

onions
Internet-Draft
Intended status: Informational
Expires: 6 April 2026

L. Dunbar, Ed.
Futurewei
Q. Sun
China Telecom
B. Wu, Ed.
Huawei
L. CONTRERASMURILLO, Ed.
Telefonica
C. Xie
China Telecom
3 October 2025

Applying Attachment Circuit and Traffic Engineering YANG Data Model to
Edge AI Use Case
draft-dunbar-onions-ac-te-applicability-00

Abstract

This document explores how existing IETF YANG data models, specifically the Attachment Circuit (AC)-as-a-Service (ACaaS) and Traffic Engineering (TE) topology data models, can be applied to support a use case involving dynamic AI model placement at Edge Cloud sites. The use case involves selecting optimal Edge locations for real-time AI inference based on end-to-end network performance between street-level cameras and Edge Cloud compute nodes. By mapping the use case across multiple network segments and applying relevant YANG data models to retrieve and request specific services objectives such as bandwidth, latency, and reliability, this document serves as a practical exercise to evaluate model applicability and identify gaps, if any.

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 6 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	3
2. Conventions used in this document	4
3. YANG Data Models and APIs to Support Dynamic AI Model Placement	5
3.1. Using AC YANG Data Model to Evaluate Access and Edge Connectivity	9
3.2. Using IETF YANG Data Models to Evaluate PE to PE Connectivity	10
3.3. Non-TE-Based PE to PE Connectivity	13
4. Potential APIs for Kubernetes to Query Network Conditions . .	13
4.1. Behavior and Semantics	14
4.2. YANG Model Integration	15
5. Dynamic UCMP Load Balancing for Periodic Inter Site AI Traffic	15
5.1. Edge Cloud Site to Local PE (Access Segment)	15
5.1.1. Time-Scoped UCMP Policy Enforcement	19
5.1.2. UCMP Enforcement for SRv6 based PE to PE segment . .	20
5.2. Cloud-Initiated UCMP Activation API	21
6. Gaps Between Existing IETF Specifications and Neotec Use Case Needs	21
6.1. Consumption Model Mismatch	22
6.2. Lack of Time-Bounded, Cloud-Triggered Policy Interfaces	22
6.3. Absence of Intent and Workload-to-Flow Mapping	22
6.4. Lack of Abstracted, Queryable Interfaces	22
6.5. Feedback and Acknowledgment Gaps	23
6.6. Security and Access Control Framework Alignment	23
7. Security Considerations	23
8. IANA Considerations	23
9. References	23
9.1. Normative References	23
9.2. Informative References	23

Acknowledgements	24
Contributors	24
Authors' Addresses	24

1. Introduction

This document explores the applicability of the Attachment Circuits YANG data model [I-D.ietf-opsawg-teas-attachment-circuit] and TE topology YANG data model, to support a simplified version of the use case described in China telecom's Neotec side meeting during IETF122(Cloud aware Network Operation for AI Services). Also, the document specifies the APIs needed to support the simplified use case.

Simplified Use Case:

Let's assume there are 10 Edge Cloud sites. An City Surveillance AI model (e.g., detecting traffic congestion or garbage classification) needs to be deployed dynamically to some of these sites in response to real time events.

High level steps for selecting Edge sites to instantiate the City Surveillance AI model:

- Step 1: A Cloud Manager needs to query the network connectivity characteristics (bandwidth, latency, topology constraints, etc.) between street cameras (or gateways, eNB that connect to those street cameras) and candidate Edge Cloud sites in order to determine the optimal locations for the City Surveillance AI model deployment.
- Step 2: Based on the information gathered, the Cloud Manager decides to deploy the City Surveillance AI model in a subset of the Edge Cloud sites (e.g., 4 sites among 10).

High level steps to support the following desired outcome:

- Suppose that the City Surveillance AI model instances in the 4 Edge Cloud sites need to exchange large volumes of data with strict performance constraints (e.g., XX Gbps bandwidth and YY ms one-way delay along the entire end-to-end paths between edge sites). This request is to a network controller to dynamically adjust UCMP (Unequal Cost Multipath) load-balancing algorithms [I-D.ietf-bess-evpn-unequal-lb][I-D.ietf-idr-link-bandwidth] on all the nodes along the paths interconnecting those 4 sites.

Disclaimer

The use of specific YANG data models (e.g., Attachment Circuit and TE topology) in this section is intended as a provisional exercise to explore how existing IETF models might address aspects of such a use case. These examples are not exclusive or exhaustive. Other data models, such as Network Slicing Service Model (NSSM) or service function chaining models, could also be relevant depending on the network architecture and service requirements. The intent is to assess the applicability and identify gaps (if any), not to pre-define the final solution set.

2. Conventions used in this document

Cloud Manager: An entity that is primarily responsible for placing service instances and managing compute resources across Edge Cloud sites. It monitors the health and scaling status of VMs, containers, or services, and makes infrastructure-level decisions, such as where to host service instances based on latency, CPU, GPU availability, or other constraints. In addition, the Cloud Manager may provide abstracted status about the cloud environment (e.g., resource availability) to external entities.

