



Security RIM for Green Computing with Remote GPU over All-Photonic Networks

Classification: APPROVED Reference Document

Confidentiality: PUBLIC

Version: 1

[GC with Remote GPU Security RIM]

October 2025



Legal

THIS DOCUMENT HAS BEEN DESIGNATED BY THE INNOVATIVE OPTICAL AND WIRELESS NETWORK GLOBAL FORUM, INC. (“IOWN GLOBAL FORUM”) AS AN APPROVED REFERENCE DOCUMENT AS SUCH TERM IS USED IN THE IOWN GLOBAL FORUM INTELLECTUAL PROPERTY RIGHTS POLICY (THIS “REFERENCE DOCUMENT”).

THIS REFERENCE DOCUMENT IS PROVIDED “AS IS” WITH NO WARRANTIES WHATSOEVER, WHETHER EXPRESS, IMPLIED, STATUTORY, OR OTHERWISE, INCLUDING WITHOUT LIMITATION ANY WARRANTY OF MERCHANTABILITY, NONINFRINGEMENT OF THIRD PARTY RIGHTS, TITLE, VALIDITY OF RIGHTS IN, FITNESS FOR ANY PARTICULAR PURPOSE, OR ANY WARRANTY OTHERWISE ARISING OUT OF ANY PROPOSAL, REFERENCE DOCUMENT, SAMPLE, OR LAW. WITHOUT LIMITATION, IOWN GLOBAL FORUM DISCLAIMS ALL LIABILITY, INCLUDING WITHOUT LIMITATION LIABILITY FOR INFRINGEMENT OF ANY PROPRIETARY RIGHTS AND PRODUCTS LIABILITY, RELATING TO USE OF THE INFORMATION IN THIS REFERENCE DOCUMENT AND TO ANY USE OF THIS REFERENCE DOCUMENT IN CONNECTION WITH THE DEVELOPMENT OF ANY PRODUCT OR SERVICE, AND IOWN GLOBAL FORUM DISCLAIMS ALL LIABILITY FOR COST OF PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES, LOST PROFITS, LOSS OF USE, LOSS OF DATA OR ANY INCIDENTAL, CONSEQUENTIAL, DIRECT, INDIRECT, PUNITIVE, EXEMPLARY, OR SPECIAL DAMAGES, WHETHER UNDER CONTRACT, TORT, WARRANTY OR OTHERWISE, ARISING IN ANY WAY OUT OF USE OR RELIANCE UPON THIS REFERENCE DOCUMENT OR ANY INFORMATION HEREIN.

EXCEPT AS EXPRESSLY SET FORTH IN THE PARAGRAPH DIRECTLY BELOW, NO LICENSE IS GRANTED HEREIN, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS OF THE IOWN GLOBAL FORUM, ANY IOWN GLOBAL FORUM MEMBER OR ANY AFFILIATE OF ANY IOWN GLOBAL FORUM MEMBER. EXCEPT AS EXPRESSLY SET FORTH IN THE PARAGRAPH DIRECTLY BELOW, ALL RIGHTS IN THIS REFERENCE DOCUMENT ARE RESERVED.

A limited, non-exclusive, non-transferable, non-assignable, non-sublicensable license is hereby granted by IOWN Global Forum to you to copy, reproduce, and use this Reference Document for internal use only. You must retain this page and all proprietary rights notices in all copies you make of this Reference Document under this license grant

THIS DOCUMENT IS AN APPROVED REFERENCE DOCUMENT AND IS SUBJECT TO THE REFERENCE DOCUMENT LICENSING COMMITMENTS OF THE MEMBERS OF THE IOWN GLOBAL FORUM PURSUANT TO THE IOWN GLOBAL FORUM INTELLECTUAL PROPERTY RIGHTS POLICY. A COPY OF THE IOWN GLOBAL FORUM INTELLECTUAL PROPERTY RIGHTS POLICY CAN BE OBTAINED BY COMPLETING THE FORM AT: www.iowngf.org/join-forum. USE OF THIS REFERENCE DOCUMENT IS SUBJECT TO THE LIMITED INTERNAL-USE ONLY LICENSE GRANTED ABOVE. IF YOU WOULD LIKE TO REQUEST A COPYRIGHT LICENSE THAT IS DIFFERENT FROM THE ONE GRANTED ABOVE (SUCH AS, BUT NOT LIMITED TO, A LICENSE TO TRANSLATE THIS REFERENCE DOCUMENT INTO ANOTHER LANGUAGE), PLEASE CONTACT US BY COMPLETING THE FORM AT: <https://iowngf.org/contact-us/>]
Copyright © 2025 Innovative Optical Wireless Network Global Forum, Inc. All rights reserved. Except for the limited internal-use only license set forth above, copying or other forms of reproduction and/or distribution of this Reference Document are strictly prohibited.

The IOWN GLOBAL FORUM mark and IOWN GLOBAL FORUM & Design logo are trademarks of Innovative Optical and Wireless Network Global Forum, Inc. in the United States and other countries. Unauthorized use is strictly prohibited. IOWN is a registered and unregistered trademark of Nippon Telegraph and Telephone Corporation in the United States, Japan, and other countries. Other names and brands appearing in this document may be claimed as the property of others.



Contents

- 1. Introduction.....4**
- 2. Review of Remote GPU over APN use case6**
 - 2.1. Target service 6
 - 2.2. Original system architecture 6
 - 2.3. Data security requirements 7
- 3. Security RIM with Today’s Hardware9**
 - 3.1. Design objectives and assumptions..... 9
 - 3.2. RIM structure10
 - 3.3. Another optional RIM regarding subsystem protecting Data in Motion 14
- 4. Performance Considerations 20**
- 5. Conclusion21**
- 6. References 22**
- History..... 24**

List of Figures

- Figure 2.2-1: Reference Implementation Model (Direct access via APN) 7
- Figure 3.1-1: Definition of subsystems 10
- Figure 3.2-1: Structure of security RIM for protection of entire data lifecycle 11
- Figure 3.3-1: Structure of optional RIM 1 for protection of data in motion..... 15
- Figure 3.3-2: Structure of security RIM 2 for protection of data in motion 17

List of Tables

- Table 3.2-1: Example implementation choices for security RIM for protection of the entire data lifecycle 13
- Table 3.3-1: Example implementation choices for security RIM 1 for protection of data in motion 16
- Table 3.3-2: Example implementation choices for security RIM 2 for protection of data in motion 18



1. Introduction

"Remote GPU over APN," a means of accessing advanced computing resources over a next-generation network, is expected to help enterprises remain competitive in the AI era, enabling them to upgrade their products and services with AI. This novel approach enables enterprises to utilize a GPU in a GPU-as-a-Service (GPUaaS) site without uploading their data in advance, in other words, with on-the-fly remote storage access from the GPUaaS site. In this use case, when sensitive data is handled, data security is a critical issue that can make or break the business viability of this use case. Considering that, the forum's use case document, "Green Computing with Remote GPU Service for Generative AI / LLM Use Case - Light Speed Data Transfer for AI Training -(GC with Remote GPU UC) [GC with Remote GPU UC]" defines a data security requirement as follows: "when training an AI/LLM using GPUs in a remote green computing center as a third-party service, the confidentiality of the training data and model must be guaranteed within the communication path and the actual computing space, including within the green data center." On the other hand, the forum's reference implementation model (RIM) and PoC reference document, "Reference Implementation Model and Proof-of-Concept Reference of Green Computing with Remote GPU (Remote GPU UC RIM PoC Ref)[Remote GPU UC RIM PoC Ref]" does not give any recommendation as to how we should accomplish the data security requirement.

Related to the above data security requirement, the forum has already developed two functional architectures, "Functional Architecture for Protection of Data in Use: IOWN Privacy Enhancing Technologies (IOWN PETs architecture) [IOWN PETs FA]" and "Functional Architecture for Protection of Data in Motion: Multi Factor Security Key Exchange and Management (MFS architecture) [MFS FA]". IOWN PETs architecture provides a data space that ensures confidentiality by seamlessly connecting environments protected by PETs through secure communication, facilitated by open interfaces and its technical guarantee to users. MFS architecture provides post-quantum encrypted communications with high crypto-agility on IOWN infrastructure.

