

RTGWG
Internet-Draft
Intended status: Standards Track
Expires: 7 January 2026

Z. Hu
Y. Zhu
J. Hu
T. Pi
China Telecom
6 July 2025

Fast Notification for Traffic Engineering and Load Balancing for tunnel-
based lossless transmission in WAN
draft-hzh-fantel-wan-tunnel-00

Abstract

With the rapid development of large language models, many emerging AI scenarios require tunnel-based lossless transmission of RDMA packets in WAN. To fulfill this requirement, WAN should support the real-time notification of network conditions to ensure high throughput, low latency, and zero packet loss data transmission. The current reactive notification solution is limited by responsiveness, coverage, and operational overhead. Therefore, we need to establish a faster and proactive notification mechanism to implement more responsive Traffic Engineering and Load Balancing.

This draft first describes typical scenarios for tunnel-based RDMA lossless transmission in WAN, then specifies the fast notification framework for implementing key TE areas (congestion control, protection, and load balance), and finally analyses the protocol implementation for fast notification.

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 7 January 2026.

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	4
3. Scenarios	4
3.1. Scenario 1: distributed model training across AIDC	4
3.2. Scenario 2: Cloud-Edge Collaborative Model Inference . .	4
4. Framework	4
4.1. Failure protection	6
4.2. Congestion control	7
4.3. Load balancing for network state changes	9
5. Solutions	10
5.1. ICMPv6-based solution	10
5.2. CNP-based solution	11
6. Security Considerations	12
7. IANA Considerations	13
8. Acknowledgments	13
9. References	13
9.1. Normative References	13
9.2. Informative References	13
Authors' Addresses	14

1. Introduction

In use cases such as distributed model training or cloud-edge collaborative model inference, WAN needs to tunnel RDMA traffic between data centers. RDMA is a widely used technology in high-performance computing or AI clusters, achieving lower latency, reduced CPU overhead, and increased network throughput. Currently, mainstream RDMA protocols (e.g., IB, RoCE) are based on the Go-Back-N mechanism, where a small number of packet losses will result in a dramatic reduction in the effective throughput. In order to achieve

lossless data transmission, WAN need to support FAsT Notification for Traffic Engineering and Load balancing (FANTEL).

ECN[RFC3168] enables a forwarding element (e.g., a router) to notify the sender for congestion control without having to drop packets. When a router detects congestion, it marks the packets with an ECN code-point in the IP header. The receiver, upon receiving marked packets, sends a Congestion Notification Packet (CNP) to the sender, which then temporarily reduces its transmission rate until the path can accommodate higher traffic. ECN still relies on end-to-end signaling, making real-time feedback challenging in long-distance WAN.

[draft-wh-rtgwg-adaptive-routing-arn] proposes a proactive notification mechanism ARN for adaptive routing, and describes the information carried in ARN to notify remote nodes for re-routing. This draft proposes a unified mechanism for congestion notifications, link failure notifications, and even to convey other relevant network events for re-routing. [draft-liu-rtgwg-adaptive-routing-notification] describes the mechanisms of delivering ARN message. This draft gives three options, each of which specifies the information carried in the ARN message and the mechanism of sending the message to specific network nodes.

Some gap analysis documents demonstrate that FANTEL achieves high-throughput, low-latency, and lossless RDMA data transmission in AI Data Centers (AIDCs). With the development of Large Language Models (LLMs), some scenarios require WANs to guarantee lossless data transmission when carrying RDMA packets over tunnels. This draft introduces the typical scenarios of distributed lossless network based on tunnels in WAN, then specifies the mechanisms of FANTEL to achieve key TE areas such as congestion control, load balancing, and failure protection, and finally analyses the protocol implementation.

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

3.1. Scenario 1: distributed model training across AIDC

The growth of computing power of a single AIDC is limited by space and power supply, making it difficult to meet the fast-growing computational demands of LLMs. Therefore, more and more enterprises are inclined to realize super-scale model training by distributed model training across multiple AIDCs.

During this scenario, AI servers across different AIDCs synchronize parameter plane data using RDMA, which requires the WAN to provide lossless transport of RDMA packets over tunnels such as SRv6 and MPLS.

3.2. Scenario 2: Cloud-Edge Collaborative Model Inference

Many enterprises achieve the deployment and application of LLMs by building on-premises computing power resources pool. However, the cost of deployment and O&M of this approach is very high, and it is difficult to meet the subsequent computing power demand of LLMs. To address this, the cloud-edge collaboration between enterprise on-premises and cloud computing resource pools provides a more efficient, agile, and cost-effective approach to realize elastic computing power scaling.

