# Draft new Recommendation ITU-T Y.NCE-DAICC

## Network capability enhancement for distributed artificial intelligent computing centers in NGNe

## 1. Scope

This draft Recommendation provides network enhancement requirements and capabilities for distributed artificial intelligent computing centers in NGNe.

- Background and motivations
- Typical scenarios
- Framework
- Network enhancement requirements
- Network enhancement capabilities
- Security considerations

## 2. References

*The following ITU-T Recommendations and other references contain provisions, which, through reference in this text, constitute provisions of this Recommendation. At the time of publication, the editions indicated were valid. All Recommendations and other references are subject to revision; users of this Recommendation are therefore encouraged to investigate the possibility of applying the most recent edition of the Recommendations and other references listed below. A list of the currently valid ITU-T Recommendations is regularly published.*

*The reference to a document within this Recommendation does not give it, as a stand-alone document, the status of a Recommendation.*

*[Editor's Note] The ongoing item ITU-T Y.NGNe-MC-reqts、ITU-T Y.NGNe-NCI-reqts focus on NGN evolution for support network and cloud interworking, may be concerned in the future.*

## 3. Definitions

## 3.1 Terms defined elsewhere

## 3.2 Terms defined in this Recommendation

 [TBD]

## 4. Abbreviations and acronyms

 [TBD]

## 5. Conventions

In this Recommendation:

The keywords "is required to" indicate a requirement which must be strictly followed and from which no deviation is permitted, if conformance to this Recommendation is to be claimed.

[TBD]

## 6.    Background and motivations

AICC (artificial intelligent computing center) is an artificial intelligent computing center built based on artificial intelligent chips. It covers the complete system of hardware infrastructure and software infrastructure. This kind of computing center is mainly used in the scenarios of intelligent computing, such as artificial intelligent algorithm model development, model training and model reasoning. As the infrastructure of AI computing power, AICC has received widespread attention. In order to improve the cognitive ability of AI models, a large number of pre-training models and high-quality large-scale data sets used to support large-scale pre-training models have emerged. The computational power level has also increased from PFLOPS to EFLOPS, even starting to enter 10 EFLOPS era. More and more intelligent computing tasks are no longer limited to a single AICC, but require collaboration from multiple AICCs. Therefore, AICC needs to be interconnected to solve the problem of insufficient capacity of a single AICC. Therefore, AICCs needs to be interconnected to solve the problem of insufficient capacity of a single AICC.
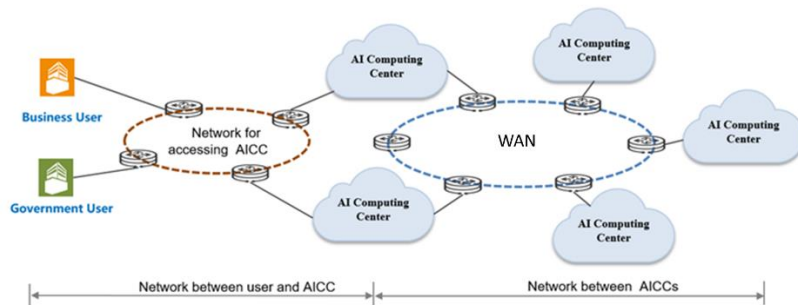


Figure 6-1: Interconnection network of multiple AICCs

The interconnection network of multiple AICCs can be divided into two parts：

(1) Network from users to the AICC: used to carry the data interaction between users and the AICC, including the transfer of AI computing data and the acquisition of AI training results, and provide the on-demand allocated network resources, so that users can safely and efficiently access to the AICC;

(2) Interconnected network between multiple AICCs: used to complete the computing tasks through the collaboration with multiple AICCs, provide QoS services such as network slicing and low latency, and ensure that network resources meet the computing requirements.

The IP network oriented to AICCs  has put forward higher requirements for transmission quality, such as the need to provide greater bandwidth, lower latency, high determinacy, and ensure efficient scheduling of data applications and algorithms. The following problems need to be solved:

(1) Slow feedback and packet loss in long distance transmission. Long distance transmission brings ultra-long link transmission delay and lagging feedback on network state, which requires the IP carrier network between AICCs to adopt different strategies according to the transmission distance step by step to solve the packet loss problem.

(2) Network utilization. When only several ultra-long distance transmission connections are working, the network utilization is very low. Network between AICCs is required to provide reasonable network utilization.

