

cats  
Internet-Draft  
Intended status: Informational  
Expires: 3 September 2026

J. Wang  
China Mobile  
J. Guo  
Inspur Computer Technology Co., Ltd.  
2 March 2026

Consideration for Computing-Power Collaboration in Computing-Aware  
Traffic Steering (CATS)  
draft-wang-cats-computing-power-collaboration-00

Abstract

This document outlines a series of challenges and considerations to explore computing-power collaboration in CATS.

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 3 September 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 . . . . .	2
2. Definition of Terms . . . . .	3
3. Challenges . . . . .	3
3.1. Architectural Heterogeneity . . . . .	3
3.2. Differences in Network Protocols . . . . .	4
3.3. Challenging Cross-domain Collaboration . . . . .	5
4. Consideration . . . . .	5
4.1. Flexible Conversion of Agreements . . . . .	6
4.2. Long-distance Low-latency Routing . . . . .	6
4.3. Exposure of Energy Status in Computing Resource Supply . . . . .	7
5. Conclusion . . . . .	8
6. Security Considerations . . . . .	8
7. IANA Considerations . . . . .	8
8. Informative References . . . . .	8
Authors' Addresses . . . . .	8

## 1. Introduction

With the continuous development and progress of the Internet, a large amount of computing resources is required to complete data processing. In order to disperse the pressure of cloud data centers, computing power gradually moves from the center to the edge, forming scattered computing resources in mobile networks. In order to make full use of scattered computing resources and provide better services, Computing-Aware Traffic Steering (CATS) is proposed to support steering the traffic among different edge sites according to both the real-time network and computing resource status as mentioned in [I-D.ietf-cats-usecases-requirements]. It requires the network to be aware of computing resource information and select a service instance based on the joint metric of computing and networking.

As artificial intelligence technology advances into the era of large models, the demand for computational power for AI training has grown exponentially. This has resulted in a continuous increase in electricity consumption. From pre-training to fine-tuning and ongoing iterative optimization, massive computing clusters operate under high loads for extended periods of time. As a result, electricity costs, power supply stability, and energy provision capacity directly impact the training efficiency, deployment scale, and iteration speed of AI models. Today, power supply is no longer just a cost issue, but it has become a critical bottleneck that hinders the further scaling and industrialization of AI technology.

Whether it's the construction of hyperscale data centers, the deployment of intelligent computing centers, or the development and application of large-scale models, all rely heavily on stable, green,

and low-cost power supply. Tight power resources, regional disparities in supply capacity, and energy consumption control policies are profoundly shaping the development pace and spatial layout of the AI industry. Today, power has become a core constraint on AI advancement—a conclusion no longer confined to individual companies' observations but a widely recognized industry consensus spanning the entire industrial chain from research and development to application.

Green requirements for AI training scenarios have now been formally incorporated into [I-D.ietf-cats-usecases-requirements]. This document outlines a series of challenges and considerations to explore computing-power collaboration in CATS.

## 2. Definition of Terms

Computing-Aware Traffic Steering (CATS): Aiming at computing and network resource optimization by steering traffic to appropriate computing resources considering not only routing metric but also computing resource metric.

Service: A monolithic functionality that is provided by an endpoint according to the specification for said service. A composite service can be built by orchestrating monolithic services.

Service instance: Running environment (e.g., a node) that makes the functionality of a service available. One service can have several instances running at different network locations.

## 3. Challenges

Computing-Power Collaboration faces a series of challenges.

### 3.1. Architectural Heterogeneity

The network architectures of Power Network and Computing-Aware Network originate from distinct design philosophies and application scenarios, exhibiting significant heterogeneity. This makes it difficult to coordinate computing and power management..

Power system networks use a "layered, partitioned, closed-loop control" architecture and primarily support grid dispatch, equipment monitoring, and fault isolation. Nodes are dispersed and mainly located at industrial sites, with network topologies dominated by trees and rings, prioritizing stability and controllability.