The above two functional architectures are very relevant to the data security requirement of the Remote GPU over APN use case. However, no document illustrates how we can accomplish the data security requirement with the following two functional architectures. The purpose of this document is to accelerate the implementation and evaluation of the Remote GPU over APN use case, taking into account data security requirements, by demonstrating how these data security requirements can be met in practice using the aforementioned two architectures.

The objective of this document is to present several RIMs using the above security architectures being discussed in the IOWN Global Forum that can satisfy those data security requirements of the Remote GPU over APN use case.

The scope of this document is limited to the following:



- This document focuses on the implementation of the components of the IOWN PETs and MFS architecture that isolate user data and the communication between these components, which directly affect performance and confidentiality assurance.
- RIMs are described to the extent that it can be achieved with current hardware. This document may be updated as hardware advances.



2. Review of Remote GPU over APN use case

2.1. Target service

The envisioned use case of the "Remote GPU over APN" is for users to generate a Large Language Mode (LLM) by performing training on GPU computing resources located in an Open APN-connected green data center, using the training data that exists at their own location. For use cases where multiple users train their respective AI models, the green data center becomes an AI factory, where the customer's data and models are received, trained, and the parameters are transferred. From a data-centric perspective, the image is "the required GPU computing resources are delivered from the green data center to the customer's data storage, and the models are trained in the customer's environment." This use case requires that the GPU computing in the green data center and the Open APNs connected to them should be available on demand for as long as needed.

When using sensitive data in this use case, data security must be considered to ensure the data owner maintains data sovereignty. The following are examples of LLM use cases where special data security considerations should be taken into account:

Medical record information

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

Drug discovery

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

For security concerns that require special consideration in these use cases, please refer to the IOWN PETs architecture [IOWN PETs FA].

2.2. Original system architecture

Remote GPU UC RIM PoC Ref [Remote GPU UC RIM PoC Ref] defines the system architecture as follows: "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 utilize storage servers for replication purposes or caching. GPUs and replication/cache servers in the green data center are intended for temporary use on demand."

A typical example of the above use case is the "Direct access (via APN)" method. In this method, the GPU accesses training data on the Office side via Open APN. This document focuses on the "Direct access (via APN)" method below.

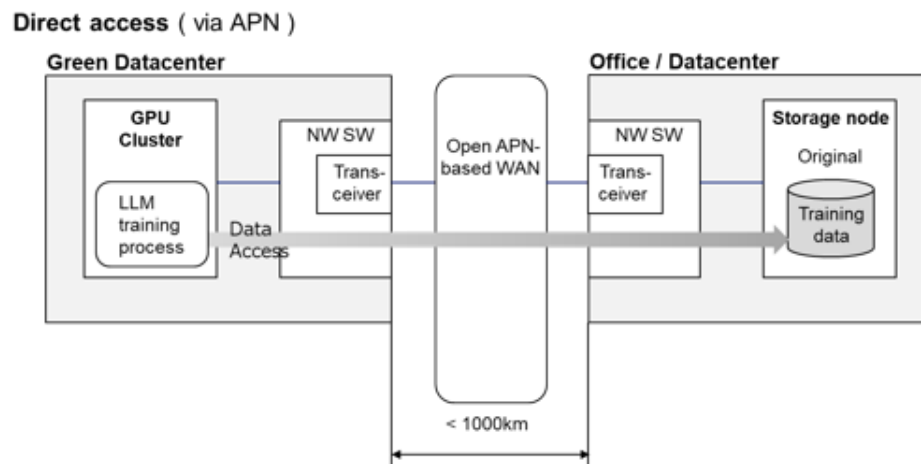


Figure 2.2-1: Reference Implementation Model (Direct access via APN)

2.3. Data security requirements

According to the Zero Trust model, isolation technologies that directly protect data access, monitoring and reporting systems that enable mitigation, early detection of incidents and auditing are required to ensure data security throughout the system.

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

Req 1: Protection of data in use in a Green Data Center

- Only authorized users can access data computed on the GPU in the green data center, taking into account internal attacks by DC operators.