(3) High-speed transmission for massive data. AICC that obtains user training data needs to distribute the data to different AICCs. So AICCs needs to have the ability of high-speed transmission for massive data.

*[Editor's Note] This work item aims to solve the problem raised by multi AICCs collaboration like low feedback and packet loss in long distance transmission. However, the background needs to be enhanced to explain why AICCs need to be interconnected.*

## 7. Typical Scenarios

Scenario 1：Mass training data transmission to different AICCs

On one hand, because some AICCs are located in remote areas, the storage cost is low, or the user data needs to be uploaded nearby, etc. The user data used for large-scale model training needs to be transferred to the nearest AICC, as shown in figure 7-1. The network from users to AICC needs to provide high-quality transmission services according to different requirements.

On the other hand, AICC that obtains user training data needs to distribute the data to different AICCs based on factors such as cost, computing resource, or parallel processing, as shown in figure 7-2. And the network between AICCs needs to have the ability of high-speed transmission for massive data.

In the scenario of meteorological data analysis, data need to be uploaded from the monitoring systems to the AICCs for reasoning, analysis and so on. In the above described process, it involves mass training data transmission to one or multiple different AICCs.
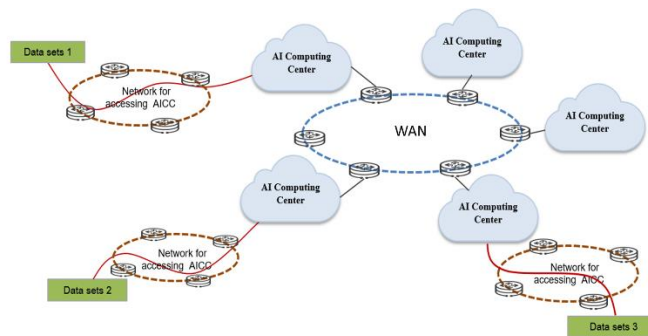


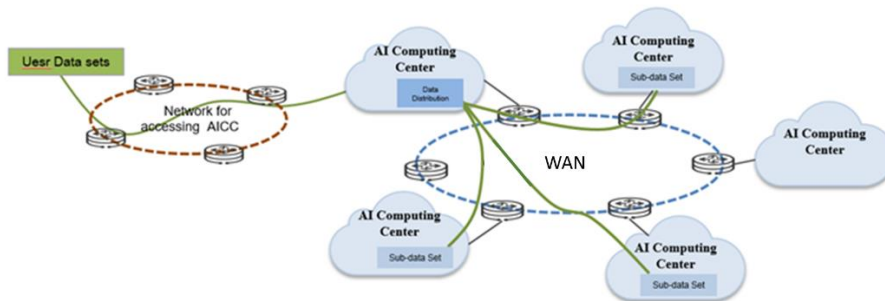Figure 7-1: Scenario of large-scale user data sets stored in different AICCs



Figure 7-2: Scenario of massive training data distribution

Scenario 2：Large-scale distributed collaborative AI model training under large computing power requirements

In the process of machine learning, artificial intelligence model finally obtains a model with relatively high accuracy through multiple parameters iteratively optimal. Large-scale model training refers to the way that multiple AICCs iterate and update part of AI model parameters based on their own data. This allows multiple AICCs to complete the parameter training of the master model without data migration.

As shown in the figure below, a user wants to conduct large-scale AI model training, and the required large number of data sets for training are stored in several different AICCs. In each AICC, the sub model parameters assigned by main model update iteratively on its own data sets, and then send the updated parameter results and possible related data sets to the AICC which host the main model, so as to complete the training task of the large model.
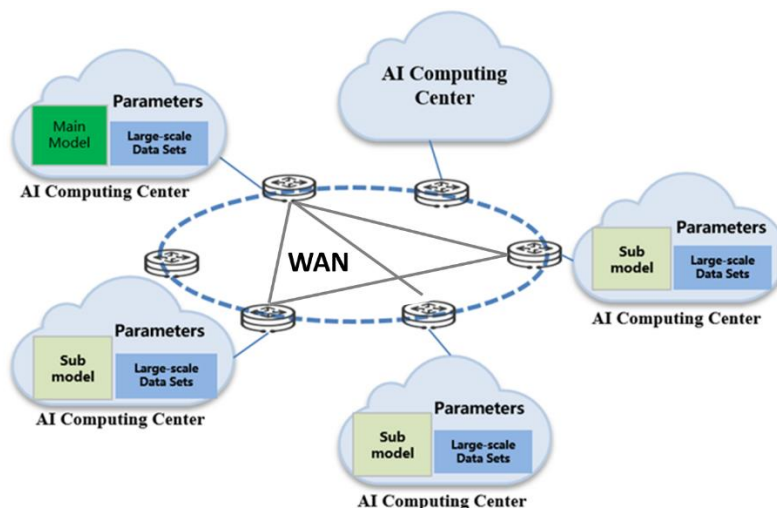
Figure 7-3: Scenario of large-scale distributed collaborative AI model training

The expansion of AICC by a single operating body will cost a lot, and the utilization rate is not so high. While multiple existing AICCs can be logically interconnected as a super virtual intelligence centers for relatively optimal utilization.

Third-party companies integrate AICCs belonging to multiple AI service operating bodies to provide agent services. In this scenario, on the one hand, the AICCs of multiple service bodies under management is secure and mutual trust by default; on the other hand, for large-scale AI computing needs, multiple AICCs can be invoked simultaneously to provide services.

When AICCs belong to different operating bodies, idle resources for AI services may be shared by security means such as blockchain or trading platform managed by one operating body.

## 8. Framework

*[Editor's Note] This clause will give the general framework.*

The overall framework can be roughly divided into three layers.

The first layer is responsible for task operation, in which the AI models, data models, algorithms, computing resources and service can be shared.

The second layer is responsible for the unified scheduling management, which provides functions such as network resources and AICC resources sensing, unified resources scheduling and management, multi-objective intelligent routing, unified authentication and authorization, and unified accounting and so on.

The third layer is infrastructure layer, which includes AICCs and Internet networks. The AICCs contains dedicated servers for AI computing and common servers. The Internet network includes the IP bearer network for user to access to AICC and the IP bearer network for Interconnection between AICCs.
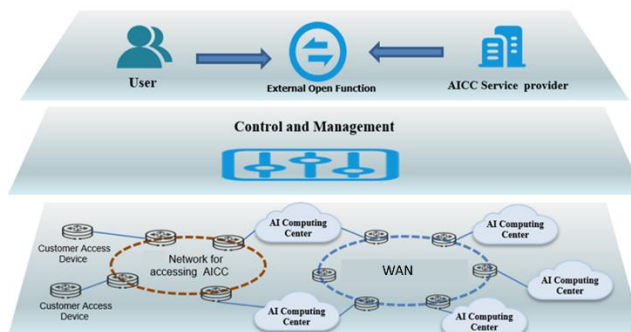
Figure 8-1: Framework of network capability enhancement for distributed AICCs in NGNe

*[Editor's Note] This figure needs to be enhanced based on NGNe architecture, and the contributions are welcomed.*

## 8.1 Working Process for user accessing to AICC

*[Editor's Note] This clause will provide working process for opening a VPN for users accessing the AICC.*

## 8.2 Working Process for collaboration of AICCs

*[Editor's Note] This clause will provide working process for service orchestration between cross-AICCs.*

## 9. Network enhancement requirements

*[Editor's Note] This clause will provide the network enhancement requirements.*

## 9.1 Interworking requirements

*[Editor's Note] This clause will provide interworking requirements.*

## 9.2 Efficient transmission requirements

*[Editor's Note] This clause will provide efficient transmission requirements.*

## 9.3 Network aggregation requirements

*[Editor's Note] This clause will provide network aggregation requirements.*

## 9.4 Other requirements

*[Editor's Note] This clause will provide other requirements.*

## 10. Network enhancement capabilities

## 10.1 Interworking capability

· Interworking capability between the customer device and the fixed-mobile integration gateway;

· Interworking capability between the fixed-mobile integration gateway and the AICCs device.

## 10.2 Efficient transmission capacity

· Provide flexible configuration and customize transmission services based on differentiated communication features and multidimensional application data features.

· Provide dedicated lines to meet QoS requirements according to user needs from user to AI Computing centre.

· Intelligent multi-target routing capability, and provides end-to-end multipath concurrent transmission capability to meet the demand for large-scale data transmission across AICC with high bandwidth

· Slicing capability;

- Lossless or low latency transmission capability, provides flexible segment transmission control capability based on the heterogeneous characteristics of underlying links in wide-area networks, to effectively alleviate the problems of lagging network status feedback caused by long-distance transmission and long transmission delay caused by packet loss and retransmission.

