCs 2330, 2679 and 2681, is the time between when the first part (e.g. the first bit) of an object (e.g. a packet) passes an observational position (e.g. where a host's network interface card connects to the wire) and the time the last part (e.g. the last bit) of that object or a related object (e.g. a response packet) passes a second (it may be the same point) observational point.

The above raises several issues including:

- How the time is synchronized if 2 observational points are used;
- Packet fragmentation issues;
- Most measurements are made by Internet hosts, which can introduce scheduling delays into the timestamps.

Delay can be measured one-way or roundtrip. One-way delays are important since today's Internet connections often use asymmetric paths, or have different quality of service in the two directions. Even symmetric paths may have different characteristics due to asymmetric queuing. Also an application may depend mostly on performance in one direction. For example the performance of a file transfer may depend more on performance in the direction in which the data flows, or in a game it may be more important to get a request to the destination before another gamer does, than it is to get the response back.

In principle the round-trip delay can be composed from the one-way delay measurements in both directions. On the other hand, it is often easier to measure round-trip delay than one-way delay since only one observation point and clock is needed. Also, many applications depend mainly on round-trip delays. RFC 2679 discusses the issues of errors and uncertainties in delay measurements related to clock accuracy and synchronization.

7.1 One-way delay measurements

One-way delay is usually measured by timestamping a packet as it enters the network and comparing that timestamp with the time the packet is received at the destination. This assumes the clocks at both ends are closely synchronized. For accurate synchronization (tens of usecs) the clocks are often synchronized with GPS. If the packet is not received at the destination within a reasonable period of time, then the one-way delay is undefined (informally, infinite according to RFC 2679), and the packet is taken to be lost.

7.2 Roundtrip delay measurements

Roundtrip delays can be composed from the individual one-way delay measurements. However, this requires making the individual measurements close to one another and being able to select the appropriate 2 measurements to add together. Since roundtrip delays are usually easier to measure than one-way delays, round-trip delays are usually measured directly.

Round-trip delay is usually measured by noting the time when the packet is sent (often this time is recorded in the packet itself), and comparing this with the time when the response packet is received back from the destination. This requires that the destination must be prepared to receive and respond (i.e. send the packet back to the source) to the test packet.

Most modern IP stacks implement an ICMP echo responder in the ICMP server [22]. Upon seeing an ICMP echo request packet, the ICMP echo responder will send a corresponding ICMP echo response packet to the sender. Thus no special server has to be installed. The ubiquitous ping tool makes use of ICMP echo and is heavily used for making round-trip delay measurements.

Another way to avoid having a server is to measure the delay between sending a TCP packet, such as a SYN packet, and then timing how long before the ACK is received, see for example [8]. The TCP stack

itself also estimates the RTT and the Web100 tool [27] allows access to various TCP RTT estimates such as minimum, maximum and smoothed RTT. It is also possible to measure the TCP RTT by passively capturing the TCP packets. The tcptrace tool [20] provides RTT reports for various passive capture tools such as tcpdump. These TCP mechanisms may be useful if, for example, ICMP echo is blocked or is suspected to be rate limited. A disadvantage of the SYN/ACK mechanism is that the frequent opening of a TCP connection via SYN packets may look like a denial of service attack.

7.3 Issues in measuring delay

7.3.1 Jitter

The variation in delay of packets as they flow from one host to another is sometimes called the “jitter.” For the characteristics hierarchy, we consider jitter to be a statistical representation of the variability in delay measurements. Other groups have defined jitter as a separate metric. While it is not a separate characteristic in this hierarchy, it is important that standards be agreed upon for common representations of variation.

As described in [9] the term jitter is used in different ways by different groups of people. The IP community therefore uses the term IP Packet Delay Variation (IPDV) defined, for a selected pair of packets in a stream of packets, as the difference in the delay of the first packet and the second of the selected packets.

Jitter, in general, is very important in sizing playout buffers for applications requiring regular delivery of packets (e.g. voice or video). Other uses include determining the dynamics of queues within a network where changes in delay variation can be linked to changes in queue length.

Given a stream of delay measurements, IPDV is easy to extract. Note that since it is the difference in delays between packet pairs, clocks do not need to be carefully synchronized. Given an IPDV probability distribution, one can also calculate statistics such as the Inter Quartile Range (IQR) to provide other estimates of jitter.