Req 2: Protection of 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.

Req 3: Protection of data at rest in the user's office

- Only authorized users can access data in storage in the user's office.

The key point in this use case is that, due to the nature of computing sensitive data in a DC operated by a third party, data isolation must not be interrupted from storage to the GPU. To eliminate the risk of data exposure, it is necessary to ensure that data is protected seamlessly between each of the above isolation measures.



This document focuses solely on isolation technologies and does not address monitoring and auditing as the appropriate level of auditability depends on individual requirements.



3. Security RIM with Today's Hardware

This section describes a Reference Implementation Model that meets the above data security requirements.

3.1. Design objectives and assumptions

The goal of this RIM is to apply the IOWN PETs architecture to this use case, which consists of components that can protect data seamlessly in E2E from Storage to GPU and also provide users with technical assurance of its confidentiality. Assuming not only external attacks but also internal attacks by operators, an isolated space (PETs Space) is formed across the office and data center, where storage and the green data center are located. Only authorized users can access the data used in this use case at any time. Before using the space, users must confirm and trust the information about its environment (configuration certification information), which consists of multiple components. In addition, the communication path between DCs should be protected by post-quantum security with high crypto-agility.

There are two key implementation points:

- The PETs Space must be composed of components that can be verified to prevent unauthorized components from entering the PETs Space.
- There should be no unprotected timing between components (even during data state transitions).

To achieve the above, this RIM consists of three major subsystems that address the aforementioned data security requirements, and these subsystems must be interconnected.

- Subsystem for the protection of data in use
This subsystem corresponds to Req .1 and protects data used by a CPU and GPUs in the Green Data Center.
- Subsystem for the protection of data in motion
This subsystem corresponds to Req .2 and protects communication from the office-side storage to the CPUs in the Green Data Center.
- Subsystem for protection of data at rest
This subsystem corresponds to Req.3 and protects stored data in office site.