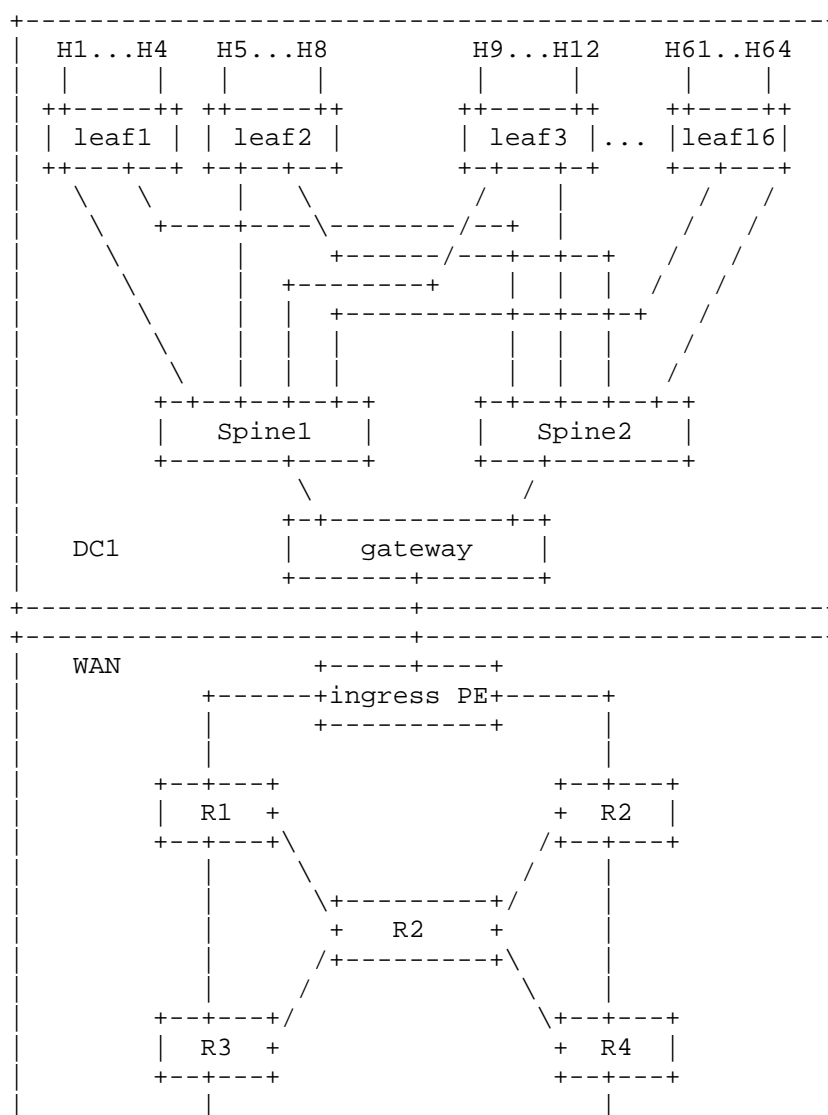
During this scenario, Parameter plane data synchronization between on-premises AI servers and cloud-based AI servers requires lossless RDMA packet transmission over WAN tunnels.

4. Framework

The framework for lossless RDMA data transmission over WAN tunnels is as follows:

- * The AI servers(H1-H64) in DC1 sends RoCEv2 packets to WAN's ingress PE.

- * At the WAN's ingress PE, the RoCEv2 packets are encapsulated according to the SRv6 tunnel protocol, Then it is sent to the path with the best load sharing across the network.
- * The WAN's P node(R1-R5) 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 AI servers in DC2.



