

(Draft) IEEE 802 Nendica Report: Intelligent Lossless Data Center Networks

1 Editor

Name	Affiliation
Guo, Liang	CAICT/ODCC
Congdon, Paul	Huawei

2 Nendica Chair

Name	Affiliation
Marks, Roger	Huawei

3 Contributors/Supporters

Name	Affiliation
Li, Jie	CAICT/ODCC
Gao, Feng	Baidu
Gu, Rong	China Mobile
Zhao, Jizhuang	China Telecom
Chen, Chuansheng	Tencent
Yin, Yue	Huawei
Song, Qingchun	Nvidia
Liu, Jun	Cisco
He, Zongying	Broadcom
Sun, Liyang	Huawei
Tang, Guangming	Meituan
Quan, Hao	Meituan
Tao, Chunlei	JD
Wang, Shaopeng	CAICT/ODCC

4

5

1 Trademarks and Disclaimers

2 *IEEE believes the information in this publication is accurate as of its publication date; such*
3 *information is subject to change without notice. IEEE is not responsible for any inadvertent errors.*

4
5 **Copyright © 2021 IEEE. All rights reserved.**

6
7 IEEE owns the copyright to this Work in all forms of media. Copyright in the content retrieved, displayed or output
8 from this Work is owned by IEEE and is protected by the copyright laws of the United States and by international
9 treaties. IEEE reserves all rights not expressly granted.

10
11 IEEE is providing the Work to you at no charge. However, the Work is not to be considered within the “Public
12 Domain,” as IEEE is, and at all times shall remain the sole copyright holder in the Work.

13
14 Except as allowed by the copyright laws of the United States of America or applicable international treaties, you
15 may not further copy, prepare, and/or distribute copies of the Work, nor significant portions of the Work, in any
16 form, without prior written permission from IEEE.

17
18 Requests for permission to reprint the Work, in whole or in part, or requests for a license to reproduce and/or
19 distribute the Work, in any form, must be submitted via email to stds-ipr@ieee.org, or in writing to:

20
21 IEEE SA Licensing and Contracts
22 445 Hoes Lane
23 Piscataway, NJ 08854

Comments on this report are welcomed by Nendica: the IEEE 802 “Network Enhancements for the
Next Decade” Industry Connections Activity: <<https://1.ieee802.org/802-nendica>>

Comment submission instructions are available at: <<https://1.ieee802.org/802-nendica/nendica-dcn>>

24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45

The Institute of Electrical and Electronics Engineers, Inc.
3 Park Avenue, New York, NY 10016-5997, USA

Copyright © 2021 by The Institute of Electrical and Electronics Engineers, Inc.
All rights reserved. Published April 2021. Printed in the United States of America.

*IEEE and 802 are registered trademarks in the U.S. Patent & Trademark Office, owned by The Institute of Electrical and
Electronics Engineers, Incorporated.*

PDF: ISBN xxx-x-xxxx-xxxx-x XXXXXXXXXX

IEEE prohibits discrimination, harassment, and bullying. For more information, visit
<http://www.ieee.org/web/aboutus/whatis/policies/p9-26.html>.

*No part of this publication may be reproduced in any form, in an electronic retrieval system, or otherwise, without the prior
written permission of the publisher.*

To order IEEE Press Publications, call 1-800-678-IEEE.
Find IEEE standards and standards-related product listings at: <http://standards.ieee.org>

1 **NOTICE AND DISCLAIMER OF LIABILITY CONCERNING THE USE OF IEEE SA** 2 **INDUSTRY CONNECTIONS DOCUMENTS**

3
4 This IEEE Standards Association (“IEEE SA”) Industry Connections publication (“Work”) is not a consensus
5 standard document. Specifically, this document is NOT AN IEEE STANDARD. Information contained in this
6 Work has been created by, or obtained from, sources believed to be reliable, and reviewed by members
7 of the IEEE SA Industry Connections activity that produced this Work. IEEE and the IEEE SA Industry
8 Connections activity members expressly disclaim all warranties (express, implied, and statutory) related
9 to this Work, including, but not limited to, the warranties of: merchantability; fitness for a particular
10 purpose; non-infringement; quality, accuracy, effectiveness, currency, or completeness of the Work or
11 content within the Work. In addition, IEEE and the IEEE SA Industry Connections activity members disclaim
12 any and all conditions relating to: results; and workmanlike effort. This IEEE SA Industry Connections
13 document is supplied “AS IS” and “WITH ALL FAULTS.”

14 Although the IEEE SA Industry Connections activity members who have created this Work believe that the
15 information and guidance given in this Work serve as an enhancement to users, all persons must rely upon
16 their own skill and judgment when making use of it. IN NO EVENT SHALL IEEE OR IEEE SA INDUSTRY
17 CONNECTIONS ACTIVITY MEMBERS BE LIABLE FOR ANY ERRORS OR OMISSIONS OR DIRECT, INDIRECT,
18 INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO:
19 PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS
20 INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT
21 LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF
22 THIS WORK, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE AND REGARDLESS OF WHETHER SUCH
23 DAMAGE WAS FORESEEABLE.

24 Further, information contained in this Work may be protected by intellectual property rights held by third
25 parties or organizations, and the use of this information may require the user to negotiate with any such
26 rights holders in order to legally acquire the rights to do so, and such rights holders may refuse to grant
27 such rights. Attention is also called to the possibility that implementation of any or all of this Work may
28 require use of subject matter covered by patent rights. By publication of this Work, no position is taken
29 by the IEEE with respect to the existence or validity of any patent rights in connection therewith. The IEEE
30 is not responsible for identifying patent rights for which a license may be required, or for conducting
31 inquiries into the legal validity or scope of patents claims. Users are expressly advised that determination
32 of the validity of any patent rights, and the risk of infringement of such rights, is entirely their own
33 responsibility. No commitment to grant licenses under patent rights on a reasonable or non-
34 discriminatory basis has been sought or received from any rights holder. The policies and procedures
35 under which this document was created can be viewed at <http://standards.ieee.org/about/sasb/iccom/>.

36 This Work is published with the understanding that IEEE and the IEEE SA Industry Connections activity
37 members are supplying information through this Work, not attempting to render engineering or other
38 professional services. If such services are required, the assistance of an appropriate professional should
39 be sought. IEEE is not responsible for the statements and opinions advanced in this Work.

TABLE OF CONTENTS

Editor	1
Nendica Chair	1
Contributors/Supporters	1
INTRODUCTION	5
Scope	5
Purpose.....	5
BRINGING THE DATA CENTER TO LIFE.....	5
A new world with data everywhere	5
EVOLVING DATA CENTER REQUIREMENTS AND TECHNOLOGY.....	7
Previous Data Center Bridging Standards	7
Requirements evolution.....	8
Characteristics of AI computing	9
Evolving technologies.....	11
CHALLENGES WITH TODAY’S DATA CENTER NETWORK	20
High throughput and low latency tradeoff	20
Deadlock free lossless network.....	21
Congestion control issues in large-scale data center networks.....	23
Configuration complexity of congestion control algorithms	25
NEW TECHNOLOGIES TO ADDRESS NEW DATA CENTER PROBLEMS.....	27
Hybrid transports for low latency and high throughput	27
PFC deadlock prevention using topology recognition	28
Improving Congestion Notification	30
Addressing configuration complexity.....	32
STANDARDIZATION CONSIDERATIONS	34
CONCLUSION	36
CITATIONS.....	37

1

1

Introduction

2

3

4 This paper is the result of a work item [1] within the IEEE 802 “Network Enhancements for the Next
5 Decade” Industry Connections Activity known as Nendica. The paper is an update to a previous
6 report, “IEEE 802 Nendica Report: The Lossless Network for Data Centers” published on August 17,
7 2018 [2]. This update provides additional background on evolving use cases in modern data centers
8 and proposes solutions to additional problems identified by this paper.

9 Scope

10 The scope of this report is the exploration of networking technologies to support the requirements
11 of modern data center networks that include support for high performance computing and artificial
12 intelligence applications. Solutions to address challenges created by evolving requirements and new
13 age technologies are proposed. Standardization considerations are identified.

14 Purpose

15 The purpose of this report is to frame high-level solutions to issues and challenges with modern
16 data center Networks. The report includes background and technical analyses of current data
17 center environments as they are applied to the evolving needs of target applications. The report
18 highlights new technologies that are changing the dynamics and operation of the data center
19 Network. The results of the analysis lead to identification and recommendation of future
20 standardization activities.

2

Bringing the data center to life

21

22 A new world with data everywhere

23 Digital transformation is driving change in both our personal and professional lives. Workflows and
24 personal interactions are turning to digital processes and automated tools that are enabled by the
25 Cloud, Mobility, and the Internet of Things. The intelligence behind the digital transformation is
26 Artificial Intelligence (AI). Data centers running AI applications with massive amounts of data are
27 recasting that data into pertinent information, automated human interactions, and refined decision
28 making. The need to interact with the data center in real-time is more important than ever in
29 today’s world where augmented reality, voice recognition, and contextual searching demand
30 immediate results. Data center networks must deliver unprecedented levels of performance, scale,
31 and reliability to meet these real-time demands.

32 Data centers in the cloud era focused on application transformation and the rapid deployment of
33 services. During the AI era, data centers are the source of information and algorithms for the real-

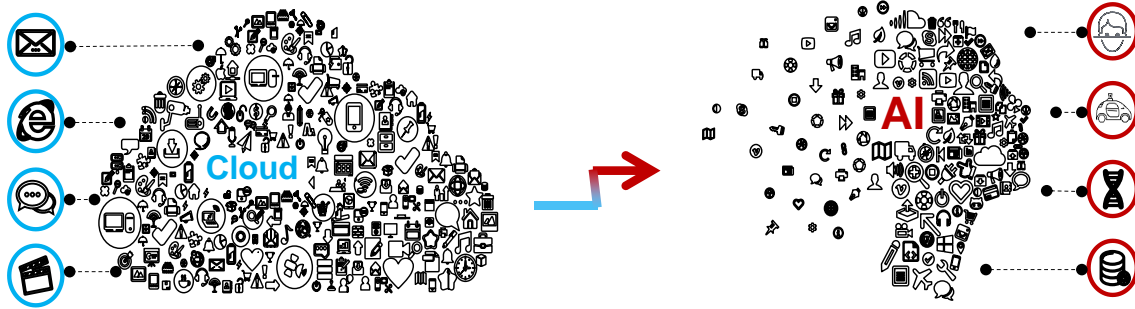


Figure 1 – Digital Transformation in the Era of AI

1 time digital transformation of our digital lives. The combination of high-speed storage and AI
 2 distributed computing render big data into fast data, access by humans, machines, and things. A
 3 high-performance, large scale data center network without packet loss is critical to the smooth
 4 operation of digital transformation.

5 For high-performance applications, such as AI, key measures of network performance include
 6 throughput, latency, and congestion. Throughput is dependent on the total capacity of the network
 7 to quickly transmit a large amount of data. Latency refers to the total delay for a transaction across
 8 the data center network. When the traffic load exceeds network capacity, congestion occurs. Packet
 9 loss is a factor that seriously affects both throughput and latency.

10 Currently, digital transformation of various industries is accelerating. It is estimated that 64% of
 11 enterprises have become the explorers and practitioners of digital transformation [3]. Among 2000
 12 multinational companies, 67% of CEOs have made digitalization the core of their corporate
 13 strategies [4]. The drive towards digital transformation in the real-time world is leading the data
 14 center network to support a ‘Data-Centric’ model of computing.

15 A large amount of data will be generated during the digitalization process, becoming a core asset,
 16 and enabling the emergence of artificial intelligence applications. The Huawei Global Industry
 17 Vision predicts that data volume will reach 180 ZB in 2025 [5]. However, data is not the “end-in-
 18 itself”. Knowledge and wisdom extracted from data are eternal values. The proportion of
 19 unstructured data (such as raw voice, video, and image data) increases continuously, and will

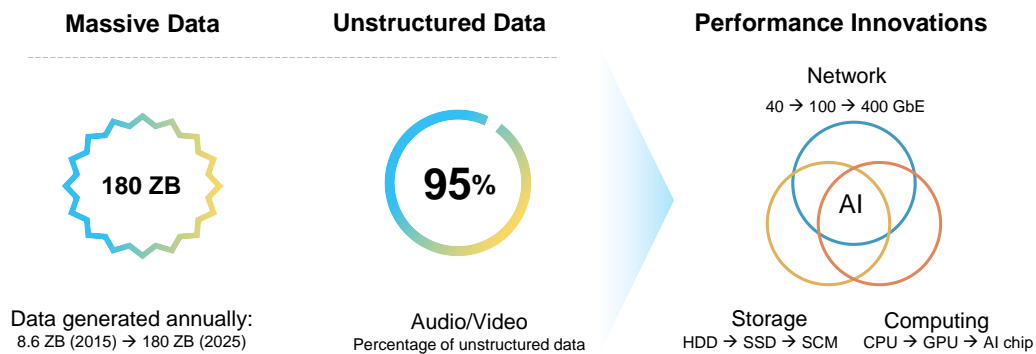


Figure 2 – Emerging Artificial Intelligence Applications

1 account for 95% of all data in the future. Current big data analytic methods are unable to keep pace
2 with the growth of data and performance innovations are needed to extract the value from the raw
3 data. An AI approach based on deep learning can filter out massive amounts of invalid data and
4 automatically extract useful information, providing more efficient decision-making and behavior
5 guidance.

6 The cloud data center architecture improved the performance and scale of applications in general.
7 The cloud platform allows rapid distribution of IT resources to create an application-centric service
8 model. In the AI era, the applications are consuming unprecedented amounts of data and the cloud
9 data center architecture is augmented with necessary performance innovations to handle the load.
10 Seamlessly introducing these innovations along with new AI applications can be tricky in an existing
11 cloud data center. Understanding how to efficiently process data based on the needs of AI
12 applications is a key focus area. Orchestrating the flow of data between the storage and computing
13 resources of the applications is a critical success factor.

3

Evolving data center requirements and technology

Previous Data Center Bridging Standards

17 During the early days of 10 Gbps Ethernet, the IEEE 802.1 Working Group developed a focus on Data
18 Center Bridging (DCB). The DCB task group defined a set of enhancements to the Ethernet, Bridges,
19 and associated protocols for use in data center environments. The use-case and focus were on
20 clustering and storage area networks, where traditionally dedicated technologies such as Infiniband
21 and Fiber Channel were used. Important objectives for Ethernet were to eliminate loss due to
22 congestion and to allocate bandwidth on links for selected traffic. The key contributions at the time
23 included the following:

- 24 • **Priority-based Flow Control (PFC):** A link level flow control mechanism that eliminates
25 packet loss and can be applied independently to each traffic class.
- 26 • **Enhanced Transmission Selection (ETS):** A queue scheduling algorithm that allows for
27 bandwidth assignments to traffic classes.
- 28 • **Congestion Notification:** A layer-2 end to end congestion management protocol that
29 detects congestion, signals across the layer-2 network to limit the transmission rate of
30 senders to avoid packet loss.
- 31 • **Data Center Bridging Capabilities Exchange Protocol (DCBX):** a discovery and capability
32 exchange protocol, working in conjunction with the Link Layer Discovery Protocol (LLDP), to
33 convey capabilities and configuration of the above features.

