Internet Engineering Task Force                             J. Rapp
Internet-Draft                                            L. Avramov
Intended status: Informational                     Cisco Systems, Inc
Expires: December 6, 2013                                 June 4, 2013

           Definitions and Metrics for Data Center Benchmarking
                          draft-dcbench-def-00

The purpose of this informational document is to establish
   definitions, discussion and measurement techniques for data center
   benchmarking. Also, it is to introduce new terminologies applicable
   to data center performance evaluations. The purpose of this document
   is not to define the test methodology, but rather establish the
   important concepts when one is interested in benchmarking network
   equipment in the data center.

Status of this Memo

This Internet-Draft is submitted in full conformance with the
   provisions of BCP 78 and BCP 79.

   Internet-Drafts are working documents of the Internet Engineering
   Task Force (IETF).  Note that other groups may also distribute
   working documents as Internet-Drafts.  The list of current Internet-
   Drafts is at http://datatracker.ietf.org/drafts/current/.

   Internet-Drafts are draft documents valid for a maximum of six months
   and may be updated, replaced, or obsoleted by other documents at any
   time.  It is inappropriate to use Internet-Drafts as reference
   material or to cite them other than as "work in progress."

   This Internet-Draft will expire on December 6, 2013.

Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the
   document authors.  All rights reserved.

   This document is subject to BCP 78 and the IETF Trust's Legal
   Provisions Relating to IETF Documents
   (http://trustee.ietf.org/license-info) in effect on the date of
   publication of this document.  Please review these documents

Table of Contents

1.  Introduction

    Traffic patterns in the data center are not uniform and are contently
    changing. They are dictated by the nature and variety of applications
    utilized in the data center. It can be largely east-west traffic
    flows in one data center and north-south in another, while some may
    combine both. Traffic patterns can be bursty in nature and contain
    many-to-one, many-to-many, or one-to-many flows. Each flow may also
    be small and latency sensitive or large and throughput sensitive
    while containing a mix of UDP and TCP traffic. All of which can
    coexist in a single cluster and flow through a single network device
    all at the same time. Benchmarking of network devices have long used
    RFC1242, RFC2432, RFC2544, RFC2889 and RFC3918. These benchmarks have
    largely been focused around various latency attributes and max
    throughput of the Device Under Test being benchmarked. These
    standards are good at measuring theoretical max throughput,
    forwarding rates and latency under testing conditions, but to not
    represent real traffic patterns that may affect these networking
    devices.


    The following defines a set of definitions, metrics and terminologies
    including congestion scenarios, switch buffer analysis and redefines
    basic definitions in order to represent a wide mix of traffic
    conditions.

## 1.1.  Requirements Language

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in RFC 2119 [6].

## 1.2. Definition format

Term to be defined. (e.g., Latency)

Definition: The specific definition for the term.

Discussion: A brief discussion about the term, it's application and any restrictions on measurement procedures.

Measurement: Methodology for the measure and units used to report measurements of this term, if applicable.

## 2.  Latency

## 2.1. Definition

Latency is a the amount of time it takes a frame to transit the DUT.

Latency can be measured with the following methods, irrespectively of the type of switching device (bit forwarding aka cut-through or store forward type of device)

FILO (First In Last Out) The time interval starting when the end of the first bit of the input frame reaches the input port and ending when the last bit of the output frame is seen on the output port

FIFO (First In First Out) The time interval starting when the end of the first bit of the input frame reaches the input port and ending when the start of the first bit of the output frame is seen on the output port

LILO (Last In Last Out) The time interval starting when the last bit of the input frame reaches the input port and the last bit of the output frame is seen on the output port

LIFO (Last In First Out) The time interval starting when the last bit of the input frame reaches the input port and ending when the first bit of the output frame is seen on the output port.

This definition replaces the previous definition of Latency defined in RFC 1242, section 3.8 and is quoted here:

For store and forward devices: The time interval starting when the last bit of the input frame reaches the input port and ending when the first bit of the output frame is seen on the output port.