7.3.2 Measurement issues

Active measurements of delay require a probe to send probe packets and a responder to receive the probe packets and possibly return response packets. The need to have a responder/server at all the destination hosts can be major drawback to deployment. For roundtrip measurements, in many cases one can avoid this difficulty by using the ICMP echo response server built into most modern hosts.

In either case there can be problems since security may block or rate limit access to the server. Rate limiting the delay probes/responses can be tricky to determine and is sometimes suspected due to anomalous results, for example, the loss rate increases as delay packets are sent at higher frequencies. One maybe able to determine whether ICMP echo rate limiting is being imposed by comparing the ICMP echo delays with those measured using TCP RTTs as mentioned above.

As is described in detail in the RFCs referenced, clock synchronization and errors are critical for one-way delay measurements, but less so for roundtrip and IPDV measurements. Many ping implementations do not provide sub millisecond timing (or in the case of Windows sub 10 millisecond). Thus delay measurements on short links such as on a Local Area Network (LAN) are often not possible using ping.

Though the IETF delay RFCs suggest using infinity for the delay if a measurement times out, several projects ignore such packets as far as delay is concerned, and simply count them as lost. The delay RFCs also tackle the problem of duplicate and out-of-order packets. It is also useful to keep track of the occurrences of such events.

Some NIC cards coalesce interrupts to reduce the interrupt load on the CPU. This can lead to aggregations of delays since packets may be held at the NIC card, until some criteria is reached, before they are delivered to the host. In such cases it may be important to use the NIC card to do the packet timing.

8 Loss Characteristics

Along a network path, packets sent out by a sender may get lost and may, in consequence, not be received by their destination. However, the loss of a single packet often has little impact on network applications, but repeated loss can have a significant effect. Therefore, the statistical properties of packet loss events are the most interesting for determining application performance. According to RFC 2680 [2], for packet loss both a singleton metric and statistically derived quantities over time series of those singleton metrics need to be taken into account. The singleton metric observes whether or not a packet sent from a source will be received by its destination, using a given protocol, at a certain point in time. In this section, we describe one-way loss, roundtrip loss, and statistical loss properties, all three are refinements of the general concept “loss”, according to the hierarchy given in Figure 1.

This section draws heavily on RFC 2680 “A one-way Packet Loss Metric for IPPM”, G. Almes, S. Kalidindi, and M. Zekauskas. See that document for further details.

As described in RFC 2680, loss is important since:

- Some applications do not perform well (or at all) if end-to-end loss between hosts is large relative to some threshold value.
- Excessive packet loss may make it difficult to support certain real-time applications (where the precise threshold of “excessive” depends on the application).
- The larger the value of packet loss, the more difficult it is for transport-layer protocols to sustain high bandwidths.
- The sensitivity of real-time applications and of transport-layer protocols to loss become especially important when very large delay-bandwidth products must be supported.

Packet loss can impact the quality of service provided by network application programs. The sensitivity to loss of individual packets, as well as to frequency and patterns of loss among longer packet sequences is strongly dependent on the application itself. For streaming media (audio/video), packet loss results in reduced quality of sound and images. For data transfers, packet loss can cause severe degradation of sustainable bandwidth. Depending on the application itself, the above-mentioned threshold values can vary significantly.

The singleton packet loss is a binary metric. The value 0 indicates successful transmission from source to destination. The value 1 indicates a lost packet.

The **singleton, one-way packet loss** from source S to destination D at time T has the value 0 if the first bit of a packet has been sent by S to D at time T , and D has received that packet. The metric has the value 1, if D did not receive that packet.

The **singleton, roundtrip packet loss** from source S to destination D at time T has the value 0 if the first bit of a packet has been sent by S to D at time T , D has received that packet, and subsequently, a reply packet has been sent by D to S that has been received by S . The metric has the value 1, if S did not receive the reply packet.

8.1 One-way loss

In the Internet, the path from a source to a destination may be different from the path from the destination back to the source (asymmetric paths). Even in the case of symmetric paths for both directions, additional traffic from other applications may cause different queuing behaviors of the two directions.

Roundtrip measurements therefore mix the properties of both path directions, possibly producing misleading results. (For example, the sustained bandwidth of a file transfer may strongly depend on the packet loss on the path from source to destination, while being largely independent of the reverse direction.)

Packet loss depends on the following parameters:

- source and destination, defining the network path to be investigated.
- protocol and protocol parameters, packets of different protocols (TCP vs. UDP vs. ICMP) may be treated differently by the intermediate routers, possibly leading to different loss probabilities and patterns. Protocol parameters like port numbers and window sizes may further influence loss properties.

