

CATS
Internet-Draft
Intended status: Standards Track
Expires: 2 February 2026

C. Li
Huawei Technologies
Z. peng
China Mobile
J. Drake
Independent
1 August 2025

Computing-Aware Traffic Steering (CATS) Using Segment Routing
draft-lbddd-cats-dp-sr-05

Abstract

This document describes a solution that adheres to the Computing-Aware Traffic Steering (CATS) framework. The solution uses anycast IP addresses as the CATS service identifier and Segment Routing (SR) as the data plane encapsulation to achieve computing-aware traffic steering among multiple services instances.

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 2 February 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. Terminology	3
3. Solution Overview	3
3.1. Realization of CATS Framework Components	3
3.1.1. CATS Identifiers	3
3.1.2. CATS Components	3
3.2. Realization of the CATS Framework Workflow	4
3.2.1. Service Announcement	4
3.2.2. Metrics Distribution	4
3.2.3. Service Demand Processing	5
3.2.4. Service Instance Affinity	7
4. Security Considerations	7
5. IANA Considerations	7
6. Acknowledgements	7
7. References	7
7.1. Normative References	7
7.2. Informative References	8
Contributors	8
Authors' Addresses	9

1. Introduction

As described in [I-D.ietf-cats-usecases-requirements], traffic steering that takes into account computing resource metrics would benefit several services, e.g., latency-sensitive service like immersive services that rely upon the use of augmented reality or virtual reality (AR/VR) techniques.

[I-D.ietf-cats-framework] defines a framework for Computing-Aware Traffic Steering (CATS). Such a framework defines an approach for making compute- and network-aware traffic steering decisions in networking environments where services are deployed in many locations.

The CATS framework is an overlay framework for the selection of the suitable service contact instance for placing a service request. The exact characterization of 'suitable' will be determined by a combination of networking and computing metrics. The CATS framework does not assume any specific data plane and control plane solutions.

This document proposes a data plane solution for the realization of CATS. The solution uses an anycast IP address as the Computing-aware Service ID (CS-ID) associated with a service. Also, the solution uses Segment Routing (SR) as the data plane encapsulation from an Ingress CATS-Router to an Egress CATS-Router.

2. Terminology

This document makes use of the terms defined in [I-D.ietf-cats-framework].

Note: Terms such as CATS Instance Selector ID (CIS-ID) may be updated to echo what will be agreed in the CATS framework [I-D.ietf-cats-framework].

3. Solution Overview

This section describes the details of realizing CATS identifiers, CATS components, and workflow.

3.1. Realization of CATS Framework Components

3.1.1. CATS Identifiers

A CATS Service ID (CS-ID) is an anycast IPv4 or IPv6 address. Such an IP address is associated with a specific service that is reachable via one or multiple service contact instances.

The CATS overlay encapsulation is established from an Ingress CATS-Router to an Egress CATS-Router connected to a service contact instance. The service contact instance is typically hosted in a service site.

Depending on the deployment requirements, CIS-IDs may be needed to indicate where to forward the packet to a specific interface pointing to a specific site in the case that multiple sites connect to the same Egress CATS-Router.

3.1.2. CATS Components

In the context of this document, CATS-Routers are required to support SR encapsulation, including SR-MPLS [RFC8660] and SRv6 [RFC8986].

The CATS Traffic Classifier (C-TC) is assumed to be running on Ingress CATS-Routers.

For each service site, one or multiple C-SMAs and C-NMAs can be implemented within the site to collect the metrics of the service instances.

3.2. Realization of the CATS Framework Workflow

3.2.1. Service Announcement

The service anycast IP address may be announced using a rendezvous service (DNS, for example). Clients can obtain the CS-ID of the service from the rendezvous service used by the application (e.g., DNS). It is out of scope of this document to provide a comprehensive list of all candidate rendezvous services.

3.2.2. Metrics Distribution

As per the CATS framework, CS-ID routes with metrics are distributed among the overlay CATS Routers. The detailed control plane solutions of metrics distribution are out of the scope of this document. However, a sample procedure is provided for the readers' convenience.

For example, BGP can be used to distribute CS-ID routes with metrics.

In the case of the C-SMA running as stand alone outside an Egress CATS-Router, the C-SMA collects the metrics of computing resource within a service site and distributes the CS-ID routes with the collected metrics to the Egress CATS-Router. Egress CATS-Routers will generate the new metrics combined with network metrics and computing-related metrics, and redistribute the CS-ID route to Ingress CATS-Routers. In the case of the C-SMA running as a logic entity on an Egress CATS-Router, the same process will be performed inside the Egress CATS-Router.

