

RTGWG
Internet-Draft
Intended status: Standards Track
Expires: 3 September 2025

Z. Hu
Y. Zhu
China Telecom
X. Geng
Huawei
J. Hu
T. Pi
China Telecom
2 March 2025

Fast congestion notification for distributed RoCEv2 network based on
SRv6
draft-hu-rtgwg-rocev2-fcn-00

Abstract

AI services (e.g. distributed model training, separated storage and model training) drive the need to transmit RDMA packets through SRv6 tunnels in WAN. RoCEv2 is the most popular open standard for achieving RDMA and network offloads over ethernet, with its congestion control based on the combination of PFC and ECN. The document defines the fast congestion notification for distributed RoCEv2 network based on SRv6 tunnels, and further extends PFC and ECN to achieve precise flow control and end-to-end congestion control.

Status of This Memo

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

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

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

This Internet-Draft will expire on 3 September 2025.

Copyright Notice

Copyright (c) 2025 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<https://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Revised BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Revised BSD License.

Table of Contents

1. Introduction	2
2. Conventions Used in This Document	3
2.1. Abbreviations	3
2.2. Requirements Language	3
3. Scenarios for distributed RoCEv2 network	3
3.1. Scenario 1: distributed model training	3
3.2. Scenario 2: separated storage and model training	4
4. Solution	4
4.1. Precise flow control	6
4.1.1. Illustration of Para-Type and Corresponding Parameter	6
4.1.1.1. Para-Type Bit 2	7
4.1.1.2. Para-Type Bit 3	7
4.1.2. Process analysis of precise flow control	7
4.2. Extend ECN to WAN	8
4.2.1. Process analysis of CNP	8
4.2.2. Fast CNP	10
5. Security Considerations	10
6. IANA Considerations	10
7. Acknowledgments	10
8. References	10
8.1. Normative References	10
8.2. Informative References	11
Authors' Addresses	11

1. Introduction

RDMA (Remote Direct Memory Access) enables direct access to memory locations on remote machines, bypassing the need for CPU involvement in data transfer processes. RDMA results in lower latency, reduced CPU overhead, and increased network throughput, making RDMA particularly beneficial for high-performance computing environments, cloud infrastructure, and storage networks.

RoCEv2 is an open standard enabling RDMA over ethernet, with its congestion control based on the combination of PFC and ECN. Priority-based Flow Control (PFC) is a data link level flow control

mechanism, which can selectively pause traffic according to its class and eliminate packet loss caused by network congestion. Explicit Congestion Notification (ECN) is an extension to network layer protocol and transport layer protocol defined in RFC3168[RFC3168], which enables the notification of network congestion.

AI services (e.g. distributed model training, separated storage and model training) drive the demand for building distributed RoCEv2 networks based on WAN. Therefore, when congestion is detected in WAN devices, it is essential to achieve fast congestion notification based on the PFC and ECN mechanisms.

2. Conventions Used in This Document

2.1. Abbreviations

AIDC: Artificial Intelligence Data Center

CNP: Congestion Notification Packet

ECN: Explicit Congestion Notification

PFC: Priority-based Flow Control

RoCEv2: RDMA over Converged Ethernet version 2

WAN: wide area network

2.2. Requirements Language

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in BCP 14 [RFC2119] [RFC8174] when, and only when, they appear in all capitals, as shown here.

3. Scenarios for distributed RoCEv2 network

3.1. Scenario 1: distributed model training

The computational power growth of a single AIDC is limited by multiple factors such as space and power consumption, making it difficult to meet the fast-growing computational demands of large models. Therefore, more and more enterprises are inclined to support super-scale model training by coordinating distributed model training across multiple AIDCs.

During distributed model training, WAN needs to carry parameter synchronization data between multiple AIDCs through tunnels, and this data is transmitted using RDMA protocols such as RoCEv2.

