# Reference Implementation Model and Proof-of-Concept Reference of Green Computing with Remote GPU

Classification: APPROVED REFERENCE DOCUMENT

Confidentiality: PUBLIC

Version 1.0

March 4, 2025

[Remote GPU Use Case RIM PoC Ref]

# Legal

# Contents

# List of Figures

# Executive Summary

This document outlines a reference implementation model and proof-of-concept (PoC) for utilizing Green Computing (GC) with remote GPU technology to enhance the time efficiency of Generative AI (GenAI) and Large Language Model (LLM) training. Through the integration of IOWN Global Forum (GF) technology and Open All-Photonic Networks (APNs), the project aims to reduce power consumption and improve the performance of data-intensive AI tasks. The PoC evaluates the effectiveness of remote GPU connections over APN compared to traditional WAN networks, focusing on training time, energy efficiency, and security. Key scenarios include direct access, data replication, and caching, with metrics collected across different configurations. The objective is to provide a flexible, cost-effective architecture for AI model development while ensuring data security through advanced isolation and monitoring. This approach supports the real-time processing of large datasets and helps reduce operational costs and energy consumption, particularly in AI training. The document sets the stage for future evaluations and demonstrates the potential of IOWN GF technology in creating more sustainable and secure AI training systems.

# 1 Introduction

## 1.1 Purpose

Since ChatGPT was launched in late 2022, the potential for AI-driven business transformation has expanded globally. Rather than simply analyzing and classifying existing data, GenAI can create entirely new things, including text, images, audio, and synthetic data. It brings breakthroughs in human creativity and productivity to business, science, and society.

GenAI applications are built by foundation models trained with complex deep learning systems at scale on vast amounts of data. Most of today's foundation models are large-scale language models (LLMs) trained on natural language processing.

The following opportunities are examples of use cases where LLMs can be applied.

1. Business meetings

   An LLM application can summarize the content of business meetings, especially extracting the points of discussion, what should be done, action items, and next steps, as well as clarifying the meeting's results.

2. Call centers

   After converting the conversation information between customers and call center staff into text, the LLM application can evaluate the content of customer's inquiries and search for similar responses to suggest potential next steps and responses.

3. Medical record information

   An LLM application can review and summarize patient examination results and other vital information such as, medications and interview responses. It can also support doctors' diagnoses based on similar symptoms and follow-up observations.

4. Demand Forecasting

   LLM can analyze historical sales data, market trends, and external factors to generate accurate demand forecasts. This helps optimize inventory levels and production schedules for manufacturing.

5. Drug discovery

   LLM enable the data analysis and prediction with vast amounts of unstructured data. These insights can lead to more accurate diagnoses, personalized treatments.

6. Digital human

   Conversational technology, particularly advancements in emotionally aware AI, can create more natural, human-centered interactions. Digital human is designed to simulate human interactions and add an emotional dimension utilizing LLM.

7. Smart Robots

   LLM enable more natural and effective interactions with robots due to the ability to understand and generate human language. And recent research, AI agent uses LLMs to automatically generate reward algorithms to train robots to accomplish complex tasks.

8. Smart Design

LLM realize to produce designs, blueprints, software code, and solutions that mimic human creativity and analytical skills.

9. Fraud detection and cybersecurity

LLM detects anomalies in customer behaviors and connections as well as patterns of accounts and behaviors that fit fraudulent characteristics. Threat detection and data generation from LLM advancements improve cybersecurity.

These large-scale calculations for GenAI / LLM consume considerable electric power. As a result, we need to use computing resources in the green data center as much as possible. Today's GPU computing services assume that computing and data storage resources are in the same data center. However, this assumption forces enterprises to make a different choice between storing confidential data in a GPU as a Service (GPUaaS) site or owning private GPU infrastructures.

We expect APN to solve this problem by enabling enterprises to use GPU in a GPUaaS site without uploading their data there in advance, in other words, with on-the-fly remote storage access from the GPUaaS site. We call this concept remote GPU over APN.

The purpose of this project is to prove that enterprises can do training for GenAI/LLM with remote GPU over APN without significantly increasing training time, power consumption, and equipment costs.

## 1.2 Project Objective

The objective of Green Computing with Remote GPU over APN for GenAI / LLM with IOWN GF technologies project is to define a reference design that will allow AI users to customize their models in a cost-effective new architecture while bringing exciting new capabilities to meet today's fine-tuned AI model development challenges.