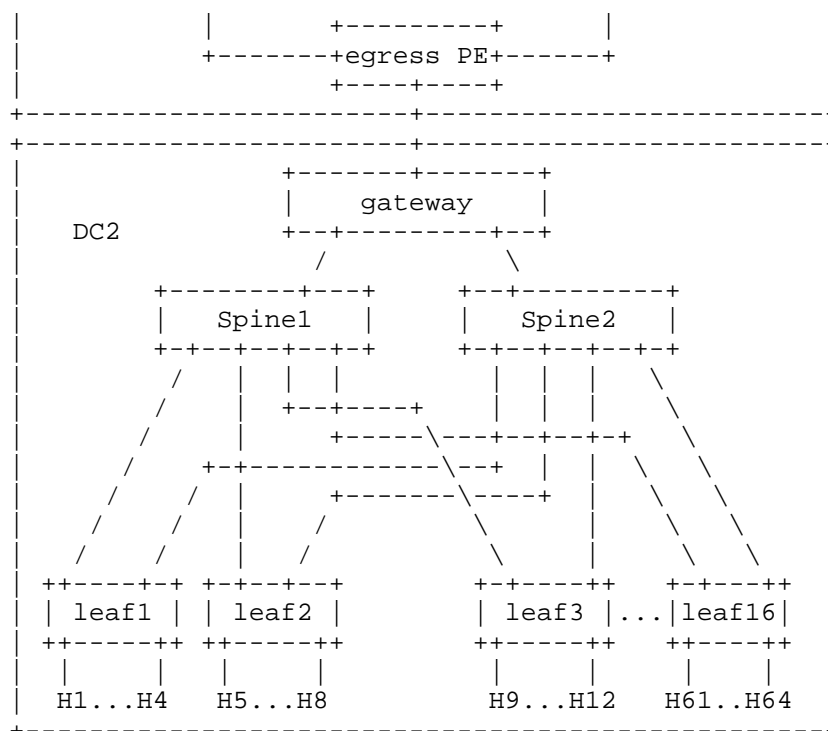


Figure 1: Network diagram

Throughout the process, the nodes in the WAN need to have a fast notification mechanism that allow ingress PE or upstream nodes take response actions to react quickly to network conditions such as link congestion, SLA violations and so on. This fast notification focuses on three key TE scenarios: failure protection, flow control, and load balancing.

4.1. Failure protection

For large-scale and dynamic networks, protection mechanisms need to ensure service continuity in case of failures. According to [draft-geng-fantel-fantel-gap-analysis], The current failure handling methodology for reactive is BFD and FRR, they lack flexibility and responsiveness in complex typologies. Therefore, The WAN needs to have fast notification of failures, allowing near-instantaneous and dynamic protection responses, minimizing user impact.

When carrying RDMA traffic based on tunneling in WAN, the ingress PE is responsible for path alignment. Therefore, in failure protection scenario, the architecture for fast notification is as follows:

- * When a node detects a local link/device failure, it collects failure information about the affected path or flow.
- * The node sends a quick notification to the ingress PE with failure information (Here not only the information about the affected flows should be carried, but also the information about the next hop affected by the failed link/device should be indicated).
- * Ingress PE receives the fast notification and excludes the path containing the fault link/device (to avoid RDMA traffic bypass), ensures minimal disruption and quick recovery by switching the affected traffic to a backup path or ECMP, etc.

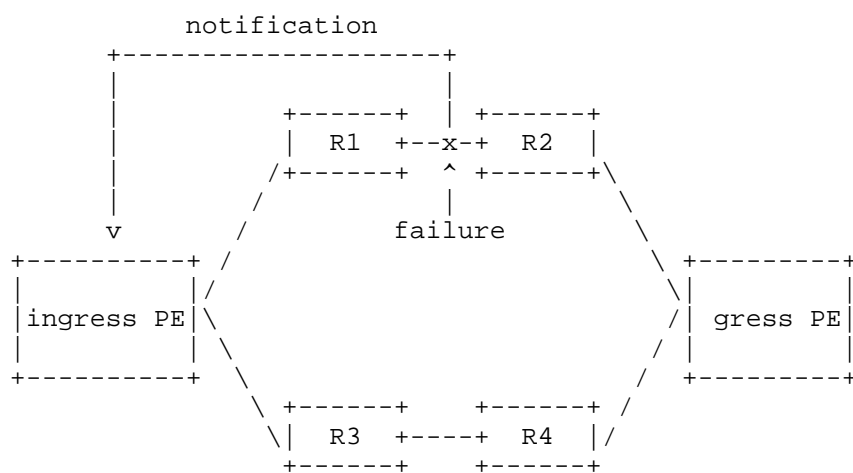


Figure 2: Failure protection procession

4.2. Congestion control

RDMA traffic is bursty and highly sensitive to packet loss, and WAN require proactive congestion control mechanisms. [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, in order to achieve ECN-based congestion control across WANs between DCs. However, [draft-geng-fantel-fantel-gap-analysis] analysis that ECN/TCP methods still relies on end-to-end signaling and lacks precise real-time feedback.

Currently, PFC is widely used in data centers to prevent data loss due to congestion. PFC uses a step-by-step back-pressure mechanism to control the upstream to stop or continue transmitting traffic. PFC achieves link-layer priority-based traffic control, but still faces problems such as queue head blocking and deadlock due to coarse control granularity.

Therefore, when carrying RDMA traffic in WAN based on tunnel, it is important to have the ability to quickly notify the ingress PE to reduce/increase the transmission rate of the corresponding paths, and also to notify the upstream to stop/continue the flow. In the congestion control scenario, the architecture of fast notification is as follows.

- * When the node detects that the buffer corresponding to a service reaches the threshold x , it will quickly notify the ingress PE of the paths or flows that need to be quenched, queue situation, and other information; when the node detects that the buffer corresponding to a service reaches the threshold y , it will quickly notify the upstream node of the paths or flows that need to be quenched, and other information.
- * After the Ingress PE receives the notification, it reduces the traffic on the congested path by switching flows to other paths; after the upstream node receives the notification, it stops send traffic to corresponding path.
- * When the node detects that the buffer corresponding to the service is lower than the threshold x , it will quickly notify the ingress PE of the path or flow that needs to be resumed; when the node detects that the buffer corresponding to the service is lower than the threshold y , it will quickly notify the upstream node of the path or flow that needs to be resumed.
- * After the Ingress PE receives the notification, it redistributes the traffic load by switching flows back to that path; after the upstream node receives the notification, it continues to send the traffic to corresponding path.

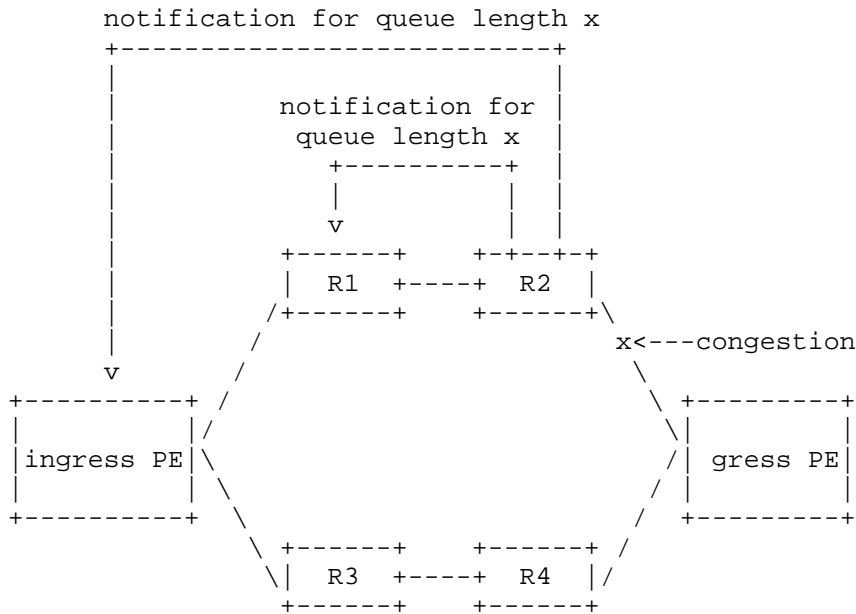


Figure 3: Congestion control procession

4.3. Load balancing for network state changes

Devices and links in WAN often carry multiple services simultaneously. In addition to failure, congestion and other conditions, dynamic load balancing based on link utilization, node load and other changes in network status can ensure efficient use of available bandwidth. It can prevent single node or link becomes overwhelmed with excessive traffic. In dynamic networks, Proper load balancing improves network performance and prevents bottlenecks.

When carrying RDMA traffic based on tunneling in a WAN, the ingress PE is responsible for path alignment. Thus in a load balancing scenario, the architecture for fast notification is as follows.

- * When a node detects the network state change, it collects the network state change information, such as link utilization, queue buildup.
- * The node sends fast notification to the ingress PE with information about the network state change.
- * Ingress PE receives the fast notification and redistributes the traffic load to paths by ECMP.