34 These contributions were important to the expansion of Ethernet into the specialized markets of
35 cluster computing and storage area networks. However, continued evolution is needed as the
36 environments and technologies have changed. Today's data centers are deployed on massive scale,
37 using Layer-3 protocols and highly orchestrated management systems. Ethernet links have
38 advanced from 10 Gbps to 400 Gbps, with active plans to increase speeds into the Tbps range. New

1 applications, such as Artificial Intelligence (AI) are placing new demands on the infrastructure and
2 driving architectural changes. Continued innovation is needed to further expand the use of Ethernet
3 in modern data center environments.

4 **Requirements evolution**

5 AI applications put pressure on the data center network. Consider AI training for self-driving cars
6 as an example. The deep learning algorithm relies heavily on massive data and high-performance
7 computing capabilities. The training data collected each day is approaching the petabyte level (1PB
8 = 1024 TB), and if traditional hard disk storage and common CPUs were used to process the data, it
9 could take at least one year to complete the training process. This is clearly impractical. To improve
10 AI data processing efficiency, revolutionary changes are needed in the storage and computing fields.
11 For example, storage performance needs to improve by an order of magnitude to achieve more
12 than 1 million input/output operations per second (IOPS) [6].

13 To meet real-time data access requirements, storage media has evolved from hard disk drives
14 (HDDs) to solid-state drives (SSDs) to storage-class memory (SCMs). This has reduced storage
15 medium latency by more than 1000 times. Without similar improvements in network latency, these
16 storage improvements cannot be realized and simply move the bottleneck from the media to the
17 network. With networked SSD drives, the communication latency accounts for more than 60% of
18 the total storage end-to-end latency. With the move to SCM drives, this percentage could increase
19 to 85% unless improvements in network performance are achieved. This creates a scenario where
20 the precious storage media is idle more than half of the time. When you consider recent
21 improvements in both storage media and AI computing processors together, the communication
22 latency accounts for more than 50% of the total latency, further hindering improvements and
23 wasting resources [7].

24 AI applications and environments are growing in scale and complexity. For example, there were 7
25 ExaFLOPS and 60 million parameters in Microsoft's Resnet of 2015. Baidu used 20 ExaFLOPS and
26 300 million parameters when training their deep speech system in 2016. In 2017, the Google NMT
27 used 105 ExaFLOPS and 8.7 billion parameters [8]. New characteristics of AI computing are requiring
28 an evolution of data center network.

29 Traditional protocols are no longer able to satisfy the requirements of new applications that serve
30 our daily lives. In a simple example, the online food take-out industry at Meitan has increased nearly
31 500% in the last four years. The number of transactions has increased from 2.149 billion to 12.36
32 billion where those transactions all occur within a few hours at peak mealtimes. The Meituan
33 Intelligent Scheduling System is responsible for orchestrating a complex multi-person, multi-point
34 real-time decision-making process for end-users, businesses and over 600,000 delivery drivers. The
35 drivers report positioning data 5 billion times a day that are used to calculate optional paths for the
36 drivers and deliver optimal solutions within 0.55 milliseconds. When the back-end servers use
37 TCP/IP protocols, the amount of data copied between kernel buffers, application buffers and NIC
38 buffers stresses the CPU and memory bus resources causing increased delay and an inability to meet
39 the application requirements. The newer Remote Direct Memory Access (RDMA) protocol
40 eliminates data copies and frees CPU resources to perform driver path and take-out order
41 calculations at scale. The improved efficiency of RDMA puts more pressure on the network, moving
42 the bottleneck to the data center network infrastructure where low-latency and lossless behavior
43 become the new critical requirements.

1 **Characteristics of AI computing**

2 Traditional data center services (web, database, and file storage) are transaction-based and the
 3 calculated results are often deterministic. For such tasks, there is little correlation or dependency
 4 between a single transaction and the associated network communication. The occurrence and
 5 duration of the traditional transactions are random. AI computing, however, is different. It is an
 6 optimization problem with iterative convergence. This causes high spatial correlation within the
 7 data sets and computing algorithms and creates temporal correlations between communication
 8 flows.

9 AI computing works on big data and demands fast data. To achieve this, it must operate in parallel
 10 to “divide-and-conquer” the problem. The computing model and input data sets are large (e.g in a
 11 100 MB node, the AI model with 10K rules requires more than 4 TB memory). A single server cannot
 12 provide enough storage capacity and processing resources to handle the problem sequentially.
 13 Concurrent AI computing and storage nodes are required to shorten the processing time. This
 14 distributed AI computing and storage requirement highlights the need for a fast, efficient, and
 15 lossless data center network that has the flexibility to support two distinct parallel modes of
 16 operation: model parallel computing and data parallel computing.

17 **Model Parallel Computing**

18 In model parallel computing, each node computes one part of the overall algorithm. Each node
 19 processes the same set of data, but with a different portion of the algorithm, resulting in an estimate
 20 for a differing set of parameters. The nodes exchange their estimates to converge upon the best
 21 estimate for all the data parameters. With model parallel computing, there is an initial distribution
 22 of the common data set to a distributed number of nodes, followed by a collection of individual
 23 parameters from each of the participating nodes. Figure 3 shows how parameters of the overall
 24 model may be distributed across computing nodes in a model parallel mode of operation.

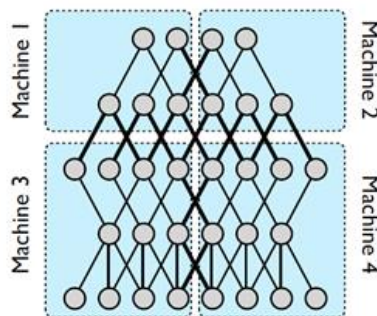


Figure 3 - Model parallel training

25 **Data Parallel Computing**

26 In data parallel computing, each node loads the entire AI algorithm model, but only processes part
 27 of the input data. Each node is trying to estimate the same set of parameters using a different view
 28 of the data. When a node completes a round of calculations, the parameters are weighted and
 29 aggregated by a common parameter server as seen in Figure 4. The weighted parameter update
 30 requires that all nodes upload the information synchronously.

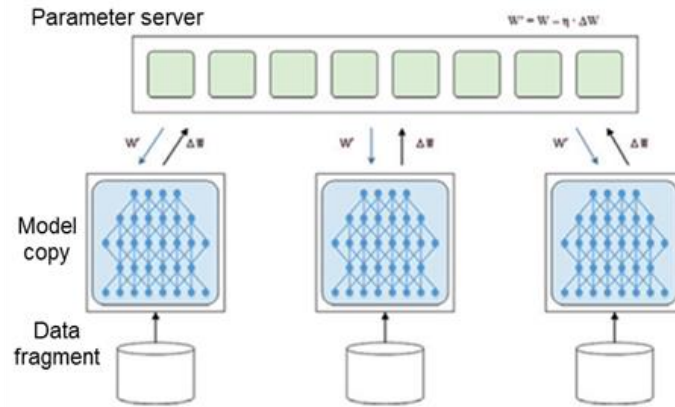


Figure 4 - Data parallel training

1 Regardless of the parallel computing approach, data center networks feel the pressure of
 2 demanding communication. When the network becomes the bottleneck, the waiting time for
 3 computing resources can exceed 50% of the job completion time [9].

4 With all AI applications, the computing model is iterative and requires a synchronization step that
 5 creates network incast congestion. Figure 5 shows how incast congestion occurs with AI training.
 6 The training process is iterative and there are many parameters synchronized on each iteration. The
 7 workers download the model and upload newly calculated results (ΔM) to a parameter server
 8 during a synchronization step. The uploading to the parameter server creates incast. When the
 9 computing time is improved by deploying new compute technology, the pressure on the network
 10 and resulting incast increases.

11 The communication between the worker nodes and the parameter server constitutes a collection
 12 of interdependent network flows. In the iteration process of distributed AI computing, many burst
 13 traffic flows are generated to distributed data to workers within milliseconds, followed by an incast
 14 event of smaller sized flows directed at the parameter server when the intermediate parameters
 15 are delivered and updated. During the exchange of these flows packet loss, congestion, and load
 16 imbalance can occur on the network. As a result, the Flow Completion Time (FCT) of some of the
 17 flows is prolonged. If a few flows are delayed, storage and computing resource can be underutilized.
 18 Consequently, the completion time of the entire application is delayed.

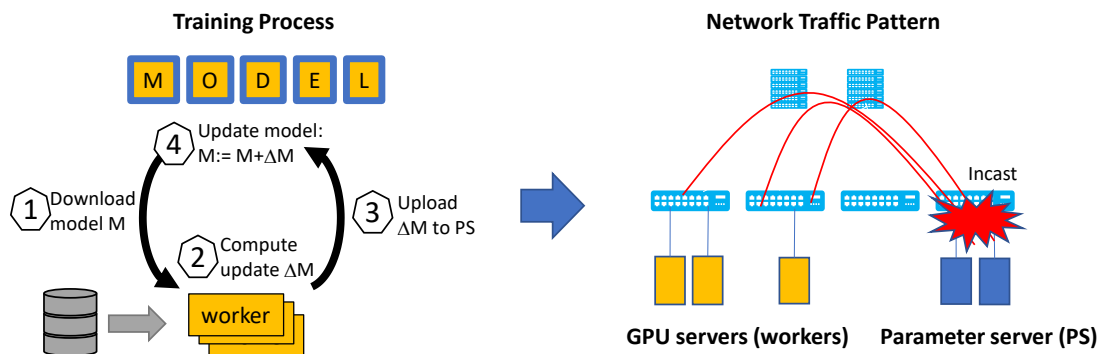


Figure 5 - Periodic incast congestion during training

1 Distributed AI computing is synchronous, and it is desirable for the jobs to have a predictable
 2 completion time. When there is no congestion, dynamic latency across the network is small
 3 allowing the average FCT to be predictable and therefor the performance of the entire application
 4 is predictable. When congestion causes dynamic latency to increase to the point of causing packet
 5 loss, FCT can be very unpredictable. Flows that complete in a time that is much greater than the
 6 average completion contributes to what is known as tail latency. Tail latency is the small percentage
 7 of response times from a system, out of all of responses to the input/output (I/O) requests it serves,
 8 that take the longest in comparison to the bulk of its response times. Reducing tail latency as much
 9 as possible is extremely critical to the success of parallel algorithms and the whole distributed
 10 computing system. To maximize the use of computing resources in the data center, tail latency
 11 should be addressed.

12 Evolving technologies

13 Progress can be seen when evolving requirements and evolving technologies harmonize. New
 14 requirements often drive the development of new technologies and new technologies often enable
 15 new use cases that lead to, yet again, a new set of requirements. Breakthroughs in networked
 16 storage, distributed computing, system architecture and network protocols are enabling the
 17 advancement of the next generation data center.

18 SSDs and NVMeoF: High throughput, low-latency network

19 In networked storage, a file is distributed to multiple storage servers for IO acceleration and
 20 redundancy. When a data center application reads a file, it accesses different parts of data from
 21 different servers concurrently. The data is aggregated through a data center switch at nearly the
 22 same time. When a data center application writes a file, the writing of data can trigger a series of
 23 storage transactions between distributed and redundant storage nodes. The entire sequence of
 24 transactions must complete before the original write action is satisfied. Figure 6 shows an example
 25 of data center communication triggered by the networked storage service model.

26 The example highlights the importance of the network enabling both high throughput and low
 27 latency simultaneously. The bulk data being written to the primary storage server is transmitted
 28 multiple times to the replicas. The small sized acknowledgments and commit messages must be
 29 sequenced and ultimately delivered to the originating client before the transaction can complete,
 30 emphasizing the need for ultra-low latency.

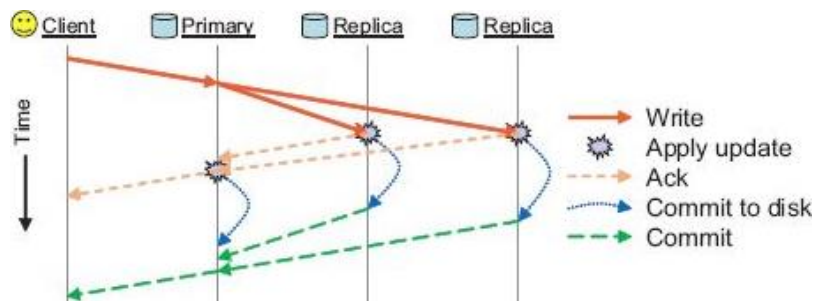


Figure 6 - Networked storage service model

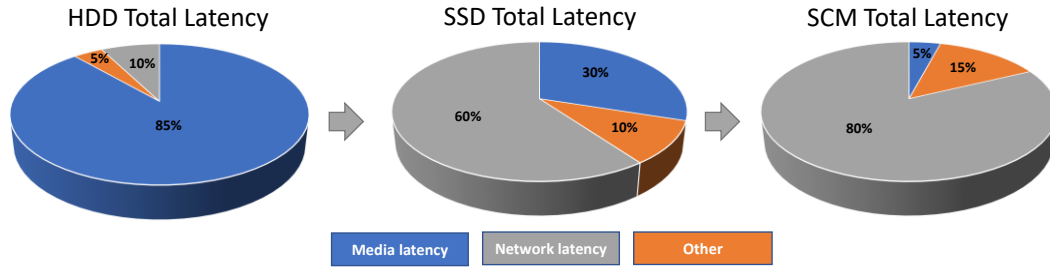


Figure 7 – End-to-end latency breakdown for HDD and SSD

1 Massive improvements in storage performance have been achieved as the technology has evolved
 2 from HDD to SSD to SCM using the Non-Volatile Memory Express (NVMe) interface specification.
 3 Accessing storage media via NVMe has decreased access time by a factor of 1000 over previous HDD
 4 technology. Sample seek times between the various technologies include: HDD = 2-5 ms, SATA SSD
 5 = 0.2 ms, and NVMe SSD = 0.02 ms. SCM is generally three to five times faster than NVMe flash
 6 SSDs.

7 NVMe-over-fabrics (NVMeoF) involves deploying NVMe for networked storage. The much faster
 8 access speed of the medium results in greater network bottlenecks and the impact of network
 9 latency becomes more significant. Figure 7 shows how network latency has become the primary
 10 bottleneck with faster NVMe based storage. Network latency was a negligible part of end-to-end
 11 networked HDD storage latency, but is poised to become a significant component of latency with
 12 networked SCM storage.. To maximize the IOPS performance of the new medium, the network
 13 latency problem must be resolved first.

