

Computing-Aware Traffic Steering
Internet-Draft
Intended status: Standards Track
Expires: 6 August 2026

Y. Kehan
China Mobile
C. Li
Huawei Technologies
L. M. Contreras
Telefonica
J. Ros-Giralt
Qualcomm Europe, Inc.
G. Zeng
Huawei Technologies
2 February 2026

CATS Metrics Definition
draft-ietf-cats-metric-definition-05

Abstract

Computing-Aware Traffic Steering (CATS) is a traffic engineering approach that optimizes the steering of traffic to a given service instance by considering the dynamic nature of computing and network resources. In order to consider the computing and network resources, a system needs to share information (metrics) that describes the state of the resources. Metrics from network domain have been in use in network systems for a long time. This document defines a set of metrics from the computing domain used for CATS.

Discussion Venues

This note is to be removed before publishing as an RFC.

Discussion of this document takes place on the Computing-Aware Traffic Steering Working Group mailing list (cats@ietf.org), which is archived at <https://mailarchive.ietf.org/arch/browse/cats/>.

Source for this draft and an issue tracker can be found at <https://github.com/VMatrix1900/draft-cats-metric-definition>.

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 August 2026.

Copyright Notice

Copyright (c) 2026 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 and Definitions	4
3. Design Principles	4
3.1. Three-Level Metrics	4
3.2. Level 0: Raw Metrics	5
3.3. Level 1: Normalized Metrics in Categories	6
3.4. Level 2: Single Normalized Metric.	7
4. CATS Metrics Framework and Specification	8
4.1. CATS Metric Fields	8
4.2. Aggregation and Normalization Functions	11
4.2.1. Aggregation	11
4.2.2. Normalization	12
4.3. On the Meaning of Scores in Heterogeneous Metrics Systems	13
4.4. Level Metric Representations	13
4.4.1. Level 0 Metrics	14
4.4.2. Level 1 Metrics	14
4.4.3. Level 2 Metrics	15
5. Comparison among Metric Levels	16
6. CATS L2 Metric Registry Entry	18
6.1. Summary	18
6.1.1. ID (Identifier)	18
6.1.2. Name	18
6.1.3. URI	18
6.1.4. Description	19

6.1.5.	Change Controller	19
6.1.6.	Version	19
6.2.	Metric Definition	19
6.2.1.	Reference Definition	19
6.2.2.	Fixed Parameters	19
6.3.	Method of Measurement	19
6.3.1.	Reference Methods	19
6.3.2.	Packet Stream Generation	20
6.3.3.	Traffic Filtering (Observation) Details	20
6.3.4.	Sampling Distribution	20
6.3.5.	Runtime Parameters and Data Format	20
6.3.6.	Roles	20
6.4.	Output	20
6.4.1.	Type	21
6.4.2.	Reference Definition	21
6.4.3.	Metric Units	21
6.4.4.	Calibration	21
6.5.	Administrative Items	21
6.5.1.	Status	21
6.5.2.	Requester	21
6.5.3.	Revision	21
6.5.4.	Revision Date	21
6.5.5.	Comments and Remarks	21
7.	Implementation Guidance on Using CATS Metrics	22
8.	Security Considerations	22
9.	IANA Considerations	22
10.	References	22
10.1.	Normative References	22
10.2.	Informative References	23
Appendix A.	Appendix A	24
A.1.	Level 0 Metric Representation Examples	24
A.1.1.	Compute Raw Metrics	24
A.1.2.	Communication Raw Metrics	24
A.1.3.	Delay Raw Metrics	25
Contributors	25
Authors' Addresses	26

1. Introduction

Service providers are deploying computing capabilities across the network for hosting applications such as distributed AI workloads, AR/VR and driverless vehicles, among others. In these deployments, multiple service instances are replicated across various sites to ensure sufficient capacity for maintaining the required Quality of Experience (QoE) expected by the application. To support the selection of these instances, a framework called Computing-Aware Traffic Steering (CATS) is introduced in [I-D.ietf-cats-framework].

CATS is a traffic engineering approach that optimizes the steering of traffic to a given service instance by considering the dynamic nature of computing and network resources. To achieve this, CATS components require performance metrics for both communication and compute resources. Since these resources are deployed by multiple providers, standardized metrics are essential to ensure interoperability and enable precise traffic steering decisions, thereby optimizing resource utilization and enhancing overall system performance.

Metrics from network domain have already been defined in previous documents, e.g., [RFC9439], [RFC8912], and [RFC8911], and been in use in network systems for a long time. This document focuses on categorizing the relevant metrics at the computing domain for CATS into three levels based on their complexity and granularity.

2. Conventions and Definitions

This document uses the following terms defined in [I-D.ietf-cats-framework]:

- * Computing-Aware Traffic Steering (CATS)
- * Service
- * Service site
- * Service contact instance
- * CATS Service Contact Instance ID (CSCI-ID)
- * CATS Service Metric Agent (C-SMA)
- * CATS Network Metric Agent (C-NMA)

3. Design Principles

3.1. Three-Level Metrics

As outlined in [I-D.ietf-cats-usecases-requirements], the resource model that defines CATS metrics **MUST** be scalable, ensuring that its implementation remains within a reasonable and sustainable cost. Additionally, it **MUST** be useful in practice. To that end, a CATS system should select the most appropriate metric(s) for instance selection, recognizing that different metrics may influence outcomes in distinct ways depending on the specific use case.

Introducing a definition of metrics requires balancing the following trade-off: if the metrics are too fine-grained, they become unscalable due to the excessive number of metrics that must be communicated through the metrics distribution protocol. (See [I-D.rcr-opsawg-operational-compute-metrics] for a discussion of metrics distribution protocols.) Conversely, if the metrics are too coarse-grained, they may not have sufficient information to enable proper operational decisions.

Conceptually, it is necessary to define at least two fundamental levels of metrics: one comprising all raw metrics, and the other representing a simplified form---consisting of a single value that encapsulates the overall capability of a service instance.

However, such a definition may, to some extent, constrain implementation flexibility across diverse CATS use cases. Implementers often seek balanced approaches that consider trade-offs among encoding complexity, accuracy, scalability, and extensibility.

To ensure scalability while providing sufficient detail for effective decision-making, this document provides a definition of metrics that incorporates three levels of abstraction:

- * ***Level 0 (L0): Raw metrics.*** These metrics are presented without abstraction, with each metric using its own unit and format as defined by the underlying resource.
- * ***Level 1 (L1): Metrics normalized within categories.*** These metrics are derived by aggregating L0 metrics into multiple categories, such as network and computing. Each category is summarized with a single L1 metric by normalizing it into a value within a defined range of scores.
- * ***Level 2 (L2): Single normalized metric.*** These metrics are derived by aggregating lower level metrics (L0 or L1) into a single L2 metric, which is then normalized into a value within a defined range of scores.

3.2. Level 0: Raw Metrics

Level 0 metrics encompass detailed, raw metrics, including but not limited to:

- * **CPU:** Base Frequency, boosted frequency, number of cores, core utilization, memory bandwidth, memory size, memory utilization, power consumption.

- * GPU: Frequency, number of render units, memory bandwidth, memory size, memory utilization, core utilization, power consumption.
- * NPU: Computing power, utilization, power consumption.
- * Network: Bandwidth, capacity, throughput, bytes transmitted, bytes received, host bus utilization.
- * Storage: Available space, read speed, write speed.
- * Delay: Time taken to process a request.

L0 metrics serve as foundational data and do not require classification. They provide basic information to support higher-level metrics, as detailed in the following sections.

L0 metrics can be encoded and exposed using an Application Programming Interface (API), such as a RESTful API, and can be solution-specific. Different resources can have their own metrics, each conveying unique information about their status. These metrics can generally have units, such as bits per second (bps) or floating point instructions per second (flops).

Regarding network-related information, [RFC8911] and [RFC8912] define various performance metrics and their registries. Additionally, in [RFC9439], the ALTO WG introduced an extended set of metrics related to network performance, such as throughput and delay. For compute metrics, [I-D.rcr-opsawg-operational-compute-metrics] lists a set of cloud resource metrics.

3.3. Level 1: Normalized Metrics in Categories

L1 metrics are organized into distinct categories, such as computing, communication, service, and composed metrics. Each L0 metric is classified into one of these categories. Within each category, a single L1 metric is computed using an `_aggregation function_` and normalized to a unitless score that represents the performance of the underlying resources according to that category. Potential categories include:

- * ***Computing:** A normalized value derived from computing-related L0 metrics, such as CPU, GPU, and NPU utilization.
- * ***Communication:** A normalized value derived from communication-related L0 metrics, such as communication throughput.
- * ***Service:** A normalized value derived from service-related L0 metrics, such as tokens per second and service availability

- * ***Composed:** A normalized value derived from an aggregation function that takes as input a combination of computing, communication and service metrics. For example, end-to-end delay computed as the sum of all delays along a path.

Editor note: detailed categories can be updated according to the CATS WG discussion.

L0 metrics, such as those defined in [RFC8911], [RFC8912], [RFC9439], and [I-D.rcr-opsawg-operational-compute-metrics], can be categorized into the aforementioned categories. Each category will employ its own aggregation function (e.g., weighted summary) to generate the normalized value. This approach allows the protocol to focus solely on the metric categories and their normalized values, thereby avoiding the need to process solution-specific detailed metrics.