Neotec: Network Operation for Telco Cloud.

Network Orchestrator (or Orchestrator): A logical entity that interfaces with the Cloud Manager to receive service requests or queries and coordinates the end-to-end connectivity across multiple network domains. It abstracts underlying domain-specific technologies (e.g., L2/L3 VPNs, SR paths, TE tunnels) and disseminates policies to individual Network Controllers, enabling seamless stitching of diverse network segments to meet service-level requirements.

Network Controller: A domain specific control entity responsible for managing and configuring network resources within a single administrative or technological domain (e.g., IP/MPLS core, access network). It receives high-level intent or service instructions from a Network Orchestrator and translates them into device-level configurations using protocols such as NETCONF, BGP, or PCEP.

UE: User Equipment

UPF: User Plane Function [TS.23.501-3GPP]

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. YANG Data Models and APIs to Support Dynamic AI Model Placement

Multiple Edge Cloud sites are located in close geographic proximity, allowing any of them to potentially host AI inference model instances for city surveillance tasks such as traffic congestion detection or garbage classification. A Cloud Manager must evaluate the end-to-end data paths between the street-level cameras and each Edge Cloud site to determine which sites offer sufficient compute resources and the required network performance (e.g., low latency and high throughput). This assessment enables dynamic placement of AI inference models in response to real-time events while ensuring reliable data delivery from the field.

This path typically spans three segments:

- The first is the access segment from the cameras local access node, such as a WiFi Access Point or an eNB, to a PE router.
- The second segment traverses the provider transport network between that Access PE and the PE router serving the candidate Edge Cloud site.
- The third segment connects that Edge PE to the Edge Cloud Gateway or compute node where the AI workload is deployed.

There are two primary types of access connectivity for the cameras: via cellular networks (e.g., through eNBs/gNBs and User Plane Functions (UPFs)) or via WiFi Access Points. In cellular based access, cameras connect to the network through eNBs, which forward user data to UPFs via GTP-U tunnels. The UPFs serve as termination points for GTP-U tunnels and are often co-located with or adjacent to Edge Cloud sites. When the Cloud Manager selects Edge Cloud sites for hosting AI inference modules, it must ensure that the corresponding nearby UPFs are assigned as the serving UPFs for the associated eNBs. This enables the establishment of GTP-U tunnels from the eNBs directly to the selected Edge Cloud locations, minimizing latency and improving efficiency for real-time AI processing.

For cameras connected through WiFi Access Points, no GTP tunneling is involved. These access points typically connect to PEs through Layer 2 or Layer 3 links. In this case, the Attachment Circuit (AC) YANG model can be directly used to represent the logical connectivity between the WiFi AP and the Access PE. The Cloud Manager or orchestrator can query the AC model to evaluate operational status, available bandwidth, latency, and packet loss, and determine whether the path is capable of supporting the target AI workload.

In both access scenarios, after determining candidate Edge Cloud sites, the Cloud Manager evaluates the transport network (Segment 2) between the access side PE (or the PE adjacent to the UPF) and the Edge side PE. This segment may be traffic engineered and can be modeled using the IETF TE topology model [RFC8795] and the TE tunnel model [RFC8776]. These data models expose metrics such as available bandwidth, TE delay, and administrative attributes, which the orchestrator can use to assess whether the underlying transport paths meet the end-to-end service constraints (captured in an SLA).

Finally, the last segment (Segment 3: from the PE to the Edge Cloud Gateway or compute node) is again modeled using the AC YANG data model. By combining insights from the AC and TE data models across all three segments, the Cloud Manager and orchestrator can select Edge Cloud sites that not only have available compute capacity but also meet the network performance requirements to support low latency, high bandwidth AI model inference.

Important Clarification:

The Attachment Circuit (AC) and Traffic Engineering (TE) Topology YANG models described above are internal network models used by the network controller and orchestrator. They are not directly exposed to external entities like the Cloud Manager. Instead, the network controller abstracts relevant information from these internal models and presents it through simplified, service-oriented APIs (e.g., onions APIs), allowing the Cloud Manager to query network performance metrics (such as latency and available bandwidth) without accessing raw topology or low-level link state.

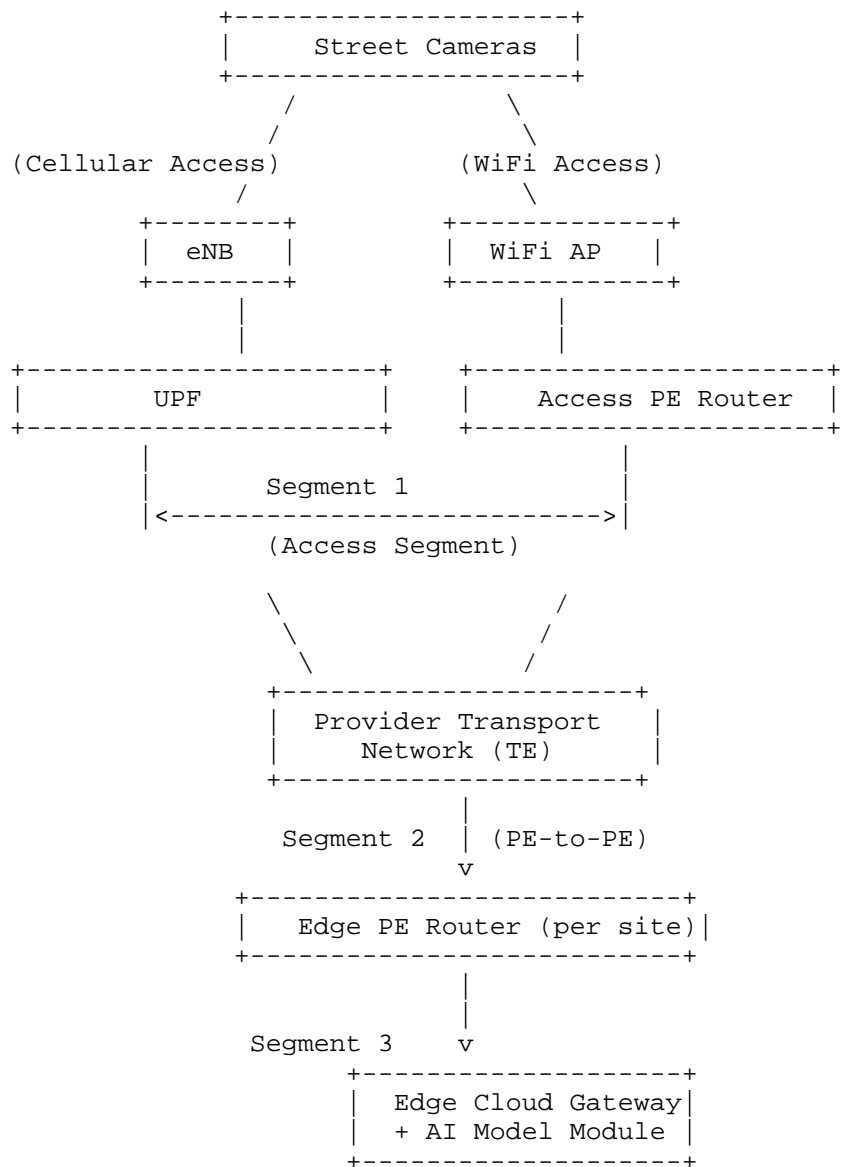


Figure 1: Network Segments

For the first and third segments, the Attachment Circuit (AC) YANG model, defined in [I-D:ietf-opsawg-teas-attachment-circuit], provides a standardized abstraction of the logical link between a service access point and the provider's routing infrastructure. This model

allows the Cloud Manager to retrieve configuration and operational status of these links, including encapsulation type, bandwidth, and link status. While the Bearer Service Model [I-D:ietf-opsawg-teas-attachment-circuit] can provide additional contextual information, such as interface type, location, and administrative details, it is not the primary focus of this document. For the purposes of real-time workload placement and path performance assessment, the emphasis here is on network performance metrics (e.g., latency, available bandwidth) exposed via AC and TE models. Integration with bearer-level metadata (e.g., location) may be considered in future enhancements or broader orchestration frameworks. By querying this information, the orchestrator can determine whether each access circuit is operational and capable of delivering real-time video streams to or from the AI inference point.

The second segment, the PE to PE transport path across the provider network, requires querying a different set of YANG models. The base network and topology models defined in [RFC8345] provide the foundation, while the Traffic Engineering Topology model [RFC8795] adds details such as TE metrics, link delay, bandwidth availability, and administrative constraints. These models are typically maintained by network controllers and made accessible via NETCONF or RESTCONF interfaces. The Cloud Manager can initiate a query to evaluate the performance characteristics of available transport paths between Access PEs and candidate Edge PEs. This query is handled by the Network Orchestrator, which interacts with the Network Controllers managing each domain. The orchestrator uses IETF defined topology and TE models, such as the TE Topology model [RFC8795] and the TE Tunnel model [RFC8776], to retrieve current path attributes, including available bandwidth and latency. In architectures that support a Path Computation Element (PCE), the orchestrator may also issue path computation requests via PCEP or a REST API to identify transport paths that meet the performance objectives required for AI inference workloads.

While L2SM [RFC8466] and L3SM [RFC8299] offer useful abstractions for provisioning Layer 2 and Layer 3 VPN services, they are primarily focused on expressing service level intent rather than evaluating real time transport performance. As such, they are not directly applicable to the dynamic, performance sensitive selection process described in the Neotec use case.

By combining the Attachment Circuit and TE topology models, the Cloud Manager can construct an end to end view of network connectivity from each camera access point to every potential Edge Cloud site. It can then select the subset of Edge Cloud locations that not only have sufficient computing resources but also meet the required network performance criteria for serving the selected set of street cameras.

This allows dynamic, network aware placement of AI inference models, ensuring efficient use of edge infrastructure while meeting service quality requirements.

3.1. Using AC YANG Data Model to Evaluate Access and Edge Connectivity

In the context of dynamic AI model placement at Edge Cloud sites, the AC data model is useful for querying the health and performance of the access (Segment 1) which connects Access Points or eNBs to Access PE routers, and edge (Segment 3) which links Edge PEs to Edge Cloud Gateways or compute nodes.

Each attachment circuit is modeled with attributes such as administrative state, operational status, encapsulation type, and configured bandwidth. When augmented with telemetry or traffic engineering extensions, the AC model also supports real time status such as current utilization, available bandwidth, latency, and even packet loss.

However, the Cloud Manager does not directly interact with the AC model. Instead, the network controller queries and interprets the AC YANG data model and exposes only the relevant service-level metrics through the onions API interface. This ensures that the Cloud Manager can assess whether a given access or edge segment supports the required performance (e.g., latency less than 10 ms, throughput greater than 500 Mbps) without needing detailed knowledge of the network.

For example, if the Cloud Manager specifies a target Edge Cloud site, the network controller may internally retrieve attachment circuit metrics and present a simple API response indicating available bandwidth and latency thresholds met or not met:

```
GET https://<controller>/restconf/data/ietf-ac:attachment-circuits  
/ac[pe-address='192.0.2.11']
```

Figure 2: Get AC Status

```
{
  "ietf-ac:ac": [
    {
      "name": "ac-to-eNB-001",
      "pe-address": "192.0.2.11",
      "ce-interface": "ge-0/0/1",
      "oper-status": "up",
      "admin-status": "enabled",
      "encapsulation-type": "dot1q",
      "bandwidth": {
        "configured": 1000,
        "available": 850,
        "unit": "Mbps"
      },
      "performance-metrics": {
        "latency": {
          "average": 5,
          "max": 10,
          "unit": "ms"
        },
        "packet-loss": {
          "percentage": 0.01
        }
      }
    }
  ]
}
```

Figure 3: Access Segments Response

In this example, the attachment circuit is operational (`oper-status: up`), with 850 Mbps of available bandwidth and average latency of 5 ms, making it a strong candidate for supporting a real time AI application. The Cloud Manager may apply similar queries to all PE addresses associated with candidate Edge Cloud sites and their corresponding access points, filtering the results to identify which circuits meet both latency and throughput thresholds.

3.2. Using IETF YANG Data Models to Evaluate PE to PE Connectivity

Segment 2 of the end-to-end path connects the Access PE to the Edge PE router serving a candidate Edge Cloud site. This segment typically traverses the provider's core or aggregation network and is crucial for ensuring that real-time AI inference traffic can meet stringent latency and bandwidth requirements.

The network controller uses internal models such as the IETF Network Topology YANG model [RFC8345] and the Traffic Engineering Topology YANG model [RFC8795] to monitor transport network characteristics. These models describe nodes, TE links, available bandwidth, delay, SRLGs, and other performance attributes. Similarly, TE tunnels can be modeled using [RFC8776] to represent established LSPs or segment routing paths.

Importantly, the Cloud Manager does not directly access these detailed TE models. Instead, the controller or orchestrator processes this information and presents it via an abstracted API interface, exposing key metrics like end-to-end latency, available bandwidth, and path health between specified endpoints (e.g., Access Node to Edge Site).

This abstraction enables the Cloud Manager to make workload placement decisions based on network performance, without having to manage or interpret the detailed topology or traffic engineering constructs.

For example, to query available transport links from an Access PE at PE1 to a set of Edge PEs, the Cloud Manager may retrieve the list of TE links and their attributes using the following RESTCONF-style request:

```
GET https://<controller>/restconf/data/ietf-te-topology:te-topologies  
/topology=default/te-node=PE1
```

Figure 4: Get PE-PE Path Status

```
{
  "te-links": [
    {
      "name": "PE1-to-PE5",
      "link-id": "link-001",
      "oper-status": "up",
      "te-default-metric": 30,
      "te-bandwidth": {
        "max-link-bandwidth": 10000,
        "available-bandwidth": 7000,
        "unit": "Mbps"
      },
      "delay": {
        "unidirectional-delay": 8,
        "unit": "ms"
      },
      "adjacent-te-node-id": "PE5"
    }
  ]
}
```

Figure 5: PE-PE Segment

In this example, the link from PE1 to PE5 offers 7 Gbps of available bandwidth with 8 ms unidirectional delay, which may satisfy a 500 Mbps, sub-10 ms latency requirement for AI data ingestion. The Cloud Manager can evaluate similar TE links or tunnels to other Edge PEs (e.g., PE6, PE7) and compare their performance characteristics.

Alternatively, if the network has PCE deployed, the Cloud Manager can issue a path computation request, using PCEP to query the end to end path metrics and validate whether a PE to PE segment can meet specified SLA constraints. If Segment Routing (SR-MPLS or SRv6) is deployed, these models can also include SR label stacks or SID lists needed for forwarding decisions.

By querying the TE topology and tunnel models, the orchestrator can build a filtered set of feasible transport segments that support the expected latency and bandwidth for the AI workload. This insight, combined with data from the Attachment Circuit models for Segments 1 and 3, allows the orchestrator to make holistic decisions about AI workload placement and optimal traffic steering across the network.

3.3. Non-TE-Based PE to PE Connectivity

The previous sections assume that the transport network between Access PEs and Edge PEs is traffic-engineered and that the orchestrator has access to detailed models, such as ietf-TE-topology, ietf-TE, or SR policy models. However, in some deployments, particularly those relying on Internet transit or best effort IP core, traffic engineering is not explicitly available.

In these cases, performance visibility and decision-making must rely on more general network telemetry. Useful data can still be obtained through IETF defined YANG models, including:

- RFC 8343 (Interface Management YANG Model) and RFC 8344 (IP Management YANG Model) for retrieving interface status and link-level attributes
- Operational telemetry models, such as ietf-system-telemetry or real time telemetry via YANG Push [RFC8641][RFC8639], to gather metrics such as interface utilization, packet loss, round trip time (RTT), and jitter, especially if measured using active probes or synthetic monitoring.

The Cloud Manager's query can trigger the network orchestrator to collect and aggregate per-hop or per-segment metrics using these models or active measurements (similar to IP SLA). The resulting data allows the Cloud Manager to estimate end-to-end path performance and make workload placement decisions accordingly, even if deterministic path selection or policy enforcement is not possible.

While this approach lacks the fine-grained path control of TE-based networks, it still enables adaptive, network-aware service placement in less structured or loosely-managed environments.

4. Potential APIs for Kubernetes to Query Network Conditions

This section outlines a potential onions API interface to enable Kubernetes (or its external workload orchestrators) to make network-aware workload placement decisions, such as selecting appropriate Edge Cloud sites for deploying latency-sensitive AI inference modules.

Cloud-native platforms such as Kubernetes do not consume YANG-modeled data directly. Instead, they rely on REST-style APIs or gRPC interfaces with JSON or protobuf-encoded payloads. To support Neotec use cases, the network infrastructure can expose an abstracted API that translates YANG-modeled topology and performance data into a form consumable by Kubernetes or its scheduling extensions.

A representative onions API endpoint could take the form:

```
GET /network-advisor/query-path-performance?  
source-node=<access-node-id>&target-node=<edge-site-id>
```

Figure 6: REST style API

This API allows Kubernetes to query the performance of the network path between a street camera's access node (e.g., a WiFi AP or eNB/UPF) and a candidate Edge Cloud site. The API response includes metrics such as available bandwidth, average path latency, and operational status. A sample response might be:

```
{  
  "source-node": "AP-235",  
  "target-node": "EdgeSite-4",  
  "path-latency-ms": 6.5,  
  "available-bandwidth-mbps": 920,  
  "path-status": "healthy"  
}
```

Figure 7: JSON Response

4.1. Behavior and Semantics

This API provides a read-only, low-latency interface for workload schedulers to assess whether a given path meets predefined service-level thresholds (e.g., latency less than 10 ms, bandwidth larger than 500 Mbps). The source node and target node identifiers correspond to access nodes (e.g., PE routers adjacent to UPFs or WiFi APs) and Edge Cloud PE routers respectively. These identifiers are assumed to be mapped within the operator domain.

The semantics of the fields are as follows:

- path-latency-ms: Derived from YANG-modeled metrics in ietf-TE-topology, representing end to end unidirectional delay.
- available-bandwidth-mbps: Aggregated from TE-bandwidth in the same topology model.
- path-status: Reflects operational state derived from oper-status fields in the TE and AC models.

4.2. YANG Model Integration

The backend implementation of this API is expected to query IETF defined YANG models using RESTCONF or NETCONF. Specifically:

- Topology and path delay metrics are sourced from the ietf-te-topology and ietf-network-topology models [RFC8795], [RFC8345].
- Access circuit status and available bandwidth can be derived from the AC model [opsawg-teas-attachment-circuit].

A shim layer exposes REST or gRPC APIs that accept requests from cloud managers (e.g., Kubernetes) and translates them into queries over IETF defined YANG data models using NETCONF or RESTCONF. This architecture enables real-time network path evaluations without requiring cloud systems to interact directly with YANG or underlying network protocols.

5. Dynamic UCMP Load Balancing for Periodic Inter Site AI Traffic

In the Neotec use case described in Section 1, AI inference modules are deployed across four Edge Cloud sites to support distributed city surveillance. These modules periodically exchange large volumes of data, for instance, during result aggregation or synchronized event analysis. These data exchanges are not continuous but are periodic and event driven, requiring guaranteed bandwidth and low latency for short time windows.

The underlying network connecting these Edge Cloud sites typically includes multiple paths between nodes and across multiple network segments. An end-to-end path between Edge Cloud Site A and B spans at least three segments:

- The first is the access segment from the Edge Cloud A to its closest PE. There could be multiple PEs to Edge Cloud A for multi-homing.
- The second segment traverses the provider transport network between the PE serving Edge Cloud A and the PE serving the Edge Cloud B.
- The third segment connects that Edge PE to the Edge Cloud B's Gateway

5.1. Edge Cloud Site to Local PE (Access Segment)

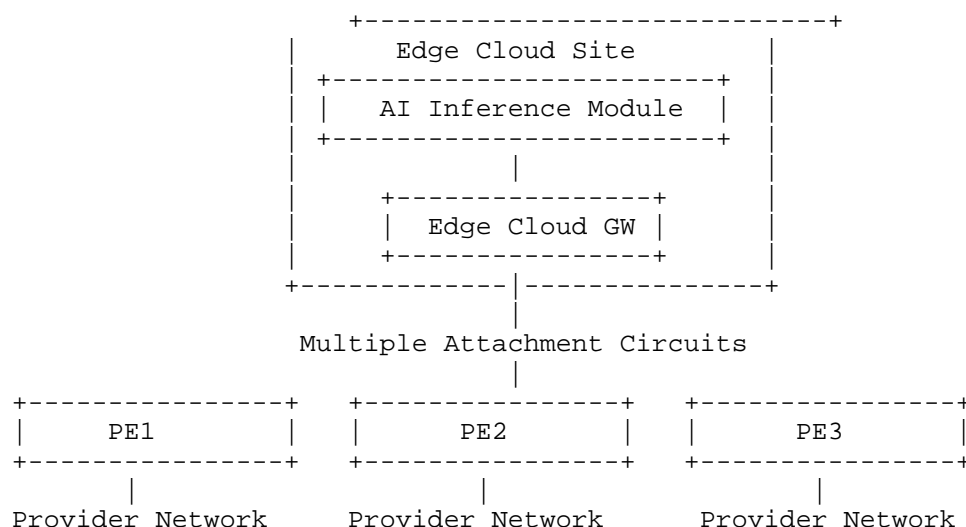


Figure 8: Edge Cloud Access Segment

The Edge Cloud Gateway has multiple logical links (attachment circuits) to a set of PEs (e.g., PE1, PE2, PE3). These links may vary in latency, bandwidth, or current load. During periodic AI data bursts, the Edge Cloud GW must push all other traffic away from PE1, PE2, and PE3, reserving their full available bandwidth for AI flows. Or it could push all other traffic to one of the PEs that has the lowest bandwidth and highest latency. This can be implemented by dynamically updating forwarding policies or QoS profiles to deprioritize or reroute non-AI traffic, ensuring that the AI inference module has uncontested access to network capacity during critical synchronization windows. Depending on the request from the Cloud Manager, the Network Orchestrator can determine what exact policies to push to the PEs and the Edge GW..

YANG Data Models for Policy Enforcement:

To support dynamic UCMP-based traffic steering across PE1, PE2, and PE3, the Network Orchestrator can utilize the following IETF YANG models:

- ietf-routing [RFC8349] enables configuration of routing instances and static routes. The data model allows per-prefix route entries with multiple weighted next hops, supporting unequal cost paths. It can be used to define static next-hop policies from the Edge GW toward PE1, PE2, and PE3.

- ietf-qos-policy [I-D.claise-opsawg-qos-packet-marking] defines traffic classification, marking, and treatment policies. It can enforce rate limits or scheduling priority for non-AI traffic routed to PE3.
- ietf-ac [I-D.ietf-opsawg-teas-attachment-circuit] can provides performance characteristics (bandwidth, latency, packet loss), which inform the decision to assign traffic classes.
- ietf-traffic-classifier / ietf-packet-policy (or OpenConfig equivalents) can be used to match traffic classes (e.g., AI vs non-AI flows), enables forwarding decisions at the Edge GW and ingress of PEs.

Policy Example: Prioritized AI Flow Assignment and Non-AI Rerouting:

Assuming PE1 has the highest available capacity and best latency toward the remote site:

The orchestrator uses ietf-routing [RFC8349] to define weighted static next hops at the Edge GW:

- PE1: 2/3 of AI traffic.
- PE2: 1/3 of AI traffic.
- PE3: no AI traffic.

```

{
  "routing:routing": {
    "routing-instance": [
      {
        "name": "ai-routing-instance",
        "routing-protocols": {
          "routing-protocol": [
            {
              "type": "static",
              "name": "ai-policy",
              "static-routes": {
                "ipv4": {
                  "route": [
                    {
                      "destination-prefix": "198.51.100.0/24",
                      "next-hop": [
                        {"outgoing-interface": "to-PE1", "weight": 66},
                        {"outgoing-interface": "to-PE2", "weight": 34}
                      ]
                    }
                  ]
                }
              }
            }
          ]
        }
      ]
    }
  }
}

```

Figure 9: UCMP policy example

Non-AI traffic is matched by classifiers and redirected entirely to PE3:

- Apply QoS policy using `ietf-qos-policy` [I-D.claise-opsawg-qos-packet-marking] to limit available bandwidth on PE3 to non-AI traffic.