8.2 Roundtrip loss

From a network-centric viewpoint, the more useful loss metric is one-way loss, as described in RFC2680. Roundtrip loss is, in fact, the combination of two one-way loss measurements. The focus of this document is, however, on the impact of network characteristics on grid applications. For this reason, roundtrip loss becomes an important characteristic of its own. Furthermore, as measurements of round-trip loss are frequently reported, it must be represented in the hierarchy.

Roundtrip loss occurs when a sender does not receive a reply for a packet, although the network (or application-level) protocol expects such a reply. Assuming statistical independence of packet loss in both directions of a network path, roundtrip loss properties may be derived from the individual metrics for each direction. However, for applications that use request/reply protocols (like RPC or HTTP), traffic in both directions is not independent, so roundtrip loss is more precisely measured directly.

8.3 Loss patterns (statistical properties)

Currently, the Internet community primarily uses only one statistically derived quantity for loss, namely the *loss average*, defined as the average of the singleton loss values over a series of sent packets. The loss average corresponds to the *loss rate* which is given as a percentage between 0 % and 100 %.

Other statistically derived quantities, like loss burstiness patterns, loss free seconds, or conditional loss probability [4], are of interest to specific transport protocols as well. Techniques are being developed to measure these types of loss, and should be represented under loss pattern, or possibly other characteristics.

[more work is needed here](#)

8.4 Issues in measuring loss

The general problem with measurements of packet loss (as with other properties) is that the measurements need to be as close as possible to application data transfer, in order to produce significant results. Since packets of different protocols are supposed to be forwarded differently by routers in the Internet, loss properties need to be derived for each relevant protocol separately.

- **ICMP (ping)**

The *ping* utility uses ICMP *ECHO_REQUEST* / *ECHO_REPLY* packets to determine both roundtrip delay and roundtrip loss. *ping* reports the roundtrip loss rate. Typical implementations of *ping* send only one packet per second, possibly causing fewer loss events than application-data transfers which typically send bursts of back-to-back packets.

Furthermore, ping's default packet size is much smaller than the MTU of existing networks. Such small packets sometimes cause much smaller loss rates than the longer packets which are typical for application data transfers.

- **TCP**

Loss information for TCP as the predominant transport protocol is the most important for applications. However, TCP packet loss can only be measured at the kernel level of a protocol stack. Retrieving this information is thus dependent on the interface provided by the operating system. TCP roundtrip loss might also be measured by SYN/ACK packet pairs. However, as these packets are used in the special case of connection establishment, the significance of the collected loss data is questionable for TCP data transfers.

The *sting* tool [24] can accurately measure loss rates, separately for both directions of a TCP connection. Sting re-implements TCP in user space to retrieve information about lost or duplicate acknowledgments in order to derive which packets have been lost. However, the implementation relies on cooperation of the operating system kernel, limiting sting's applicability and portability.

- **UDP**

UDP lends itself for application-level loss measurements, both one-way and roundtrip. In such measurements, a series of packets containing timestamps are sent from source to destination. The destination records for each packet the event of receiving. The receiver needs to define a timeout value for each packet after which the packet is supposed to be lost. Besides the measured UDP packets, both parties need to communicate across a (TCP) channel with enforced packet delivery to exchange control information, e.g., for determining a suitable timeout value.

9 Topology Characteristics

There are two approaches to characterizing network topology: physical and functional. The physical approach is to determine the physical links that connect the network together. By determining the connections between links, along with their capacities, queuing algorithms, and traffic load, the network can be modeled and its behavior analyzed or predicted. The functional approach differs in that it makes use of end-to-end information, under the assumption that such observations are more readily available and usable than modeling low-level network behavior. Functional topology representations attempt to group and arrange network sites according to their perceived closeness determined by traffic performance, rather than according to the actual connections of physical links. There are advantages and disadvantages to each approach, and combinations of the two approaches are used. We will not attempt to evaluate the merits of each approach here.

The topology characteristics are chosen so that both approaches share part of the same tree, while being distinguished from one another. In this way, an application could be written to make use of information from either technique, or to deliberately require a specific technique.

The topology hierarchy is shown in Figure 2. The top layer divides topology characteristics into links, nodes, and routes. Links are the edges used to build a topology graph. A system describing a set of links should provide links to form a connected graph that other systems can process. Routes are descriptions of end-to-end paths used by data. They may be, but are not necessarily associated with links used to form a topology graph.

