

MPLS  
Internet-Draft  
Intended status: Informational  
Expires: 23 April 2026

A. Mahale  
Cerebras Systems  
K. Kompella  
V. P. Beeram  
Juniper Networks  
D. Patel  
AMD  
20 October 2025

MPLS for AIDC Probing  
draft-amahale-mpls-for-aidc-00

## Abstract

This document describes a method for using Multi-Protocol Label Switching (MPLS) encapsulation to perform scalable and vendor-agnostic network probing within AI/ML data centers. The goal is to detect and isolate gray failures—non-deterministic hardware and software faults—in large-scale lossless networks. The approach enables targeted probing at per-link and per-node granularity, independent of IP/BGP control plane operation, and is extensible to Various CLOS topologies.

## 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 23 April 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. Problem Statement: Gray Failures . . . . .	3
3. Network Probing Overview . . . . .	3
3.1. End-to-End Probing . . . . .	3
3.2. Direct Network probing . . . . .	3
4. MPLS Operations for Probing in AIDC . . . . .	4
4.1. MPLS Refresher . . . . .	4
4.2. MPLS Operations and Label Stack Use . . . . .	4
4.3. MPLS Myths and Clarifications . . . . .	4
5. MPLS for AI Cluster Network Probing . . . . .	5
5.1. Topology Considerations . . . . .	5
5.2. Scaling and Label Allocation . . . . .	7
5.3. Failure Correlation and Control Plane Independence . . . . .	7
6. Security Considerations . . . . .	7
7. IANA Considerations . . . . .	7
8. Acknowledgments . . . . .	7
9. References . . . . .	7
Authors' Addresses . . . . .	8

## 1. Introduction

The advent of large-scale AI/ML data centers and the adoption of lossless networking paradigms have increased the operational risk of gray failures—partial, intermittent, or non-deterministic faults in network components. These failures are notoriously difficult to detect and isolate, especially in high-performance environments that rely on congestion control mechanisms such as PFC, ECN, and DCQCN.

This document proposes a vendor-agnostic probing mechanism leveraging MPLS encapsulation to detect gray failures. The technique provides deterministic path visibility and decouples the probing infrastructure from the routing control plane.

## 2. Problem Statement: Gray Failures

Gray failures refer to partial or intermittent faults in network devices such as switches, routers, NICs, optics, and cables that do not manifest as complete outages. These failures often evade monitoring systems and may take hours or days to isolate.

Traditional white-box monitoring approaches that rely on ASIC error counters or register captures are insufficient because the failure mode is not deterministic. ASICs may not be capable of self-detection when specific functional blocks wedge or stall.

Proactive network probing provides a means of external validation of data-plane health by exercising network paths in controlled ways.

## 3. Network Probing Overview

Network probing supplements standard telemetry and monitoring systems by introducing synthetic traffic to verify forwarding correctness. There are two primary categories of network probing mechanisms:

- \* End-to-End Probing
- \* Direct Network Probing

Generally prober machines are co-located with other AI compute and attached to the leaf switches.

### 3.1. End-to-End Probing

End-to-end probing systems, such as [Pingmesh], rely on hosts sending probes to all other hosts to verify reachability. While effective at identifying connectivity and control-plane issues, this approach lacks sufficient entropy to exercise all paths in large CLOS networks and does not directly map probe failures to specific components.

### 3.2. Direct Network probing

Direct network probing targets intermediate hops and links directly. Implementations such as [NetNorad] generate probes toward every network node or link from each source. This provides superior fault localization but presents challenges in targeting granularity, particularly in pure IP networks where multiple encapsulations may be required.

MPLS provides a clean mechanism to express path and target semantics through label stacking.

## 4. MPLS Operations for Probing in AIDC

### 4.1. MPLS Refresher

MPLS is a 4-byte shim header inserted between Ethernet and IP headers. It contains:

- \* 20-bit Label field
- \* 3-bit QoS (Traffic Class)
- \* 1-bit Bottom of Stack (BOS)
- \* 8-bit TTL

MPLS allows multiple labels to be stacked to represent a sequence of hops or links, enabling source routing and fine-grained path control.

