

AI Computing Network

Lily Lyu (Huawei)

Jieyu Li (China Mobile)

July Plenary 2024

AICN study item

AICN study item website: <https://1.ieee802.org/nendica-aicn/>

IEEE 802 Nendica Initiating Motion (2024-03-14)

- To initiate a Nendica Study Item on AI computing network

Progress:

Contributions discussion:

- Contributed text: Load balancing requirements and challenges
 - Experiment data to support the presented view
- Contributed text: Scale requirements and challenges
 - Data and some concepts clarification(AZ, convergent points)
- Contributed text: Availability requirements and challenges
 - Clarify scope of availability

Report draft discussion:

- Draft 0.1 (802.1-24-0022) integrates the contents of the previous study item proposal to form the framework of the document
 - “Discussion addressed points that could be addressed in progressing the draft toward a more complete report. (page number/reference list/some figures, concepts clarification)”
- Draft (802.1-24-0022) R0, R1, R2 incorporates contributed texts and makes update according received comments
 - “Discussion included suggestions to clarify some information and to ensure that assertions are well supported citation or other means.”

AICN report draft

Background information

- ▲ Introduction
 - Scope
 - Purpose
 - Abbreviation
- ▲ Stepping into the Large-Scale AI era
 - ChatGPT ignites enthusiasm for large-scale AI models
 - Large-scale AI models show emergent abilities
- ▲ Large-scale AI model Training
 - AI training process
 - Distributed AI system and parallelism
- ▲ Communication characteristics in AI training
 - Sparsity of traffic in space
 - Sparsity of traffic in time
 - Huge amount of traffic for communication
- AI computing networks
- ▲ Requirements and Challenges of AI computing Networks
 - Scale
 - Efficiency
 - Availability
 - Future technologies
- Standard considerations
- References

Remaining work ①

Remaining work ②

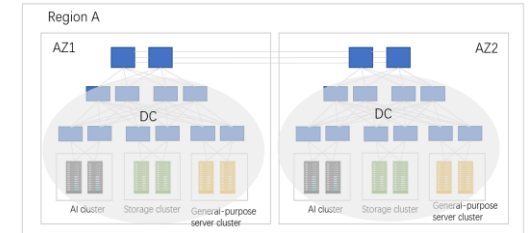
clarify data amount requirement of parallelisms, and model development figure

Total compute of distributed AI system = single AI accelerator compute * Scale * Efficiency * Availability

Scale

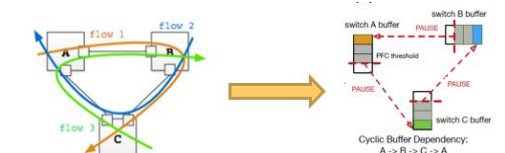
Distributed AI training across locations to address growing power consumption

- Challenges:
 - limited bandwidth of long-distance link
 - unpredictable transmission on long-distance links



Optimized network topology to address growing network cost

- Challenges:
 - irregularity of paths increases risk of deadlock causing by link-level flow control



Efficiency

Traffic management to address communication contention in order to increase AI accelerators utilization

- Path optimization --- load balancing
 - ECMP issue
 - Packet spray issue
- Data rate optimization --- flow control/congestion control
 - Intra-job flows
 - Inter-job flows

Work Item Proposal

Work item: AI computing Network (AICN)

Purpose:

- Understand the requirement of network for AI computing.
- Look for potential standardization opportunity in IEEE802.

Scope:

- Study main factors (parallelism, collective communication) in AI training which impact traffic.
- Analyze the major challenges for the network.
- Investigate future network technologies.
- Identify potential standard work.

Deliverables:

- A complete AICN report, including
 - Background/Use cases
 - AI computing network requirements and challenges
 - Potential technologies
 - Standardization considerations

Schedule:

- 2 months to draft a complete version of AICN
- 2 months to circulate the report for comments and start comment resolution

Co-editor: Lily Lyu (Huawei), Jieyu Li (China mobile)

Motion Text

To initiate a Nendica work item on AI computing network

Proposed:

Second: