

RTGWG Working Group  
Internet Draft  
Intended status: Informational  
Expires: August 24, 2025

P. Huo  
G. Chen  
ByteDance  
C. Lin  
New H3C Technologies  
H. Dai  
ByteDance  
February 24, 2025

A OSF Framework for Artificial Intelligence (AI) Network  
draft-hcl-rtgwg-osf-framework-02

## Abstract

This document describes a framework for Artificial Intelligence (AI) network. Particularly, the document identifies a set of AI network components, describes their interactions, and exemplifies the workflow of the control and data planes.

## 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 August 24, 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 (<http://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 Simplified BSD License text as described in

## Table of Contents

1. Introduction.....	2
1.1. Requirements Language.....	3
1.2. Terminology.....	3
2. OSF Framework and Components.....	3
2.1. Framework Overview.....	3
2.2. OSF Functional Components.....	5
2.3. OSF-TM.....	5
2.4. OSF-Ingress.....	5
2.5. OSF-Egress.....	5
2.6. OSF-Forwarder.....	6
2.7. OSF-CFC.....	6
3. Deployment Considerations.....	6
4. OSF Framework Workflow.....	6
4.1. OSF Topology Manage.....	6
4.2. Load Balancing in OSF Packet Transmission.....	7
4.3. OSF Congestion Control.....	8
4.3.1. Credit-based Flow Control.....	9
4.3.2. Congestion Control Based on Link Quality Detection..	10
4.4. Rapid Link Failure Switchover in OSF.....	11
5. Security Considerations.....	11
6. IANA Considerations.....	11
7. References.....	11
7.1. Normative References.....	11
7.2. Informative References.....	11
Authors' Addresses.....	12

## 1. Introduction

With the widespread application of Artificial Intelligence (AI), the demand for AI networks is increasing. As described in [I-D.draft-hcl-ai-network-problem-00], with the development of AI networks, the model parameters for AI training are becoming increasingly large. In order to meet the demands of large-scale AI training, AI training networks typically adopt a distributed cluster approach, which presents the following new requirements for the network:

- o AI training networks need a new load balancing method to mitigate the impact of uneven loads caused by burst traffic and achieve as much load balancing as possible.

- o AI training networks require a new congestion control mechanism that can quickly detect congestion when it occurs locally, communicate the congestion state, and then perform global congestion control. This is more efficient than performing congestion control locally, and thus necessitates a global end-to-end congestion control mechanism.
- o AI training networks need to have the ability for fast fault recovery.

This document proposes an AI network architecture to meet the new requirements for AI networks.

### 1.1. 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.

### 1.2. Terminology

AI: Artificial Intelligence

OSF: Open Scheduled Fabric

OSF-TM: OSF Topo Manager

OSF-Ingress: OSF Ingress router

OSF-Egress: OSF Egress router

OSF-Forwarder: OSF Forwarder Router

OSF-CFC: OSF Credit-based Flow Control

## 2. OSF Framework and Components

### 2.1. Framework Overview

A high-level view of the OSF framework, without expanding the functional entities in the network, is illustrated in Figure 1.

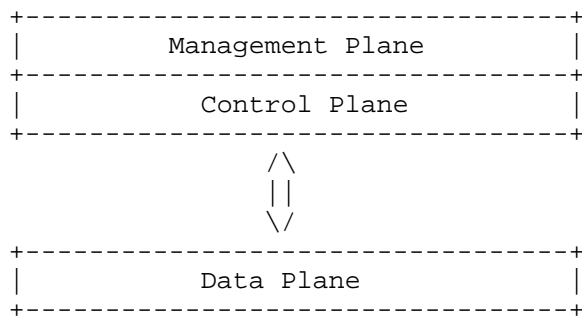
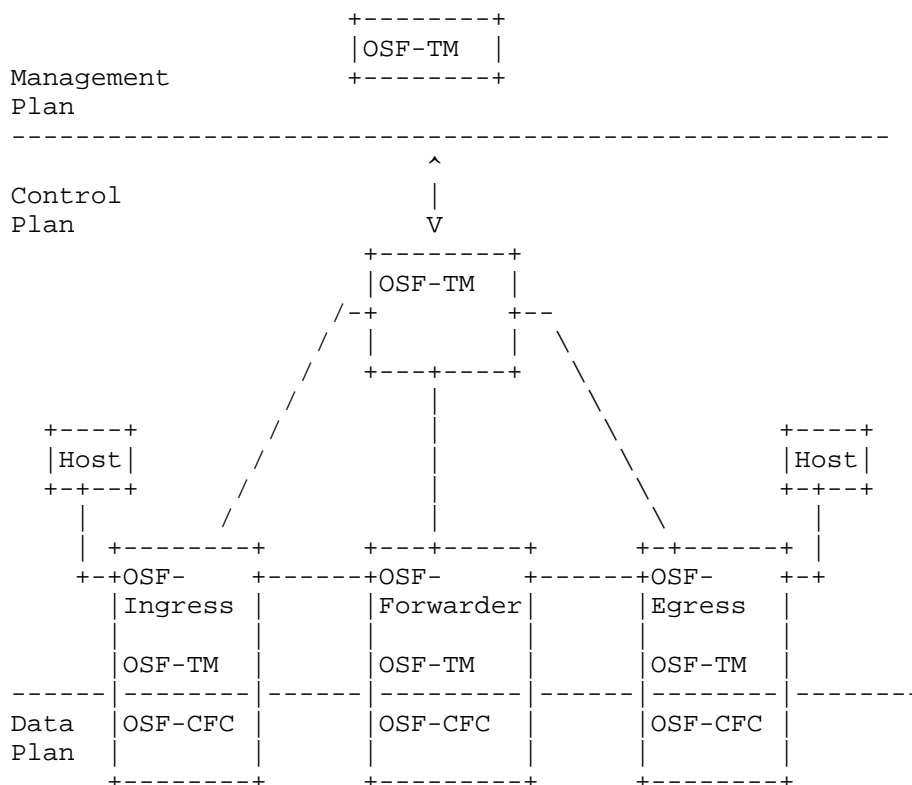


Figure 1: OSF Interactions

For the OSF network, define the following layers:

- o OSF Management Layer: Mainly responsible for monitoring, configuring, and maintaining OSF devices.
- o OSF Control Layer: Mainly responsible for maintaining OSF network topology information, congestion detection, and fast switching.
- o OSF Data Layer: Mainly responsible for encapsulating, forwarding, and decapsulating packets based on the routing information issued by the control layer, and sending packets to the authorization system to achieve congestion control.

## 2.2. OSF Functional Components



## 2.3. OSF-TM

Discovery and Maintenance of OSF Network Topology involves collecting node information, internal connection information between each interface, and external interface information for each node to generate the OSF topology.

OSF-TM is also responsible for maintaining the quality of all links, in order to select the best path for packet transmission based on link quality and avoid network congestion.

## 2.4. OSF-Ingress

The entry point for OSF data packets, where path selection and load balancing are performed based on OSF-MS, encapsulating the packets and sending them towards the exit.

## 2.5. OSF-Egress

The exit point for OSF data packets, where the packets are decapsulated, reordered, and delivered to the recipient.

## 2.6. OSF-Forwarder

The forwarder between OSF entrance and OSF exit forwards based on destination interface information, disregarding the content of the packets.

## 2.7. OSF-CFC

OSF-CFC operates at the data layer and is used for congestion control during OSF packet forwarding. For details on the specifics of congestion control, refer to section 4.3.

## 3. Deployment Considerations

The OSF-TM and OSF-MS components at the control layer can operate in a centralized processing mode, a distributed processing mode, or a hybrid mode.

In the centralized mode, all topology information, network quality information, and metric information are mOSFtOSFed centrally.

In the distributed mode, all topology information and network quality information are maintained in a distributed manner across devices and synchronized.

In the hybrid mode, stable information such as topology information is maintained in a distributed manner, while information that changes frequently is maintained centrally to reduce flooding in the network, such as network quality information. Ultimately, network metric information is generated based on network quality information and maintained centrally.

## 4. OSF Framework Workflow

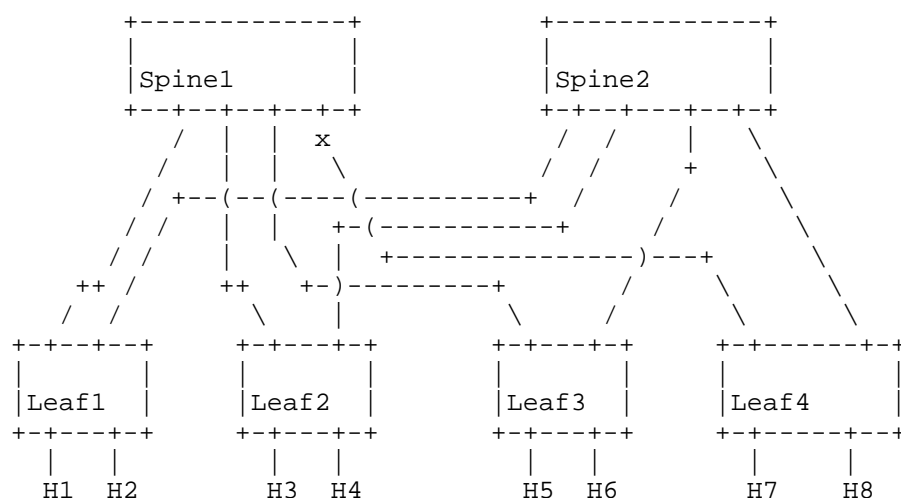
### 4.1. OSF Topology Manage

In an OSF network, the OSF topology can be generated through a topology discovery protocol for use in load balancing across multiple paths during data forwarding.

During load balancing, it is necessary to exclude paths with poor link quality. OSF-TM is responsible for maintaining quality information for each link. Link quality information is detected and reported by the switches, and link state synchronization is achieved through the link state protocol.

In the event of link failure, the switches near the failure point need to be the first to detect and quickly notify the head end, which then performs global traffic scheduling.

OFP-TM can be deployed in a distributed manner or in a centralized manner.



As shown in the figure, topology discovery is performed between devices to dynamically maintain the network topology.

Each device maintains the link quality with its neighbors, and the overall link quality for all links is ultimately maintained by OSF-TM.

When OSF-Ingress sends packets to OSF-Egress, path selection is based on the topology and link quality information.

For example, in the diagram, when H1 sends a packet to H7 and network congestion occurs between Spine1 and Leaf4, Spine1 detects the congestion and notifies Leaf1. Leaf1 then reselects the path, changing the forwarding path to H1->Leaf1->Spine2->Leaf4->H7 based on the new information.

#### 4.2. Load Balancing in OSF Packet Transmission

Traditional load balancing typically involves hashing based on the five-tuple of packets. However, for AI networks, the small amount of traffic and the large load per flow can lead to imbalanced loads.

OSF-Ingress needs to dynamically calculate the bandwidth for each path based on the path bandwidth maintained in OSF-TM, while disregarding paths with excessively high congestion levels.

OSF-Ingress performs load balancing for packets destined for the same interface, allowing packet aggregation. For oversized packets, they can be fragmented into smaller segments for transmission to ensure a more balanced load. After packet aggregation or fragmentation, the packets need to be sorted and numbered, then sent sequentially through ECMP links.

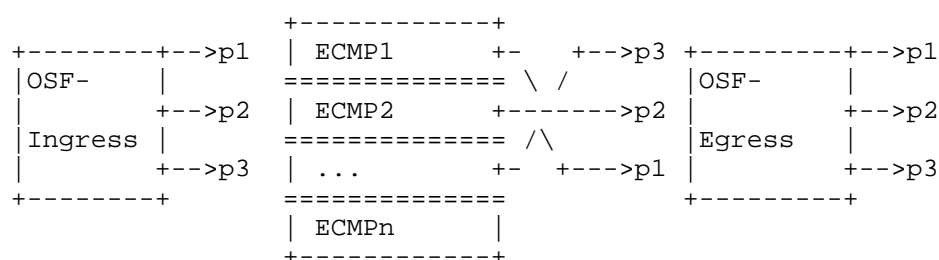
The intermediate OSF-Forwarder devices forward the packets based on the final destination interface information, ultimately sending the packets to OSF-Egress. There is no need to identify the content of the packets, handle packet reordering, or process packet fragmentation and aggregation.

The OSF-Forwarder needs to process authorization requests through the OSF-CFC component. For details about the specific handling process of OSF-CFC, refer to section 4.3.1.

After the packets are sent to OSF-Egress, it is necessary to ensure that the packets are delivered sequentially before handing them over to the receiver. If the received packet is composed of multiple original packets aggregated together, it should be separated into original packets before delivery. If the packet is composed of fragments of a large packet, they need to be reassembled into a complete packet before delivery.

As shown in the diagram below, for three packets p1, p2, and p3 of the same flow, OSF-Ingress no longer performs hash-based route selection. Instead, it sequentially selects the optimal ECMP paths for all packets, ensuring the maximum utilization of the bandwidth across all paths. To ensure sequential delivery at the receiving end, packets need to be numbered and sorted by OSF-Ingress before sending. The OSF-Forwarder devices along the path are responsible for forwarding the packets to OSF-Egress. At OSF-Egress, to ensure that the receiver can receive packets in the original order, the received packets need to be sorted before being delivered to the receiver. In the diagram, the order of the packets received by OSF-Egress is p3->p2->p1. After sorting by OSF-Egress, the packets are delivered to the receiver in the original order p1->p2->p3.

If the packets have undergone fragmentation or aggregation at OSF-Ingress, they also need to be reconstructed into the original packets by OSF-Egress before being delivered to the receiver. This document does not specify the specific format and encapsulation of packet numbering.



#### 4.3. OSF Congestion Control

If network congestion occurs, network performance will severely deteriorate. Therefore, we need to ensure that network congestion is



Congestion control includes both active and passive control.

Passive congestion control typically involves testing the current network state and providing feedback so that the sending end can quickly react to congestion. It controls network congestion by allocating rates based on measurement information.

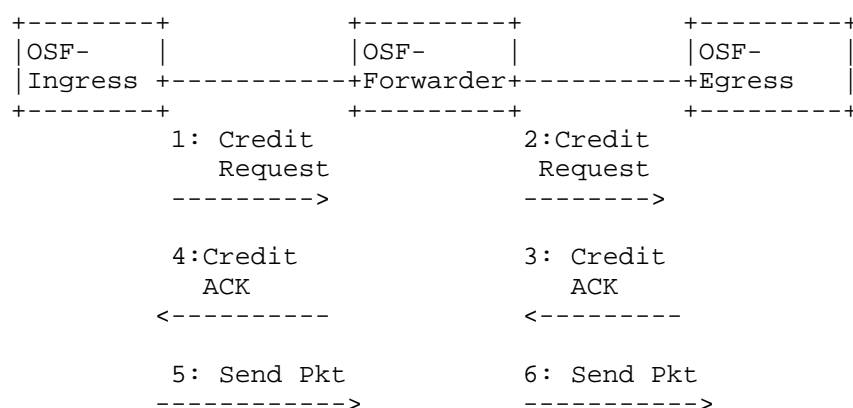
The active congestion control scheme aims to prevent congestion by only sending data when the network has sufficient available bandwidth. In this context, we will mainly introduce the implementation of active congestion control.