9.1 Links

The building block of a topology graph are the links between nodes. Links are unidirectional. Each link identifies a source and destination, which can be endpoints, routers, switches, or autonomous systems.

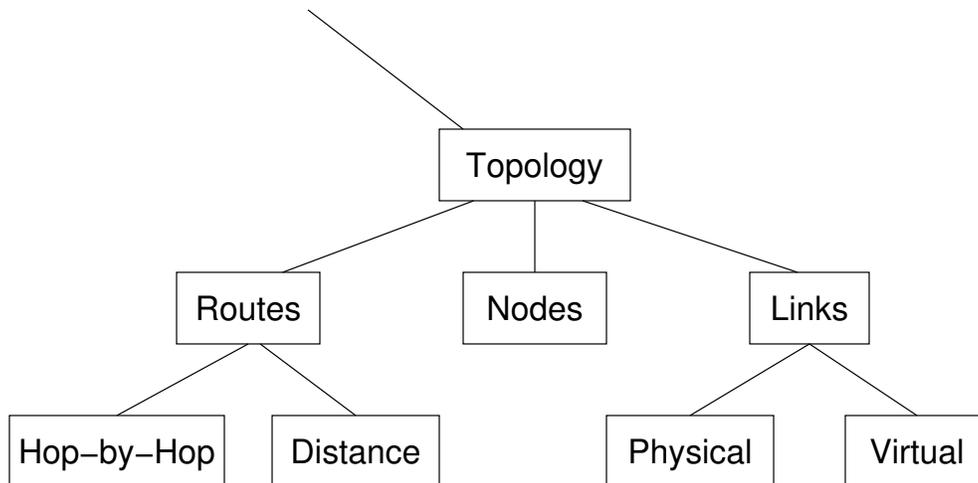


Figure 2: Hierarchy of topology characteristics, refining Fig. 1

Links are further subdivided into physical links, intended to represent actual connections between nodes, and virtual links, which represent a derived view of behavior from end-to-end application measurements.

A link may be annotated with a “Valid-From” field, indicating the node or domain for which it is meaningful. This information might be used for a single-source network view built of virtual links or when describing a physical link that exists only in a particular VPN.

Bandwidth and delay information for the link may also be available with the same source and destination. It may be beneficial to associate this information directly in the link information.

9.1.1 Physical

Links are reported at the finest level available to the reporting system, i.e. links should not be aggregated before being reported as links. However, not all observers are capable of resolving networks at the same resolution. For example, many observers may only have partial knowledge, such as observing a single link between routers when there is really an ATM cloud between them. Furthermore, fluctuating routes can quickly invalidate any observed topology. As with all network information, the characteristics described here should be regarded as best-effort.

The primary distinguishing factor of a physical link is not that it is guaranteed to indicate a single electrical or optical connection between switches, but that it is part of a graph that indicates how network traffic actually moves or may move between nodes of the network.

Because routers use different IP addresses for each port, and those IP addresses typically resolve to different names, it is frequently difficult to identify the individual routers, in particular when using several tools. For example, traceroute identifies only the receiving interface of each router. Measurements based on traceroute, therefore, do not provide a useful source-destination pair, but only a series of destination interfaces, which each indicate a single link. Assuming the router’s switching fabric is not itself a bottleneck, this detail may be unimportant. However, if traceroute measurements are being combined with other tools, difficulty may be experienced combining the information because different names or IP addresses may be used to refer to the same router. Resolving these issues is beyond the scope of describing the characteristics, but it is worth noting that comparing destinations of links may be more useful than sources due to the way traceroute works [13].

9.1.2 Virtual

Not all topology graphs are designed to represent the physical connections between nodes or the flow of traffic across the network. Some are designed by measuring bandwidth between sites and identifying groups that appear to be well connected. This approach may be taken across a variety of sites distributed around the Internet, or using a single-source tree [6, 16, 25, 28].

Virtual links should be used when a system has determined a relationship between nodes that is represented by a graph, but where there is not a direct attempt to correlate the graph to actual network connections.

9.2 Nodes

In wired networks, queue overflow in routers is the predominant reason for packet loss. To model network behavior, analytical models require information about the static and dynamic properties of the queues. Information about queues is described under this characteristic.

Queue information can be either raw or derived. Raw information, such as that obtained by SNMP, can indicate the precise drop rate, length, and queuing algorithm used by the router. Derived information can be obtained by sending traffic through the router and deducing its behavior through analytical means [17]. However the information is obtained, it does not effect the characteristic it attempts to determine, although measurements should indicate how the information was derived.