For bit forwarding devices: The time interval starting when the end of the first bit of the input frame reaches the input port and ending when the start of the first bit of the output frame is seen on the output port.

2.2 Discussion

FILO is the most important measuring method. Any type of switches MUST be measured with the FILO mechanism: FILO will include the latency of the switch and the latency of the frame as well as the serialization delay. It is a picture of the 'whole' latency going through the DUT. For applications, which are latency sensitive and can function with initial bytes of the frame, FIFO MAY be an additional type of measuring to supplement FILO.

LIFO mechanism can be used with store forward type of switches but not with cut-through type of switches, as it will provide negative latency values for larger packet sizes. Therefore this mechanism MUST NOT be used when comparing latencies of two different DUTs.

2.3 Measurement

The measuring methods to use for benchmarking purposes are as follow:

1) FILO MUST be used as a measuring method, as this will include the latency of the packet; and today the application commonly need to read the whole packet to process the information and take an action.

2) FIFO MAY be used for certain applications able to proceed data as the first bits arrive (FPGA for example)

3) LIFO MUST not be used, because it subtracts the latency of the packet; unlike all the other methods.


3 Jitter

3.1 Definition

The definition of Jitter is covered extensively in RFC 3393. This definition is not meant to replace that definition, but it is meant

to provide guidance of use for data center network devices.

The use of Jitter is in according with the variation delay definition from RFC 3393:

The second meaning has to do with the variation of a metric (e.g., delay) with respect to some reference metric (e.g., average delay or minimum delay). This meaning is frequently used by computer scientists and frequently (but not always) refers to variation in delay.

## 3.2 Discussion

Jitter can be measured in different scenarios:-packet to packet delay variation-delta between min and max packet delay variation for all packets sent.

## 3.3 Measurement

The jitter MUST be measured when sending packets of the same size. Jitter MUST be measured as packet to packet delay variation and delta between min and max packet delay variation of all packets sent. A histogram MAY be provided as a population of packets measured per latency or latency buckets.


## 4 Physical Layer Calibration

## 4.1 Definition

The calibration of the physical layer consists of defining and measuring the latency of the physical devices used to perform test on the DUT.

It includes the list of all physical layer components used as listed here after:

-type of device used to generate traffic / measure traffic

-type of line cards used on the traffic generator

-type of transceivers on traffic generator

-type of transceivers on DUT

-type of cables

-length of cables

-software name, and version of traffic generator and DUT

-list of enabled features on DUT MAY be provided and is recommended [especially the control plane protocols such as LLDP, Spanning-Tree etc.]. A comprehensive configuration file MAY be provided to this effect.

## 4.2 Discussion

Physical layer calibration is part of the end to end latency, which should be taken into acknowledgment while evaluating the DUT. Small variations of the physical components of the test may impact the latency being measure so they MUST be described when presenting results.

## 4.3 Measurement

It is RECOMMENDED to use all cables of : the same type, the same length, when possible using the same vendor. It is a MUST to document the cables specifications on section [4.1s] along with the test results. The test report MUST specify if the cable latency has been removed from the test measures or not. The accuracy of the traffic generator measure MUST be provided [this is usually a value in the 20ns range for current test equipments].

## 5 Line rate

## 5.1 Definition

The transmit timing, or maximum transmitted data rate is controlled by the "transmit clock" in the DUT.  The receive timing (maximum ingress data rate) is derived from the transmit clock of the connected interface.

The line rate or physical layer frame rate is the maximum capacity to send frames of a specific size at the transmit clock frequency of the DUT.

The frequency ("clock rate") of the transmit clock in any two connected interfaces will never be precisely the same, therefore a tolerance is needed, this will be expressed by Parts Per Million (PPM) value. The IEEE standards allow a specific +/- variance in the transmit clock rate, and Ethernet is designed to allow for small, normal variations between the two clock rates. This results in a tolerance of the line rate value when traffic is generated from a testing equipment to a DUT.