- PE3 may be tagged "degraded" in the `ietf-ac` model to signal its backup nature.

This ensures that the AI traffic is prioritized through the most capable paths (PE1 and PE2) with precise weight. All non-critical traffic is offloaded to the least-desired path (PE3), reducing congestion. Orchestrator needs to dynamically adapt this logic per service request.

Fallback Logic: If available capacity on PE1 and PE2 is not sufficient, AI traffic still receives best possible routing via PE1 and PE2, and Non-AI traffic may be allowed limited access to PE2 with reduced weight (e.g., 10%) while still primarily routed via PE3.

This tiered policy enforcement ensures strict adherence to SLA goals for AI inference while maintaining service continuity for other traffic classes using standard YANG-based configuration interfaces.

5.1.1.1. Time-Scoped UCMP Policy Enforcement

One critical limitation in existing IETF YANG models is the lack of native support for time-scoped UCMP policy activation. To support the bursty nature of AI traffic, the following strategies can be applied:

- Define policy activation via external triggers from the Cloud Manager using an API call to the orchestrator.
- Orchestrator maintains a mapping between the requested activation time window and the temporary configuration to be applied.
- Extend existing YANG models (e.g., `ietf-routing` or `ietf-sr-policy`) with optional augmentation for: `start-time`, `duration`, and `expiration-action`

Example augmentation (conceptual):

```
augment "/routing:routing/routing-instance/static-routes/route" {
  leaf burst-policy-start-time {
    type yang:date-and-time;
  }
  leaf burst-policy-duration-sec {
    type uint32;
  }
  leaf expiration-action {
    type enumeration {
      enum revert-to-default;
      enum retain;
    }
  }
}
```

Figure 10: Example Augmentation

In the absence of standard YANG support, this behavior need to be implemented in the orchestrator logic by maintaining policy lifecycle state and timers, pushing temporary configuration via NETCONF/RESTCONF at start-time, and reverting after duration-sec.

5.1.2. UCMP Enforcement for SRv6 based PE to PE segment

Assuming an SRv6 underlay among the PEs, the network controller can use the ietf-sr-policy YANG model to update the traffic distribution weights across pre-established paths. For example, if three SRv6 paths exist between EdgeSite-A and EdgeSite-C, the controller can push the following configuration to the ingress node:

```
sr-policy {  
  color 4001;  
  endpoint "2001:db8:100::1";  
  candidate-paths {  
    preference 100;  
    path {  
      weight 70;  
      sid-list [2001:db8:10::1, 2001:db8:11::2];  
    }  
    path {  
      weight 20;  
      sid-list [2001:db8:20::1, 2001:db8:21::2];  
    }  
    path {  
      weight 10;  
      sid-list [2001:db8:30::1, 2001:db8:31::2];  
    }  
  }  
}
```

Figure 11: Using SR Policy

This UCMP configuration tells the network to distribute traffic unequally across the three paths based on their capability. The underlying topology and metrics are derived from ietf-TE-topology and ietf-TE models, which expose bandwidth, latency, and available resources for each link.

Similar UCMP behavior can also be implemented over SR-MPLS, MPLS-TE, or enhanced IP networks, using the corresponding IETF YANG models (ietf-TE, ietf-routing, etc.). The key point is that the network paths are preexisting, and the only dynamic action is adjusting how traffic is forwarded among them in response to a cloud service request.

5.2. Cloud-Initiated UCMP Activation API

A simplified example of a cloud-initiated API call to the network controller might look like:

```
POST /network-policy/ucmp-activation
{
  "source-sites": ["EdgeSite-A", "EdgeSite-B"],
  "dest-sites": ["EdgeSite-C", "EdgeSite-D"],
  "start-time": "2025-05-01T10:00:00Z",
  "duration-sec": 300,
  "min-bandwidth-mbps": 5000,
  "max-latency-ms": 10
}
```

Figure 12: Burst Network Request

This request informs the network controller that a high-volume, low-latency data exchange will occur and that UCMP forwarding policies should be applied to optimize transport between the specified sites for the specified duration.

6. Gaps Between Existing IETF Specifications and Neotec Use Case Needs

The Neotec use case, supporting real time, event driven placement and coordination of AI inference workloads across Edge Cloud sites, requires close interaction between cloud orchestration platforms and programmable transport networks. While IETF has standardized robust YANG models for traffic engineering (e.g., *ietf-te-topology*, *ietf-sr-policy*, and *ietf-ac*), these models are network-internal, and fall short in addressing cloud-driven, time-scoped network adaptation requirements.

This document evaluates two core capabilities against existing IETF YANG models:

- Network-aware workload placement at Edge sites
- Dynamic UCMP policy activation for inter-site AI data exchange

From these exercises, the following gaps have been identified:

6.1. Consumption Model Mismatch

Most IETF YANG models are accessed via NETCONF/RESTCONF, and are designed for network operator tools. In contrast, cloud-native environments rely on REST/gRPC APIs, JSON payloads, and declarative interfaces (e.g., Kubernetes CRDs). There is no standardized translation layer that exposes network topology or path performance in a form consumable by external cloud orchestrators or AI services.