3.4. Level 2: Single Normalized Metric.

The L2 metric is a single score value derived from the lower level metrics (L0 or L1) using an aggregation function. Different implementations may employ different aggregation functions to characterize the overall performance of the underlying compute and communication resources. The definition of the L2 metric simplifies the complexity of collecting and distributing numerous lower-level metrics by consolidating them into a single, unified score.

TODO: Some implementations may support the configuration of Ingress CATS-Forwarders with the metric normalizing method so that it can decode the information from the L1 or L0 metrics.

Figure 1 provides a summary of the logical relationships between metrics across the three levels of abstraction.

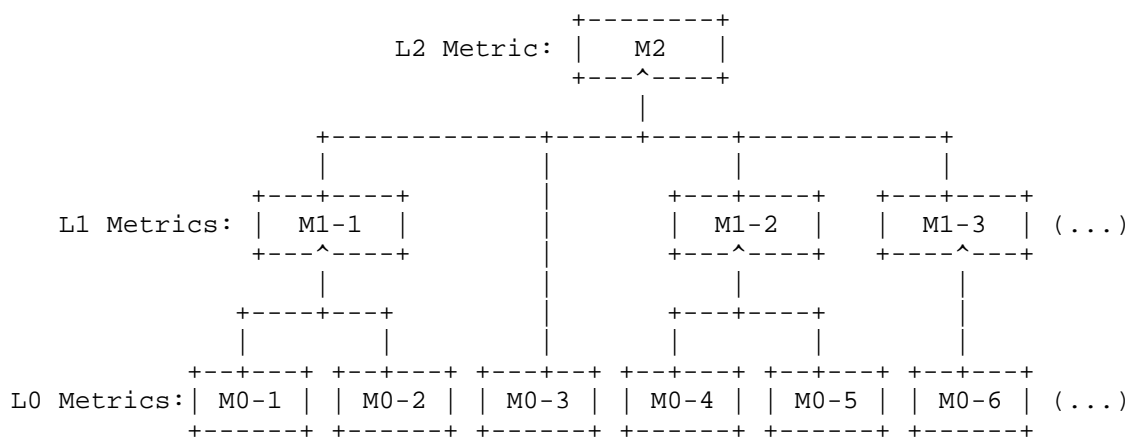


Figure 1: Logic of CATS Metrics in levels

4. CATS Metrics Framework and Specification

The CATS metrics framework is a key component of the CATS architecture. It defines how metrics are encoded and transmitted over the network. The representation should be flexible enough to accommodate various types of metrics along with their respective units and precision levels, yet simple enough to enable easy implementation and deployment across heterogeneous edge environments.

4.1. CATS Metric Fields

This section defines the detailed structure used to represent CATS metrics. The design follows principles established in related IETF specifications, such as the network performance metrics outlined in [RFC9439].

Each CATS metric is expressed as a structured set of fields, with each field describing a specific property of the metric. The following definition introduces the fields used in the CATS metric representations.

- Cats_metric:
 - Metric_type:
 - The type of the CATS metric.
 - Examples: compute_cpu, storage_disk_size, network_bw, compute_delay, network_delay, compute_norm, storage_norm, network_norm, delay_norm.
 - Format:
 - The encoding format of the metric.
 - Examples: int, float.
 - Format_std (optional):
 - The standard used to encode and decode the value field according to the format field.
 - Example: ieee_754, ascii.
 - Length:
 - The size of the value field measured in octets.
 - Examples: 2, 4, 8, 16, 32, 64.
 - Unit:
 - The unit of this metric.
 - Examples: mhz, ghz, byte, kbyte, mbyte, gbyte, bps, kbps, mbps, gbps, tbps, tflops, none.
 - Source (optional):
 - The source of information used to obtain the value field.
 - Examples: nominal, estimation, normalization, aggregation.
 - Statistics(optional):
 - The statistical function used to obtain the value field.
 - Examples: max, min, mean, cur.
 - Level:
 - The level this metric belongs to.
 - Examples: L0, L1, L2.
 - Value:
 - The value of this metric.
 - Examples: 12, 3.2.

Figure 2: CATS Metric Fields

Next, we describe each field in more detail:

- * *Metric_Type (type)*: This field specifies the category or kind of CATS metric being measured, such as computational resources, storage capacity, or network bandwidth. It acts as a label that enables network devices to identify the purpose of the metric.
- * *Format (format)*: This field indicates the data encoding format of the metric, such as whether the value is represented as an integer, a floating-point number, or has no specific format.