The 20-bit label space (1,048,576 possible labels) is sufficient for even the largest data centers. MPLS labels may also be used to exercise lossless and lossy queues by mapping QoS bits to specific hardware queues.

### 4.2. MPLS Operations and Label Stack Use

The POP operation is central to MPLS probing. Each network element receiving a packet with a label stack removes the top label and forwards the remaining packet based on the next label or IP header.

This enables hierarchical targeting where each label represents a link or node. A single probe packet can traverse multiple layers of the CLOS fabric.

### 4.3. MPLS Myths and Clarifications

MPLS needs an additional protocol to function.

MPLS is an encapsulation, not a routing protocol. Labels MAY be distributed using dynamic protocols such as LDP, RSVP, or BGP, but they MAY also be configured statically. Static assignment is sufficient for probing use cases.

MPLS is complex to implement in ASICs.

MPLS lookup is simpler than IPv4/IPv6 LPM lookups, as it is an exact match on a fixed-width 20-bit key. Modern ASICs implement MPLS forwarding efficiently using hash-based tables.

Network ASICs cannot handle multiple labels.

Most modern data center ASICs support 8 or more MPLS labels in the stack, sufficient for multi-stage Clos topologies.

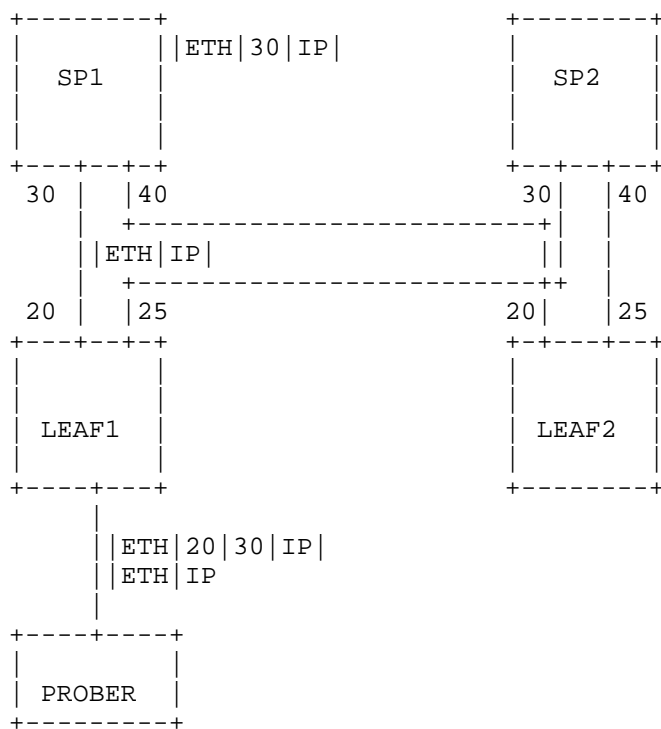
\_SRv6 obsoletes MPLS.\_

While SRv6 provides similar functionality, it requires IPv6 and adds considerable packet overhead and implementation complexity. MPLS offers a lightweight and control-plane-agnostic alternative suitable for IPv4-only environments.