5.2 Discussion

For a transmit clock source, most Ethernet switches use "clock modules" (also called "oscillator modules") that are sealed, internally temperature-compensated, and very accurate. The output frequency of these modules is not adjustable because it is not necessary.  Many test sets, however, offer a software-controlled adjustment of the transmit clock rate, which should be used to compensate the test equipment to not send more than line rate of the DUT.

To allow for the minor variations typically found in the clock rate of commercially-available clock modules and other crystal-based oscillators, Ethernet standards specify the maximum transmit clock rate variation to be not more than +/- 100 PPM (parts per million) from a calculated center frequency. Therefore a DUT must be able to accept frames at a rate within +/- 100 PPM to comply with the standards.

Very few clock circuits are precisely +/- 0.0 PPM because:

1.The Ethernet standards allow a maximum of +/- 100 PPM (parts per million) variance over time. Therefore it is normal for the frequency of the oscillator circuits to experience variation over time and over a wide temperature range, among external factors.

2.The crystals or clock modules, usually have a specific  +/- PPM variance that is significantly better than +/- 100 PPM. Often times this is +/- 30 PPM or better in order to be considered a "certification instrument".

When testing an Ethernet switch throughput at "line rate", any specific switch will have a clock rate variance. If a test set is running +1 PPM faster than a switch under test, and a sustained line rate test is performed,  a gradual increase in latency and eventually packet drops as buffers fill and overflow in the switch can be observed. Depending on how much clock variance there is between the two connected systems, the effect may be seen after the traffic stream has been running for a few hundred microseconds, a few milliseconds, or seconds. The same low latency and no-packet-loss can be demonstrated by setting the test set link occupancy to slightly less than 100 percent link occupancy. Typically 99 percent link occupancy produces excellent low-latency and no packet loss. No Ethernet switch or router will have a transmit clock rate of exactly +/- 0.0 PPM. Very few (if any) test sets have a clock rate that is precisely +/- 0.0 PPM.

Test set equipment manufacturers are well-aware of the standards, and

allows a software-controlled +/- 100 PPM "offset" (clock-rate adjustment) to compensate for normal variations in the clock speed of "devices under test". This offset adjustment allows engineers to determine the approximate speed the connected device is operating, and verify that it is within parameters allowed by standards.

5.3 Measurement

"Line Rate" CAN be measured in terms of "Frame Rate":

Frame Rate = Transmit-Clock-Frequency / (Frame-Length*8 + Minimum_Gap + Preamble + Start-Frame Delimiter)

Example for 1 GB Ethernet speed with 64-byte frames: Frame Rate = 1,000,000,000 /(64*8 + 96 + 56 + 8) Frame Rate = 1,000,000,000 / 672 Frame Rate = 1,488,095.2 frames per second.

Considering the allowance of +/- 100 PPM, a switch may "legally" transmit traffic at a frame rate between 1,487,946.4 FPS and 1,488,244 FPS.  Each 1 PPM variation in clock rate will translate to a 1.488 frame-per-second frame rate increase or decrease.

In a production network, it is very unlikely to see precise line rate over a very brief period. There is no observable difference between dropping packets at 99% of line rate and 100% of line rate.

-Line rate CAN measured at 100% of line rate with a -100PPM adjustment.

-Line rate SHOULD be measured at 99,98% with 0 PPM adjustment.

6   Buffering

6.1 Buffer

6.1.1 Definition

Buffer Size: the term buffer size, represents the total amount of frame buffering memory available on a DUT. This size is expressed in Byte; KB (kilobytes), MB (megabytes) or GB (gigabyte). When the buffer size is expressed it SHOULD be defined by a size metric defined above. When the buffer size is expressed, an indication of the frame MTU used for that measurement is also necessary as well as the cos or dscp value set; as often times the buffers are carved by quality of service implementation.

Example: Buffer Size of DUT when sending 1518 bytes frames is 18 Mb.

Port Buffer Size: the port buffer size is the amount of buffer a single ingress port, egress port or combination of ingress and egress buffering location for a single port. The reason of mentioning the three locations for the port buffer is, that the DUT buffering scheme can be unknown or untested, and therefore the indication of where the buffer is located helps understand the buffer architecture and therefore the total buffer size. The Port Buffer Size is an informational value that MAY be provided from the DUT vendor. It is not a value that is tested by benchmarking. Benchmarking will be done using the Maximum Port Buffer Size or Maximum Buffer Size methodology.