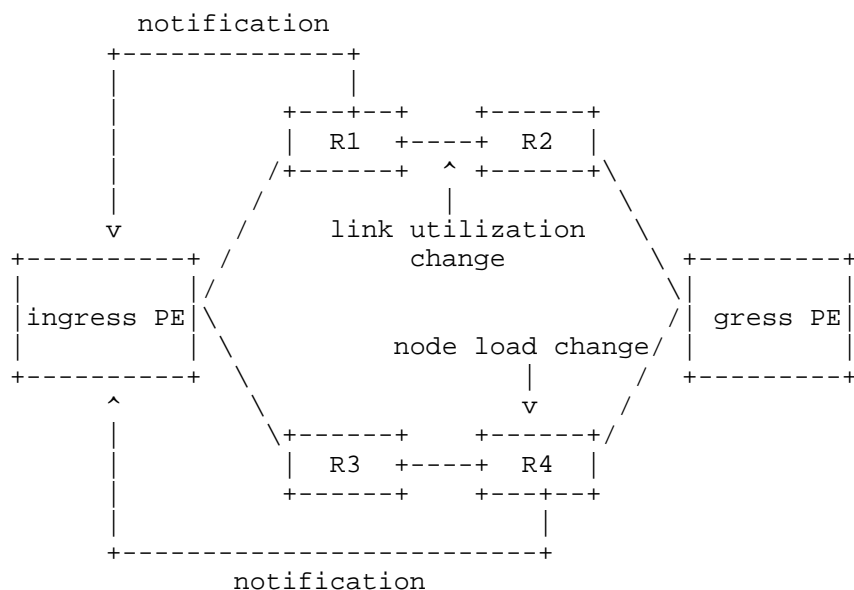


Figure 4: Load balancing for network state changes

5. Solutions

Based on the framework analysis of fast notification in key TE areas, a unified protocol implementation for fast notification should be established, with explicit forwarding procedures to realize tunnel-based lossless transmission of RDMA packets in WAN.

5.1. ICMPv6-based solution

The source quench mechanism has been deprecated in ICMPv6 because TCP's built-in congestion avoidance algorithms are more efficient, and source quench may interfere with their normal operation. However, when transmitting RDMA data over WAN tunnels, the source quench notification is confined within the WAN domain (this message is used by WAN devices such as Ingress PE or transit node for traffic engineering) and does not affect transport layer congestion control.

This document specifies a new ICMP message to realize rapid notification in key traffic engineering areas including failure protection, congestion control, and load balancing. The message format is defined as follows:

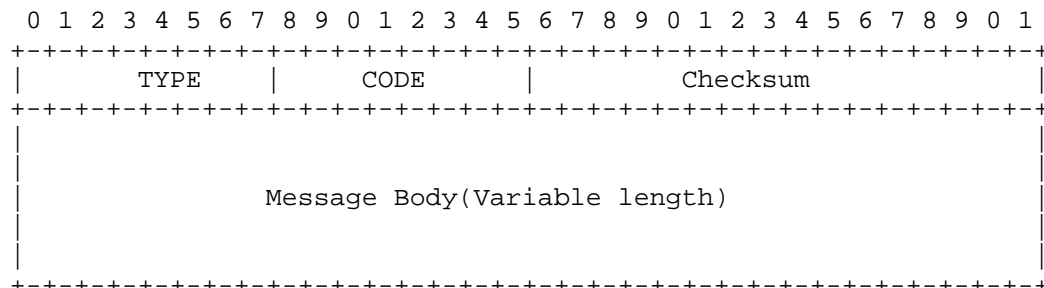


Figure 5: new ICMP message for fast notification

```
*TYPE:8-bit identifier for fast notification ICMP message (to be reserved).
```

*CODE: 8-bit identifier for TE areas, When it is set to 1, it indicates the fast notification for failure protection; When it is set to 2, it indicates the fast notification for failure elimination; When it is set to 3, it indicates the fast notification for congestion control; When it is set to 4, it indicates the fast notification for congestion elimination; When it is set to 5, it indicates the fast notification for load balancing. Other bits are not defined.

*Checksum: Used for error-checking the packet.

*Message Body: It carries notification information specific to each areas: for failure protection, it includes path, five-tuple of flow, and failure cause; for congestion control, it contains path and buffer status; for load balancing, it comprises link utilization and device load. This field format need to be designed with extensibility, while subsequent refinements and specific packet forwarding mechanisms(TBD).