3.2. Scenario 2: separated storage and model training

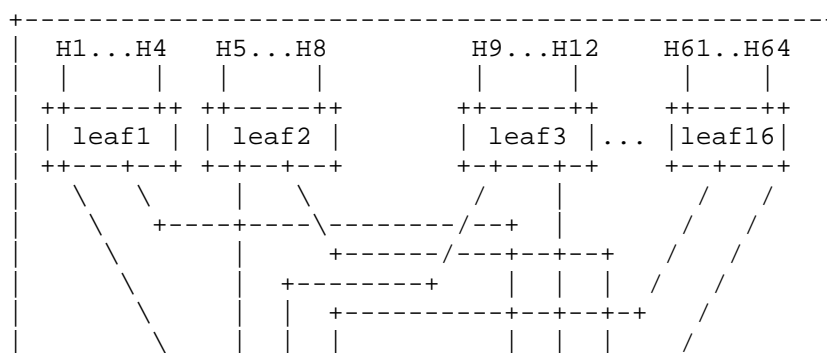
In scenarios such as computational power rental, computing clusters and storage clusters are usually located in different DCs. To prevent sensitive data from being stored in the compute cluster and causing data leakage, clients request that this data be directly transmitted to the memory of the compute cluster for model training.

During separated storage and model training, WAN needs to carry sample data from storage clusters to compute clusters through tunnels, and this data is transmitted using RDMA protocols such as RoCEv2.

4. Solution

The scenarios mentioned in Chapter 3 can be summarized as the network diagram shown in Figure 1. In these scenarios, the sender in DC needs to transmit RoCEv2 packets to the receiver in another DC through SRv6 tunnels in WAN, the process is as follows:

- * The sender in DC sends RoCEv2 packets to WAN's ingress PE through the leaf, spine, and gateway devices.
- * At the WAN's ingress PE, the RoCEv2 packets are encapsulated according to the SRv6 tunnel protocol.
- * The WAN's P node transits the payload from ingress PE to egress PE through SRv6 tunnels.
- * At the WAN's egress PE, the payload are decapsulated to RoCEv2 packets and transmitted to the receiver in DC.