Maximum Port Buffer Size: this is in most cases the same as the Port Buffer Size. In certain switch architecture called SoC (switch on chip), there is a concept of port buffer and shared buffer pool available for all ports. Maximum Port Buffer, defines the scenario of a SoC buffer, where this amount in B (byte), KB (kilobyte), MB (megabyte) or GB (gigabyte) would represent the sum of the port buffer along with the maximum value of shared buffer this given port can take. The Maximum Port Buffer Size needs to be expressed along with the frame MTU used for the measurement and the cos or dscp bit value set for the test.

Example: a DUT has been measured to have 3KB of port buffer for 1518 frame size packets and a total of 4.7 MB of maximum port buffer for 1518 frame size packets and a cos of 0.

Maximum DUT Buffer Size: this is the total size of Buffer a DUT can be measured to have. It is most likely different than the Maximum Port Buffer Size. It CAN also be different from the sum of Maximum Port Buffer Size. The Maximum Buffer Size needs to be expressed along with the frame MTU used for the measurement and along with the cos or dscp value set during the test.

Example: a DUT has been measured to have 3KB of port buffer for 1518 frame size packets and a total of 4.7 MB of maximum port buffer for 1518 frame size packets. The DUT has a Maximum Buffer Size of 18 MB at 1500 bytes and a cos of 0.

Burst: The burst is a fixed number of packets sent over a percentage of linerate of a defined port speed. The amount of frames sent are evenly distributed across the interval T. A constant C, can be defined to provide the average time between two consecutive packets evenly spaced.

Microburst: it is a burst. A microburst is when packet drops occur

when there is not sustained or noticeable congestion upon a link or device. A characterization of microburst is when the Burst is not evenly distributed over T, and is less than the constant C [C= average time between two consecutive packets evenly spaced out].

Intensity of Microburst: this is a percentage, representing the level of microburst between 1 and 100%. The higher the number the higher the microburst is. $I=[1-[ (TP2-Tp1)+(Tp3-Tp2)+....(TpN-Tp(n-1) ] / Sum(packets)]]*100$

6.1.3 Discussion

When measuring buffering on a DUT, it is important to understand what the behavior is for each port, and also for all ports as this will provide an evidence of the total amount of buffering available on the switch. The terms of buffer efficiency here helps one understand what is the optimum packet size for the buffer to be used, or what is the real volume of buffer available for a specific packet size. This section does not discuss how to conduct the test methodology, it rather explains the buffer definitions and what metrics should be provided for a comprehensive data center device buffering benchmarking.

6.1.3 Measurement

When Buffer is measured:
-the buffer size MUST be measured
-the port buffer size MAY be provided for each port
-the maximum port buffer size MUST be measured
-the maximum DUT buffer size MUST be measured
-the intensity of microburst MAY be mentioned when a microburst test is performed
-the cos or dscp value set during the test SHOULD be provided


6.2 Incast
6.2.1 Definition

The term Incast, very commonly utilized in the data center, refers to the traffic pattern of many-to-one or many-to-many conversations. Typically in the data center it would refer to many different ingress server ports(many), sending traffic to a common uplink (one), or multiple uplinks (many). This pattern is generalized for any network as many incoming ports sending traffic to one or few uplinks. It can also be found in many-to-many traffic patterns.

Synchronous arrival time: When two, or more, frames of respective sizes L1 and L2 arrive at their respective one or multiple ingress

ports, and there is an overlap of the arrival time for any of the
bits on the DUT, then the frames L1 and L2 have a synchronous arrival
times. This is called incast.

Asynchronous arrival time: Any condition not defined by synchronous.

Percentage of synchronization: this defines the level of overlap
[amount of bits] between the frames L1,L2..Ln.