14 There are two distinct types of latency; static latency and dynamic latency. Static latency includes
 15 serial data latency, device forwarding latency, and optical/electrical transmission latency. This type
 16 of latency is determined by the capability of the switching hardware and the transmission distance
 17 of the data. It usually is fixed and very predictable. Figure 8 shows that current industry
 18 measurements for static latency are generally at the nanosecond (10^{-9} second) or sub-microsecond
 19 (10^{-6}) level, and account for less than 1% of the total end-to-end network delay.
 20

21 Dynamic latency plays a much greater role in total end-to-end network delay and is greatly affected
 22 by the conditions within the communication environment. Dynamic latency is created from delays
 23 introduced by internal queuing and packet retransmission, which are caused by network congestion
 24 and packet loss. Parallel AI computing models create unique traffic patterns that result in heavy
 25 network congestion. The key to low end-to-end network latency is to address dynamic latency and
 26 the key to addressing dynamic latency is mitigating congestion.
 27

28 The major component of dynamic latency is the delay from packet retransmission when packets are
 29 dropped within the network. Packet loss latency is an order magnitude greater than queuing delay
 30 and has proven to have a severe impact on applications. Packet loss occurs when switch buffers are
 31 overrun because of congestion (NOTE: we ignore packet loss due low-probability bit errors during
 32 transmission). There are two key types of congestion that lead to packet loss: in-network
 33 congestion and incast congestion. In-network congestion occurs on switch-to-switch links within
 34 the network fabric when the links become overloaded, perhaps due to ineffective load balancing.
 35 Incast congestion occurs at the edge of the network when many sources are sending to a common

End-to-end Network Latency Breakdown

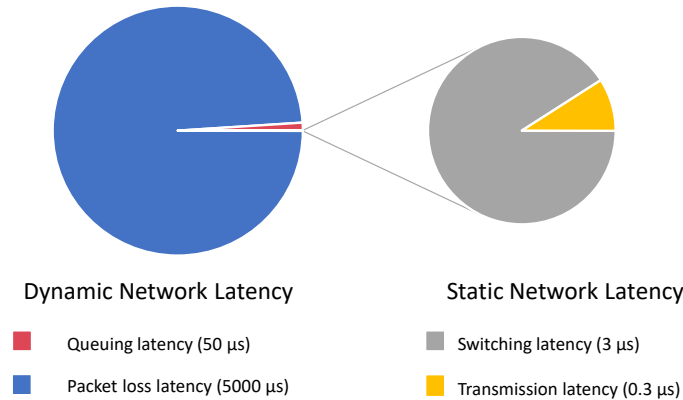


Figure 8 – Network Latency Breakdown

1 destination at the same time. AI computing models inherently have a phase when data is
 2 aggregated after a processing iteration from which incast congestion (many-to-one) easily occurs.

3 GPUs: Ultra-low latency network for parallel computing

4 Today’s AI computing architecture includes a hybrid mix of Central Processing Units (CPUs) and
 5 Graphics Processing Units (GPUs). GPUs, originally invented to help render video games at
 6 exceptional speeds, have found a new home in the data center. The GPU is a processor with
 7 thousands of cores capable of performing millions of mathematical operations in parallel. All AI
 8 learning algorithms perform complex statistical computations and deal with a huge number of
 9 matrix multiplication operations – perfectly suited for a GPU. However, to scale the AI computing
 10 architecture to meet the needs of today’s AI applications in a data center, the GPUs must be
 11 distributed and networked. This places stringent requirements on communication volume and
 12 performance.

13 Facebook recently tested the distributed machine learning platform Caffe2, in which the latest
 14 multi-GPU servers are used for parallel acceleration. In the test, computing tasks on eight servers
 15 resulted in underutilized resources on the 100 Gbit/s InfiniBand network. The presence of the

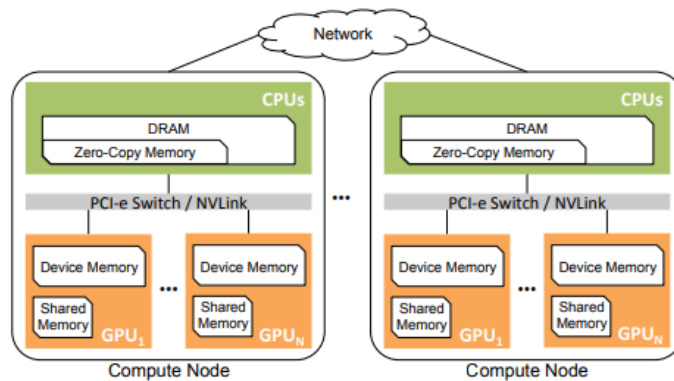


Figure 9 – Distributed AI Computing Architecture

1 network and network contention reduced the performance of the solution to less than linear scale
2 [10]. Consequently, network performance greatly restricts horizontal extension of the AI system.

3 GPUs provide much higher memory bandwidth than today's CPU architectures. Nodes with multiple
4 GPUs are now commonly used in high-performance computing because of their power efficiency
5 and hardware parallelism. Figure 9 illustrates an architecture with multi-GPU nodes, each of which
6 consists of a host (CPUs) and several GPU devices connected by a PCI-e switch or NVLink. Each GPU
7 can directly access its local relatively large device memory, much smaller and faster shared memory,
8 and a small, pinned area of the host node's DRAM, called zero-copy memory [11].

9 SmartNICs

10 Over the years there have been periods of time when performance improvements in CPU speeds
11 and Ethernet links have eclipsed one another. Figure 10 shows the relative historical performance
12 gains with Ethernet link speeds [12] and benchmark improvements for CPU performance [13].
13 During some historical periods, the processing capability of a traditional CPU was more than enough
14 to handle the load of an Ethernet link and the cost savings of a simplified network interface card
15 (NIC) along with the flexibility of handling the entire networking stack in software was a clear
16 benefit. During other periods, the jump in link speed from the next iteration of IEEE 802.3 standards
17 was too much for the processor to handle and a more expensive and complex SmartNIC with
18 specialized hardware offloads became necessary to utilize the Ethernet link. As time goes on and
19 the SmartNIC offloads mature, some of them become standard and included in the base features of
20 what is now considered a common NIC. This phenomenon was seen with the advent of the TCP
21 Offload Engine (TOE) which supported TCP checksum offloading, large segment sending and receive
22 side scaling.

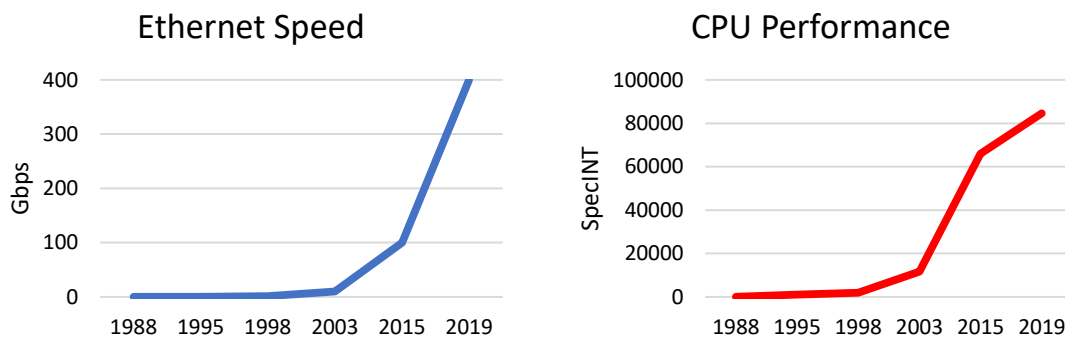


Figure 10 – Historical Performance Comparison

23 In today's world, there are signs of Moore's law fading while Ethernet link speeds continue to soar.
24 The latest iteration of IEEE 802.3 standards is achieving 400 Gbps. Couple this divergence with the
25 added complexity of software-defined networking, virtualization, storage, message passing and
26 security protocols in the modern data center, and there is a strong argument that the SmartNIC
27 architecture is here to stay. So, what exactly is a data center SmartNIC today?

28 Figure 11 shows a data center server architecture including a SmartNIC. The SmartNIC includes all
29 the typical NIC functions, but also includes key offloads to help accelerate applications running on
30 the server CPU and GPU. The SmartNIC does not replace the CPU or the GPU but rather
31 complements them with networking offloads. Some of the key offloads include virtual machine

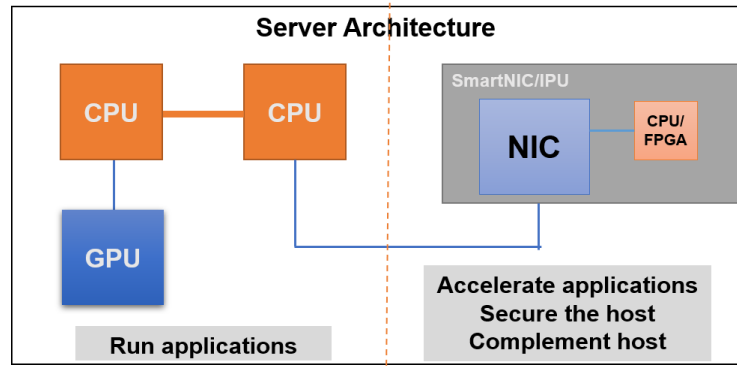


Figure 11 – Server Architecture with SmartNIC

1 interface support, flexible match-action processing of packets, overlay tunnel termination and
 2 origination, encryption, traffic metering, shaping and per-flow statistics. Additionally, SmartNICs
 3 often include entire protocol offloads and direct data placement to support RDMA and NVMe-oF
 4 storage interfaces.

5 One new critical component of today's SmartNIC is programmability. A criticism of SmartNICs in
 6 the past was their inability to keep pace with the rapidly changing networking environment. The
 7 early cloud data center environments favored using the CPU for most networking functions because
 8 the required feature set for the NIC was evolving faster than the development cycle of the
 9 hardware. Today's SmartNICs however have an open and flexible programming environment. They
 10 are essentially a computer in front of the computer with an open source development environment
 11 based on Linux and other software-defined networking tools such as Open vSwitch [14]. It is
 12 essential that SmartNICs integrate seamlessly into the open source ecosystem to enable rapid
 13 feature development and leverage.

14 SmartNICs in the data center increase the overall utilization and load on the network. They can
 15 exacerbate the effects of congestion by fully and rapidly saturating a network link. At the same
 16 time, they can respond quickly to congestion signals from the network to alleviate intermittent
 17 impact and avoid packet loss. The programmability of the SmartNIC allows it to adapt to new
 18 protocols that can coordinate with the network to avoid conditions such as incast.

19 Remote Direct Memory Access (RDMA)

20 RDMA is a new technology designed to solve the high latency problem of server-side data processing
 21 in network applications. With RDMA data transfers directly from one computer's memory to
 22 another without the intervention of either's operating system. This allows for high bandwidth, low
 23 latency network communication and is particularly suitable for use in massively parallel computer
 24 environments. Figure 12 shows the principles of the RDMA protocol.

25 There are three different transports for the RDMA protocol: Infiniband, iWarp and RoCEv1/RoCEv2.

26 *Infiniband*

27 In 2000, the InfiniBand Trade Association (IBTA) released the initial Infiniband specification with
 28 support for RDMA. Infiniband is tailored for an efficient hardware design that ensures reliable data

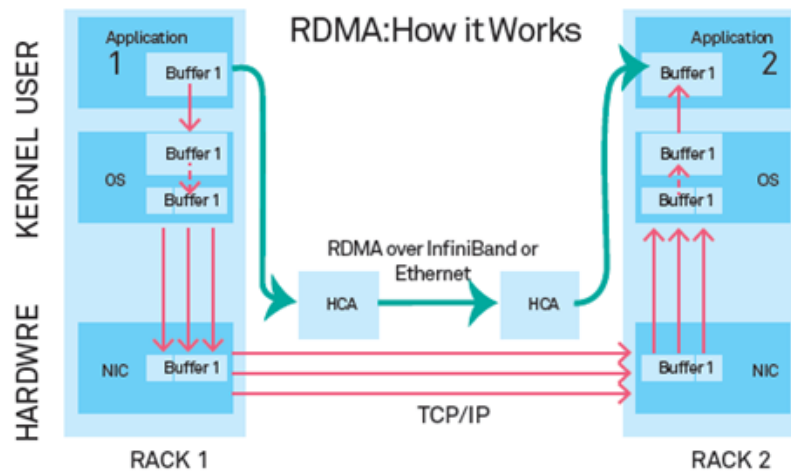


Figure 12 - Working principle of RDMA

1 transmission and direct access to the memory of remote nodes. Infiniband is a unique network
 2 solution requiring specific Infiniband switches and Infiniband interface cards.

3 *iWarp*

4 iWarp is an RDMA protocol, defined by the IETF in 2014 to run over TCP. Using TCP as a transport
 5 allows iWarp to traverse the Internet and wide area as well as a standard Ethernet network and
 6 within a data center. While iWarp can be implemented in software, to obtain the desired
 7 performance within the data center specialized iWarp NICs card are used.

8 *RDMA over Converged Ethernet (RoCE)*

9 In April 2010, the IBTA released the RoCEv1 specification, which augments the Infiniband
 10 Architecture Specification with the capability of supporting Infiniband over Ethernet (IBoE). The
 11 RoCEv1 standard specifies an Infiniband network layer directly on top of the Ethernet link layer.
 12 Consequently, the RoCEv1 specification does not support IP routing. Since Infiniband relies on a
 13 lossless physical transport, the RoCEv1 specification depends on a lossless Ethernet environment.

14 Modern data centers tend to use Layer-3 technologies to support large scale and greater traffic
 15 control. The RoCEv1 specification required an end-to-end layer-2 Ethernet transport and did not
 16 operate effectively in a layer-3 network. In 2014, the IBTA published RoCEv2, which extended
 17 RoCEv1 by replacing the Infiniband Global Routing Header (GRH) with an IP and UDP header. Now
 18 that RoCE is routable, it is easily integrated into the preferred data center environment. However,
 19 to obtain the desired RDMA performance, the RoCE protocol is offloaded to special network
 20 interface cards. These network cards implement the entire RoCEv2 protocol, including the UDP
 21 stack, congestion control and any retransmission mechanisms. While UDP is lighter weight than
 22 TCP, the additional support required to make RoCEv2 reliable adds complication to the network
 23 card implementation. RoCEv2 still depends upon the Infiniband Transport Protocol, which was
 24 designed to operate in a lossless Infiniband environment, so RoCEv2 still benefits from a lossless
 25 Ethernet environment.

Technology	Data Rates (Gbit/s)	Latency	Key Technology	Advantage	Disadvantage
TCP/IP over Ethernet	10, 25, 40, 50, 56, 100, or 200	500-1000 ns	TCP/IP Socket programming interface	Wide application scope, low price, and good compatibility	Low network usage, poor average performance, and unstable link transmission rate
Infiniband	40, 56, 100, or 200	300-500 ns	InfiniBand network protocol and architecture Verbs programming interface	Good performance	Large-scale networks not supported, and specific NICs and switches required
RoCE/RoCEv2	40, 56, 100, or 200	300-500 ns	InfiniBand network layer or transport layer and Ethernet link layer Verbs programming interface	Compatibility with traditional Ethernet technologies, cost-effectiveness, and good performance	Specific NICs required Still have many challenges to
Omni-Path	100	100 ns	OPA network architecture Verbs programming interface	Good performance	Single manufacturer and specific NICs and switches required

Table 1 – Comparison of RDMA Network Technologies

1 Figure 13 shows the most common RDMA protocol stacks and their associated standards bodies.
 2 Table 1 compares the details of different implementations. RDMA has become the protocol of
 3 choice for high-speed storage, AI and Machine Learning applications in large scale cloud data
 4 centers. There are real world examples of tens of thousands of servers running RDMA in production.
 5 Applications have reported impressive performance improvements by adopting RDMA [15]. For
 6 example, distributed training for machine learning has accelerated more than 100 times and the
 7 I/O speed of networked SSD storage has improved more than 50 times using RDMA for
 8 communications as opposed to TCP/IP. These improvements majorly stem from the hardware
 9 offloading characteristic of RDMA.