- * ***Format standard (format_std, optional)*:** This optional field indicates the standard used to encode and decode the value field according to the format field. It is only required if the value field is encoded using a specific standard, and knowing this standard is necessary to decode the value field. Examples of format standards include `ieee_754` and `ascii`. This field ensures that the value can be accurately interpreted by specifying the encoding method used.
- * ***Length (length)*:** This field indicates the size of the value field measured in octets (bytes). It specifies how many bytes are used to store the value of the metric. Examples include 4, 8, 16, 32, and 64. The length field is important for memory allocation and data handling, ensuring that the value is stored and retrieved correctly.
- * ***Unit (unit)*:** This field defines the measurement units for the metric, such as frequency, data size, or data transfer rate. It is usually associated with the metric to provide context for the value.
- * ***Source (source, optional)*:** This field describes the origin of the information used to obtain the metric. It may include one or more of the following non-mutually exclusive values:
 - **'nominal'.** Similar to [RFC9439], "a 'nominal' metric indicates that the metric value is statically configured by the underlying devices. For example, bandwidth can indicate the maximum transmission rate of the involved device.
 - **'estimation'.** The 'estimation' source indicates that the metric value is computed through an estimation process.
 - **'directly measured'.** This source indicates that the metric can be obtained directly from the underlying device and it does not need to be estimated.
 - **'normalization'.** The 'normalization' source indicates that the metric value was normalized. For instance, a metric could be normalized to take a value from 0 to 1, from 0 to 10, or to take a percentage value. This type of metrics do not have units.
 - **'aggregation'.** This source indicates that the metric value was obtained by using an aggregation function.

Nominal metrics have inherent physical meanings and specific units without any additional processing. Aggregated metrics may or may not have physical meanings, but they retain their significance relative to the directly measured metrics. Normalized metrics, on the other hand, might have physical meanings but lack units.

- * ***Statistics (statistics, optional)*:** This field provides additional details about the metrics, particularly if there is any pre-computation performed on the metrics before they are collected. It is useful for services that require specific statistics for service instance selection.
 - 'max'. The maximum value of the data collected over intervals.
 - 'min'. The minimum value of the data collected over intervals.
 - 'mean'. The average value of the data collected over intervals.
 - 'cur'. The current value of the data collected.
- * ***Level (level)*:** This field specifies the level at which the metric is measured. It is used to categorize the metric based on its granularity and scope. Examples include L0, L1, and L2. The level field helps in understanding the level of detail and specificity of the metric being measured.
- * ***Value (value)*:** This field represents the actual numerical value of the metric being measured. It provides the specific data point for the metric in question.

4.2. Aggregation and Normalization Functions

In the context of CATS metric processing, aggregation and normalization are two fundamental operations that transform raw and derived metrics into forms suitable for decision-making and comparison across heterogeneous systems.

4.2.1. Aggregation

Aggregation functions combine multiple metric values into a single representative value. This is particularly useful when metrics are collected from multiple sources or over time intervals. For example, CPU usage metrics from multiple service instances may be aggregated to produce a single load indicator for a service. Common aggregation functions include:

- * **Mean average:** Computes the arithmetic average of a set of values.

- * Minimum/maximum: Selects the lowest or highest value from a set.
- * Weighted average: Applies weights to values based on relevance or priority.

The output of an aggregation function is typically a Level 2 metric, derived from multiple Level 0 metrics, or a level 2 metric, derived from multiple Level 0 or 1 metrics.

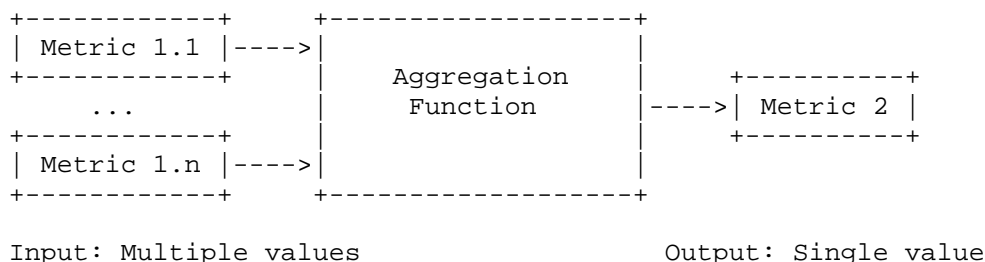


Figure 3: Aggregation function

4.2.2. Normalization

Normalization functions convert metric values with or without units into unitless scores, enabling comparison across different types of metrics and systems. This is essential when combining metrics from a heterogeneous set of resources (e.g, latency measured in milliseconds with CPU usage measured in percentage) into a unified decision model.

Normalization functions often map values into a bounded range, such as integers from 0, to 5, or real numbers from 0 to 1, using techniques like:

- * Sigmoid function: Smoothly maps input values to a bounded range.
- * Min-max scaling: Rescales values based on known minimum and maximum bounds.
- * Z-score normalization: Standardizes values based on statistical distribution.