In contrast, Computing-Aware Network are centered around data centers and intelligent computing clusters, utilizing a "flat, highly aggregated" architecture. The focus is on supporting large-scale data transmission, computing power scheduling, and distributed computing. The nodes are highly concentrated and the topologies are primarily based on spine-leaf architectures, prioritizing bandwidth and transmission efficiency.

However, the significant differences in design objectives, topology structures, and node characteristics between these two architectures present challenges for efficient interconnection during computing power coordination. The exchange of data and scheduling of resources must overcome architectural barriers, resulting in increased network deployment and modification costs, as well as issues such as data transmission delays and resource scheduling disconnects. These factors greatly hinder the overall efficiency of computing-power coordination.

### 3.2. Differences in Network Protocols

The Power Network and Computing-Aware Network have historically developed independently, each with its own closed network protocol system. However, this has resulted in protocol inconsistencies that pose challenges for data exchange and command transmission during power-computing coordination. This has become a major technical bottleneck.

The Power Network primarily uses specialized industrial protocols such as IEC 61850 and DL/T 860, which are designed for real-time control and equipment monitoring in the grid. These protocols prioritize low-latency and high-reliability transmission of small data packets, meeting the communication needs of power equipment. However, they struggle to seamlessly integrate with general-purpose network protocols, leading to compatibility issues.

On the other hand, the Computing-Aware Network mainly relies on the TCP/IP protocol suite, which includes general-purpose protocols like HTTP, FTP, and RDMA. These protocols prioritize the transmission of large data volumes and high bandwidth, catering to scenarios such as computing resource scheduling and data exchange. However, they lack the ability to adapt to the real-time control commands of power systems.

As a result, direct interaction between power system operational data/control commands and computing system workload/scheduling demands is hindered, requiring the deployment of additional protocol conversion devices. This increases system complexity and operational costs, while also introducing extra transmission delays and the risk of data packet loss. Ultimately, this compromises the real-time performance and reliability of computing-power collaboration.

### 3.3. Challenging Cross-domain Collaboration

The core requirement for cross-domain computing-power coordination is to achieve dynamic matching and real-time scheduling of computing resources and power resources. This process imposes extremely high demands on network latency, presenting a critical challenge that constrains coordination effectiveness. In cross-regional computing-power coordination scenarios, such as ultra-long-distance interconnection, cross-regional virtual power plant coordination, and intelligent computing cluster scheduling, real-time collection of power system load data, renewable energy output data, and computing system workload/energy consumption data is essential.

This data must be transmitted via networks to the coordination dispatch center, where it undergoes analysis and decision-making before dispatch instructions are relayed back to terminal nodes, forming a closed-loop “collection-analysis-decision-execution” process. This closed-loop process imposes extremely stringent latency requirements. End-to-end latency for power control commands must be kept below 10ms, with certain critical scenarios demanding 5ms, while jitter must be 1ms to ensure synchronized execution of dispatch commands.

However, current wide-area networks suffer from complex transmission links, multiple cross-domain routing hops, and network congestion, making it difficult to consistently meet these latency requirements. Exceeding latency thresholds may cause computational power scheduling delays, untimely power load transfers, and even risks such as grid frequency fluctuations or computational cluster overloads, severely compromising the safety and effectiveness of computing-power coordination.

## 4. Consideration

To better achieve computing-power coordination, it is necessary to enable flexible protocol conversion; long-distance low-latency routing; and exposure of energy status for computing resource supply.

#### 4.1. Flexible Conversion of Agreements

Flexible protocol conversion serves as the core enabler for resolving protocol incompatibilities between power and computing systems, achieving power-computing synergy, and overcoming the challenges of data isolation and command transmission difficulties.

To address compatibility challenges between specialized industrial protocols for power systems (e.g., IEC 61850, DL/T 860) and general-purpose protocols for computing systems (e.g., TCP/IP protocol suite, RDMA), a flexible and efficient protocol conversion framework must be established—not merely deploying single conversion devices. A modular, scalable conversion architecture should be adopted to support real-time parsing, adaptation, and conversion of multiple protocols. This enables both the transformation of power system control commands and operational data into computing system protocols for efficient transmission to the collaborative dispatch center, and the conversion of computing system dispatch commands into specialized protocols recognizable by power equipment to ensure precise execution.