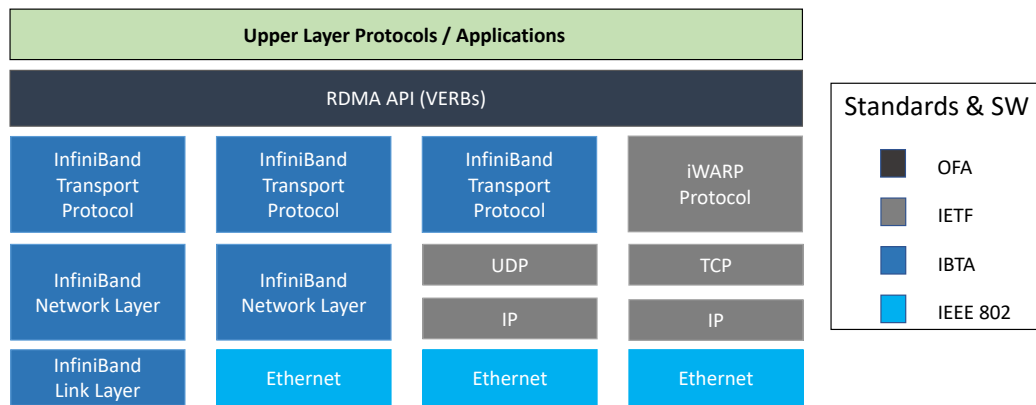


Figure 13 – RDMA protocol stacks and standards

10 **GPU DirectRDMA**

11 Combining two good ideas can often create a breakthrough idea. GPU DirectRDMA comprises the
 12 PeerDirect technology of PCIe and the RDMA technology of the network to deliver data directly to

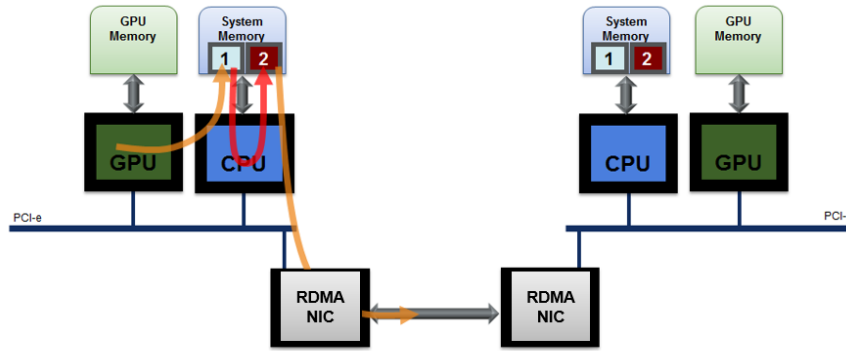


Figure 14 - The Data Transfer Before GPU DirectRDMA

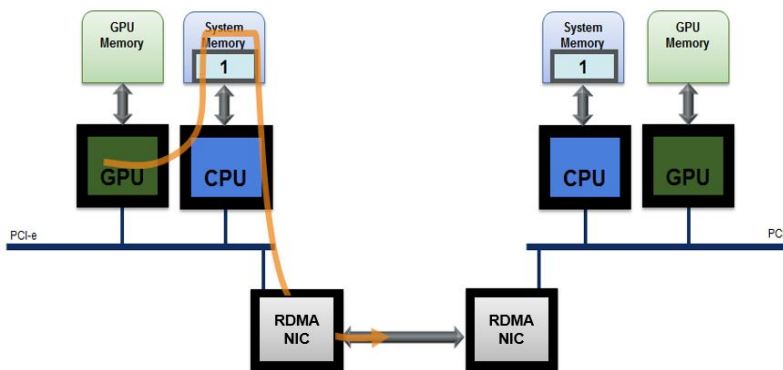
1 a GPU memory. This technology can be supported by any PCIe peer which provides access to its
 2 memory, such as NVIDIA GPU, XEON PHI, AMD GPU, FPGA, and so on.

3 GPU communications uses “pinned” buffers for data movement. A SmartNIC may also use “pinned”
 4 memory to communicate with a remote “pinned” memory across the network. These two types of
 5 “pinned” memory are separate sections of host memory that are dedicated to the GPU and the
 6 SmartNIC.

7 Before GPU DirectRDMA, when one GPU transferred data to another GPU in a remote server, the
 8 source GPU needed to copy the data from GPU memory to CPU memory which was pinned by the
 9 GPU. Then the host CPU copied the data from the GPU pinned memory to memory pinned by the
 10 SmartNIC. Next, the SmartNIC used RDMA to transmit the data to the remote server across the
 11 network. On the remote server side, the reverse process took place. The data arrived at the memory
 12 pinned by the SmartNIC, then the CPU copied the data to the memory pinned by the GPU, and
 13 eventually the data arrived at the remote GPU memory. Figure 14 shows the GPU-to-GPU data copy
 14 process before the existence of GPU DirectRDMA.

15 While the cost of copying data between the GPU and CPU is much lower than the cost of using TCP
 16 to pass the data between GPUs, it still suffers from several issues:

- 17 1. Consumption of CPU resources. The CPU may become a bottleneck during the data copy.



Picture 15 - The Data Transfer Using GPU Direct

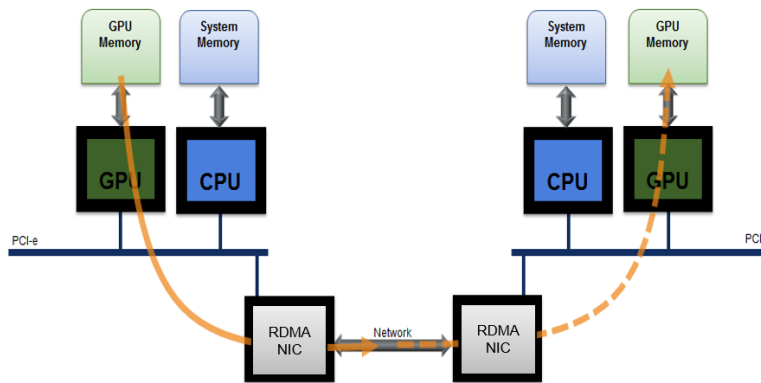


Figure 16 - The Data Transfer Using GPU DirectRDMA

- 1 2. Increased latency and reduced bandwidth. The additional memory copies take time and
- 2 reduce I/O bandwidth.
- 3 3. Host memory consumption. Multiple sets of pinned buffers reduce available host memory
- 4 which impacts application performance and increases system TCO.

- 5 Optimizations such as write-combining and overlapping GPU computation with data transfer allow
- 6 the network and the GPU to share “pinned” buffers. This eliminates the need to make a redundant
- 7 copy of the data in host memory and allows the data to be directly transferred via RDMA. On the
- 8 receiver side the data is directly written to the GPU pinned host buffer after arriving via RDMA. This
- 9 technique eliminates buffer copies between the CPU and the GPU and is known as GPU Direct
- 10 technology.

- 11 A further optimization is to create an RDMA channel between the local GPU memory and the
- 12 remote GPU memory to eliminate CPU bandwidth and latency bottlenecks. This results in
- 13 significantly improved communication efficiency between GPUs in remote nodes. For this
- 14 optimization to work, the CPU coordinates RDMA communication tasks for the GPU and SmartNIC.
- 15 The SmartNIC directly accesses GPU memory to send and receive data to a remote GPU memory.
- 16 This technique is known as GPU DirectRDMA technology.

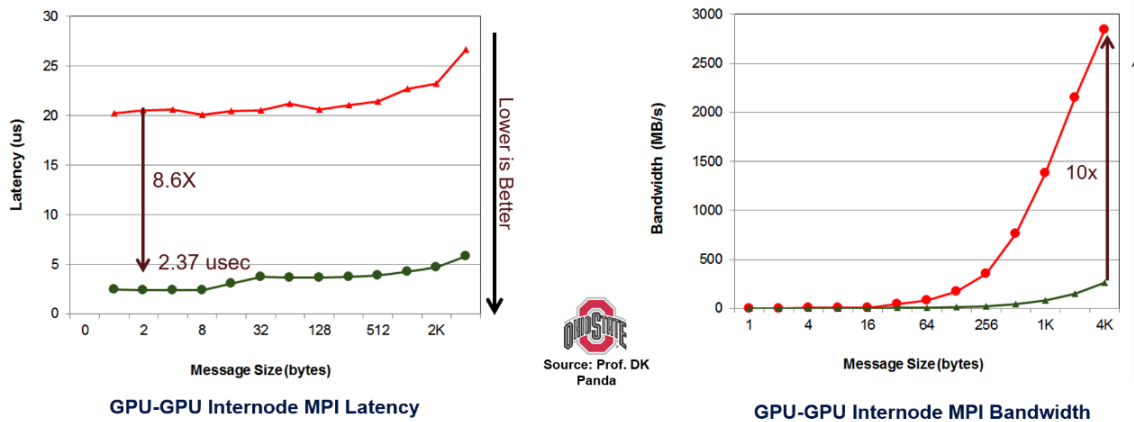


Figure 17 - GPU DirectRDMA Performance (From OSU)

1 Figure 17 shows how GPU DirectRDMA technology improves GPU communication performance by
 2 a factor of 10 over the traditional approach. These improvements have made GPU DirectRDMA
 3 technology a mandatory component of HPC and AI applications, improving both performance and
 4 scalability.

4

Challenges with today's data center network

High throughput and low latency tradeoff

8 Simultaneously achieving both low latency and high throughput in a large-scale data center is
 9 difficult. To achieve low latency, it is necessary to allow flows to begin transferring at line rate while
 10 at the same time maintaining near empty switch queues. Aggressively starting flows at line rate
 11 allows them to consume all available network bandwidth instantly and can lead to extreme
 12 congestion at convergence points. Deep switch buffers absorb temporary congestion to avoid
 13 packet loss but delay the delivery of latency sensitive packets. While deep switch buffers provide
 14 more resources for balancing the tradeoff between low latency and high throughput it is
 15 increasingly difficult to build switches with deep buffers. Switch capacity continues to increase with
 16 link speeds and higher port density, but the buffer size of commodity switching chips cannot keep
 17 pace. Figure 18 shows hardware trends for top-of-the-line data center switches chips manufactured
 18 by Broadcom [16].

19 Using a low ECN marking threshold can help slow aggressive flows and keep switch queue levels
 20 empty, but this reduces throughput. High throughput flows benefit from larger switch queues and
 21 higher ECN marking thresholds to prevent overreacting to temporary congestion and slowing down
 22 unnecessarily.

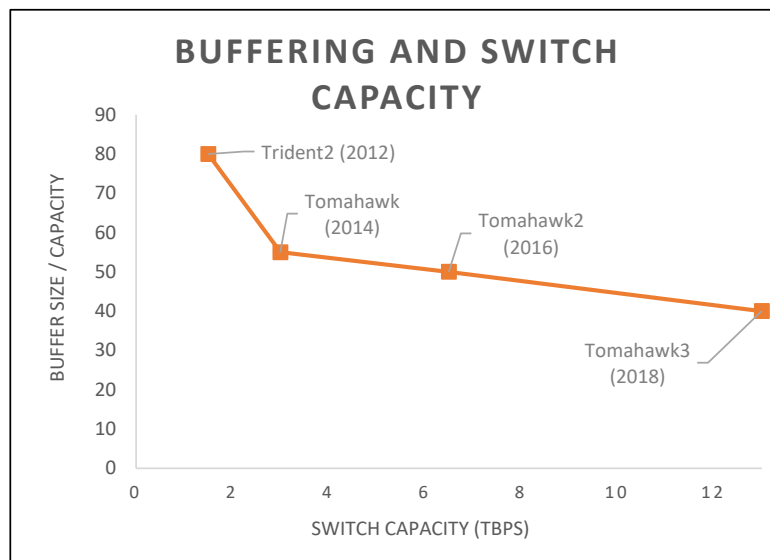


Figure 18 – Switch Chip Buffer Trends

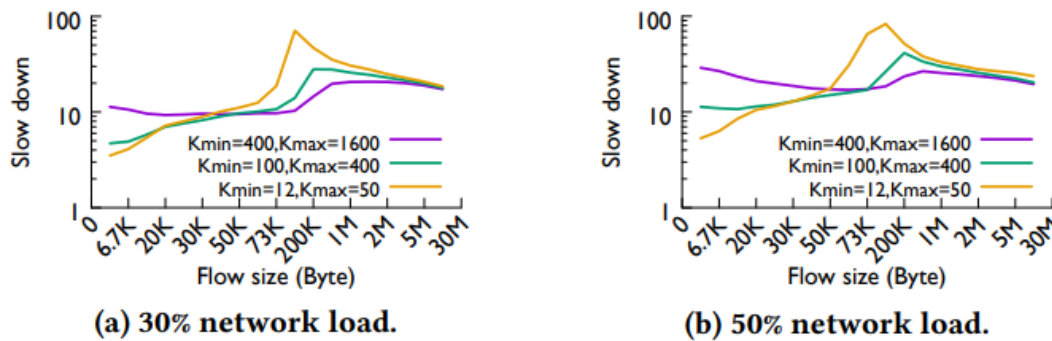


Figure 19 – FCT slowdown distribution with different ECN thresholds, using WebSearch

1 Experimentation shows the tradeoff between high throughput and low latency exists after varying
 2 algorithms, parameters, traffic patterns and link loads [15]. Figure 19 from [15] shows how flow
 3 completion times (FCT) are delayed beyond their theoretical minimum FCT when using different
 4 ECN marking thresholds (Kmin, Kmax) during a controlled experiment using a public RDMA
 5 WebSearch traffic workload as the input. Lower values for Kmin and Kmax will cause ECN markings
 6 to occur more quickly and force flows to slow down more aggressively. As seen in the figure, when
 7 using low ECN thresholds, small flows which are latency-sensitive have lower FCT slowdown, while
 8 big flows which are typically bandwidth-hungry suffer from larger FCT slowdown. The trend is more
 9 obvious when the network load is higher (Figure 19-b when the average link load is 50%).

10 Deadlock free lossless network

11 RDMA advantages over TCP include low latency, high throughput, and low CPU usage. However,
 12 unlike TCP, RDMA needs a lossless network; i.e. there should be no packet loss due to buffer
 13 overflow at the switches [17]. The RoCE protocol runs on top of UDP with a go-back N retransmission
 14 strategy that severely impacts performance when retransmission is invoked. As such, RoCE requires
 15 Priority-based Flow Control (IEEE Std 802.1Q-2018, Clause 36 [18]) to ensure that no packet loss
 16 occurs in the data center network. Figure 20 from [19] shows how the RoCE service throughput
 17 decreases rapidly with increasing packet loss rate. Losing as little as one in one thousand packets
 18 decreases RoCE service performance by roughly 30%.