As described in Section 3.4 of [I-D.ietf-cats-framework], CATS can be deployed in a distributed model, centralized model, or a hybrid model. In a centralized model or hybrid model, the routes with metrics may be collected by centralized controllers. BGP-LS may be a candidate solution to collect the route with metrics from CATS-Routers to controllers; the use of BGP-LS is however out of the scope of this document.

A centralized controller may also install the forwarding policy on Ingress CATS-Routers to steer the traffic; how these policies are communicated to the routers is out of the scope of this document.

3.2.3. Service Demand Processing

Two SR [RFC8402] data plane approaches are supported: SRv6 [RFC8986] and SR-MPLS [RFC8660]. This section introduces a solution based upon SRv6 and SR-MPLS as data planes for CATS purposes.

An Ingress CATS-Router generates SRv6/SR-MPLS encapsulations from itself to Egress CATS-Routers according to the SR policy received from a controller. An Ingress CATS-Router receives service routes with network and computing-related metrics from Egress CATS-Routers. An C-PS will select the best service site according to the received service routes and routing policies. Once the best service site is selected, the associated Egress CATS-Router can be determined and the appropriate SR encapsulation from an Ingress CATS-Router to the C-PS-computed Egress CATS-Router can be selected.

When a service demand is received by an Ingress CATS-Router, it is classified by the C-TC component. When a matching classification entry is found for this demand, the Ingress CATS-Router encapsulates and forwards it to the C-PS selected Egress CATS-Router via the matching SR tunnel.

3.2.3.1. SRv6

As shown in Figure 1, SRv6 tunnels are established from Ingress CATS-Routers to Egress CATS-Routers.

There may be multiple encapsulations from a single Ingress CATS-Router to different Egress CATS-Routers so that the ingress can choose the best Egress CATS-Router connected to the target site.

Furthermore, there may be multiple tunnels from a single Ingress CATS-Router to a single Egress CATS-Router, e.g., to provide different connectivity performance guarantees.

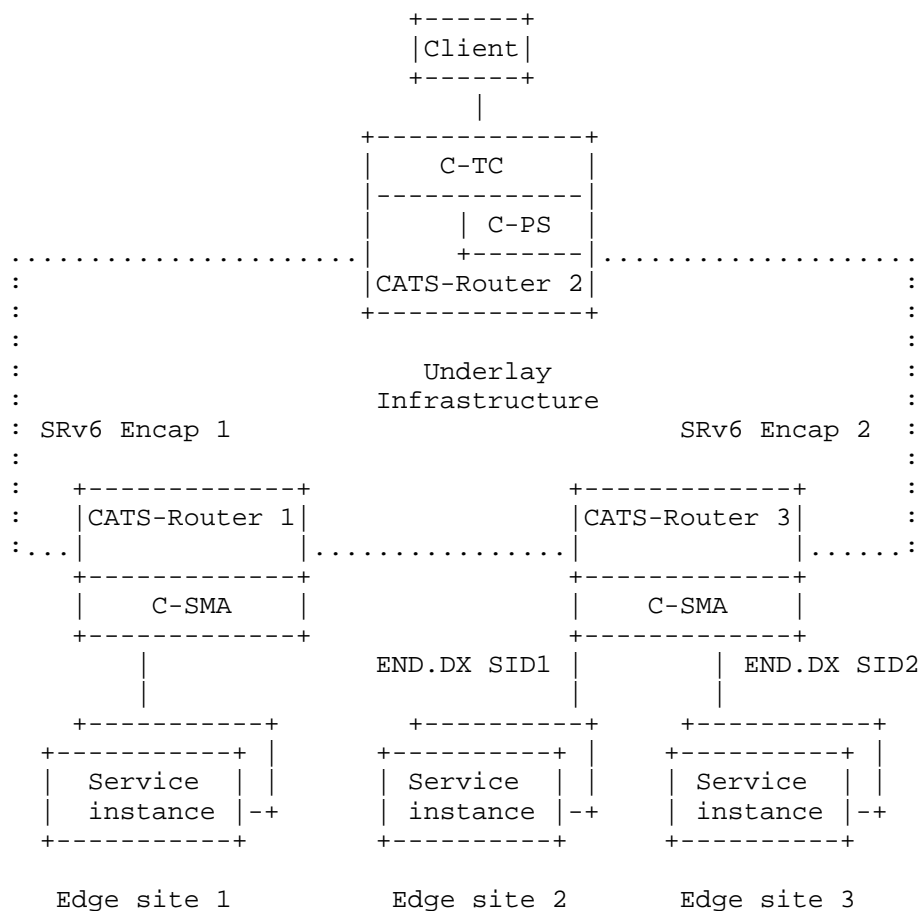


Figure 1: Using SRv6 in CATS

In some cases, multiple service sites may be connected to a single Egress CATS-Router. To demux these sites, a specific attachment circuit must be provided to indicate the specific target service. In order to explicitly indicate the interface towards a site, an END.DX [RFC8986] is encoded as the last segment in the SRv6 encapsulation. The associated END.DX is learned from the control plane.