A full list of the information to be stored here needs to be determined. It should include static drop algorithm, capacity, etc, as well as dynamic length and drop information. Information must be maintained on a per-port basis, and needs to be connected to the link information in some way. IP address is OK for IP-level routers, but switches are more complex...

9.3 Routes

A route describes the path followed by traffic between two endpoints. Routes can either be described using a set of links from a particular topology, or they can be described using arbitrary cost characteristics. Two types of routes are supported, hop-by-hop and distance routes. Hop-by-hop routes indicate the links of a topology graph followed by traffic between two endpoints. Information about the path may be derived from information about the links. Distance routes, on the other hand, are not necessarily based on links, but may be based on a combination of other characteristics that are measured or calculated separately regarding the end-to-end behavior of the network between the particular two end-routes.

A hop-by-hop route can consist of either physical or virtual links (possibly both, although we doubt that a single system would generate such a route). A single system may generate both a hop-by-hop route and a distance route, if it produces both hop-by-hop information, as well as end-to-end information that is not well-specified by the links in the topology or by the other characteristics.

9.3.1 Hop-by-hop

A route is a sequence of *links* used by traffic from a source to a destination (unidirectional). The set of all routes does not necessarily utilize every link in a topology, and a valid or shortest path across links in the topology does not indicate that such a path is used for a route. For example, many universities and research labs have connections to multiple WANs, generally at least one commercial and one research network. Although a local company connected to the commercial WAN may have a one-hop path to the university's research network, the service contracts do not allow their traffic to be routed to the research network.

In general, the only way to determine routes based on physical links is to run traceroute or examine the routing tables directly. In the event that routes are oscillating, or traffic is split across different links, there

may be multiple valid routes between the same source-destination pair. Frequently this behavior is observed through traceroutes collected at periodic intervals, therefore a time series of routes can be just as interesting as a time series of available bandwidth measurements.

Most of the difficulties with determining routes are restricted to routes over physical links. Routes over virtual links [6, 16, 25, 28] are not generally subject to the same issues because the links themselves are determined by the overall end-to-end behavior of the system.

9.3.2 Distance

Distance route characteristics are derived about specific routes, and frequently used to calculate the topology represented by virtual links. A single available bandwidth or latency between two endpoint would be better represented through one of those particular characteristics, however for systems that make use of a derived characteristic using both of those values, that quantity is best represented here.

There are an infinite number of distance characteristics below this node in the hierarchy. All such characteristics meet the requirement that they produce a single quantity that is indicative of the “distance” or “closeness” of two nodes in a network [11]. For example, closeness is a function combining available bandwidth with roundtrip delay. In principle, an application could make a choice between two data providers by using any single distance characteristic that measures each possibility. In practice, it might be better to chose a particular characteristic that weights bandwidth, latency, long transfers, or short transfers according to the application’s need.

10 Full Copyright Notice

Copyright (C) Global Grid Forum (2002). All Rights Reserved.

This document and translations of it may be copied and furnished to others, and derivative works that comment on or otherwise explain it or assist in its implementation may be prepared, copied, published and distributed, in whole or in part, without restriction of any kind, provided that the above copyright notice and this paragraph are included on all such copies and derivative works. However, this document itself may not be modified in any way, such as by removing the copyright notice or references to the GGF or other organizations, except as needed for the purpose of developing Grid Recommendations in which case the procedures for copyrights defined in the GGF Document process must be followed, or as required to translate it into languages other than English.

The limited permissions granted above are perpetual and will not be revoked by the GGF or its successors or assigns.

This document and the information contained herein is provided on an “AS IS” basis and THE GLOBAL GRID FORUM DISCLAIMS ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

11 Intellectual Property Rights

The GGF takes no position regarding the validity or scope of any intellectual property or other rights that might be claimed to pertain to the implementation or use of the technology described in this document or the extent to which any license under such rights might or might not be available; neither does it represent that it has made any effort to identify any such rights. Copies of claims of rights made available for publication and any assurances of licenses to be made available, or the result of an attempt made to obtain a general

license or permission for the use of such proprietary rights by implementers or users of this specification can be obtained from the GGF Secretariat.

The GGF invites any interested party to bring to its attention any copyrights, patents or patent applications, or other proprietary rights which may cover technology that may be required to practice this recommendation. Please address the information to the GGF Executive Director.