Normalized metrics facilitate composite scoring and ranking, and can be used to produce Level 1 and Level 2 metrics.

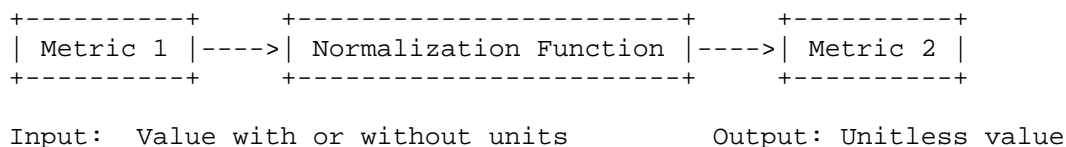


Figure 4: Normalization function

4.3. On the Meaning of Scores in Heterogeneous Metrics Systems

In a system like CATS, where metrics originate from heterogeneous resources---such as compute, communication, and storage---the interpretation of scores requires careful consideration. While normalization functions can convert raw metrics into unitless scores to enable comparison, these scores may not be directly comparable across different implementations. For example, a score of 4 on a scale from 1 to 5 may represent a high-quality resource in one implementation, but only an average one in another.

This ambiguity arises because different implementations may apply distinct normalization strategies, scaling methods, or semantic interpretations. As a result, relying solely on unitless scores for decision-making can lead to inconsistent or suboptimal outcomes, especially when metrics are aggregated from multiple sources.

To mitigate this, implementors of CATS metrics SHOULD provide clear and precise definitions of their metrics--particularly for unitless scores---and explain how these scores should be interpreted. This documentation should be designed to support operators in making informed decisions, even when comparing metrics from different implementations.

Similarly, operators SHOULD exercise caution when making potentially impactful decisions based on unitless metrics whose definitions are unclear or underspecified. In such cases, especially when decisions are critical or sensitive, operators MAY choose to rely on Level 0 (L0) metrics with units, which typically offer a more direct and unambiguous understanding of resource conditions.

4.4. Level Metric Representations

4.4.1. Level 0 Metrics

Several definitions have been developed within the compute and communication industries, as well as through various standardization efforts---such as those by the [DMTF]---that can serve as L0 metrics. L0 metrics contain all raw metrics which are not considered to be standardized in this document, considering about their diversity and many other existing work.

See Appendix A for examples of L0 metrics.

4.4.2. Level 1 Metrics

L1 metrics are normalized from L0 metrics. Although they don't have units, they can still be classified into types such as compute, communication, service and composed metrics. This classification is useful because it makes L1 metrics semantically meaningful.

The sources of L1 metrics is normalization. Based on L0 metrics, service providers design their own algorithms to normalize metrics. For example, assigning different cost values to each raw metric and do weighted summation. L1 metrics do not need further statistical values.

4.4.2.1. Normalized Compute Metrics

The metric type of normalized compute metrics is "compute_norm", and its format is unsigned integer. It has no unit. It will occupy an octet. Example:

Basic fields:

```
Metric type: compute_norm
Level: L1
Format: unsigned integer
Length: one octet
Value: 5
```

Source:

```
normalization
```

Metric Type	Level	Format	Length	Value	Source
8bits	2bits	1bit	3bits	8bits	3bits

Figure 5: Example of a normalized L1 compute metric

4.4.2.2. Normalized Communication Metrics

The metric type of normalized communication metrics is "communication_norm", and its format is unsigned integer. It has no unit. It will occupy an octet. Example:

Basic fields:

Metric type: communication_norm

Level: L1

Format: unsigned integer

Length: one octet

Value: 1

Source:

normalization

Metric Type	Level	Format	Length	Value	Source
8bits	2bits	1bit	3bits	8bits	3bits

Figure 6: Example of a normalized L1 communication metric

4.4.2.3. Normalized Composed Metrics

The metric type of normalized composed metrics is "delay_norm", and its format is unsigned integer. It has no unit. It will occupy an octet. Example:

Basic fields:

Metric type: composed_norm

Level: L1

Format: unsigned integer

Length: an octet

Value: 8

Source:

normalization

Metric Type	Level	Format	Length	Value	Source
8bits	2bits	1bit	3bits	8bits	3bits

Figure 7: Example of a normalized L1 composed metric

4.4.3. Level 2 Metrics

A Level 2 metric is a single-value, normalized metric that does not carry any inherent physical unit or meaning. While each provider may employ its own internal methods to compute this value, all providers must adhere to the representation guidelines defined in this section to ensure consistency and interoperability of the normalized output.