When the traffic reaches the Egress CATS-Router, the SRv6 packet is decapsulated and the traffic is forwarded to the service contact instance. How the packet is handled beyond that point is out of the scope.

3.2.3.2. SR-MPLS

Similarly, SR-MPLS can be used as the overlay CATS encapsulation. The forwarding path is encoded as an MPLS label stack, and a potential VPN label can be included as the last label to indicate to steer the traffic through a specific interface to a target service contact instance in the case multiple service sites connect to the same Egress CATS-Router.

3.2.4. Service Instance Affinity

As per [I-D.ietf-cats-framework], different services may have different notions of what constitutes a 'flow' and may thus identify a flow differently. Typically, a flow is identified by the 5-tuple transport coordinates (source and destination addresses, source and destination port numbers, and protocol).

Note: This section will be updated to reflect the discussion in the WG about affinity.

4. Security Considerations

This document specifies a CATS solution using anycast IP addresses as CS-IDs and SR as data plane. It does not introduce further security threats considering to the existing ones in [RFC8402], [RFC8660], [RFC8986] and [I-D.ietf-cats-framework].

Anycast-related security considerations are discussed in Section 4.4 of [RFC7094].

5. IANA Considerations

This document makes no requests for IANA action.

6. Acknowledgements

Many thanks to Mohamed Boucadair for his valuable help on this document.

7. References

7.1. Normative References

`[I-D.ietf-cats-framework]`

Li, C., Du, Z., Boucadair, M., Contreras, L. M., and J. Drake, "A Framework for Computing-Aware Traffic Steering (CATS)", Work in Progress, Internet-Draft, draft-ietf-cats-framework-11, 7 July 2025, <<https://datatracker.ietf.org/doc/html/draft-ietf-cats-framework-11>>.

[RFC8402] Filsfils, C., Ed., Previdi, S., Ed., Ginsberg, L., Decraene, B., Litkowski, S., and R. Shakir, "Segment Routing Architecture", RFC 8402, DOI 10.17487/RFC8402, July 2018, <<https://www.rfc-editor.org/info/rfc8402>>.

[RFC8660] Bashandy, A., Ed., Filsfils, C., Ed., Previdi, S., Decraene, B., Litkowski, S., and R. Shakir, "Segment Routing with the MPLS Data Plane", RFC 8660, DOI 10.17487/RFC8660, December 2019, <<https://www.rfc-editor.org/info/rfc8660>>.

[RFC8986] Filsfils, C., Ed., Camarillo, P., Ed., Leddy, J., Voyer, D., Matsushima, S., and Z. Li, "Segment Routing over IPv6 (SRv6) Network Programming", RFC 8986, DOI 10.17487/RFC8986, February 2021, <<https://www.rfc-editor.org/info/rfc8986>>.

7.2. Informative References

`[I-D.ietf-cats-usecases-requirements]`

Yao, K., Contreras, L. M., Shi, H., Zhang, S., and Q. An, "Computing-Aware Traffic Steering (CATS) Problem Statement, Use Cases, and Requirements", Work in Progress, Internet-Draft, draft-ietf-cats-usecases-requirements-07, 10 June 2025, <<https://datatracker.ietf.org/doc/html/draft-ietf-cats-usecases-requirements-07>>.

[RFC7094] McPherson, D., Oran, D., Thaler, D., and E. Osterweil, "Architectural Considerations of IP Anycast", RFC 7094, DOI 10.17487/RFC7094, January 2014, <<https://www.rfc-editor.org/info/rfc7094>>.

Contributors

Dirk Trossen
Huawei Technologies
Email: dirk.trossen@huawei.com

Luigi Iannone
Huawei Technologies
Email: luigi.iannone@huawei.com

Yizhou Li
Huawei Technologies
Email: liyizhou@huawei.com

Hang Shi
Huawei Technologies
Email: shihang9@huawei.com

Authors' Addresses

Cheng Li
Huawei Technologies
China
Email: c.l@huawei.com

Zongpeng Du
China Mobile
China
Email: duzongpeng@chinamobile.com

John E Drake
Independent
United States of America
Email: je_drake@yahoo.com