Simultaneously, the conversion process must balance low latency with high reliability, avoiding additional transmission delays and data packet loss. By optimizing conversion algorithms and streamlining conversion workflows, millisecond-level response times for protocol conversion can be achieved. This ensures seamless data interaction and command transmission during computing-power coordination, laying the foundation for cross-system collaborative dispatch.

#### 4.2. Long-distance Low-latency Routing

Long-distance low-latency routing design is a critical measure for meeting the high-latency requirements of cross-domain computing-power coordination and enabling dynamic cross-regional matching of computing and power resources. To address current challenges such as complex wide-area network links, multiple cross-domain routing hops, and frequent congestion, a dedicated wide-area routing system for computing-power coordination must be established, balancing long-distance transmission with low-latency requirements.

On one hand, optimized routing planning algorithms should dynamically select optimal transmission paths based on computing-power coordination service priorities, minimizing routing hops and avoiding congested network segments to ensure the shortest paths and lowest latency for cross-domain data transmission and command delivery.

On the other hand, implementing deterministic networking technologies such as TSN+SRv6 through techniques like time slot scheduling and path reservation ensures stability and predictable latency for long-distance transmission. This results in a cross-domain end-to-end latency of less than 10ms, reduces core scenarios to under 5ms, and maintains jitter within 1ms.

Simultaneously, deploying multi-path redundant routing addresses link failures, enabling sub-second fault self-healing. This ensures uninterrupted, low-latency transmission over long distances, supporting the efficient implementation of cross-domain computing-power coordination scenarios like “East Data, West Computing” and “West Power, East Transmission.”

#### 4.3. Exposure of Energy Status in Computing Resource Supply

Exposing the energy status of computing resources is a crucial prerequisite for achieving dynamic coordination between computing and power systems, enhancing resource utilization efficiency, and breaking down barriers to matching computing and power resources. Currently, computing systems and power systems operate independently. The energy consumption, energy efficiency, and power supply requirements of computing clusters are not effectively exposed to power dispatch systems.

Similarly, information such as the power supply capacity, green power generation output, and load fluctuations of the power system is not synchronized to computing dispatch systems, leading to a disconnect in resource matching between the two. Therefore, it is necessary to establish a unified computing-power status perception system to promote the comprehensive exposure and sharing of energy status information for computing resource supply. Specifically, deploy energy monitoring devices within computing clusters to collect real-time data on energy consumption, power demands, and energy efficiency from computing nodes. This information should be transmitted via standardized interfaces to the computing-power coordination dispatch center.

Simultaneously, synchronize real-time power system data—including supply load, green power generation output, and electricity price fluctuations—to the computing dispatch system. This achieves bidirectional transparency between computing energy status and power supply conditions. Through this state exposure, the coordination center can precisely grasp the resource status of both parties, enabling deep synergy between computing power scheduling and power dispatch. This optimizes resource allocation, enhances green electricity consumption rates, and improves computing power operational efficiency.

## 5. Conclusion

This document highlights the challenges and considerations for Computing-Power Collaboration in CATS.

## 6. Security Considerations

TBD.

## 7. IANA Considerations

TBD.

## 8. 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-14,  
2 February 2026, <[https://datatracker.ietf.org/doc/html/  
draft-ietf-cats-usecases-requirements-14](https://datatracker.ietf.org/doc/html/draft-ietf-cats-usecases-requirements-14)>.

## Authors' Addresses

Jing Wang  
China Mobile  
No.32 XuanWuMen West Street  
Beijing  
100053  
China  
Email: wangjingjc@chinamobile.com

Jianchao Guo  
Inspur Computer Technology Co., Ltd.  
Beijing  
China  
Email: guojianchao01@inspur.com