References

- [1] G Almes, S Kalidindi, and M Zekauskas. A one-way delay metric for IPPM. RFC2679, September 1999.
- [2] G Almes, S Kalidindi, and M Zekauskas. A one-way packet loss metric for IPPM. RFC2680, September 1999.
- [3] G Almes, S Kalidindi, and M Zekauskas. A round-trip delay metric for IPPM. RFC2681, September 1999.
- [4] J. C. Bolot. Characterizing end-to-end packet delay and loss in the internet. *High-Speed Networks*, 2(3):305–323, 1993.
- [5] Ann Chervenak and GGF: Data Replication Research Group. An architecture for replica management in grid computing environments. <http://www.sdsc.edu/GridForum/RemoteData/Papers/ggf1replica.pdf>.
- [6] Mark Coates, A O Hero III, Robert Nowak, and Bin Yu. Internet tomography. *IEEE Signal Processing Magazine*, 19(3):47–65, May 2002.
- [7] R. L. Cottrell and Connie Logg. A new high performance network and application monitoring infrastructure. Technical Report SLAC-PUB-9202, SLAC, 2002.
- [8] R. L. Cottrell and M. Shah. Measuring rtt by using syn/acks instead of pings. <http://www-iepm.slac.stanford.edu/monitoring/limit/limiting.html#synack>, December 1999.
- [9] C. Demichelis and P. Chimento. Ip packet delay variation metric for ippm. <http://www.ietf.org/internet-drafts/draft-ietf-ippm-ipdv-09.txt>, April 2002.
- [10] Constantinos Dovrolis, Parameswaran Ramanathan, and David Moore. What do packet dispersion techniques measure? In *IEEE INFOCOM 2001*, 2001.
- [11] Tiziana Ferrari and Francesco Giacomini. Network monitoring for grid performance optimization. *Computer Communications*, 2002. submitted for publication.
- [12] GGF. Grid scheduling area. <http://www.mcs.anl.gov/~schopf/ggf-sched>.
- [13] Ramesh Govindan and Hongsuda Tangmunarunkit. Heuristics for internet map discovery. In *IEEE INFOCOM 2000*, Tel Aviv, Israel, March 2000.
- [14] T. J. Hacker and B. D. Athey. The end-to-end performance effects of parallel tcp sockets on a lossy wide-area network.
- [15] Internet2 end-to-end performance initiative. <http://www.internet2.edu/e2eperf/>.

- [16] Sugih Jamin, Cheng Jin, Yixin Jin, Danny Raz, Yuval Shavitt, and Lixia Zhang. On the placement of internet instrumentation. In *IEEE INFOCOM 2000*, Tel Aviv, Israel, March 2000.
- [17] Jun Liu and Mark Crovella. Using loss pairs to discover network properties. In *ACM SIGCOMM Internet Measurement Workshop 2001 (IMW2001)*, November 2001.
- [18] M. Mathis and M. Allman. A framework for defining empirical bulk transfer capacity metrics. RFC3148, July 2001.
- [19] Net100. <http://www.net100.org>.
- [20] S. Ostermann. tcptrace. <http://irg.cs.ohiou.edu/software/tcptrace/tcptrace.html>.
- [21] V. Paxson, G. Almes, J. Mahdavi, and M. Mathis. Framework for IP performance metrics. RFC2330, May 1998.
- [22] J. Postel. Internet control message protocol. RFC792, September 1981.
- [23] J. W. Tukey R. McGill and W. A. Larsen. Variations of box plots. *The American Statistician*, 32(1):12–16, February 1978.
- [24] Stefan Savage. Sting: a tcp-based network measurement tool. In *USENIX Symposium on Internet Technologies and Systems*, pages 71–79, Boulder, CO, October 1999.
- [25] Gary Shao, Fran Berman, and Rich Wolski. Using effective network views to promote distributed application performance. In *Proceedings of the 1999 International Conference on Parallel and Distributed Processing Techniques and Applications (PDPTA'99)*, 1999.
- [26] M. Swamy and R. Wolski. Representing dynamic performance information in grid environments with the network weather service. 2nd IEEE International Symposium on Cluster Computing and the Grid (to appear), May 2002.
- [27] Web100 Project team. Web100 project. <http://www.web100.org>.
- [28] Wolfgang Theilmann and Kurt Rothermel. Dynamic distance maps of the internet. In *IEEE INFOCOM 2000*, Tel Aviv, Israel, March 2000.
- [29] R. Wolski. Dynamically forecasting network performance using the network weather service. *Cluster Computing*, 1:119–132, January 1998.