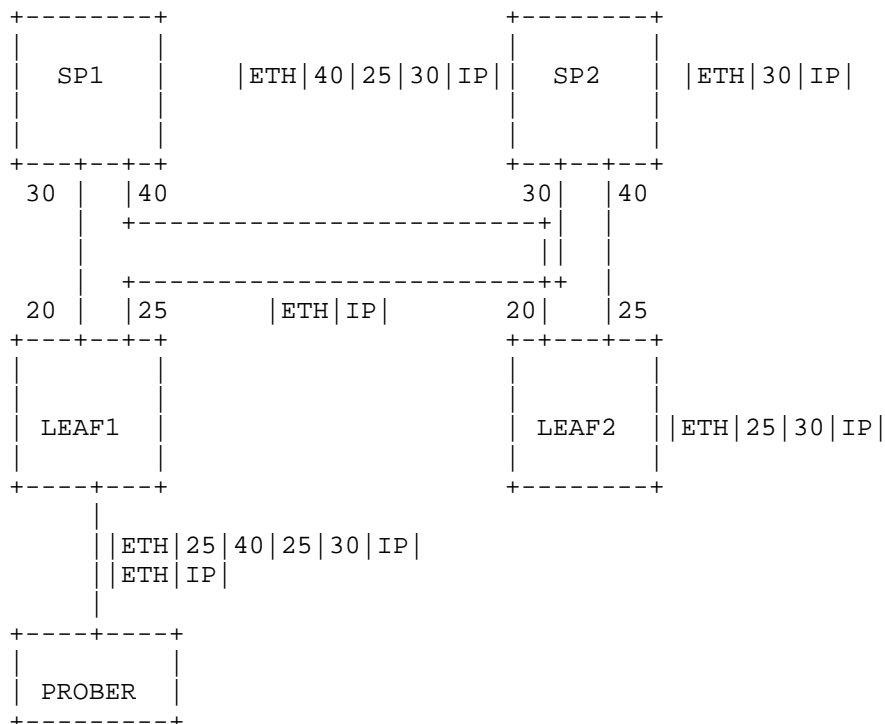
## 5. MPLS for AI Cluster Network Probing

### 5.1. Topology Considerations

Prober to Leaf1 to SP1 to Leaf1 to Prober



Prober to Leaf1 to SP2 to Leaf2 to SP2 to Leaf1 to Prober



In a three-stage CLOS topology, MPLS labels are provisioned statically for each leaf and spine link. A probe server connected to each leaf switch sends labeled packets that describe both the forward and return paths.

For example: \* Stage-1 probes target directly connected leafspine links. \* Stage-2 probes target remote leafspine links.

Each probe's IP destination is set to the local probe server, enabling the packet to return once all labels are popped.

Labels have local significance and may be reused across the topology, simplifying configuration.

Example: Label 20: LF-1 to SP-1 link Label 30: SP-1 to LF-1 link

The prober sends a two-label packet for Stage-1 probing. LF-1 pops label 20 and forwards label 30 to SP-1. SP-1 pops label 30 and returns the packet toward LF-1, completing the loop.

Stage-2 probing uses four-label stacks to test multi-hop paths. A complete mesh of connectivity can be established across the entire data center fabric.

## 5.2. Scaling and Label Allocation

MPLS provides a 20-bit label space (1 million+ labels). Even large topologies with thousands of leaf-spine interfaces consume only a small fraction of the space.

Label provisioning may be static or automated via controller software. A consistent label pattern can be repeated across leaves and spines, simplifying operational overhead.

An Example label allocation scheme can use a formula like:  $2 * (m + M * n) + 100$  Where  $m$  is leaf index and  $n$  is spine index and  $M$  is number of leaf switches.

## 5.3. Failure Correlation and Control Plane Independence

Because MPLS probing operates entirely in the data plane, it continues to function even if the IP or BGP control plane is unavailable. End-to-end probing shares fate with routing protocols, whereas MPLS probing isolates data-plane verification.

This separation provides improved resilience and faster failure localization in large AI/ML clusters.

## 6. Security Considerations

MPLS probing MUST be isolated from tenant or production traffic. Probes SHOULD be rate-limited and authenticated where feasible. Incorrect label provisioning MAY cause unintended forwarding loops or leakage into production paths. Also MPLS TTL will prevent forwarding loop packets from looping indefinitely.

## 7. IANA Considerations

To be added.

## 8. Acknowledgments

## 9. References

[RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.

[NetNorad] Facebook Engineering, "NetNorad: End-to-End and Direct Probing in Production Networks," 2016.

[Pingmesh] Microsoft Research, "Pingmesh: A Large-Scale System for Data Center Network Latency Measurement," 2015.

#### Authors' Addresses

Aditya Mahale  
Cerebras Systems  
Email: aditya.ietf@gmail.com

Kireeti Kompella  
Juniper Networks  
Email: kireeti.ietf@gmail.com

Vishnu Pavan Beeram  
Juniper Networks  
Email: vbeeram@juniper.net

Devang Patel  
AMD  
Email: devang.patel@amd.com