19 Priority-based Flow Control (PFC) prevents packet loss due to buffer overflow by pausing the
 20 upstream sending device when the receiving device input buffer occupancy exceeds a specified
 21 threshold. While this provides the necessary lossless environment for RoCE, there are problems
 22 with the large-scale use of PFC. One such problem is the possibility of a PFC deadlock.

23 Deadlocks in lossless networks using PFC style backpressure have been studied for many years [20,
 24 21, 22]. A PFC deadlock occurs when there is a cyclic buffer dependency (CBD) among switches in
 25 the data center network. The CBD is created when a dependent switch, in a sequence of switches,
 26 is waiting for the availability of buffers in other switches before transmitting a packet. If the switches
 27 involved in the CBD are using PFC and the sequence of switches are physically connected in a loop,
 28 a PFC deadlock can occur. RDMA flows in a Clos data center network are distributed across multiple
 29 equal cost paths to achieve the highest possible throughput and lowest latency. While there are no

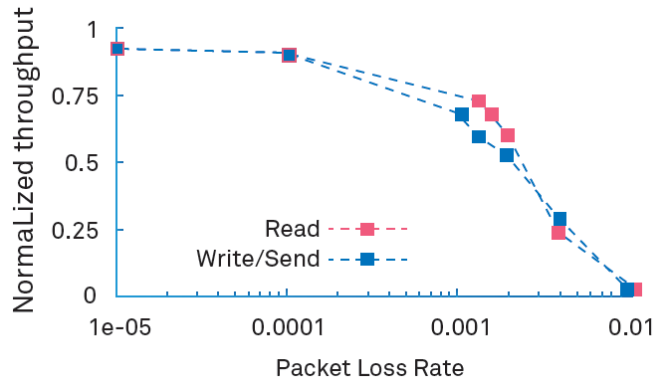


Figure 20 – Impact of packet loss on RDMA throughput

1 loops in the logical topology, these paths naturally contain loops in the physical topology. A PFC
 2 deadlock in the network can completely halt network traffic.

3 Consider the example in Figure 21. The figure shows four phases of PFC deadlock creation. In phase
 4 1, four flows are equally load balanced across the Clos fabric and the network is running smoothly.
 5 In phase 2, the red cross indicates a transient or permanent fault in the topology, such as link failure,
 6 port failure, or route failure. Due to the failure, in the example, traffic between H1 and H7 (green
 7 and yellow lines) is re-routed. The re-routing pushes more traffic through leaves 2 and 3 causing a
 8 potential overflow in spine 1 and spine 2 as shown in phase 2. In the example we assume the
 9 pressure on spine 1 occurs first. To avoid loss, the spine 1 switch issues PFC towards leaf 3, shown
 10 in phase 3. Traffic in leaf 3 now backs up, causing further backups around the topology and a
 11 cascade of PFC messages along the loop backward towards the original point of congestion. Phase
 12 4 shows the resulting PFC deadlock.

13 When the network size is small, the probability of PFC deadlock is low. However, at larger scale and
 14 with the high-performance requirements of the RoCE protocol, the probability of PFC deadlock
 15 increases significantly. Achieving larger scale and optimal performance is a key objective of the

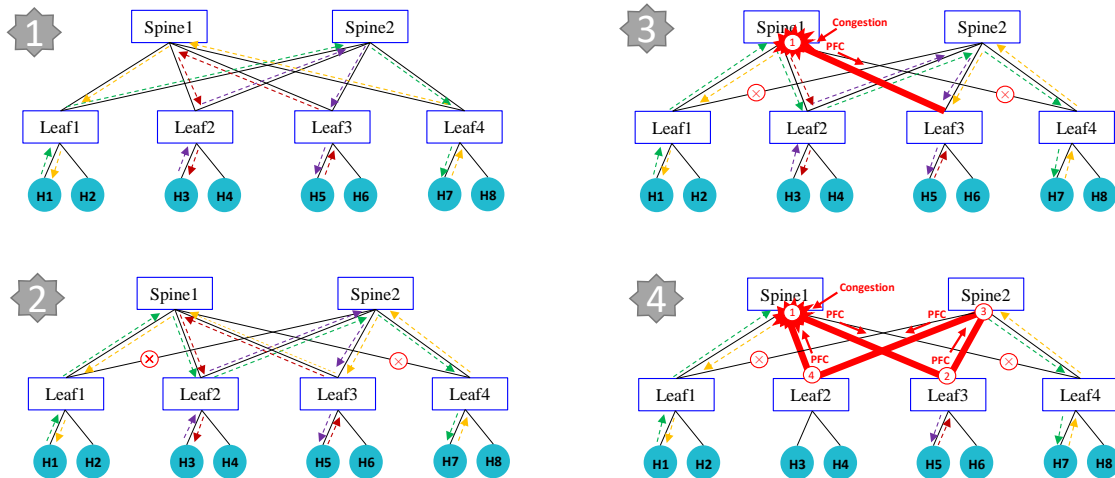


Figure 21 – Example PFC Deadlock

1 Intelligent Lossless Data Center Network of the future. Section 5 discusses a possible new
 2 technology for PFC deadlock prevention.

3 **Congestion control issues in large-scale data center networks**

4 RDMA technology was initially used by customers in constrained, conservative, small scale
 5 environments such as high-performance cluster computing or targeted storage networks. Tuning
 6 the resources required for the dedicated environment was manageable by the network operator,
 7 at least to some degree. However, the performance advantages of RDMA have proven useful in
 8 many application environments and there is a strong desire to use RDMA in a large-scale. Figure 22
 9 shows an example of a large-scale RoCE network. In the example, the entire data center network is
 10 based on Ethernet. The computing cluster and storage cluster use the RDMA protocol while the X86
 11 server cluster uses traditional TCP/IP.

12 In the large-scale data center network scenario, TCP and RoCE traffic can traverse common parts of
 13 the network for several different reasons. Traditional web-based applications using high-speed
 14 storage backends mix end-user TCP requests with RDMA storage requests to read and write data.
 15 The management and software-defined control plane of RDMA devices is typically based on TCP
 16 while using RoCE for data communications. AI/ML applications use RoCE to interconnect GPUs and
 17 CPUs, but still may be using TCP-based storage solutions. This leads to multiple combinations of
 18 TCP and RoCE between computing-and-computing, storage-and-storage, and computing-and-
 19 storage systems.

20 In theory, separating TCP and RoCE traffic within the network should be easy. IEEE Std 802.1Q
 21 defines 8 classes of service that can map to 8 queues with differing queue scheduling algorithms.
 22 Different switch queues can be used to isolate the different traffic types. While the queues and the
 23 buffer management are implemented in hardware on the switch chip, there is a performance and
 24 cost tradeoff problem with the memory. Allocating sufficient dedicated memory to each queue on
 25 each port to absorb microbursts of traffic without incurring packet loss can be too expensive and

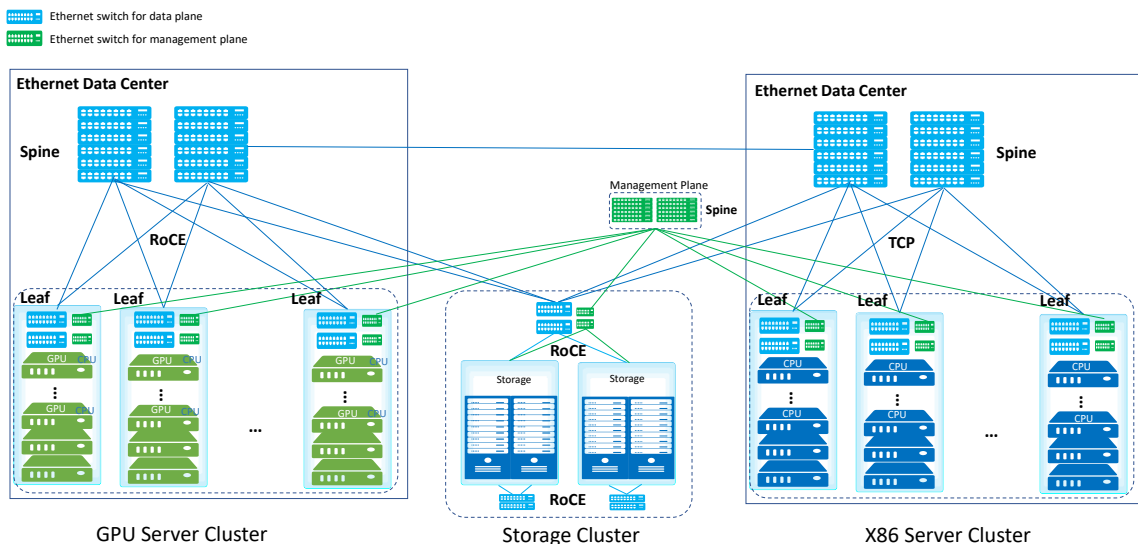


Figure 22 – RoCE application in large-scale data center networks

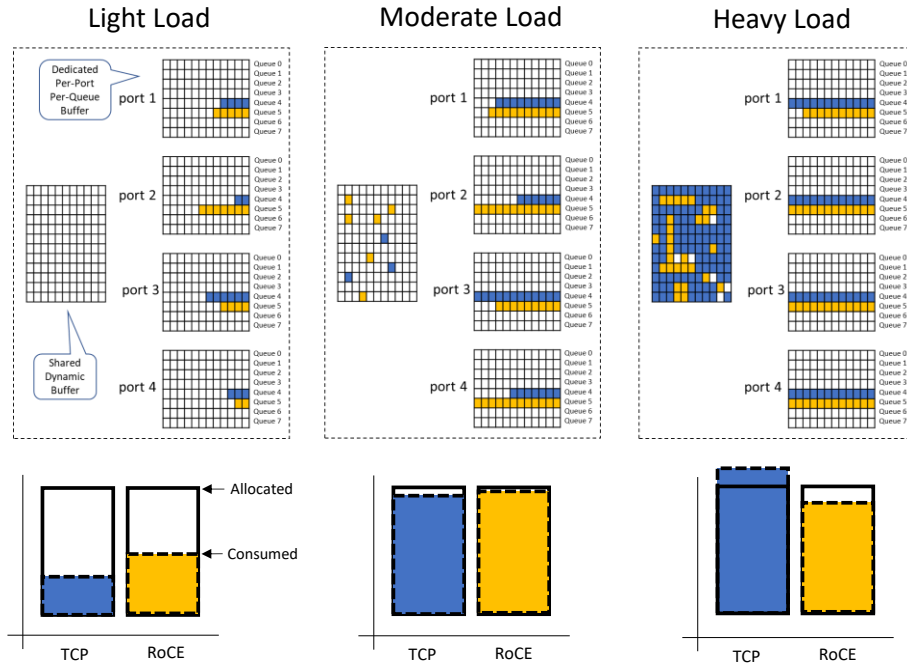


Figure 23 – TCP and RoCE coexistence with smart buffering.

1 technically challenging as the number of ports per switch chip goes up. To address this tradeoff,
 2 switch chip vendors implement a smart buffering mechanism that allows for a hybrid of fixed and
 3 shared buffers.

4 The core idea of smart buffering is the creation of a dynamic shared buffer. The goal is to optimize
 5 buffer utilization and burst absorption by reducing the number of dedicated buffers while providing
 6 a dynamic and self-tuning shared pool across all ports to handle temporary bursts [23].

7 An example smart buffer architecture, as shown in Figure 23. Each port has a fixed number of
 8 dedicated buffers for each of its queues and a common dynamic pool of centralized surplus buffers.
 9 The approach considers that congestion in a typical data center environment is localized to a subset
 10 of egress ports at any given point in time and rarely occurs on all ports simultaneously. This
 11 assumption allows the centralized on-chip buffer to “right-size” the memory usage for overall cost
 12 and power consumption while still providing resources to congested ports exactly when needed by
 13 deploying self-tuning thresholds.

14 Contrasted with static per-port buffer allocation schemes found in other switch architectures, the
 15 smart buffer approach significantly improves buffer utilization and enables better performance for
 16 data center applications. However, the shared dynamic pool has consequences on traffic class
 17 isolation in congested situations. TCP and RoCE flows may impact one another when they traverse
 18 common links, even if they are using separate traffic classes on those links. TCP and RoCE use
 19 different congestion control mechanisms, different re-transmission strategies and different traffic
 20 class configuration (lossless verse lossy). The algorithms and configurations can lead to unfair
 21 sharing of the common resource. Figure 23 shows the problem when the switch is under heavy
 22 load. Network operators allocate the network bandwidth to different traffic classes based on the
 23 service requirements of the network, but over time and during periods of congestion the bandwidth

1 allocations cannot be met. The different congestion control methods create different traffic
2 behavior that impacts the smart buffering mechanism's ability to fairly allocated the dynamic shared
3 buffer pool. In this case, TCP preempts RoCE bandwidth, even when it is assigned to separate traffic
4 classes. The RoCE flow completion delay has been seen to increase by 100 times. ODCC conducted
5 several tests to verify the problem of traffic coexistence [24].

6 **Configuration complexity of congestion control algorithms**

7 Historically, HPC data center networks were small in scale and optimized through manual
8 configuration. However, a goal of the Intelligent Lossless Data Center Network is to enable HPC and
9 AI data centers to grow to cloud scale and be provisioned through automation. Manual
10 configuration and hand tuning parameters are not possible at cloud scale, but the proper operation
11 of the HPC data center requires network wide consistent configuration of several attributes. Some
12 of the key attributes include:

- 13 • Consistent mapping of network priorities to switch traffic classes (i.e. switch queues).
- 14 • Consistent assignment of application traffic to network priorities.
- 15 • Consistent enablement of PFC on lossless traffic classes.
- 16 • Bandwidth allocations for traffic classes using Enhanced Transmission Scheduling (ETS).
- 17 • Buffer threshold settings for PFC and ensuring there is enough headroom to avoid loss.
- 18 • Buffer threshold settings for ECN marking.