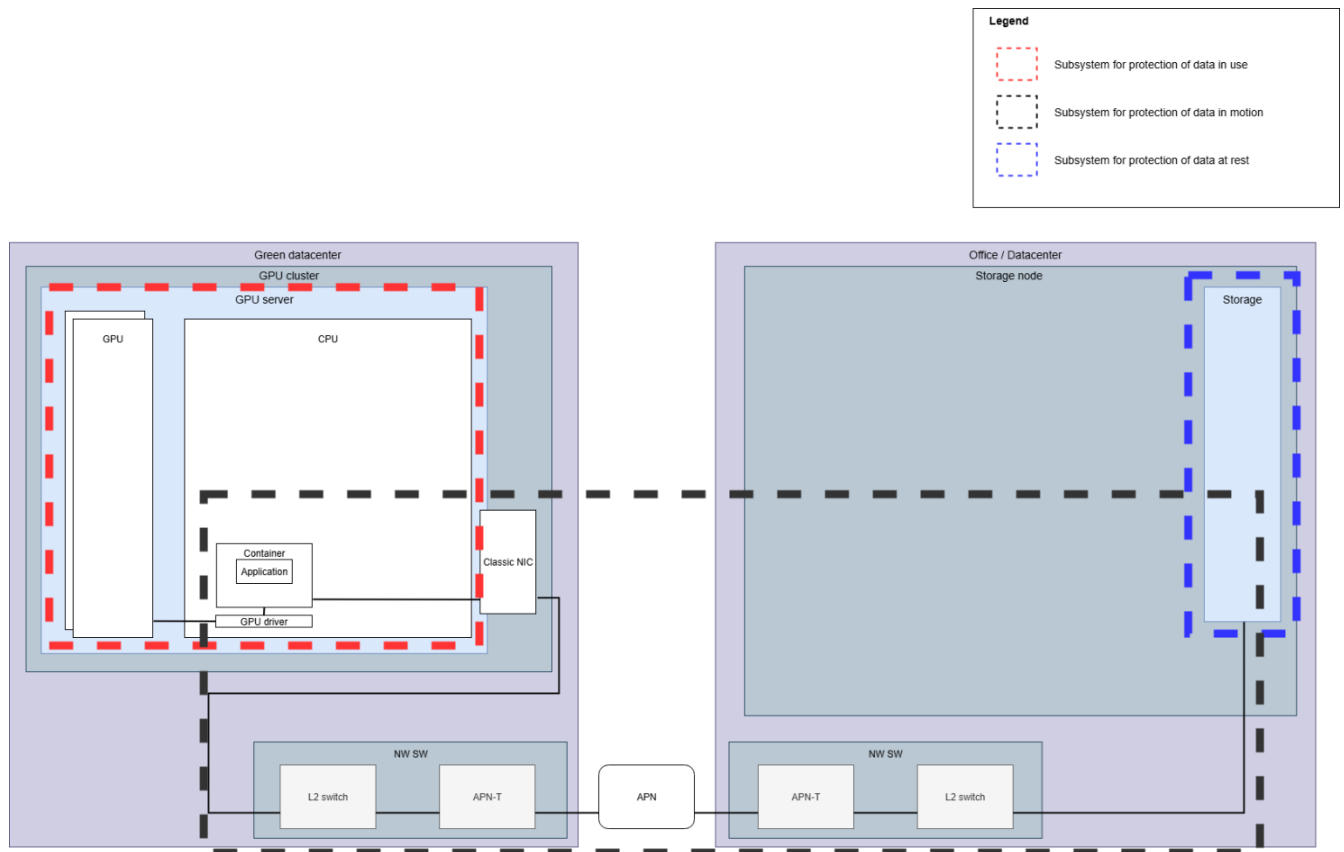


Figure 3.1-1: Definition of subsystems

3.2. RIM structure

This subsection outlines how IOWN PETs architecture using MFS architecture to achieve end-to-end data confidentiality from storage to GPU can be constructed.

Figure 3.2-1 shows the RIM structure centering secure seamless connection of Confidential Virtual Machines (CVMs) using TEE on a GPU cluster and storage node with MACsec/IPsec using the MFS architecture that meets the data security requirements of the aforementioned Remote GPU use case.

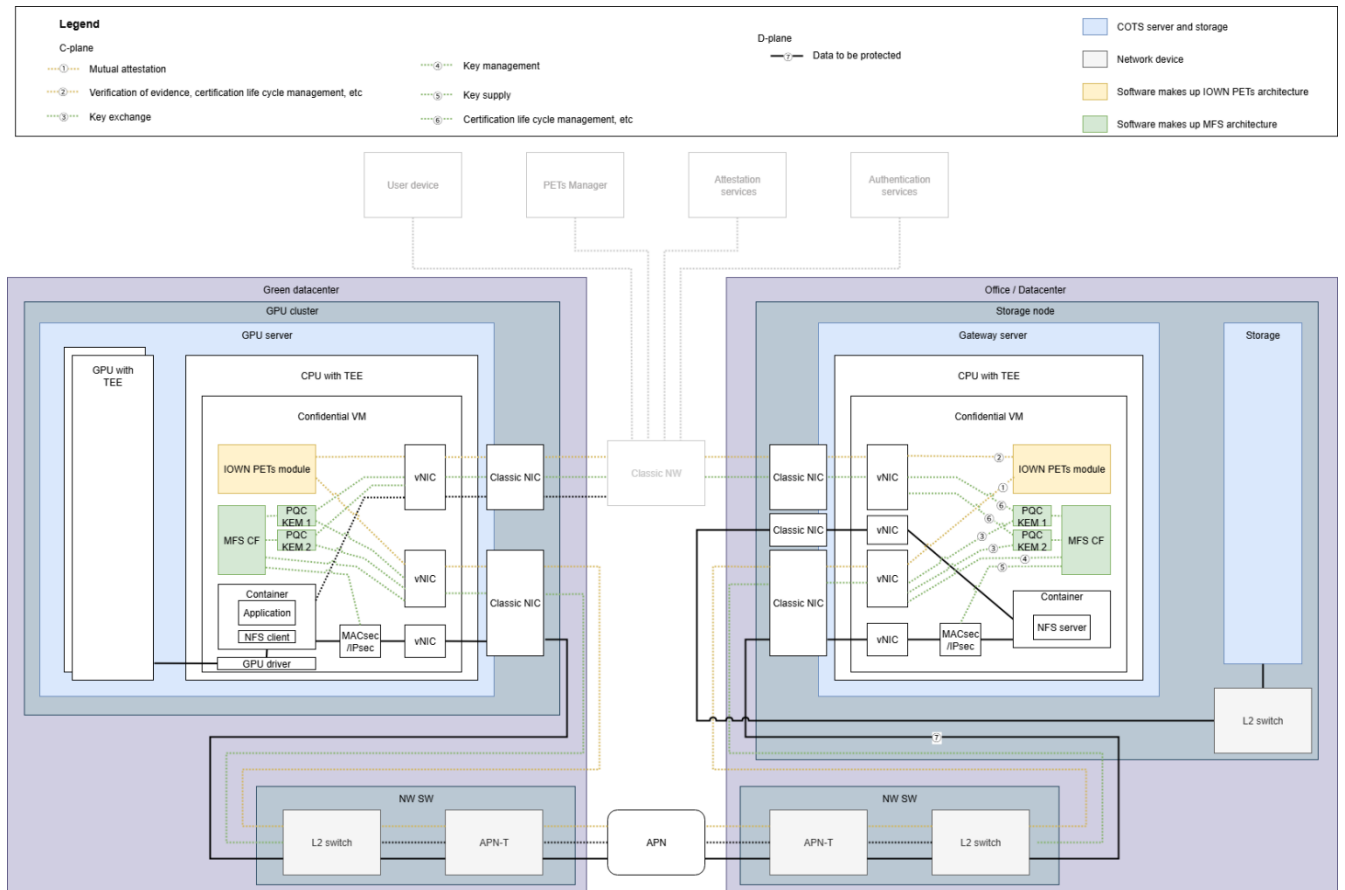


Figure 3.2-1: Structure of security RIM for protection of entire data lifecycle

3.2.1. Challenges of existing approaches

Although cloud vendors have offered confidential computing services utilizing Trusted Execution Environment (TEE) in recent years, there is no best practice for a confidential data space that combines multiple TEEs:

- There was no way to establish trust between the components forming the confidential data space, preventing untrusted components from entering the data space.
- The connection between arbitrary confidential computing spaces has traditionally been left up to the application, so the platform cannot guarantee end-to-end data confidentiality between them, nor can it take full advantage of the capabilities of APN.

In addition, in the quantum computer era, the risk of a sudden compromise of cryptographic algorithms is expected to increase; therefore, it is desirable that a cryptographic system is not only quantum-safe but also has high crypto-agility.



3.2.2. Design overview

In this configuration, the confidential data space is formed solely by the TEE and devices that can be securely connected to the TEE. Components mutually verify their configurations through IOWN PETs modules, thereby eliminating contamination from untrusted components. In addition, CVMs at both locations are connected to each other via L2/L3 (MACsec/IPsec) encrypted communication, utilizing the MFS architecture between virtual NICs (vNICs) within the CVMs or between containers within the CVMs. The MFS concept is realized by combining PQC algorithms based on different mathematical problems. This configuration allows applications in containers to communicate securely between CVMs without being aware of encrypted communication, taking full advantage of APN performance. Another option for encrypted communication between CVMs is TLS from within the containers. GPU TEE is used to achieve confidentiality on the GPU. GPU TEE refers to a TEE built on a GPU by modifying the GPU hardware design. The modified GPU can ensure its own computing security without relying on CPU TEE [GPU TEE]. While it describes using multiple NICs to clearly represent network segments, this configuration is not mandatory.

Data stored in the Storage should be encrypted, and communication between the Storage and the CVM located in the Gateway server should also be encrypted with encryption protocols such as IPsec and TLS. Furthermore, all switching of encrypted communication between GPU-CVM, CVM-CVM, and Storage-CVM is performed within TEE, which guarantees confidentiality, thereby achieving end-to-end data confidentiality from storage to GPU. Data transfer from Storage to the CVM on the GPU server can use NFS, for example, as described in the figure. For the yellow and green dashed lines in Figure 3.2-1, refer to the interface information described in IOWN PETs architecture [IOWN PETs FA] and MFS architecture [MFS FA], respectively.