Example: two 64 bytes frames, of length L1 and L2, arrive to ingress
port 1 and port 2 of the DUT. There is an overlap of 6.4 bytes
between the two where L1 and L2 were at the same time on the
respective ingress ports. Therefore the percentage of synchronization
is 10%.

6.2.2 Discussion


In this scenario, buffers are solicited on the DUT. In a ingress
buffering mechanism, the ingress port buffers would be solicited
along with Virtual Output Queues, when available; whereas in an
egress buffer mechanism, the egress buffer of the one outgoing port
would be used.

In either cases, regardless of where the buffer memory is located on
the switch architecture; the Incast creates buffer utilization.

When one or more frames having synchronous arrival times at the DUT
they are considered forming an incast.


6.2.3 Measurement

It is a MUST to measure the number of ingress and egress ports. It is
a MUST to have a non null percentage of synchronization, which MUST
be specified.



7 Application Throughput: Data Center Goodput

7.1. Definition

In Data Center Networking, a balanced network is a function of
maximal throughput 'and' minimal loss at any given time. This is
defined by the Goodput. Goodput is the application-level throughput.
It is measured in bytes / second. Goodput is the measurement of the
actual payload of the packet being sent.

7.2. Discussion

In data center benchmarking, the Goodput is a value that SHOULD be measured. It provides a realistic idea of the usage of the available bandwidth. A goal in data center environments is to maximize the Goodput while minimizing the loss.

7.3. Measurement

When S is the total bytes received from all senders [not inclusive of packet headers or TCP headers - it's the payload] and Ft is the Finishing Time of the last sender; the Goodput G is then measured by the following formula: G= S / Ft  bytes per second

Example: a TCP file transfer over HTTP protocol on a 10Gb/s media. The file cannot be transferred over Ethernet as a single continuous stream. It must be broken down into individual frames of 1500 bytes when the standard MTU [Maximum Transmission Unit] is used. Each packet requires 20 bytes of IP header information and 20 bytes of TCP header information, therefore 1460 byte are available per packet for the file transfer. Linux based systems are further limited to 1448 bytes as they also carry a 12 byte timestamp. Finally, the date is transmitted in this example over Ethernet which adds a 26 byte overhead per packet.

G= 1460/1526 x 10 Gbit/s which is 9.567 Gbit/s or 1.196 Gigabytes per second.

Please note: this example does not take into consideration additional Ethernet overhead, such as the interframe gap (a minimum of 96 bit times), nor collisions (which have a variable impact, depending on the network load).

When conducting Goodput measurements please document in addition to the 4.1 section:

-the TCP Stack used

-OS Versions

-NIC firmware version and model

For example, Windows TCP stacks and different Linux versions can influence TCP based tests results.


8.  References

8.1.  Normative References

   [1]    Bradner, S. "Benchmarking Terminology for Network
          Interconnection Devices", RFC 1242, July 1991.

   [2]    Bradner, S. and J. McQuaid, "Benchmarking Methodology for
          Network Interconnect Devices", RFC 2544, March 1999.

8.2.  Informative References

   [3]    Mandeville R. and Perser J., "Benchmarking Methodology for LAN
          Switching Devices", RFC 2889, August 2000.

   [4]    Stopp D. and Hickman B., "Methodology for IP Multicast
          Benchmarking", BCP 26, RFC 3918, October 2004.

8.3.  URL References

   [5]  Yanpei Chen, Rean Griffith, Junda Liu, Randy H. Katz, Anthony D.
        Joseph, "Understanding TCP Incast Throughput Collapse in
        Datacenter Networks",
        http://www.eecs.berkeley.edu/~ychen2/professional/TCPIncastWREN2009.pdf"
.

8.4.  Acknowledgments

Authors' Addresses

        Jacob Rapp
        Cisco Systems
        170 West Tasman Drive
        San Jose, CA 95134
        United States
        Phone: +1 408 853 2970
        Email: jarapp@cisco.com


        Lucien Avramov
        Cisco Systems
        170 West Tasman drive
        San Jose, CA 95134
        United States
        Phone: +1 408 526 7686
        Email: lavramov@cisco.com