The scope of this project is to:

1. Describe the Green Computing with Remote GPU over APN for GenAI / LLM and its key requirements

2. Define the Technology Evaluation Criteria, which include reference cases and key benchmarks

3. Develop the Reference Implementation Model (RIM), which provides a practical implementation of IOWN Technologies as a reference model to realize the use case

4. Define the Proof of Concept (PoC) Reference, which provides guidelines for conducting PoCs for the use case to evaluate the RIM with the defined Technology Evaluation Criteria

5. Run and evaluate the PoC based on the PoC Reference

## 1.3 Document Scope

We have completed the first step of this activity, covering Steps 1 and 2. Please refer to the document "Green Computing with Remote GPU Service for Generative AI / LLM Use Case" for further information. This document is the second step of this activity, which covers Steps 3 and 4 to engage early adopters in the GenAI / LLM industry.

The following sections of this document describe one of the Reference Implementation Models of remote GPU service Infrastructure and define its system validation testing and performance benchmark as Proof-of-Concept references.

There are several approaches to building remote GPU service infrastructure, such as confidential computing technologies; these will be addressed in future versions of this document.

# 2 Reference Implementation Models

This section describes the initial reference implementation models (hereafter RIM) for each use cases.

## 2.1 Basic Strategy of System Architecture and Design

We use IOWN networking technology as a design strategy for systems that transfer large amounts of data for GenAI/LLM training between GPUs in green data centers and data storage in user offices and private clouds (Figure 1).



*Figure 1 Training Data Transfer on the Open APN to GPU Clusters for LLMs*

It is valuable to reduce connectivity costs from traditional technologies, with Open APN of dynamic path allocation functionality and on-demand cost model. It also adopts deterministic networking capabilities that guarantee fixed delays between POPs by data-centric infrastructure capabilities. Today, high-performance GPUs are increasingly used as GPU-as-a-service. With Open APN on demand, users get a data-centric infrastructure without copying any data. This means you can connect GPU servers to your storage servers on-demand via Open APN for AI customization.

For these use cases, traffic types such as low latency and bandwidth must be managed based on the AI model size, data volume, training duration, and checkpoint iterations.

For GenAI/LLM customization that uses a large amount of data in a relatively short training period of a week, Direct Access with Open APN can help reduce operational costs and increase data confidentiality.

## 2.2 Reference Implementation Model for Remote GPU

The envisioned use case is for users to generate an LLM by training on GPU computing located on Open APN-connected green data centers using the training data that exists at their location.

### 2.2.1 System Architecture

The use case assumes that the physical distance between the green data center and the user offices is approximately 1000 km or less. It also assumes that the green data center has at least one GPU server and that the user office has at least one storage server that can be connected through the appropriate network interface. Green data centers may have storage servers for replication or caching. GPUs and replication/cache servers in the green data center are intended for temporary use on demand.

*Figure 2 GPU/Storage - three storage access methods*

For sensitive data use cases, data confidentiality must be ensured throughout the system so that the data owner can maintain data sovereignty. According to the Zero Trust model requires isolation technologies that directly protect data access, and monitoring and reporting systems that enable mitigation, early detection of incidents, and auditing.

The requirements for isolation in this system can be broken down into the following elements.

**Protection for data in use and data at rest**

- Only authorized users can access data computed or stored on the GPUs and replication/cache servers in the green data center, taking into account internal attacks by DC operators.
- Only authorized users can access data in storage in the user's office.

**Protection for data in motion**

- Data in communication paths between storage in the user office and GPUs in the green data center, as well as within data centers and servers, including bus communication, should be adequately protected, taking into account internal attacks by DC operators. Post-quantum security should also be considered.

It is also necessary to ensure that the data is protected seamlessly between each of the above isolation measures.

Regarding monitoring and reporting, the system should be auditable to prove that data has not been compromised, especially on the GPUs and replication/cache servers in green data centers, however the appropriate level of auditability depends on the implementation and requirements and does not affect the system's performance, so detailed requirements are not discussed here.

## 2.2.2 Initial Design for Each Functional Node

### 2.2.2.1 GPU Cluster

GPU cluster is a group of computers that have GPUs on every node. Multiple GPUs provide accelerated computing power for heavy computational tasks such as LLM training.

### 2.2.2.2 Storage node

Storage nodes, consisting of dedicated storage equipment that operates as storage servers, provide file storage services that store large-scale training data for GPU computing. Storage nodes are also used to store checkpoints during model training processing.

The storage performance in this PoC is measured by the processing time of AI training. In addition, the fundamental storage performance is measured by throughput and latency using tools (such as the nfsiostat or iostat command).

### 2.2.2.3 Replication / Cache node

Replication nodes and cache nodes are effective for quickly reading large amounts of training data when the latency of the Internet or Open APN is significant (such as long-distance connections).

Replicating or caching training data close to the GPU allows high-speed model training without affecting network latency.

### 2.2.3 Data access sequence of the AI training process

This pattern shows a general system sequence diagram for the AI training process on GPU computing nodes and storage.

- Prerequisites:
    - Model Data and training data are stored on local or remote storage.
    - GPU computing and storage are connected via Ethernet.
    - The storage has a file access service and is mounted from a GPU server.
- Description:
    - In initialization phase, the model data is loaded from storage to GPU server.
    - The training data is also loaded into the memory of the GPU server in the Initialization phase.
    - In the loop phase, any training data that could not be loaded into memory in the Initialization phase will be loaded into the GPU server.
    - After training phase, checkpoint data is saved to storage.
- System sequence diagram:

*Figure 3 general system sequence diagram on AI training*

## 2.2.4 The choices of System design

This section describes the design choice of data access for model training. The sequence diagram focuses on the data communication between GPU computing and remote storage and shows a series of processes from start to finish.

For the data communication use case with remote locations, the following three choices of system design are defined. For each system design, a sequence diagram and the prerequisites for defining the sequence diagram are specified.

- System design #1: Direct access
- System design #2: Replication
- System design #3: Cache on Internet / Open APN

This document also describes an implementation model using an NFS-based storage server widely used in enterprise AI data infrastructures. The system design of PoC should be chosen from one or more of these three designs.

## 2.2.4.1 System design #1: Direct Access

This section describes a data sequence where GPU computing reads and writes data directly to and from Remote storage.

- Hardware conditions
  - ○ GPU Cluster
  - ○ Storage node
  - ○ Network switch for 100G or 10G connection
  - ○ Open APN for 100G or 10G connection [remote site]



*Figure 4 System diagram of Direct access at same site*



*Figure 5 System diagram of Direct access at remote sites*

- Recommended conditions (if using an NFS-based storage server):
  - ○ Remote storage has an NFS server service running and supports NFS v4 and NFS over RDMA.
  - ○ GPU Computing requires NFS client software to be installed and support NFS over RDMA.
  - ○ If a firewall is configured, NFS communication must be permitted (NFSv4=TCP/2049 port, NFS over RDMA=TCP&UDP/20049 port).
  - ○ A common NFSv4 domain ID and User ID must be configured for GPU computing and Remote storage.
  - ○ Training data is stored in an NFS shared directory and can be read and written from GPU computing.
- Description:
  - ○ Remote storage handles large amounts of data communication (reading and writing), so it is desirable for it to have sufficient storage performance. In general, the following is recommended (but not required):
    - ▪ RAID configuration using multiple SSDs for high-speed read and write.
    - ▪ Scalable storage that can handle increases in data volume
- Data sequence diagram
  - ○ Data read/write sequence

*Figure 6 Direct access - Data sequence for read/write*

In this PoC, the usefulness of these technologies is evaluated by comparing the process execution performance of GPUs and storage at different locations under Open APN with the performance of the local environments. File servers and high-performance GPU servers, which are conventionally located at the same site, are placed at remote sites connected by Open APN. Large amounts of training data and models for generated AI such as LLM and others, are transferred to these sites.

- KR1.1
    - Increase of training time for direct connections at remote sites over Open APN is less than 10% from direct connections at same sites.
    - The training time is measured by adding a command to record the timestamp in a log file after processing the Read Time, GPU Time, and Save Time in the training loop of a training script written in Python or another language (such as trainer.py).

In this direct access use case, AI training can be performed without replicating data beyond a security domain such as the country. Personal information and other confidential information in training data should be stored only in the storage node and not stored in the GPU node from the perspective of data sovereignty. It is therefore also essential to ensure that no training data remains on the GPU nodes.

## 2.2.4.2 System design #2: Replication

In this pattern, all training data is replicated from remote storage to local storage in advance, and GPU computing only reads and writes data to and from local storage.

- Hardware conditions
    - Replication node for training data
    - Internet connection (from 1G to 100G)
    - Open APN for 100G or 10G connection
    - Other hardware conditions are the same as in section 2.2.4.1

*Figure 7 System diagram of the Internet with Replication*



*Figure 8 System diagram of Open APN with Replication*

- Recommended conditions (if using an NFS-based storage server):
  - GPU computing and local storage are connected via Ethernet at the site.
  - Local and remote storage are connected via IOWN infrastructure (including the Open APN).
  - Local storage has an NFS server service running and supports NFS v4 and NFS over RDMA.
  - GPU Computing requires NFS client software to be installed and support NFS over RDMA.
  - If a firewall is set between local storage and remote storage, replication communication must be allowed.
  - A common NFSv4 domain ID and user ID are set for GPU computing and local storage.
  - Training data is stored in an NFS-shared directory on local storage and can be read and written from GPU computing.
- Description:

- ○ Local storage should meet the same requirements as the remote storage described in "2.2.4.1. System design #1: Direct Access - Description".
- System sequence diagram
  - ○ The system sequence diagram describes "1. data communication between GPU computing and local storage", and "2. data replication communication between local storage and remote storage" in separate sequence diagrams.
    - ▪ Data read/write sequence
      - ▪ The following is a sequence diagram of how GPU computing reads and writes data to local storage.



*Figure 9 Replication - Data sequence for read/write (GPU computing to Local storage)*

- ▪ Replication sequence
  - ▪ The sequence diagram of data replication between storage devices is shown below.

*Figure 10 Replication sequence (Local storage to Remote storage)*

- KR2.1
  - Comparing the training time with the Internet connection and Open APN connection.
  - In the case the replication process is needed before the access to the replication node, The replication processing time should be added to the Read Time.
  - The measurement for the training time (Read Time, GPU Time, and Save Time) is the same as KR1.1 in 2.2.4.1.

In this replication use case, the AI training performance excluding the replication process is nearly the same as direct access on the same site. Depending on the size of training data and the number of pretraining loops, this replication design can be advantageous over another designs.

## 2.2.4.3  System design #3: Cache on the Internet / Open APN

This pattern shows a data sequence diagram when GPU computing uses Cache storage in conjunction with Remote storage for data communication. By using Cache storage with Remote storage, data can be read and written at high speed without being affected by the latency between Remote storage and GPU computing.

The network between sites is connected via the Internet / Open APN, and the first-time data is read, it takes longer because the data is not cached in cache storage. Data is also written to the cache storage, so write performance is also fast.

In the case Open APN is used as a network between sites, significantly faster processing speeds can be expected compared to the case the Internet is used.

- Hardware conditions
  - Cache node for training data
  - Internet connection (from 1G to 100G)

○ Open APN for 100G or 10G connection

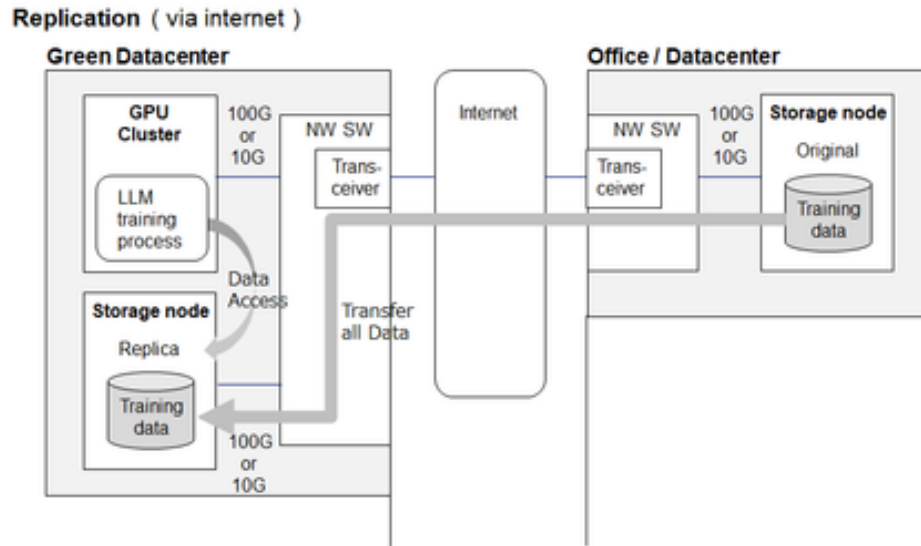○ Other hardware conditions are the same as in section 2.2.4.1.
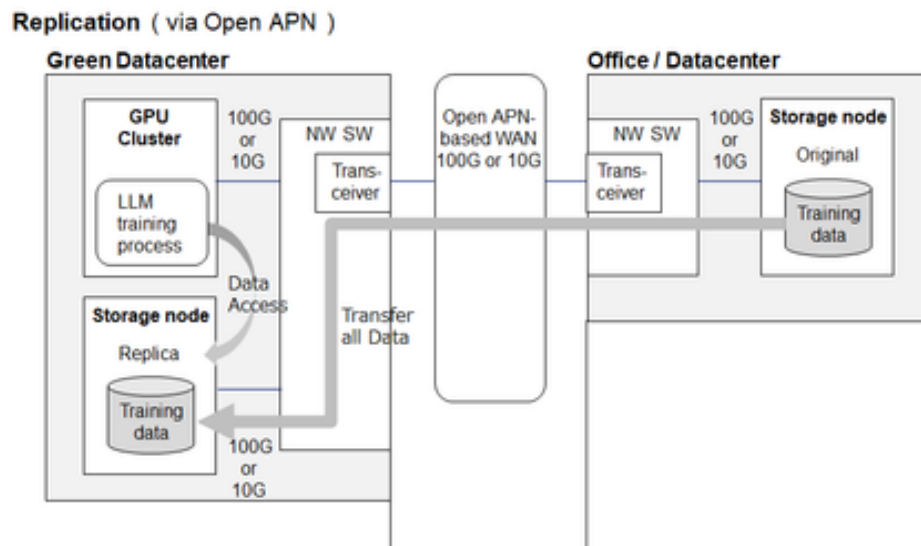


*Figure 11 System diagram of Internet with Cache*



*Figure 12 System diagram of Open APN with Cache*

● Recommended conditions (if using an NFS-based Cache node):

○ GPU computing and local storage are connected via Ethernet at the site.

○ Cache node and remote Storage node are connected via the Internet / Open APN.

○ Cache node has an NFS server service running and supports NFS v4 and NFS over RDMA.

○ GPU Computing requires NFS client software to be installed and support NFS over RDMA.

○ Data caching communication must be allowed if a firewall is set between local and remote storage.

○ A common NFSv4 domain ID and user ID are set for GPU computing and local storage.

- ○ Training data is stored in an NFS-shared directory on local storage and can be read and written from GPU computing.
- Description:
  - ○ Local storage should meet the same requirements as the remote storage described in "2.2.4.1. System design #1: Direct Access - Description".
- System sequence diagram
  - ○ Since the system sequences for reading and writing data are different, sequence diagrams are written separately for each.
    - Data read sequence
      - The sequence diagram for GPU computing reading data is shown below.



*Figure 13 Cache - Data sequence for read*

- Write sequence
  - The sequence diagram for GPU computing writing data is shown below.



*Figure 14 Cache - Data sequence for write*

- KR3.1
  - ○ Comparing the training time with the Internet connection and Open APN connection.
  - ○ In the case certain caching process is needed before the access to cache node, The cache processing time should be added to the Read Time.

- ○ The measurement for the training time (Read Time, GPU Time, and Save Time) is the same as KR1.1 in 2.2.4.1.

In this cache use case, since the general storage handles caching automatically, and in many cases good performance is expected, if the cost of the cache nodes is acceptable.

# 3 Technology Evaluation Criteria

This section defines PoC References use case of green computing with remote GPU over Open APN. This PoC reference section describes the methods for validating the key requirements and evaluating metrics for benchmarking. The key requirements are the objective criteria that need to be met. The PoC is performed with more desirable implementation for users. To benchmark the performance of green computing with remote GPU over Open APN, we compare the performance between direct access at the same site and direct access at remote sites over Open APN. In the use case of Replication, we compare the performance over Open APN and an Internet connection including the initial data replication. In the use case of Cache, we compare the performance over Open APN and an Internet connection.

## 3.1 System design #1: Direct Access

- **Validation Testing**
  - **Requirements for Evaluation:**
    - The container for LLM environment is deployed on the local storage of GPU server.
    - The training data for LLM is stored on a storage Server.
    - Open APN network.
    - The directory of training data on a storage server, such as an NFS server etc. is mounted from the GPU server.
  - **Evaluation Methods:**
    - Measuring the increase of training time for the remote site environment over Open APN and the same site environment.
  - **Metrics to be Collected:**
    - Read time, GPU time and Save time of training time (See data access sequence in section 2.2.3 and Initialization process time is not included). The training phase is repeated multiple times depending on the number of configurations. Specifically, the time is measured by adding a command to record the timestamp to log file after processing the Read Time, GPU Time, and Save Time in the training loop of a training script written in Python or another language (such as trainer.py). If isolation technologies for data security are applied to the PoC system, performance data must be reported both with and without such security technologies in place.
      - Training time [Read time]
      - Training time [GPU time]
      - Training time [Save time]
  - **Logs to be collected:**
    - The logs of LLM training (information logs and error logs)
  - **Points of Verify:**
    - Normal processing continues without error and with the completion of AI training.
  - **Security Considerations to Report**
    - Measures to protect data in use and data at rest in a Green Data Center
    - Measures to protect data in motion
    - Measures to protect data at rest in the user office

Note: The PoC requirement does not specify the threats to be addressed or the specific measures to be used, but the report must include how the above three points were considered in the PoC. The reports that do not apply or only partially consider data security measures are also acceptable.

- **Benchmark Testing**
  - Comparing the total training time for the remote environment over Open APN and the same site environment and measuring the increased rate of remote connection.
  - Storage access speed (Measured by recording the output of nfsiostat command etc. See section 2.2.2). The storage performance for the remote site environment over Open APN should be almost same as that for the same site environment.
    - Throughput
    - Latency
  - Power consumption of each piece of equipment (GPU cluster, Network switches and Storage node excluding the Open APN network part). The method of calculation should be ether of the following. And also, total power consumption per site is used to compare the difference between green and standard data centers. Power consumption for remote site environment over Open APN should be almost the same as that for same site environment. Generally, a green data center's PUE is expected to be 10-30 percent better than a standard data center, and the total power consumption is expected to be better than that of typical data center/office single site.
    - The estimation of the power consumption based on the rated power consumption from the device specifications
    - The actual monitoring of power consumption during the training procedure
  - The cost of each piece of equipment (GPU cluster, Network switches and Storage node excluding the Open APN network part). The cost should be calculated from the market price.

# 3.2. System design #2: Replication

- **Validation Testing**
  - **Requirements for Evaluation:**
    - The container for LLM environment is deployed on the local storage of GPU server.
    - The training data for LLM is stored on storage Server.
    - Open APN network. And Internet connection
    - The replicated directory of training data on replicated storage node which can be mounted from GPU server.
  - **Evaluation Methods:**
    - Measuring each training time for Open APN connection with Replication node and Internet connection with Replication node.
  - **Metrics to be Collected:**
    - Measured metrics are the same as in section 3.1 System design #1: Direct Access.
  - **Logs to be collected:**
    - The logs of LLM training (information logs and error logs)
  - **Points of Verify:**
    - Normal processing continuation without error and completion of AI training.
  - **Security Considerations to Report**

- ▪ Measures to protect data in use and data at rest in a Green Data Center

- ▪ Measures to protect data in motion

- ▪ Measures to protect data at rest in the user office

Note: The PoC requirement does not specify the threats to be addressed or the specific measures to be used, but the report must include how the above three points were considered in the PoC. The reports that do not apply or only partially consider data security measures are also acceptable.

- ● **Benchmark Testing**
  - ○ Comparing the total training time for the Open APN connection with Replication and the internet connection with Replication node.
  - ○ Storage access speed of replicated storage node (Measured by recording the output of nfsiostat command etc. See section 2.2.2).
    - ▪ Throughput
    - ▪ Latency
  - ○ Power consumption of each equipment (GPU Cluster, replicated storage node, network switches and storage node for original data. Excluding the Open APN network and internet network part). The method of calculation should be either of the following, and also total power consumption per site to compare the difference between green and standard data centers:
    - ▪ The estimation of the power consumption based on the rated power consumption from the device specifications
    - ▪ The actual monitoring of power consumption during the training procedure
  - ○ The cost of each equipment (GPU Cluster, replicated storage node, network switches and storage node for original data. Excluding the Open APN network and internet network part). The total cost should be calculated from the market price of each equipment.

# 3.3. System design #3: Cache

- ● **Validation Testing**
  - ○ **Requirements for Evaluation:**
    - ▪ The container for LLM environment is deployed on the local storage of GPU server.
    - ▪ The training data for LLM is stored on storage Server.
    - ▪ Open APN network.
    - ▪ The cache node for training data on storage server which can be accessed from GPU server.
  - ○ **Evaluation Methods:**
    - ▪ Measuring each training time for Open APN connection with Cache node and Internet connection with Cache node.
  - ○ **Metrics to be Collected:**
    - ▪ Measured metrics are the same as in section 3.1 System design #1: Direct Access.
  - ○ **Logs to be collected:**
    - ▪ The logs of LLM training (information logs and error logs)
  - ○ **Points of Verify:**
    - ▪ Normal processing continuation without error and completion of AI training.
  - ○ **Security Considerations to Report:**

- Measures to protect data in use and data at rest in a Green Data Center
- Measures to protect data in motion
- Measures to protect data at rest in the user office

Note: The PoC requirement does not specify the threats to be addressed or the specific measures to be used, but the report must include how the above three points were considered in the PoC. The reports that do not apply or only partially consider data security measures are also acceptable.

- **Benchmark Testing**
  - Comparing the total training time for the Open APN connection with Cache and the internet connection with Cache node.
  - Storage access speed of cache node (Measured by recording the output of nfsiostat command etc. See section 2.2.2 ).
    - Throughput
    - Latency
  - The power consumption of each piece of equipment (GPU cluster, Cache node, network switches and storage node for original data. Excluding the Open APN network and internet network part). The method of calculation should be ether or the following. And the total power consumption per site are also used to compare the difference between green and standard data centers.
    - The estimation of the power consumption based on the rated power consumption from the device specifications
    - The actual monitoring of power consumption during the training procedure
  - The cost of each piece of equipment (GPU cluster, Cache node, network switches and storage node for original data excluding the Open APN network and internet network part). The cost should be calculated from market price of each equipment.

# 4  Conclusion

This document explains the basic strategy of system architectures and designs that utilize IOWN GF Technology to efficiently realize Green Computing with Remote GPU Service for GenAI / LLM Use Cases.

As explained in section 1.2, this document covered items 3 and 4 described in Section 1.1. the Objective to engage early adopters in customization of the GenAI/LLM. This document will be followed by proof-of-concept demonstrations and evaluation to prove the validity of IOWN GF Technology and its effectiveness in your own customized GenAI/LLM.

# 5 Appendix

Figure 15 indicates a typical LLM training pipeline with the timing of basic checkpoint saving based on an example of system workloads in GPT3 on OSS Megatron-LM. This section provides a detailed explanation of 2.2.3 Data access sequence on the AI training process. There are initialization and training processes between the GPU server and the Storage server, but after obtaining the train data, the system enters the iterative training phase.

In the typical pattern of LLM learning, the training data is repeatedly forwarded and backward during training. This process is called iteration. At the end of each iteration, it is time to save the completed training data using a checkpoint. At this point, the I/O bandwidth increases. The checkpoint mainly involves updating the parameters. When saving a checkpoint, the GPU core utilization rate drops until the saving is complete. After the checkpoint is complete, the iteration starts again. It is possible to reduce the time it takes to complete LLM training by reducing the time to save checkpoints and starting forward without waiting for the checkpoint to finish. This is called distributed checkpointing.
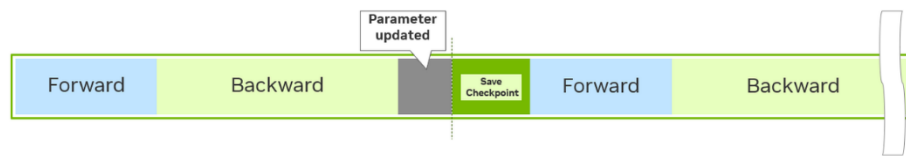


*Figure 15 Typical LLM training pipeline.*

Further information about the training pipeline of LLM and distributed checkpoint.

# References

- [Distributed Checkpoints — NVIDIA NeMo Framework User Guide](#)
- [Release NVIDIA Megatron Core 0.7.0 · NVIDIA/Megatron-LM](#)

# Acknowledgments

The management team is listed here, followed by contributors in alphabetical order.

- Fumiaki Kudoh
- Hideaki Tagami
- Joao Kluck Gomes
- Kei Karasawa
- Nobu Sasaki
- Yoshinobu Nakayama

# History

| Revision | Release Date | Summary of Changes |
| --- | --- | --- |
| 1.0 | March 4, 2025 | Initial Release |