OSF adopts a passive congestion control mechanism at the control layer, which involves detecting the link state and adjusting the forwarding path to control network congestion.

At the data layer, OSF adopts an active congestion control mechanism by deploying an authorization mechanism internally. Before sending packets, a scheduling process is initiated. The packet is only sent after successful scheduling, thereby implementing proactive congestion control and achieving deterministic and precise congestion control to strike the optimal balance between congestion and link load.

#### 4.3.1. Credit-based Flow Control

As shown in the diagram, before OSF-Ingress at the data layer sends packets to OSF-Egress, it first sends Credit Request messages along the path to request credit. Upon receiving the Credit Request, the OSF-Forwarder reserves bandwidth on the receiving interface and then continues to send Credit Request messages downstream to request credit, until all devices on the packet transmission path receive the Credit Request and respond with Credit ACK, reserving the receiving bandwidth.



When transmitting within the OSF network, congestion control can be autonomously conducted by the data layer.

It is recommended to deploy an authorization mechanism internally, where packets are scheduled before being sent. This proactive congestion control approach aims to achieve deterministic and precise congestion control by maintaining the optimal balance between congestion and link loads.

If network congestion occurs in the network, network performance will significantly degrade. Therefore, efforts should be made to prevent network congestion. OSF-CFC implements end-to-end network congestion control to minimize incidents of network congestion.

The workflow of OSF-CFC is as follows:

Step 1: The OSF-Ingress sender initiates an authorization request when sending a packet, requesting the required queue bandwidth resources along the transmission path.

Step 2: The next-hop node receives the authorization request, reserves queue bandwidth resources on the local egress interface. If resources are insufficient, an authorization rejection is issued, leading to Step 6.

Step 3: If the destination interface is not local, resources are reserved on the egress interface and the authorization request is further forwarded to the next-hop node.

Step 4: The final authorization request reaches the destination node, which responds with an authorization reply.

Step 5: Upon receiving the authorization reply, the initiator sends the packet to the next-hop node.

Step 6: Upon receiving an authorization rejection, the reserved resources on the local egress interface are released. If the device is not the initiator of the authorization request, the rejection is forwarded to the initiator.

Step 7: Upon receiving an authorization rejection, the initiator releases the reserved resources on the local egress interface and notifies OSF-MS of the network congestion message.

#### 4.3.2. Congestion Control Based on Link Quality Detection

The control layer precisely tests and explicitly feeds back network link state information using all devices, and selects paths on OSF-Ingress based on the network state information, excluding paths with high congestion characteristics.

The process involves link state detection, link state announcement, and path selection. All of these processes are completed within the OSF-TM component.

#### 4.4. Rapid Link Failure Switchover in OSF

As the size of the AI network grows, the number of network cards and optical modules increases, leading to a corresponding rise in the probability of failures.

Current failure switchover mechanisms typically involve nearby devices handling path switching, but lack a global path scheduling method for rapid failure switchover. For OSF networks, it is essential for the control plane to incorporate a fast failure detection and notification mechanism, enabling nodes near the point of failure to swiftly detect issues and promptly notify the OSF-MS component. This allows the OSF-MS component to quickly adjust paths on a global scale in response to the detected failure.

#### 5. Security Considerations

TBD.

#### 6. IANA Considerations

This document does not request any IANA allocations.

#### 7. References

##### 7.1. Normative References

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

[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>>.

##### 7.2. Informative References

TBD

PengFei Huo  
ByteDance  
China  
Email: huopengfei@bytedance.com

Gang Chen  
ByteDance  
China  
Email: chengang.gary@bytedance.com

Changwang Lin  
New H3C Technologies  
China  
Email: linchangwang.04414@h3c.com

Huichen Dai  
ByteDance  
China  
Email: daihuichen@bytedance.com



Expires August 24, 2025

[Page 12]