## 10.3 Network aggregation capabilities

- Provide in-network aggregation capability to meet the demand for high utilization of network resources between AICCs and to alleviate end-to-end communication bottlenecks.

- Provide end-network coordination capability to meet the demand for flexible coordination between end-system transmission services and in-network aggregation services.

## 10.4 Other capabilities

- Unified billing capability taking resource usage into account;

- Unified account authentication and authorization capability；

- Unified operation and maintenance analysis capability；

[TBD]

## 11. Security considerations

*[Editor's Note] This clause will provide security considerations.*

**Annex A**

**A.1 justification for proposed Draft new Recommendation ITU-T Y.NCE-DAICC: "Network capability enhancement for distributed artificial intelligent computing centers in NGNe"**

| Question: | Q2/13 | **Proposed new ITU-T Recommendation** | Geneva, 13-24 March 2023 | | |
|---|---|---|---|---|---|
| **Reference and title:** | ITU-T Y.NCE-DAICC: "Network capability enhancement for distributed artificial intelligent computing centers in NGNe" | | | | |
| **Base text:** | TD225-R1/WP3 | | **Timing:** | 2025-03 | |
| **Editor(s):** | Jianfei Li, China Unicom lijf299@chinaunicom.cn Ran Pang, China Unicom pangran@chinaunicom.cn Qianying Zhao, China Telecom zhaoqy50@chinatelecom.cn Long Luo, University of Electronic Science and Technology of China llong@uestc.edu.cn Chuxuan Zeng, China Unicom zengchuxuan@chinaunicom.cn | | **Approval process:** | AAP | |

**Scope** (defines the intent or object of the Recommendation and the aspects covered, thereby indicating the limits of its applicability):

This draft Recommendation provides network enhancement requirements and capabilities for distributed artificial intelligent computing centers in NGNe.

- Background and motivations
- Typical scenarios
- Framework
- Network enhancement requirements
- Network enhancement capabilities
- Security considerations

**Summary** (provides a brief overview of the purpose and contents of the Recommendation, thus permitting readers to judge its usefulness for their work):

Intelligent computing is a kind of computing mode based on artificial intelligence technology, mainly used to process complex, high-dimensional, dynamic, and unstructured data and issues, which plays an important role in the field of meteorological data analysis, intelligent security, autonomous driving, and so on. Artificial Intelligence Computing Center (AICC) has received widespread attention as an infrastructure for intelligent computing, more and more intelligent computing tasks are no longer limited to a single AICC, but require collaboration from multiple AICCs. Therefore, AICC needs to be interconnected to solve the problem of insufficient capacity of a single AICC.

The above scenarios put forward higher capabilities and requirements for transmission networks to serve the AICCs. The network is required to provide greater bandwidth, lower latency, high determinacy, more flexible adaptive collaboration capability with the demand of the intelligent computing service, and also to solve issues like slow feedback and packet loss in long distance transmission.

Therefore, this recommendation analyses typical scenarios, provides framework, network enhancement requirements and capabilities for distributed artificial intelligent computing centers in NGNe.

**Relations to ITU-T Recommendations or to other standards** (approved or under development):

(1) Relations to Y.NGNe-MC-reqts

The ongoing draft Recommendation Y.NGNe-MC-reqts in Q2/SG13 studies requirements and capabilities of NGN evolution to support multi-connection for network and cloud interworking.

This draft is specific to the network requirements of AICCs.

(2) Relations to Y.NGNe-NCI-reqts

The ongoing draft Recommendation Y.NGNe- NCI -reqts in Q2/SG13 studies the general requirement to support network and cloud interworking.

And this draft recommendation concentrates on the detailed requirements and capabilities on network between AICCs.

(3) Relations to IETF standards

The CATS (Computing-Aware Traffic Steering) working group is chartered to consider the problem of how the network edge can steer traffic between clients of a service and sites offering the service. There are a number of drafts describing scenarios in CATS.

This draft recommendation concentrates on the network capability enhancement for AICCs in NGNe.

**Liaisons with other study groups or with other standards bodies:**

IETF CATS

**Supporting members that are committing to contributing actively to the work item:**

China Unicom, China Telecom, University of Electronic Science and Technology of China

———————————