19
20 The IEEE 802.1 Working Group defined the Data Center Bridging eXchange protocol (DCBX) to
21 automate the discovery, configuration, and misconfiguration detection of many of the data center
22 network configuration attributes. DCBX leverages the Link Layer Discovery Protocol (LLDP) to
23 exchange a subset of configuration attributes with a network peer, and if the peer is 'willing' to
24 accept recommended settings, the two peers can create a consistent configuration. This consistent
25 configuration can propagate across the entire data center network if all devices are running DCBX.
26 The protocol, however, does not exchange all key attributes for a data center network. In particular,
27 it does not enable the automatic setting of buffer thresholds, which can be quite complex to
28 determine and critical to the proper operation of the network.
29

30 **Adaptive PFC Headroom Calculation**

31 The PFC buffer threshold determines when pause frames are sent as seen in Figure 24. If the
32 receiver's buffer fills past the XOFF threshold, the receiver sends a pause frame. When the buffer
33 drains and empties below the XON threshold, the receiver may send an un-pause frame canceling
34 the previous pause or it may simply let the original pause timeout. The XOFF threshold must be set
35 in such a way to allow in-flight frames to be received. The buffer memory available beyond the
36 XOFF threshold is often called headroom and must be available to ensure lossless operation. Finding
37 the best XON/XOFF thresholds can be tricky. Overestimating the threshold is not practical because
38 it wastes precious switch memory and reduces the number of lossless traffic classes that can be
39 supported. Underestimating the threshold leads to packet loss and poor performance for protocols
40 such as RoCE. Finding the optimal setting is difficult because it requires the complex calculation of
41 many obscure parameters [25].
42
43
44

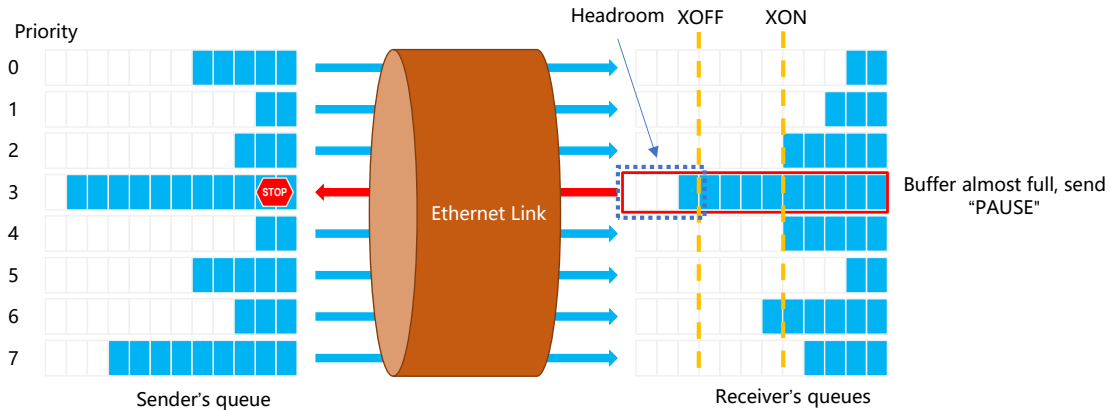


Figure 24 – Priority-based Flow Control (PFC)

- 1 Some of these obscure parameters include:
- 2 • Maximum frame size on the network
- 3 • Speed of the link
- 4 • The length of the cable
- 5 • Internal switch and transceiver latency
- 6 • Response time of sender
- 7 • Internal memory cell size of the receiver's buffer architecture

8

9 Clearly these parameters are not something a network operator can easily obtain. Many are

10 internal to the switch implementation and will differ from vendor to vendor. In addition, the

11 propagation delay, which includes the product of the link speed and cable length, can vary on every

12 port of the network. With thousands of ports to configure, a network operator will benefit from an

13 automated solution that configures PFC headroom.

14 **Dynamic ECN Threshold Setting**

15 The threshold for marking Explicit Congestion Notification (ECN) bits in congested packets is another

16 important configuration setting for the smooth operation of the network. As shown in Figure 18

17 above, setting the ECN threshold low helps achieve low latency, but at the cost of high throughput

18 for larger flows. Setting a high ECN threshold has better performance for throughput-oriented traffic

19 but slows down flow completion time for latency-sensitive smaller flows. As workloads change

20 within the data center network an ideal solution is to dynamically adjust the ECN threshold to

21 balance the tradeoff between high throughput and low latency.

22

23 The congestion control algorithms enabled by ECN involve collaboration between network adapters

24 and network switches. The ECN thresholds in switches and rate reduction and response parameters

25 on NICs and protocol stacks on end stations need to be coordinated as the workload changes. This

26 coordination can result in an untenable set of configuration parameters that need to be updated in

27 real-time. Many network operators only use a recommended static configuration based on the

28 experience of engineers over time. However, the static configuration does not adapt to real-time

29 changes in network traffic that are driven by measurable fluctuations in an application's I/O and

30 communication profile. Different static settings can result in different service performance for the

31 same application and using the same settings for different applications can result in sub-optimal

1 performance for the aggregate of applications on the data center network. Measuring the
2 characteristics of network traffic for the set of the application I/O and communication profiles can
3 lead to a predictive algorithm that dynamically adjusts the ECN threshold in switches and the rate
4 reduction and response parameters at end-stations.

5

6 **New technologies to address new data** 7 **center problems**

8 **Hybrid transports for low latency and high throughput**

9 Traditional data center transport protocols, such as DCTCP [26] and RoCEv2 with DCQCN [19] are
10 sender driven. They attempt to measure and match the instantaneous bandwidth available along
11 the path by pushing data into the channel and awaiting feedback or measurements from the
12 receiver. They continue to push more and more data into the channel until congestion is
13 experienced, at which point they reduce their sending rate to avoid packet loss. There can be many
14 methods of determining when congestion is experienced and how to adjust the sending rate, but
15 the basic premise of sender driven transports is the same – continue to adjust the sending rate up
16 or down based upon an estimation of the available channel bandwidth. This is a very well-known
17 and mature approach to transport congestion control that has been shown to be successful in highly
18 diverse networks such as the Internet. Accurately estimating of the available bandwidth depends,
19 not only, on detecting congestion, but on creating it. Congestion signal delays and untimely
20 adjustments to the sending rate can cause fluctuations to queue depths, leading to variance in
21 throughput and latency. Large buffers in routers and switches can absorb these fluctuations to
22 avoid packet loss.

23 A receiver driven transport, such as ExpressPass [27], can be used to avoid fluctuations in queue
24 depths and minimize buffering along the path from sender to receiver. With receiver driven
25 transports, the sender's transmissions are paced by the receiver's schedule. A request-grant or
26 credit-based protocol is used to pace the sender and avoid congestion while fully utilizing network
27 bandwidth. The approach is especially good at handling incast congestion where the receiver is
28 overrun by multiple simultaneous senders. The challenge with receiver driven transports is that the
29 receiver must now estimate the available bandwidth along the path. Similar techniques for
30 congestion detection can be used and the receiver driven approach as the advantage of receiving
31 those congestion signals first. Perhaps a more significant challenge with receiver driven transports
32 is the inherent delay built into the initial buffer request by the sender. The initial request-grant
33 exchange penalizes small flows which, in most cases, are latency sensitive and constitute the
34 majority of flows in the data center network.

35 A hybrid driven transport, such as NDP [28] or Homa [29], attempts to borrow the best qualities
36 from sender driven and receiver driven transports to reduce latency and increase throughput by
37 avoiding congestion. A hybrid approach allows the sender to transmit a certain amount of
38 unscheduled traffic into the network without waiting for a buffer grant by the receiver, but then it

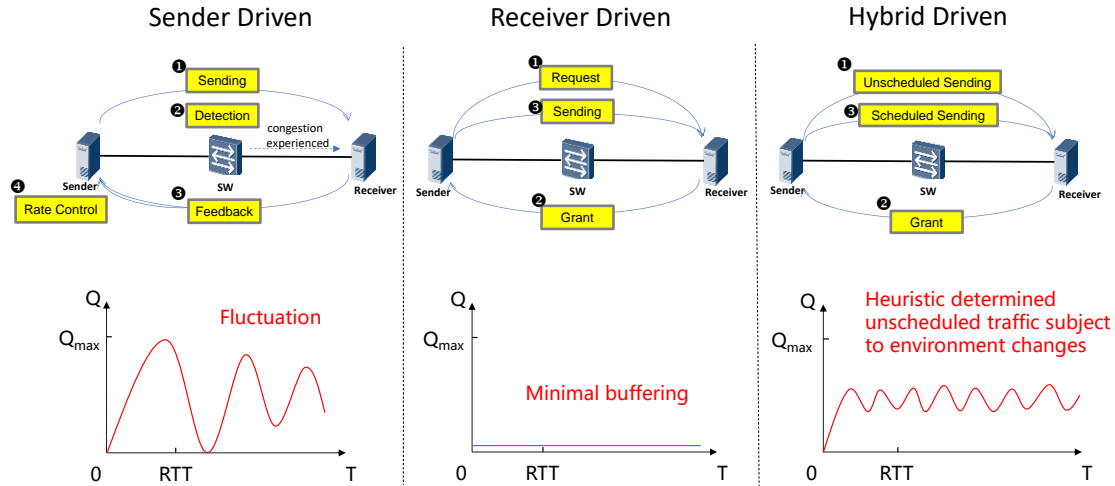


Figure 25 – Transport styles with conceptual network buffering implications

1 must transition to a scheduled receiver driven approach after the unscheduled traffic is sent. The
 2 unscheduled traffic has no additional latency penalties and benefits small flows but can create
 3 minor fluctuations in buffer occupancy which can lead to moderate packet loss. Since the amount
 4 of unscheduled traffic is small, the overall buffer occupancy remains low which leads to more
 5 bounded latency and low packet loss. Adjusting the amount of unscheduled traffic based on
 6 heuristics helps tune the network for high throughput and low latency while maintaining low buffer
 7 utilization. Figure 25 shows the high-level approach to each of the different transport types and a
 8 conceptual graph of buffer utilization over time.

9 PFC deadlock prevention using topology recognition

10 Traffic on a well-balanced Clos networks is loop free and typically flows from uplink to downlink at
 11 ingress and downlink to uplink at egress. However, rerouting occurs when transient link faults are
 12 detected, and traffic may flow from uplink to uplink as shown in Figure 21. According to [22], the
 13 probability of rerouted traffic is approximately 10^{-5} . While 10^{-5} is not a high probability, given the
 14 large traffic volume and the large scale of data center networks the chance of a deadlock occurring
 15 is possible and even the slightest probability of a deadlock can have dramatic consequences. PFC
 16 deadlocks are real! The larger the scale, the higher the probability of PFC deadlock, and the lower
 17 the service availability from this critical resource.

18 A mechanism to prevent PFC deadlock involves discovering and avoiding CBD loops. The core idea
 19 of the deadlock-free algorithm is to break the circular dependency by identifying traffic flows that
 20 create it. The first step in achieving this is to discover the topology and understand the port
 21 orientation of every switch port in the network. An innovative distributed topology and role auto-
 22 discovery protocol is used to identify network locations and roles of across the data center network.

23 The topology and role discovery protocol automatically determines a device's level within the
 24 topology and the orientation of each of the device's ports. The level within the topology is defined
 25 as the number of hops from the edge of the network. For example, a server or storage endpoint is
 26 at level 0 and the top-of-rack switch connected to that server or storage endpoint is at level 1. The

Discovery protocol exchange automatically determines:

1. Topology level of devices in network
 - 0 = End-station or server edge
 - 1 = Leaf
 - n+1 = Spine
2. Port orientation for each link
 - Uplink
 - Downlink
 - Crosslink

HINT: Servers are always at level 0 with uplinks.

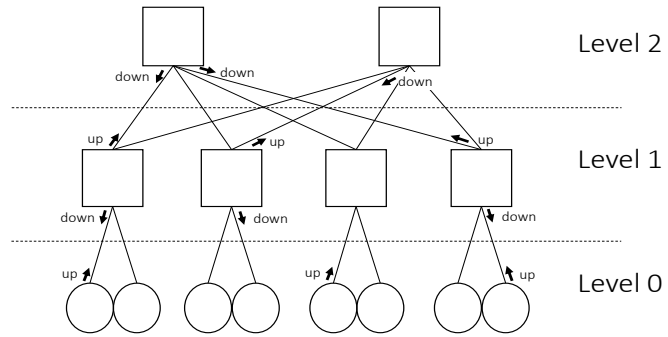


Figure 26 – Topology and Role Discovery

1 port orientation of a port can be either an uplink, downlink or a crosslink. An uplink orientation, for
 2 example, is determined for a port of a device that is connected to another device at a higher level.

3 The protocol starts out by recognizing known conditions. Servers and storage endpoints are always
 4 at level 0 and their port orientation is always an uplink. Switches are initialized without any
 5 knowledge of their level or port orientation, but as the information is propagated by a discovery
 6 protocol, the algorithm converges upon an accurate view. Figure 26 shows the resulting topology
 7 and role discovery in a simple Clos network.

8 Once the protocol has recognized the topology and port roles, the deadlock free mechanism can
 9 identify potential CBD points in the network and adapt the forwarding plane to break the buffer
 10 dependencies. Figure 27 shows how potential CBD points in the topology can be recognized. In a
 11 properly operating Clos network, there is no CBD and flows will typically traverse a switch ingress
 12 and egress port pair that has three of four possible port orientation combinations. The flow may
 13 pass from a port oriented as a downlink to a port oriented as an uplink. In the spine of the network,
 14 the flow may pass from a port oriented as a downlink to another port oriented as a downlink.
 15 Finally, as the flow reaches its destination, the flow may pass from a port oriented as an uplink to a
 16 port oriented as a downlink. A CBD may exist in the case where a flow has been rerouted and now
 17 passes from a port oriented as an uplink to another port oriented as an uplink.

18 After recognizing the CBD point, the forwarding plane is responsible for breaking the CBD. The CBD
 19 exists because a set of flows are using the same traffic class and are traversing a series of switches
 20 that now form a loop due to the flow rerouting. The buffer dependency is the shared buffer memory

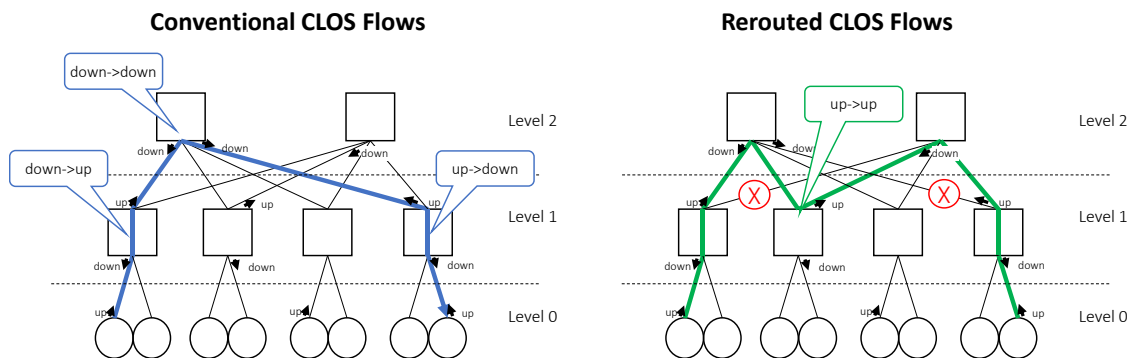


Figure 27 – Identifying CBD points in rerouted flows.

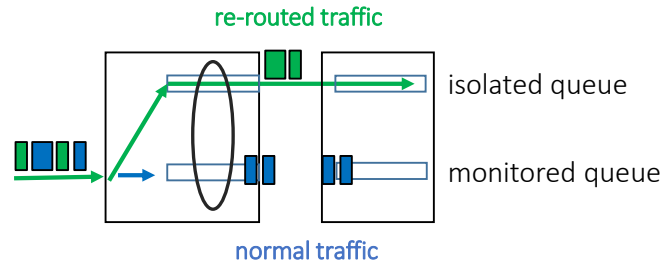


Figure 28 – Queue switch according to CBD reroute flow recognition.