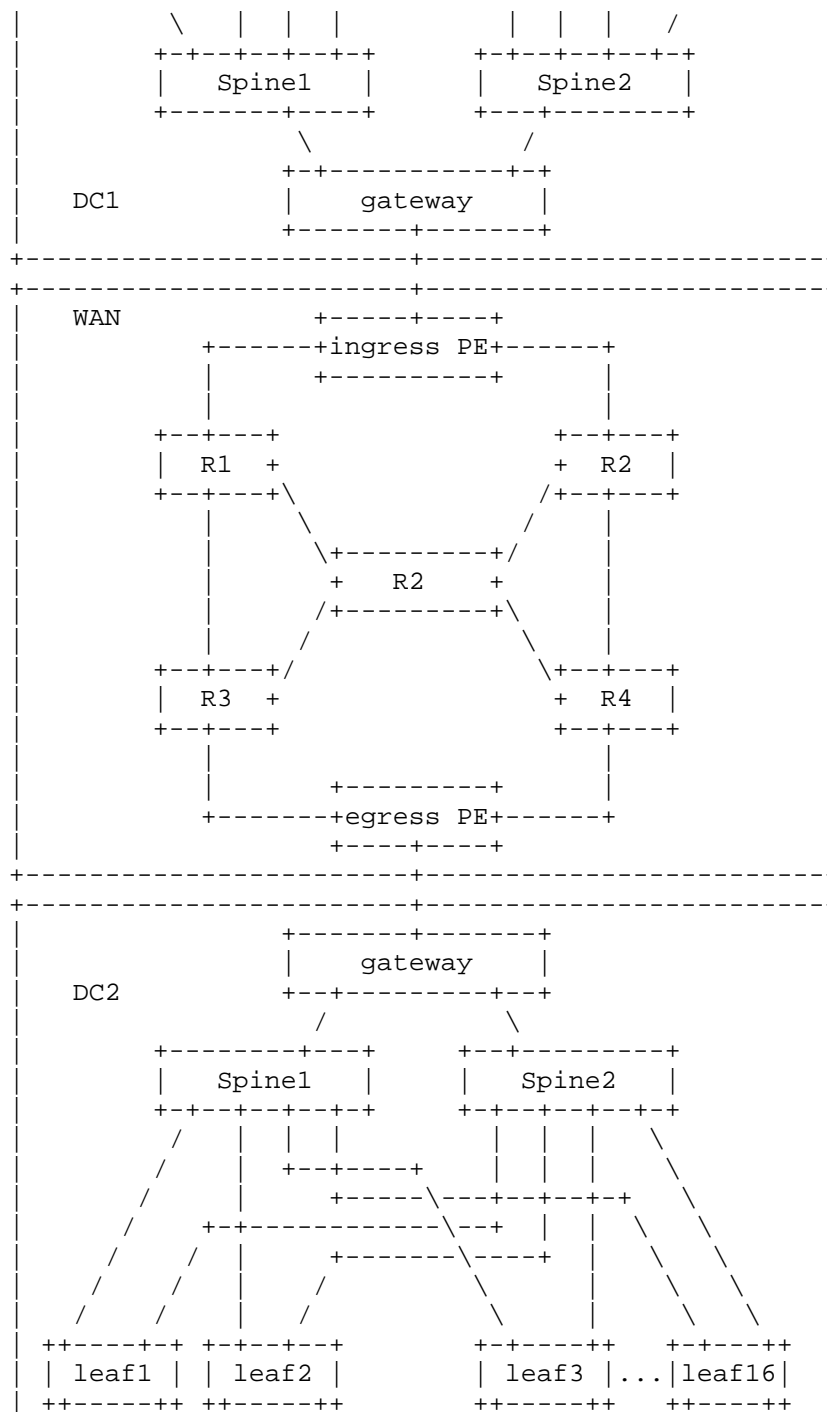




Figure 1: Network diagram

Within DC, RoCEv2 implements congestion control based on PFC and ECN. In WAN, it is necessary to extend PFC and ECN so that fast notifications can be achieved when congestion is detected by WAN devices.

4.1. Precise flow control

PFC is the feature used in ethernet to prevent data loss due to congestion. In PFC, traffic is classified into 8 priorities according to the 802.1Q protocol, and each priority maintains its queue. When the queue length of the router's receive port exceeds a certain threshold, the router sends a PAUSE frame to upstream to stop transmitting traffic. When the queue length of the router's receive port drops below a certain threshold, the router sends a RESUME frame to upstream to continue transmitting traffic.

As WAN devices often carry multiple services simultaneously, if PFC is triggered due to the congestion in a specific service, it may impact other services on the port and pose security risks. Therefore, when carrying RoCEv2 packets through tunnel in WAN, precise flow control needs to be implemented based on traffic information such as inner IP header, outer IP header and priority. This document proposes several parameters for the ARN mechanism ([I-D. draft-wh-rtgwg-adaptive-routing-arn][I-D.wh-rtgwg-adaptive-routing-arn]), enabling fast notification for congested flow in WAN.

4.1.1. Illustration of Para-Type and Corresponding Parameter

[I-D. draft-wh-rtgwg-adaptive-routing-arn][I-D.wh-rtgwg-adaptive-routing-arn] propose the ARN packet format as follows, as well as two parameters: the Para-Type Bit 0 for flow identifier and the Para-Type Bit 1 for path identifier. Specifically, flow identifier is based on the five-tuple from packet header to indicate the affected flow, path identifier is based on the 32-bit path ID to uniquely identify the affected path in the network.

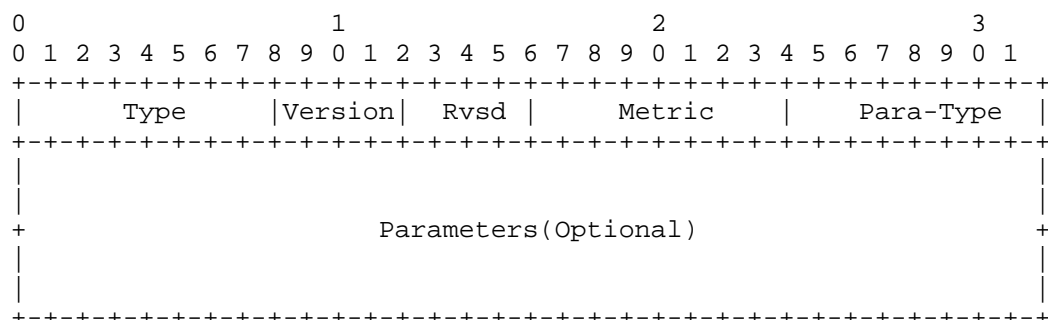
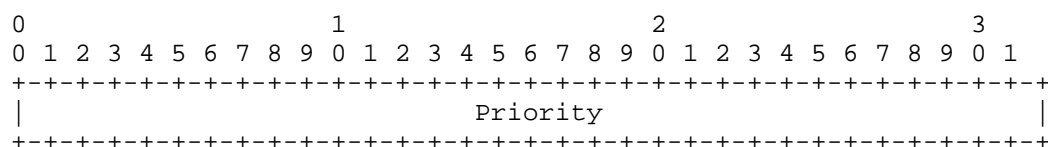


Figure 2: ARN format

This document further defines the Para-Type Bit 2 for priority identifier and the Para-Type Bit 3 for compression time.

4.1.1.1. Para-Type Bit 2

When bit2 of Para-Type is 1, the following parameter is concluded in Parameters to indicate the priority of affected flows:

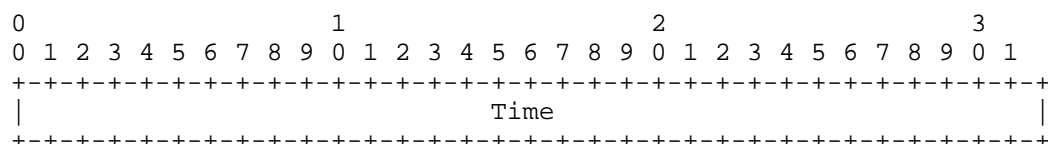


Priority ID: The 32-bit field is used to identify the priority of affected flows.

Figure 3: ARN format

4.1.1.2. Para-Type Bit 3

When bit3 of Para-Type is 1, the following parameter is concluded in Parameters to indicate the compression time of affected flows:



Time: The 32-bit field is used to uniquely identify the compression time of affected flows.

Figure 4: Compression time

4.1.2. Process analysis of precise flow control

Taking Figure 1 as an example, within DC1 and DC2, PFC is still used to implement the flow control based on port priority.

Within WAN, precise flow control can be used to achieve congestion control at three levels (service level, priority level, and flow level), effectively avoiding impact on other services. Specifically, when the receiving port of device detects congestion, it will notify upstream devices to pause the corresponding flow:

- * When configured to stop the service in case of congestion, the router will pass on ARN packet carrying flow identifier (identify the corresponding service based on the five-tuple of outer IP header) and compression time to upstream devices.
- * When configured to stop specific flow in case of congestion, the router will pass on ARN packet carrying flow identifier (identify the corresponding flow based on the five-tuple of inner IP header) and compression time to upstream devices.
- * When configured to stop flow based on priority, the router will pass on ARN packet carrying priority identifier and compression time to upstream devices.

Upon receiving the ARN packet, upstream devices will halt the corresponding flow based on the configured mode.

Between WAN and DC1 or DC2, considering the capabilities of the gateway devices, congestion control can be achieved by PFC or precise flow control.

4.2. Extend ECN to WAN

4.2.1. Process analysis of CNP

ECN enables a forwarding element (e.g., a router) to notify the sender for congestion control without having to drop packets [RFC3168[RFC3168]]. When a router detects congestion, instead of dropping packets as a signal of congestion, it marks the packets with an ECN codepoint in the IP header. The marked packets alert the receiver that packet loss is imminent, and then the receiver alerts the sender by sending Congestion Notification Packet (CNP). After receiving the CNP, the sender knows to slow down the transmission rate temporarily until the flow path is ready to handle a higher rate of traffic.

Within WAN, devices should be able to notify the sender to reduce the transmission rate when congestion is detected. [RFC6040[RFC6040]] redefines how the explicit congestion notification (ECN) field of the IP header should be constructed on entry to and exit from any IP-in-IP tunnel (including the SRv6 tunnel). It defines the rules for ingress encapsulation and egress decapsulation as follows:

Incoming Header (also equal to departing Inner Header)	Departing Outer Header	
	Compatibility Mode	Normal Mode
Not-ECT	Not-ECT	Not-ECT
ECT(0)	Not-ECT	ECT(0)
ECT(1)	Not-ECT	ECT(1)
CE	Not-ECT	CE

Figure 5: Tunnel ingress encapsulation behaviors

Arriving Inner Header	Arriving Outer Header			
	Not-ECT	ECT(0)	ECT(1)	CE
Not-ECT	Not-ECT	Not-ECT(!!!)	Not-ECT(!!!)	<drop>(!!!)
ECT(0)	ECT(0)	ECT(0)	ECT(1)	CE
ECT(1)	ECT(1)	ECT(1) (!)	ECT(1)	CE
CE	CE	CE	CE(!!!)	CE

Currently unused combinations are indicated by '(!!!)' or '(!)'