Metric type is "norm_fi". The format of the value is unsigned integer. It has no unit. It will occupy an octet. Example:

Basic fields:

```
Metric type: norm_fi
Level: L2
Format: unsigned integer
Length: an octet
Value: 1
```

Source:

```
normalization
```

Metric Type	Level	Format	Length	Value	Source
8bits	2bits	1bit	3bits	8bits	3bits

Figure 8: Example of a normalized L2 metric

The single normalized value also facilitates aggregation across multiple service instances. When each instance provides its own normalized value, no additional statistical processing is required at the instance level. Instead, aggregation can be performed externally using standardized methods, enabling scalable and consistent interpretation of metrics across distributed environments.

5. Comparison among Metric Levels

Metrics are progressively consolidated from L0 to L1 to L2, with each level offering a different degree of abstraction to address the diverse requirements of various services. Table 1 provides a comparative overview of these metric levels.

Level	Encoding Complexity	Extensibility	Stability	Accuracy
0	High	Low	Low	High
1	Medium	Medium	Medium	Medium
2	Low	High	High	Medium

Table 1: Comparison among Metrics Levels

Since Level 0 metrics are raw and service-specific, different services may define their own sets---potentially resulting in hundreds or even thousands of unique metrics. This diversity introduces significant complexity in protocol encoding and standardization. Consequently, L0 metrics are generally confined to bespoke implementations tailored to specific service needs, rather than being standardized for broad protocol use. In contrast, Level 1 metrics organize raw data into standardized categories, each normalized into a single value. This structure makes them more suitable for protocol encoding and standardization. Level 2 metrics take simplification a step further by consolidating all relevant information into a single normalized value, making them the easiest to encode, transmit, and standardize.

Therefore, from the perspective of encoding complexity, Level 1 and Level 2 metrics are recommended.

When considering extensibility, Level 0 metrics allow new services to define their own custom metrics. However, this flexibility requires corresponding protocol extensions, and the proliferation of metric types can introduce significant overhead, ultimately reducing the protocol's extensibility. In contrast, Level 1 metrics introduce only a limited set of standardized categories, making protocol extensions more manageable. Level 2 metrics go even further by consolidating all information into a single normalized value, placing the least burden on the protocol.

Therefore, from an extensibility standpoint, Level 1 and Level 2 metrics are recommended.

Regarding stability, Level 0 raw metrics may require frequent protocol extensions as new metrics are introduced, leading to an unstable and evolving protocol format. For this reason, standardizing L0 metrics within the protocol is not recommended. In contrast, Level 1 metrics involve only a limited set of predefined categories, and Level 2 metrics rely on a single consolidated value, both of which contribute to a more stable and maintainable protocol design.

Therefore, from a stability standpoint, Level 1 and Level 2 metrics are preferred.

In conclusion, for CATS, Level 2 metrics are recommended due to their simplicity and minimal protocol overhead. If more advanced scheduling capabilities are required, Level 1 metrics offer a balanced approach with manageable complexity. While Level 0 metrics are the most detailed and dynamic, their high overhead makes them unsuitable for direct transmission to network devices and thus not recommended for standard protocol integration.

6. CATS L2 Metric Registry Entry

This section gives an initial Registry Entry for the CATS L2 metric.

6.1. Summary

This category includes multiple indexes to the Registry Entry: the element ID, Metric Name, URI, Metric Description, Metric Controller, and Metric Version.

6.1.1. ID (Identifier)

IANA has allocated the Identifier 1 for the Named Metric Entry in Section 5. See Section 5.1.2 for mapping to Names.

6.1.2. Name

Norm_Passive_CATS-L2_RFCXXXXsecY_Unitless_Singleton

Naming Rule Explanation

- * Norm: Metric type (Normalized Score)
- * Passive: Measurement method
- * CATS-L2: Metric level (CATS Metric Framework Level 2)
- * RFCXXXXsecY: Specification reference (To-be-assigned RFC number and section number)
- * Unitless: Metric has not units
- * Singleton: Metric is a single value

6.1.3. URI

To-be-assigned.

6.1.4. Description

This metric represents a single normalized score used within CATS. It is derived by aggregating one or more CATS L0 and/or L1 metrics, followed by a normalization process that produces a unitless value. The resulting score provides a concise assessment of the overall capability of a service instance, enabling rapid comparison across instances and supporting efficient traffic steering decisions.

6.1.5. Change Controller

IETF

6.1.6. Version

1.0

6.2. Metric Definition

6.2.1. Reference Definition

[I-D.ietf-cats-metric-definition] Core referenced sections:
Section 3.4 (L2 Level Metric Definition), Section 4.2 (Aggregation and Normalization Functions)

6.2.2. Fixed Parameters

- * Normalization score range: 0-10 (0 indicates the poorest capability, 10 indicates the optimal capability)
- * Data precision: decimal number (unsigned integer)

6.3. Method of Measurement

This category includes columns for references to relevant sections of the RFC(s) and any supplemental information needed to ensure an unambiguous method for implementations.

6.3.1. Reference Methods

Raw Metrics collection: Collect L0 service and compute raw metrics using platform-specific management protocols or tools (e.g., Prometheus [Prometheus] in Kubernetes). Collect L0 network performance raw metrics using existing standardized protocols (e.g., NETCONF [RFC6241], IPFIX [RFC7011]).

Aggregation logic: Refer to [I-D.ietf-cats-metric-definition] Section 4.2.1 (e.g., Weighted Average Aggregation).

Normalization logic: Refer to [I-D.ietf-cats-metric-definition] Section 4.2.2 (e.g., Sigmoid Normalization).

The reference method aggregates and normalizes L0 metrics to generate L1 metrics in different categories, and further calculates a L2 singleton score for full normalization.

6.3.2. Packet Stream Generation

N/A

6.3.3. Traffic Filtering (Observation) Details

N/A

6.3.4. Sampling Distribution

Sampling method: Continuous sampling (e.g., collect L0 metrics every 10 seconds)

6.3.5. Runtime Parameters and Data Format

CATS Service Contact Instance ID (CSCI-ID): an identifier of CATS service contact instance. According to [I-D.ietf-cats-framework], a unicast IP address can be an example of identifier. (format: ipv4-address-no-zone or ipv6-address-no-zone, complying with [RFC6991])

Service_Instance_IP: Service instance IP address (format: ipv4-address-no-zone or ipv6-address-no-zone, complying with [RFC6991])

Measurement_Window: Metric measurement time window (Units: seconds, milliseconds; Format: uint64; Default: 10 seconds)

6.3.6. Roles

C-SMA: Collects L0 service and compute raw metrics, and optionally calculates L1 and L2 metrics according to service-specific strategies.

C-NMA: Collects L0 network performance raw metrics, and optionally calculates L1 and L2 metrics according to service-specific strategies.

6.4. Output

This category specifies all details of the output of measurements using the metric.

6.4.1. Type

Singleton value

6.4.2. Reference Definition

Output format: Refer to [I-D.ietf-cats-metric-definition]
Section 4.4.3

Score semantics: 0-3 (Low capability, not recommended for steering),
4-7 (Medium capability, optional for steering), 8-10 (High
capability, priority for steering)

6.4.3. Metric Units

Unitless

6.4.4. Calibration

Calibration method: Conduct benchmark calibration based on standard
test sets (fixed workload) to ensure the output score deviation of
C-SMA and C-NMA is lower than 0.1 (one abnormal score in every ten
test rounds).

6.5. Administrative Items

6.5.1. Status

Current

6.5.2. Requester

To-be-assgined

6.5.3. Revision

1.0

6.5.4. Revision Date

2026-01-20

6.5.5. Comments and Remarks

None

7. Implementation Guidance on Using CATS Metrics

<Authors' Note: This section has been moved to [I-D.ietf-cats-framework] at the suggestion of the chairs, since this document focuses primarily on metric definitions rather than implementation details.>

8. Security Considerations

TBD

9. IANA Considerations

TBD

10. References

10.1. 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-19, 20 November 2025, <<https://datatracker.ietf.org/doc/html/draft-ietf-cats-framework-19>>.

[I-D.ietf-cats-metric-definition]

Yao, K., Li, C., Contreras, L. M., Ros-Giralt, J., and H. Shi, "CATS Metrics Definition", Work in Progress, Internet-Draft, draft-ietf-cats-metric-definition-04, 20 October 2025, <<https://datatracker.ietf.org/doc/html/draft-ietf-cats-metric-definition-04>>.

[RFC6241] Enns, R., Ed., Bjorklund, M., Ed., Schoenwaelder, J., Ed., and A. Bierman, Ed., "Network Configuration Protocol (NETCONF)", RFC 6241, DOI 10.17487/RFC6241, June 2011, <<https://www.rfc-editor.org/rfc/rfc6241>>.

[RFC6991] Schoenwaelder, J., Ed., "Common YANG Data Types", RFC 6991, DOI 10.17487/RFC6991, July 2013, <<https://www.rfc-editor.org/rfc/rfc6991>>.

[RFC7011] Claise, B., Ed., Trammell, B., Ed., and P. Aitken, "Specification of the IP Flow Information Export (IPFIX) Protocol for the Exchange of Flow Information", STD 77, RFC 7011, DOI 10.17487/RFC7011, September 2013, <<https://www.rfc-editor.org/rfc/rfc7011>>.

- [RFC8911] Bagnulo, M., Claise, B., Eardley, P., Morton, A., and A. Akhter, "Registry for Performance Metrics", RFC 8911, DOI 10.17487/RFC8911, November 2021, <<https://www.rfc-editor.org/rfc/rfc8911>>.
- [RFC8912] Morton, A., Bagnulo, M., Eardley, P., and K. D'Souza, "Initial Performance Metrics Registry Entries", RFC 8912, DOI 10.17487/RFC8912, November 2021, <<https://www.rfc-editor.org/rfc/rfc8912>>.
- [RFC9439] Wu, Q., Yang, Y., Lee, Y., Dhody, D., Randriamasy, S., and L. Contreras, "Application-Layer Traffic Optimization (ALTO) Performance Cost Metrics", RFC 9439, DOI 10.17487/RFC9439, August 2023, <<https://www.rfc-editor.org/rfc/rfc9439>>.

10.2. Informative References

- [DMTF] "DMTF", n.d., <<https://www.dmtf.org/>>.
- [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-13, 28 January 2026, <<https://datatracker.ietf.org/doc/html/draft-ietf-cats-usecases-requirements-13>>.
- [I-D.rcr-opsawg-operational-compute-metrics]
Randriamasy, S., Contreras, L. M., Ros-Giralt, J., and R. Schott, "Joint Exposure of Network and Compute Information for Infrastructure-Aware Service Deployment", Work in Progress, Internet-Draft, draft-rcr-opsawg-operational-compute-metrics-08, 21 October 2024, <<https://datatracker.ietf.org/doc/html/draft-rcr-opsawg-operational-compute-metrics-08>>.
- [performance-metrics]
"performance-metrics", n.d., <<https://www.iana.org/assignments/performance-metrics/performance-metrics.xhtml>>.
- [Prometheus]
"Prometheus", n.d., <<https://prometheus.io/>>.

Appendix A. Appendix A

A.1. Level 0 Metric Representation Examples

Several definitions have been developed within the compute and communication industries, as well as through various standardization efforts---such as those by the [DMTF]---that can serve as L0 metrics. This section provides illustrative examples.

A.1.1. Compute Raw Metrics

This section uses CPU frequency as an example to illustrate the representation of raw compute metrics. The metric type is labeled as `compute_CPU_frequency`, with the unit specified in GHz. The format should support both unsigned integers and floating-point values. The corresponding metric fields are defined as follows:

Basic fields:

```
Metric Type: compute_CPU_frequency
Level: L0
Format: unsigned integer, floating point
Unit: GHz
Length: four octets
Value: 2.2
```

Source:

```
nominal
```

Metric Type	Level	Format	Unit	Length	Value	Source
8bits	2bits	1bit	4bits	3bits	32bits	3bits

Figure 9: An Example for Compute Raw Metrics

A.1.2. Communication Raw Metrics

This section takes the total transmitted bytes (TxBytes) as an example to show the representation of communication raw metrics. TxBytes are named as "communication type_TxBytes". The unit is Mega Bytes (MB). Format is unsigned integer or floating point. It will occupy 4 octets. The source of the metric is "Directly measured" and the statistics is "mean". Example:


```

Basic fields:
  Metric type: "communication type_TXBytes"
  Level: L0
  Format: unsigned integer, floating point
  Unit: MB
  Length: four octets
  Value: 100
Source:
  Directly measured
Statistics:
  mean

```

Metric Type	Level	Format	Unit	Length	Value	Source	Statistics
8bits	2bits	1bit	4bits	3bits	32bits	3bits	2bits

Figure 10: An Example for Communication Raw Metrics

A.1.3. Delay Raw Metrics

Delay is a kind of synthesized metric which is influenced by computing, storage access, and network transmission. Usually delay refers to the overall processing duration between the arrival time of a specific service request and the departure time of the corresponding service response. It is named as "delay_raw". The format should support both unsigned integer or floating point. Its unit is microseconds, and it occupies 4 octets. For example:

```

Basic fields:
  Metric type: "delay_raw"
  Level: L0
  Format: unsigned integer, floating point
  Unit: Microsecond(us)
  Length: four octets
  Value: 231.5
Source:
  aggregation
Statistics:
  max

```

Metric Type	Level	Format	Unit	Length	Value	Source	Statistics
8bits	2bits	1bit	4bits	3bits	32bits	3bits	2bits

Figure 11: An Example for Delay Raw Metrics

Contributors

Mohamed Boucadair
Orange

Email: mohamed.boucadair@orange.com

Zongpeng Du
China Mobile
Email: duzongpeng@chinamobile.com

Hang Shi
Huawei
Email: shihang9@huawei.com

Authors' Addresses

Kehan Yao
China Mobile
China
Email: yaokehan@chinamobile.com

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

L. M. Contreras
Telefonica
Email: luismiguel.contrerasmurillo@telefonica.com

Jordi Ros-Giralt
Qualcomm Europe, Inc.
Email: jros@qti.qualcomm.com

Guanming Zeng
Huawei Technologies
China
Email: zengguanming@huawei.com