5.2. CNP-based solution

[RFC 7514] introduces the Fast CNP mechanism, which enables intermediate nodes to directly send Fast CNP packets to sender. The CNP packet format is as follows:

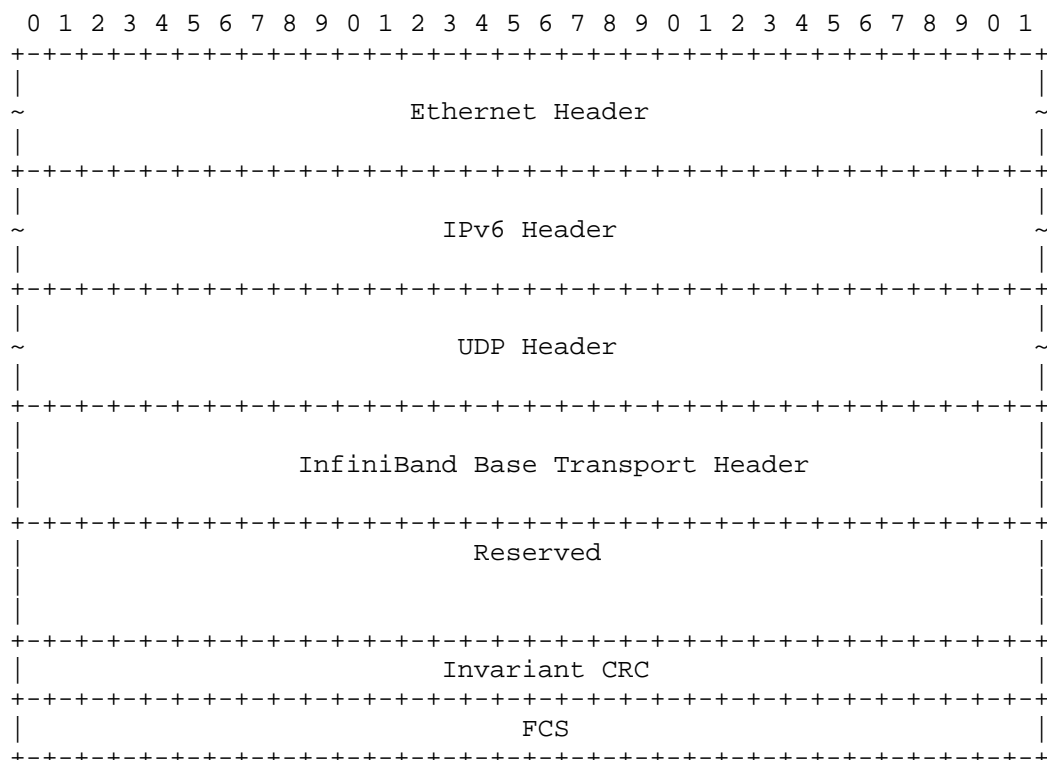


Figure 6: new ICMP message for fast notification

The sender differentiates CNP packets from data packets using the Opcode in the InfiniBand Base Transport Header (BTH) and reduce the transmission rate at which it sends data packets based on the destination QP in BTH. WAN routers do not need to parse BTH information. This document specifies a new UDP port number (to be reserved) to enable routers to perform fast notification processing. The ARN message format defined in [draft-wh-rtgwg-adaptive-routing-arn] could be considered as a potential implementation in this case. However, as it primarily targets congestion scenarios, further specification is needed for both the notification information applicable to other TE scenarios and the corresponding packet forwarding mechanisms (TBD).

6. Security Considerations

TBD

7. IANA Considerations

TBD

8. Acknowledgments

TBD

9. References

9.1. 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>>.

9.2. 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>>.

[I-D.liu-rtgwg-adaptive-routing-notification]

Liu, Y., lihesong, and W. Duan, "Adaptive Routing Notification for Load-balancing", Work in Progress, Internet-Draft, draft-liu-rtgwg-adaptive-routing-notification-02, 12 June 2025, <<https://datatracker.ietf.org/doc/html/draft-liu-rtgwg-adaptive-routing-notification-02>>.

[I-D.xiao-rtgwg-rocev2-fast-cnp]

Min, X. and lihesong, "Fast Congestion Notification Packet (CNP) in RoCEv2 Networks", Work in Progress, Internet-Draft, draft-xiao-rtgwg-rocev2-fast-cnp-03, 9 June 2025, <<https://datatracker.ietf.org/doc/html/draft-xiao-rtgwg-rocev2-fast-cnp-03>>.

Authors' Addresses

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

Yongqing Zhu
China Telecom
Guangzhou
China
Email: zhuyq8@chinatelecom.cn

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

Tanxin Pi
China Telecom
Guangzhou
China
Email: pitx1@chinatelecom.cn