Figure 6: Tunnel egress decapsulation behaviors

At the ingress PE, there are two encapsulation modes: a REQUIRED 'normal mode' and a 'compatibility mode', which is for backward compatibility with tunnel egress do not understand ECN. In normal mode, the ingress PE constructs the outer encapsulating IP header by copying the two-bit ECN field of the incoming IP header. In compatibility mode, it clears the ECN field in the outer header to the Not-ECT codepoint.

At the egress PE, to decapsulate the inner header at the tunnel egress, The ECN field in the outgoing header is set to the codepoint at the intersection of the appropriate arriving inner header (row) and arriving outer header (column) in Figure 4, or the packet is dropped where indicated.

To enable congestion control for WAN devices, the ingress PE MUST use the normal mode to construct the outer encapsulating IP header. When a WAN device detects congestion (e.g. R2), it sets the ECN field of the outer IP header to CE. Then, the egress PE sets the ECN field in the outgoing IP header to CE during decapsulation. When receiver receives the marked packet, it sends CNP warning to the sender to slow down the transmission rate of the corresponding flow, thereby preventing congestion.

4.2.2. Fast CNP

[RFC7514[RFC7514]] defines Really Explicit Congestion Notification (RECN), also known as Fast Congestion Notification Packet (Fast CNP). By extending the RoCEv2 CNP, Fast CNP can be sent by the intermediate router directly to the sender, advising the sender to reduce the transmission rate at which it sends the flow of RoCEv2 data traffic.

When transporting RoCEv2 packets through SRv6 tunnels in WAN, intermediate devices are unable to send Fast CNP back to the ingress PE based on the information in RoCEv2 packets. Additionally, the ingress PE is unable to recognize the Fast CNP and therefore cannot forward them to the sender. Therefore, this document proposes a Congestion SID (END.CON) to address aforementioned issues.

END.CON is a 128-bit value configured on PE and disseminated to other routers via IGP. The endpoint action associated with END.CON is to decapsulate the packet and then look up the routing table for traffic forwarding. The workflow of END.CON is as follows:

- * When the ingress PE passes RoCEv2 packets through SRv6 tunnels, END.CON is encapsulated as the source address of outer IP header.
- * When WAN devices detect congestion, they encapsulate an outer IP header (destination address is set to END.CON) outside the Fast CNP and send it back to the ingress PE.
- * When the ingress PE detects that the destination address of the received packet hits END.CON, it will decapsulate the Fast CNP and pass it to the sender in DC.

5. Security Considerations

TBD

6. IANA Considerations

TBD

7. Acknowledgments

TBD

8. References

8.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC3688] Mealling, M., "The IETF XML Registry", BCP 81, RFC 3688, DOI 10.17487/RFC3688, January 2004, <<https://www.rfc-editor.org/info/rfc3688>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.

8.2. Informative References

- [RFC3168] Ramakrishnan, K., Floyd, S., and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", RFC 3168, DOI 10.17487/RFC3168, September 2001, <<https://www.rfc-editor.org/info/rfc3168>>.
- [RFC6040] Briscoe, B., "Tunnelling of Explicit Congestion Notification", RFC 6040, DOI 10.17487/RFC6040, November 2010, <<https://www.rfc-editor.org/info/rfc6040>>.
- [RFC7514] Luckie, M., "Really Explicit Congestion Notification (RECN)", RFC 7514, DOI 10.17487/RFC7514, April 2015, <<https://www.rfc-editor.org/info/rfc7514>>.
- [I-D.wh-rtgwg-adaptive-routing-arn]
Wang, H., Huang, H., Geng, X., Xu, X., and Y. Xia,
"Adaptive Routing Notification", Work in Progress,
Internet-Draft, draft-wh-rtgwg-adaptive-routing-arn-03, 13
September 2024, <<https://datatracker.ietf.org/doc/html/draft-wh-rtgwg-adaptive-routing-arn-03>>.

Authors' Addresses

Zehua Hu
China Telecom
Guangzhou
China
Email: huzh2@chinatelecom.cn

Yongqing Zhu
China Telecom
Guangzhou
China

Email: zhuyq8@chinatelecom.cn

Xuesong Geng
Huawei
China
Email: gengxuesong@huawei.com

Jiayuan Hu
China Telecom
Guangzhou
China
Email: hujy5@chinatelecom.cn

Tanxin Pi
China Telecom
Guangzhou
China
Email: pitxl@chinatelecom.cn