6.2. Lack of Time-Bounded, Cloud-Triggered Policy Interfaces

In the UCMP use case, the cloud controller must be able to request changes in traffic distribution policy across existing network paths for a specific time window, such as when AI model instances begin inter-site synchronization. Current IETF models (e.g., ietf-sr-policy) allow weighted path configuration but do not support time-scoping, activation triggers, or scheduling. These functions are essential for on-demand, just-in-time optimization and must be added.

6.3. Absence of Intent and Workload-to-Flow Mapping

There is no YANG model or API that allows the cloud controller to associate a workload (e.g., "AI inference service X") with network traffic that should receive enhanced treatment. While SR policies can be assigned to colors or DSCPs, there is no abstracted intent interface for service-aware flow classification or network SLA expression based on application context.

6.4. Lack of Abstracted, Queryable Interfaces

Even for read-only functions like path selection, there is no IETF-standardized API to answer high-level queries such as:

- "Which Edge sites have less than 10ms latency from these access nodes?"
- "What is the bandwidth-latency profile of the path from Access PE A to Edge PE B?"

Operators must instead manually interpret TE link state and construct custom tooling. A standardized API to query the aggregated path metrics between logical service endpoints is missing.

6.5. Feedback and Acknowledgment Gaps

There is no defined mechanism for the network controller to confirm whether a policy request (e.g., UCMP activation) was accepted or enforced, nor to notify the cloud when SLA targets are not being met during the activation window. A bi-directional feedback channel is required to support closed-loop coordination between the cloud and the network.

6.6. Security and Access Control Framework Alignment

While the IETF has defined robust identity and access control mechanisms (e.g., OAuth2, RPKI, TLS), integrating these mechanisms with cloud-native identity systems (such as Kubernetes RBAC, SPIFFE/SPIRE, or cloud IAM services) is still ad hoc. A standard framework for mutual trust establishment and token exchange between cloud and network domains is needed to support secure onions APIs under Zero Trust principles.

7. Security Considerations

To be added

8. IANA Considerations

None

9. References

9.1. Normative References

9.2. Informative References

[Neotec-Zerotrust-Access]

L. Dunbar and H. Chihi, "Neotec-Zerotrust-Access", December 2024, <<https://datatracker.ietf.org/doc/draft-dunchihi-neotec-zerotrust-access-dm/>>.

[opsawg-teas-attachment-circuit]

M. Boucadair, et al, "opsawg-teas-attachment-circuit", January 2025, <<https://datatracker.ietf.org/doc/draft-ietf-opsawg-teas-attachment-circuit/>>.

[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>>.

- [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>>.
- [RFC8345] Clemm, A., Medved, J., Varga, R., Bahadur, N., Ananthakrishnan, H., and X. Liu, "A YANG Data Model for Network Topologies", RFC 8345, DOI 10.17487/RFC8345, March 2018, <<https://www.rfc-editor.org/info/rfc8345>>.
- [RFC8639] Voit, E., Clemm, A., Gonzalez Prieto, A., Nilsen-Nygaard, E., and A. Tripathy, "Subscription to YANG Notifications", RFC 8639, DOI 10.17487/RFC8639, September 2019, <<https://www.rfc-editor.org/info/rfc8639>>.
- [RFC8641] Clemm, A. and E. Voit, "Subscription to YANG Notifications for Datastore Updates", RFC 8641, DOI 10.17487/RFC8641, September 2019, <<https://www.rfc-editor.org/info/rfc8641>>.
- [RFC8776] Saad, T., Gandhi, R., Liu, X., Beeram, V., and I. Bryskin, "Common YANG Data Types for Traffic Engineering", RFC 8776, DOI 10.17487/RFC8776, June 2020, <<https://www.rfc-editor.org/info/rfc8776>>.
- [RFC8795] Liu, X., Bryskin, I., Beeram, V., Saad, T., Shah, H., and O. Gonzalez de Dios, "YANG Data Model for Traffic Engineering (TE) Topologies", RFC 8795, DOI 10.17487/RFC8795, August 2020, <<https://www.rfc-editor.org/info/rfc8795>>.
- [TS.23.501-3GPP] 3rd Generation Partnership Project (3GPP), "System Architecture for 5G System; Stage 2, 3GPP TS 23.501 v2.0.1", December 2017.

Acknowledgements

The authors would like to thank for following for discussions and providing input to this document: xxx.

Contributors

Authors' Addresses

Linda Dunbar (editor)
Futurewei
United States of America
Email: ldunbar@futurewei.com

Qiong
China Telecom
China
Email: sunqiong@chinatelecom.cn

Wu Bo (editor)
Huawei
China
Email: lane.wubo@huawei.com

LUIS MIGUEL CONTRERAS MURILLO (editor)
Telefonica
Spain
Email: luismiguel.contrerasmurillo@telefonica.com

ChongFeng Xie
China Telecom
China
Email: xiechf@chinatelecom.cn