1 of the common traffic class (i.e. switch queue). To break the CBD, packets of the rerouted flow need
 2 to be forwarded to a separate queue. These packets can be identified because they are flowing from
 3 a port oriented as an uplink to another port oriented as an uplink. Figure 28 illustrates the process
 4 of queue remapping within the switch. In the example, the remapping of the green flow to an
 5 isolated queue will lead the elimination of PFC deadlock. The different flows can safely pass-through
 6 different queues at the point of a potential CBD.

7 ODCC, in participation with many network vendors, conducted tests to verify the deadlock free
 8 algorithm [24].

9 Improving Congestion Notification

10 A state-of-the-art congestion control mechanism for the RoCEv2 protocols in today's data centers
 11 is Data Center Quantized Congestion Notification (DCQCN) [19]. DCQCN combines the use of ECN
 12 and PFC to enable a large-scale lossless data center network. Figure 29 shows the three key
 13 components of DCQCN; a reaction point (RP), a congestion point (CP) and a notification point (NP).

14 Reaction Point (RP)

15 The RP is responsible for regulating the injection rate of packets into the network. It is typically
 16 implemented on the sending NIC and responds to Congestion Notification Packets (CNP) sent by the
 17 NP when congestion is detected in the network. When a CNP is received, the RP will decrease the
 18 current rate of injection. If the RP does not receive a CNP within a specified period, it will increase
 19 the transmit rate using a quantized algorithm specified by DCQCN.

20 Congestion Point (CP)

21 A CP is included in the switches along the path between the transmitter (RP) and the receiver (NP).
 22 The CP is responsible for marking packets with ECN when congestion is detected at an egress queue.

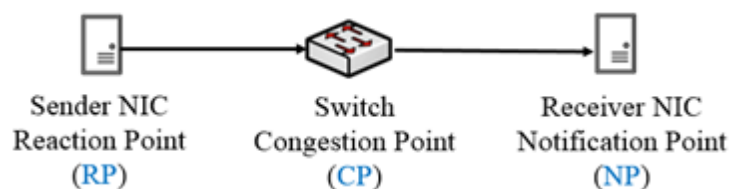


Figure 29 – Three parts of RoCE congestion control using DCQCN

1 Congestion is determined by looking at the egress queue length and evaluating it against
 2 configurable thresholds (K_{min} and K_{max}). When the queue length is less than K_{min} , traffic is not
 3 marked. When the queue length is greater than K_{max} , all packets passing through the queue are
 4 marked. When the queue length is between K_{min} and K_{max} , the marking probability increases
 5 according to the extent of the queue length, as specified by DCQCN.

6 Notification Point (NP)

7 The NP is responsible for informing the RP that congestion has been experienced by packets of a
 8 flow while traversing the network. When a data packet with an ECN flag arrives at a receiver, the
 9 NP sends a CNP packet back to the RP if one has not already been sent in the past N microseconds.
 10 It is possible to set N to 0 such that the NP will send a CNP for each packet with an ECN flag set.

11 As data center networks scale to larger sizes and support an increased number of simultaneous
 12 flows, the average bandwidth allocated to each flow can become small. Flows experiencing
 13 congestion in this environment may have their packets delayed, causing the arrival of ECN markings
 14 at the NP to also be delayed. If the rate of arrival of ECN marked packets is greater than the interval
 15 the RP uses to increase the rate of injection a problem may occur. The problem is that the RP will
 16 begin increasing the rate of injection when it should decrease the rate since the flow is congested
 17 and the missing CNP messages have simply been delayed. In this case, the end-to-end congestion
 18 control loop is not functioning correctly.

19 The impact of end-to-end congestion control loop failure in a lossless network can lead to
 20 congestion spreading. The unwanted congestion causes an increase PFC messages and an increase
 21 in the amount of time links are paused. These PFC messages further delay the propagation of ECN
 22 marked packets and only make the problem worse. In this scenario, the combination of PFC and
 23 ECN become ineffective.

24 One possible solution to this problem is for the network to intelligently supplement the CNP packets
 25 sent by the NP. The intelligence involves considering the congestion level at the egress port, the
 26 interval of the received ECN marked packets, and the interval of the DCQCN rate increase by the RP.
 27 After receiving an ECN marked packet, the CP keeps track of the frequency of received ECN marked
 28 packets as well as the packet sequence number for a congested flow. When the CP egress queue is
 29 congested and the received flow has experienced congestion further upstream, the CP may
 30 proactively supplement the CNP depending upon the rate of received ECN marked packets and the
 31 interval of the DCQCN rate increase at the RP. The CP is aware that ECN marked packets are delay

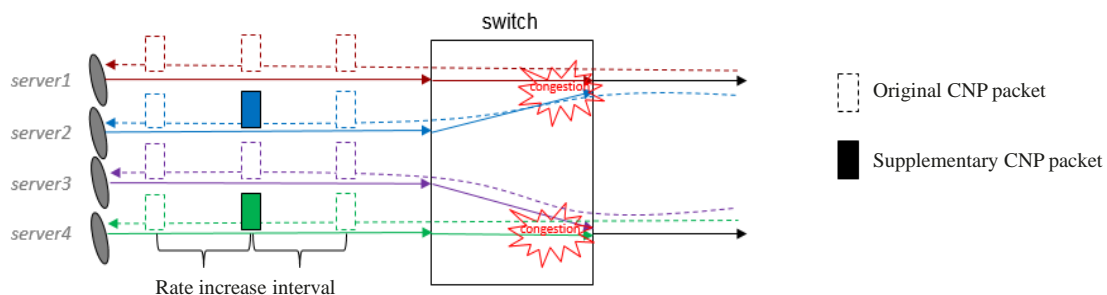


Figure 30 – Intelligent Supplemental CNP

1 and that subsequent CNP packets from the NP will be further delayed, so the supplemental CNP
2 messages prevent the end-to-end congestion control loop failure. The supplemental CNP operation
3 is performed only when the CP egress queue is severely congested, thus latency and throughput are
4 not affected when DCQCN is operating in a normal non-congested state. The solution is shown in
5 Figure 30.

6 ODCC tested the enhanced congestion control mechanism and the effect is beneficial [30].
7 According to the test results, performance is improved by more than 30% (TCP:RoCE = 9:1 scenario).

8 **Addressing configuration complexity of congestion control algorithms**

9 With thousands of switches and tens of thousands of ports to configure, network operators need
10 automated solutions to properly configure the parameters responsible for managing congestion
11 control in the data center network. The Data Center Exchange Protocol (DCBX) defined by IEEE
12 802.1 made great strides in simplifying some of the configuration and error detection, however,
13 more is needed. Automated solutions for setting and adapting switch buffer thresholds are needed.

14 **Buffer optimization to reduce the complexity of PFC headroom configuration**

15 The key to successful PFC XOFF threshold setting is assuring there is enough headroom to absorb
16 the in-flight data once the pause frame has been issued. There is a natural delay between the time
17 when the pause frame can be sent, and the sender stops transmitting data. The headroom must
18 provide enough buffer to receive data during this delay, but the calculation for the amount of
19 memory needed can be quite complex. Annex N of IEEE Std 802.1Q-2018 [18] provides the technical
20 details of this calculation. Many of the components of the delay are internal to the switch
21 implementations and remain relatively static. For example, interface and higher-layer delay do not
22 vary for a particular configuration and implementation. These static components of delay can be
23 communicated between peers on the network, but currently there is no standard protocol that
24 allows this. The propagation delay for the medium is dependent upon the transmission speed and
25 the length of the cable. To accurately obtain this component of delay a measurement is required.

26 The Time Sensitive Networking (TSN) Task Group of the IEEE 802.1 Working Group has defined IEEE
27 Std 802.1AS-2020 Timing and Synchronization for Time-Sensitive Application [31]. A small subset
28 of this specification, along with optional time-stamping support in IEEE Std 802.3 can be used to
29 measure cable delay between two peers on a point-to-point link. IEEE Std 802.1AS, however,
30 targets time-sensitive applications in constrained environments such as audio/visual, industrial, and
31 automotive networks. Its primary focus is to enable a Precision Time Protocol (PTP) used to
32 synchronize clocks throughout the computer network. While a fine-grained synchronized clock
33 could be valuable in a data center, the burden for supporting the complete set of IEEE Std 802.1AS
34 functions in data center switching silicon could be onerous. The delay measurement facilities of IEEE
35 Std 802.1AS, on the other hand, are useful in the data center to assist in the auto-configuration of
36 PFC thresholds. Having the ability to discover and communicate this capability between peers,
37 along with other DCBX attributes, would be necessary for full automation of the configuration
38 settings.

39

40

1 Intelligent ECN threshold optimization

2 The ECN threshold determines how aggressively a switch will indicate that packets are experiencing
3 congestion and subsequently how frequently the sending station may need to adjust transmission
4 rate. The optimal threshold setting depends on the current state of the network and the types of
5 communication flows that are competing for common resources. As previously discussed, a low
6 threshold setting can benefit latency-sensitive smaller flows and a high threshold setting can have
7 better performance for throughput-sensitive larger traffic flows. The mix of these flows and their
8 communication patterns is constantly changing but has been shown to be predictable using
9 machine learning techniques that model application traffic behavior [32] [33] [34]. A machine
10 learning model that predicts data center network traffic patterns could be used to dynamically
11 adjust ECN thresholds to optimize the trade-off between low-latency and high-throughput. The
12 unfair sharing of the dynamic pool of memory in the smart buffering scheme can also be address by
13 dynamically adjusting the ECN threshold differently for TCP and RoCE traffic.

14 To train a model of network traffic patterns in the data center an AI/ML system needs an abundance
15 of real-time data from the network. The data acquisition system needs to capture the temporal
16 relationships between network devices across the data center at large scale. Traditional network
17 monitoring systems based on SNMP and/or NetConf use polling to “pull” data from the devices.
18 This approach has scaling issues, increases network traffic and makes it more difficult to correlate
19 the collected data. What is needed is a telemetry stream of essential parameters flowing directly
20 from the network devices. Telemetry is a network monitoring technology developed to collect
21 performance data quickly from physical or virtual devices. Telemetry differs from traditional
22 network monitoring technologies as it enables network devices to “push” high-precision
23 performance data to a data repository in real time and at high speeds. This improves the utilization
24 of device and network resources during data collection.

25 Using the telemetry stream of data from network devices, an AI/ML system can build a model that
26 monitors the congestion status of all queues on the entire network. The stream of parameters can
27 be used to train and retrain the network model, allowing inference engines on the network devices
28 to predict changes in the data center environment and self-adjust their ECN threshold. Inputs to the
29 model can extend well beyond the existing counters obtained by traditional network monitoring
30 systems. Essential input parameters might include:

- 31 • A snapshot of the incast ratio (N:1) at an egress port
- 32 • The mix of mice and elephants flows at an ingress port
- 33 • The rate change in switch buffer occupancy

34 Other more traditional network metrics might include:

- 35 • Port-level information
 - 36 ○ Sent and received bytes
 - 37 ○ Sent and received packets
 - 38 ○ Discarded packets in the transmit and receive directions
 - 39 ○ Received unicast packets, multicast packets, and broadcast packets
 - 40 ○ Sent unicast packets, multicast packets, and broadcast packets
 - 41 ○ Sent and received error packets
 - 42 ○ Ingress port bandwidth usage and egress port bandwidth usage
 - 43 ○ ECN packets

- 1 • Queue-level information
- 2 ○ Egress queue buffer utilization
- 3 ○ Headroom buffer utilization
- 4 ○ Received PFC frames
- 5 ○ Sent PFC frames

6 Another type of telemetry, known as in-band telemetry, provides real-time information about an
7 individual packet's experience as it traverses the network. The information is collected and
8 embedded into the packet headers by the switch data plane without involving the control plane.
9 The amount of information collected is more limited than traditional telemetry because it must be
10 included with the contents of the original data packet, which has a finite size. However, the
11 information within the packet is directly related to the network state that the packet observed
12 during its existence within the network. Each hop along the path can be instructed to insert local
13 data representing the switch hop's state. Essential information might contain:

- 14 • Ingress and egress port numbers
- 15 • Local timestamps at ingress and egress
- 16 • Egress link utilization
- 17 • Egress queue buffer utilization

18 An AI model that takes real-time telemetry input from the local device can predict the adjustments
19 needed to the ECN threshold for the desired balance between low-latency and high-throughput.
20 The in-band telemetry signals can be examined with an objective of rapidly communicating
21 appropriate congestion signals to the sending sources thus avoiding packet loss and long tail latency
22 with flow completion times.

23 6

24 **Standardization Considerations**

25 Two important standards development organizations for the future technologies discussed above
26 are the IEEE 802 LAN/MAN Standards Committee and the Internet Engineering Task Force (IETF).

27 The IEEE 802 LAN/MAN Standards Committee develops and maintains networking standards and
28 recommended practices for local, metropolitan, and other area networks, using an open and
29 accredited process, and advocates them on a global basis. The most widely used and relevant
30 standards to this report are for Ethernet, Bridging, Virtual Bridged LANs and Time Sensitive
31 Networking. The IEEE 802.1 Working Group provides the focus for Bridging, Virtual Bridged LANs
32 and Time Sensitive Networking.

33 The Internet Engineering Task Force (IETF) is the premier Internet standards body, developing open
34 standards through open processes. The IETF is a large open international community of network
35 designers, operators, vendors, and researchers concerned with the evolution of the Internet
36 architecture and the smooth operation of the Internet. The technical work of the IETF is done in
37 Working Groups, which are organized by topic into several Areas. The most relevant IETF Areas for
38 the future technologies discussed above are likely the Internet Area (int), the Routing Area (rgt) and

1 the Transport Area (tsv). A parallel organization to the IETF is the Internet Research Task Force (IRTF)
2 which focuses on longer term research issues related to the Internet. The IRTF is comprised of
3 several focused and long-term Research Groups, of which the most relevant for this report are the
4 Internet Congestion Control Research Group (iccrgr) and the Computing in the Network Research
5 Group (coinrg).

6 The IEEE 802 and IETF/IRTF have a long history of working together on developing inter-related
7 standards and technology. A standing coordination function between the Internet Architecture
8 Board (IAB) of the IETF and the leadership of the IEEE 802 Working Groups is currently place [35].
9 Traditionally these two organizations were aligned by layers of the ISO stack, where IEEE 802
10 focused on layer 2 and IETF on layer 3 and above. The lines have blurred over the years, but the
11 two organizations have continued to work together, sharing information, and developing unique
12 and valuable standards.

13 Transport protocols are typically the domain of the IETF, however providing signals from the
14 network could be a result of specifications from IEEE 802.1. A new hybrid transport that optimizes
15 the tradeoff between low latency and high throughput could likely be investigated in the IRTF's
16 Internet Congestion Control Research Group (iccrgr). A proposed standard from this research would
17 most likely be developed by the IETF's Transport Area (tsv). The key to success for the hybrid
18 transport is knowing how to best estimate the amount of unscheduled traffic for the channel and
19 how to rate control the senders in an incast scenario. Congestion signals and resource status along
20 the communication path could be provided by the network switches themselves. In-band telemetry
21 or enhanced ECN signaling by the network switches could provide the needed information and
22 represents an opportunity for specification by the IEEE 802.1 Working Group.

