# ZTR-RTT Congestion Control Algorithm Overview

# Table of contents

# Overview

NVIDIA Zero Touch RoCE (ZTR) enables data centers to seamlessly deploy RDMA over Converged Ethernet (RoCE) without requiring any special switch configuration. Built according to the InfiniBand Trade Association (IBTA) RDMA standard and fully compliant with the RoCE specifications, ZTR enables seamless deployment of RoCE. ZTR also boasts performance equivalent to traditional switch-enabled RoCE and is significantly better than traditional TCP-based memory access. Moreover, with ZTR, RoCE network transport services operate side-by-side with non-RoCE communications in ordinary TCP/IP environments.

The new NVIDIA Congestion Control algorithm, Round-Trip Time Congestion Control (RTTCC) allows ZTR to scale to thousands of servers without compromising performance. Using ZTR and RTTCC allows data center operators to enjoy ease-of-deployment and operations together with the superb performance of Remote Direct Memory Access (RDMA) at a massive scale, without any switch configuration.

The new NVIDIA congestion control algorithm, RTTCC, actively monitors network RTT to proactively detect and adapt to the onset of congestion before dropping packets. RTTCC enables dynamic congestion control using a hardware-based feedback loop that provides dramatically superior performance compared to software-based congestion control algorithms. RTTCC also supports faster transmission rates and can deploy ZTR at a larger scale.
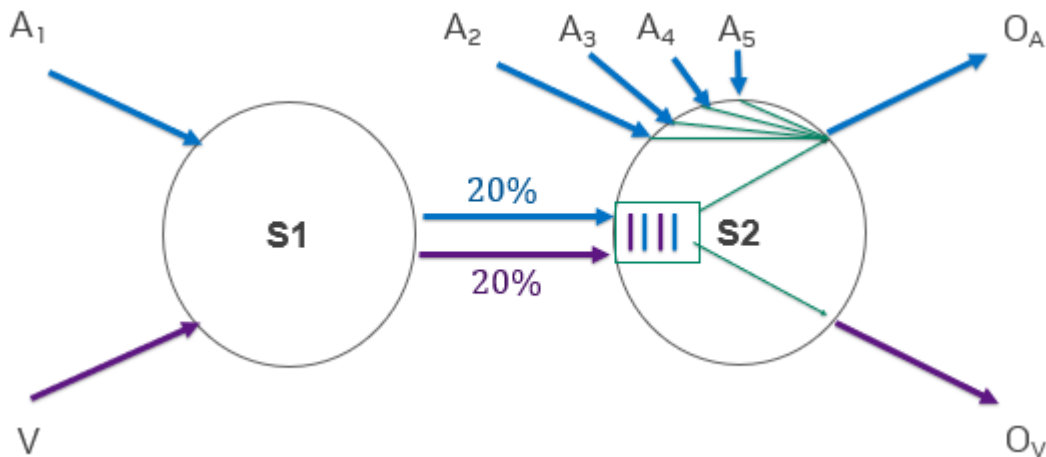
Main ZTR-RTT CC algorithm's characters are :

- Implemented on top of DPA (Data Path Accelerator)

- RTT-based congestion control

- Current default CC algorithm for RoCE

- Demonstrates better performance than DCQCN on HPC and AI workloads

- Maintain DCQCN good performance on storage workload

# Congestion Control Algorithm

Congestion Control provides performance isolation when multiple applications running on the same cluster. Additionally, it prevents congestion spreading when there is a slow receiver, reduce latency in the cluster, improves fairness, prevents parking-lot effects and packet's drop in lossy networks.

ZTR_RTTCC is NVIDIA's default Congestion Control algorithm.

The diagram below shows an example of head of the line blocking scenario.



**Head of the Line Blocking Scenario**

## Datacenter Congestion Control Challenges

The following are the Datacenter Congestion Control challenges:

- Several μ-sec of latency with hundreds of Gbps of bandwidth

    - Congestion buildup is fast, so the congestion loop should be short

- A wide variety of traffic types, topologies and applications

    - Hard to develop an algorithm that suits all

- Congestion Control algorithms are constantly being introduced with new congestion indications

- Hardware implementation is not robust enough

- Software implementation reacts too slow

# ZTR RTTCC Infrastructure



- Efficient and flexible platform to implement CC algorithms
- The code runs on in-data-path microprocessors that interact quickly with the NIC's send and receive pipes
- CNP and RTT measurement flow are included
- Rate limiting and data collection are done by the hardware
- Stateful decisions per destination
- Total decision time of ~1-1.5$\mu sec$

**ZTR RTTCC Infrastructure**

# RTT Measurement Flow



**Sender**

CC Algorithm | Transmit pipe | Receive pipe

RTT req

Wait till after the next burst of the flow

RTT req sent

Create RTT pkt and add req send timestamp

RTT event

**Receiver**

Receive pipe | Transmit pipe

Add req receive timestamp

Reverse Path

Add resp send timestamp

# ZTR RTTCC Algorithm

Congestion indications
  - RTT
  - CNP

Rate update scheme

Evaluate congestion state by comparing RTT to $\frac{1}{\sqrt{rate}}$

Additive Increase

Multiplicative Decrease

Algorithmic windowing

Mimic window behavior by adjusting the rate when RTT changed

Fast reaction on first congestion

# ZTR RTTCC Algorithm

# Document Revision History

| Revision | Date | Description |
|---|---|---|
| 1.0 | September 25, 2024 | Initial release |

NVIDIA Corporation in the U.S. and other countries. Other company and product names may be trademarks of the respective companies with which they are associated.<br/><br/>