23 PFC deadlock prevention requires an awareness of the network topology and an ability to break a
24 CBD caused by re-routed flows. P802.1Qcz Congestion Isolation has specified a mechanism using
25 LLDP to automatically recognize the level of a switch within the topology as well as the orientation
26 of each port (e.g. uplink, downlink, crosslink). A missing specification is how to recognize flows that
27 are at risk of creating a CBD and how mitigate the CBD. The mechanism specified by P802.1Qcz
28 to adjust the priority of a congesting flow could be used to adjust the priority of a flow at risk of
29 creating a CBD. Further specification of how to use these mechanisms for PFC deadlock prevention
30 could be done by the IEEE 802.1 WG.

31 Network supplemented CNPs, discussed above, augment the DCQCN protocol which currently has
32 no formal specification. Since DCQCN works in conjunction with RoCEv2, the Infiniband Trade
33 Association (IBTA) would be the natural standards organization to complete these enhancements.
34 The general idea of network supplemented CNPs could also be applied to a new IETF hybrid
35 transport protocol and most likely would be investigated by ICCRG and TSVWG. A third alternative
36 is to consider updating the mechanism defined by IEEE Std 802.1Q-2018 in Clause 30 through Clause
37 32 – Congestion Notification. To make Congestion Notification relevant in today's modern data
38 centers the Congestion Notification Messages (CNM) would need to be Layer-3 and routable.

39 Automatically setting the PFC XON/XOFF thresholds requires an accurate measurement of the
40 delays between two ends of a link in the data center. An adaptive PFC headroom algorithm could
41 be defined by the IEEE 802.1 Working Group using or augmenting the facilities already defined by
42 IEEE Std 802.3 for timestamping and IEEE Std 802.1AS for path delay measurements. A solution for
43 the data center is needed to reduce the overhead of lossless mode configuration and the associated

1 chance of error. A mechanism to communicate this capability between peers and an update to the
2 current description of how to manually calculate headroom are excellent candidates for an
3 amendment to IEEE Std 802.1Q.

4 Adjusting the ECN threshold automatically is dependent on recognizing and predicting the current
5 congestive state of the data center network. A rapid response to changing congestion status is
6 needed, but traditional network management approaches can not react quickly enough. Network
7 devices that are armed with an AI model to assist in this prediction rely on the model being well
8 trained from an accurate set of real-time data. Network telemetry can provide a new view of the
9 network state, whether that telemetry data is in-band or streamed from the network devices
10 themselves. Standards for telemetry at Layer-3 and above have historically been specified by the
11 IETF. Currently the IP Performance Measurement (ippm) group within the IETF TSV area is defining
12 In-situ Operations, Administration, and Maintenance (IOAM) [36] that provides in-band telemetry
13 at Layer-3. There are also other related and competing specifications for Layer-3 in-band telemetry
14 [37] [38]. In some environments a Layer-2 solution working in conjunction with Layer-3 may be
15 more appropriate and require standard specifications to support interoperability. This in-band
16 telemetry could be defined by the IEEE 802.1 Working Group. The Operations and Management
17 Area Working Group (opsawg) in the IETF Operations and Management area (ops) is working on a
18 framework for Network Telemetry [39]. This framework is looking at various techniques for remote
19 data collection, correlation, and consumption. For the framework to be successful it is necessary to
20 specify the information that can be extracted from the network. Supporting specifications by Layer-
21 2 devices will be needed from the IEEE 802 and specifications for Layer-3 and above devices will be
22 needed by the IETF.



Conclusion

23

24

25 Data center networks must continue to scale and innovate with new technologies to keep pace with
26 the evolving needs of high-speed computing and storage used for AI and Machine Learning
27 applications. This paper expanded upon the previous report [2] with the exploration of technical
28 challenges and potential new solutions for today's cloud scale high performance computing data
29 centers. We discussed new hybrid transport protocols that better balance the needs of both high
30 throughput and low latency communications for AI and Machine Learning. We described a solution
31 for PFC deadlock prevention using a topology recognition algorithm that leverages the existing and
32 widely deployed Link Layer Discovery Protocol (LLDP). We explored ways to reduce the feedback
33 cycle for congestion notification messages by allows the switches to supplement congestion
34 signaling. We also described approaches to reduce the complexity of switch buffer threshold
35 configuration using automated protocols and artificial intelligence models developed from advance
36 telemetry systems. Together these innovations, with the commitment to openness and
37 standardization, can advance the use of Ethernet as the premier network fabric for modern cloud
38 scale high performance data centers.

39

1
2

Citations

- [1] IEEE, "Nendica Work Item: Data Center Networks," [Online]. Available: <https://1.ieee802.org/nendica-DCN/>. [Accessed 14 05 2020].
- [2] IEEE, "IEEE 802 Nendica Report: The Lossless Network for Data Centers," 17 8 2018. [Online]. Available: <https://xploreqa.ieee.org/servlet/opac?punumber=8462817>. [Accessed 13 05 2020].
- [3] Orange, "Finding the competitive edge with digital transformation," 03 June 2015. [Online]. Available: <https://www.orange-business.com/en/magazine/finding-the-competitive-edge-with-digital-transformation>. [Accessed 1 09 2020].
- [4] J. Wiles, "Mobilize Every Function in the Organization for Digitalization," Gartner, 03 December 2018. [Online]. Available: <https://www.gartner.com/smarterwithgartner/mobilize-every-function-in-the-organization-for-digitalization/>. [Accessed 10 June 2020].
- [5] Huawei, "Huawei Predicts 10 Megatrends for 2025," Huawei, 08 August 2019. [Online]. Available: <https://www.huawei.com/en/press-events/news/2019/8/huawei-predicts-10-megatrends-2025>. [Accessed 10 June 2020].
- [6] J. Handy and T. Coughlin, "Survey: Users Share Their Storage," 12 2014. [Online]. Available: <https://www.snia.org/sites/default/files/SNIA%20IOPS%20Survey%20White%20Paper.pdf>. [Accessed 14 05 2020].
- [7] Huawei, "AI, This Is the Intelligent and Lossless Data Center Network You Want!," 13 March 2019. [Online]. Available: <https://www.cio.com/article/3347337/ai-this-is-the-intelligent-and-lossless-data-center-network-you-want.html>. [Accessed 14 05 2020].
- [8] E. K. Karuppiah, "Real World Problem Simplification Using Deep Learning / AI," 2 November 2017. [Online]. Available: https://www.fujitsu.com/sg/Images/8.3.2%20FAC2017Track3_EttikanKaruppiah_RealWorldProblemSimplificationUsingDeepLearningAI%20.pdf. [Accessed 14 05 2020].
- [9] O. Cardona, "Towards Hyperscale High Performance Computing with RDMA," 12 June 2019. [Online]. Available: https://pc.nanog.org/static/published/meetings/NANOG76/1999/20190612_Cardona_Towards_Hyperscale_High_v1.pdf. [Accessed 14 05 2020].
- [10] T. P. Morgan, "Machine Learning Gets An Infiniband Boost With Caffe2," 19 April 2017. [Online]. Available: <https://www.nextplatform.com/2017/04/19/machine-learning-gets-infiniband-boost-caffe2/>. [Accessed 14 05 2020].

- [11] Z. Jai, Y. Kwon, G. Shipman, P. McCormick, M. Erez and A. Aiken, "A distributed multi-GPU system for fast graph processing," in *VLDB Endowment*, 2017.
- [12] Wikipedia, "IEEE 802.3," 5 June 2020. [Online]. Available: https://en.wikipedia.org/wiki/IEEE_802.3. [Accessed 22 July 2020].
- [13] K. Rupp, "42 Years of Microprocessor Trend Data," February 2018. [Online]. Available: <https://www.karlrupp.net/2018/02/42-years-of-microprocessor-trend-data/>. [Accessed 22 July 2020].
- [14] The Linux Foundation, "Open vSwitch," 2016. [Online]. Available: <https://www.openvswitch.org/>. [Accessed 23 July 2020].
- [15] Y. Li, R. Miao, H. H. Liu, Y. Zhuang, F. Feng, L. Tang, Z. Cao, M. Zhang, F. Kelly, M. Alizadeh and M. Yu, "HPCC: high precision congestion control," in *Proceedings of the ACM Special Interest Group on Data Communication (SIGCOMM '19)*, New York, NY, USA, 2019.
- [16] P. Goyal, P. Shah, N. Sharma, M. Alizadeh and T. Anderson, "Backpressure Flow Control," in *Proceedings of the 2019 Workshop on Buffer Sizing (BS '19)*, New York, NY, USA, 2019.
- [17] C. Guo, H. Wu, Z. Deng, G. Soni, J. Ye, J. Padhye and M. Lipshteyn, "RDMA over Commodity Ethernet at Scale," in *In Proceedings of the 2016 ACM SIGCOMM Conference (SIGCOMM '16)*, 2016.
- [18] IEEE, IEEE Std 802.1Q-2018, IEEE Standard for Local and Metropolitan Area Networks — Bridges and Bridged Networks, IEEE Computer Society, 2018.
- [19] Y. Zhu, H. Eran, D. Firestone, C. L. M. Guo, Y. Liron, J. Padhye, S. Raindel, M. H. Yahia and M. Zhang, "Congestion Control for Large-Scale RDMA Deployments," in *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication (SIGCOMM '15)*, London, United Kingdom, 2015.
- [20] M. Karok, J. Golestani and D. Lee, "Prevention of deadlocks and livelocks in lossless backpressured packet networks," *IEEE/ACM Transactions on Networking*, vol. 11, no. 6, p. 11, 2003.
- [21] S. Hu, Y. Zhu, P. Cheng, C. Guo, K. Yan, J. Padhye and K. Chen, "Deadlocks in datacenter networks: Why do they form, and how to avoid them," in *Proceedings of the 15th ACM Workshop on Hot Topics in Networks*, 2016.
- [22] S. Hu, Y. Zhu, P. Cheng, C. Guo, K. Tan, J. Padhye and K. Chen, "Tagger: Practical PFC Deadlock Prevention in Data Center Networks," in *In Proceedings of the 13th International Conference on emerging Networking EXperiments and Technologies (CoNEXT '17)*, 2017.
- [23] S. Das and R. Sankar, "Broadcom Smart-Buffer Technology in Data Center Switches for Cost-Effective Performance Scaling of Cloud Applications," April 2012. [Online]. Available:

- <https://docs.broadcom.com/docs-and-downloads/collateral/etp/SBT-ETP100.pdf>.
[Accessed 24 June 2020].
- [24] ODCC, "ODCC lossless network test report (final draft)," 02 September 2020. [Online]. Available: <http://www.odcc.org.cn/auth/v-1300974311558307841.html>. [Accessed 03 09 2020].
- [25] Cisco Systems, Inc, "Priority Flow Control: Build Reliable Layer 2 Infrastructure," 2009. [Online]. Available: https://www.cisco.com/c/en/us/products/collateral/switches/nexus-7000-series-switches/white_paper_c11-542809.pdf. [Accessed 15 12 2020].
- [26] M. Alizadeh, A. Greenberg, D. Maltz, J. Padhye, P. Patel, B. Prabhakar, S. Sengupta and M. Sridharan, "Data center TCP (DCTCP)," in *ACM SIGCOMM 2010 conference (SIGCOMM '10)*, New York, 2010.
- [27] I. Cho, K. Jang and D. Han, "Credit-Scheduled Delay-Bounded Congestion Control for Datacenters," in *Proceedings of the Conference of the ACM Special Interest Group on Data Communication (SIGCOMM '17)*, New York, 2017.
- [28] M. Handley, C. Raiciu, A. Agache, A. Voinescu, A. W. A. G. Moore and M. Wojcik, "Re-architecting datacenter networks and stacks for low latency and high performance," in *SIGCOMM '17*, Los Angeles, 2017.
- [29] B. Montazeri, Y. Li, M. Alizadeh and J. Ousterhout, "Homa: A Receiver-Driven Low-Latency Transport Protocol Using Network Priorities," 26 03 2018. [Online]. Available: <https://arxiv.org/abs/1803.09615v1>. [Accessed 22 05 2018].
- [30] ODCC, "Lossless Network Test Specifications," 03 September 2019. [Online]. Available: <http://www.odcc.org.cn/download/p-1169553273830920194.html>. [Accessed 01 09 2020].
- [31] IEEE, IEEE Std 802.1AS-2020, IEEE Standard for Local and Metropolitan Area Networks — Timing and Synchronization for Time-Sensitive Applications, IEEE Computer Society, 2020.
- [32] L. Nie, D. Jiang, L. Guo, S. Yu and H. Song, "Traffic Matrix Prediction and Estimation Based on Deep Learning for Data Center Networks," in *2016 IEEE Globecom Workshops (GC Wkshps)*, Washington, DC, 2016.
- [33] X. Cao, Y. Zhong, Y. Zhou, J. Wang, C. Zhu and W. Zhang, "Interactive Temporal Recurrent Convolution Network for Traffic Prediction in Data Centers," *IEEE Access*, vol. 6, pp. 5276-5289, 2018.
- [34] A. Mozo, B. Ordozgoiti and S. Gomez-Canaval, "Forecasting short-term data center network traffic load with convolutional neural networks," *PLoS ONE*, vol. 13(2), no. e0191939. <https://doi.org/10.1371/journal.pone.0191939>, 2018.

- [35] IETF, "IEEE 802 and IETF Coordination Guide," 6 7 2017. [Online]. Available: <https://trac.ietf.org/trac/iesg/wiki/IEEE802andIETFCoordinationGuide>. [Accessed 1 2 2018].
- [36] IETF, "Data Fields for In-situ OAM," 17 12 2020. [Online]. Available: <https://datatracker.ietf.org/doc/draft-ietf-ippm-ioam-data/>. [Accessed 06 01 2021].
- [37] Alibaba; Arista; Barefoot Networks; Dell; Intel; Marvell; Netronome; VMware, "In-band Network Telemetry (INT) Dataplane Specification," 20 04 2018. [Online]. Available: <https://github.com/p4lang/p4-applications/blob/e5d0c4f4c9fe548e83ad91adbd38847c7dce6cfe/docs/INT.pdf>. [Accessed 06 01 2021].
- [38] J. Kumar, S. Anubolu, J. Lemon, R. Manur, H. Holbrook, A. Ghanwani, D. Cai, H. Ou, Y. Li and X. Wang, "Inband Flow Analyzer," 24 04 2020. [Online]. Available: <https://tools.ietf.org/html/draft-kumar-ippm-ifa-02>. [Accessed 06 01 2021].
- [39] IETF, "Network Telemetry Framework," 17 12 2020. [Online]. Available: <https://datatracker.ietf.org/doc/draft-ietf-opsawg-ntf/>. [Accessed 06 01 2021].

1

2