Access control to Confidential space is achieved at the VM level (TEE) or by container orchestrators on the CVM, etc. For access control to Storage, the data transfer protocol should be able to control authentication and authorization per data access unit, such as per file.

Overview for building a software security architecture

Confidential VM

- Confidential VM is an X86 system-level TEE that combines the isolation mechanism of virtualization technology to build a TEE with virtual machines as nodes [X86 system-level TEE].

IOWN PETs module

- Software to achieve the I/Fs defined in the "IOWN PETs architecture". In the scope of this document, this software provides an Inter-component I/F for authentication and attestation between components.

MFS CF

- Software, defined in the "MFS architecture", that operates multiple key exchange method modules to derive encryption/decryption keys and provides the keys to a MACsec / IPsec device. It is assumed that either two different types of PQC KEMs are performed, as described in section 3.3.2, or a hybrid of PQC KEMs and QKD is performed, as described in section 3.3.3.

PQC KEM 1/2

- Software implementation of PQC based Key Encapsulation Mechanism (KEM).

MACsec / IPsec

- Software implementation of MACsec/IPsec executable within a confidential VM.

External entity

- This document provides an overview of the external entities that appear in this configuration, although this is outside the scope of this document.

PETs user device

- The device used by a PETs user to input and output data to/from the PETs Space, or to operate and monitor the PETs Space.

PETs manager

- Entities that provide instructions for the lifecycle management of PET components, such as CVM and GPU TEE.

Attestation service

- A service performed by an entity that appraises the validity of Evidence (a set of claimed configuration information) about the platform of TEE and produces Attestation Results to be used by a PETs user.

Authentication service

- A service that issues authentication credentials such as PKI digital certificates and Verifiable Credentials.

3.2.3. Example implementation choices

Table 3.2-1: Example implementation choices for security RIM for protection of the entire data lifecycle

#	Hardware diagram block	Example implementation choice
1	CPU in GPU server and Gate-way server	CPU capable of VM-based TEE.

#	Hardware diagram block	Example implementation choice
2	GPU in GPU server	TEE-enabled GPU which has hardware-based isolation, encrypted communication between GPU and CVM, remote attestation, etc.
3	Classic NIC	10/25/100GbE NIC
4	Storage	<ul style="list-style-type: none"> • 10/25/100GbE NIC • NFSv4

3.3. Another optional RIM regarding subsystem protecting Data in Motion

Ideally, data isolation in this use case should be end-to-end, encompassing the compute space, to ensure the entire data lifecycle is seamlessly protected. However, there may be many cases in which the aforementioned implementation to protect the entire data lifecycle is not always possible due to cost, performance or service usage constraints. Security measures should be considered in light of the risk assessment results, taking into account the overall security requirements and constraints. So, there would be many cases where risk can be eliminated or accepted without implementing all of the aforementioned RIM subsystems. Therefore, in this section, as an optional RIM for a subsystem to protect data in motion, two example RIMs that only protect the communication paths between DCs with post-quantum security using the MFS described as a reference, under the assumption that threats in each DC can be eliminated or accepted by other means.

3.3.1. Design objectives and assumptions

This sub section describes a RIM for subsystem that protects only specific sections of the communication path (i.e., between near endpoints) using MFS architecture for the case where threats within each site or server can be eliminated by a combination of measures such as strict entry/exit/access control, monitoring capabilities, secure boot, etc. or risks of them can be acceptable. In that sense, the RIMs described in this section are not specific to this use case.

The two key points of implementation are as follows:

- Multiple arbitrary post-quantum key exchange algorithms should be able to be used in combination.
- Arbitrary keys generated by the MFS architecture can be provided to a low-layer encryption application without compromising the performance of the APN as much as possible.

3.3.2. Optional RIM structure 1

Figure 3.3-1 shows post-quantum encrypted communication with MFS architecture between Flexible Bridges (See Annex A.2 of [Open APN FA]) installed in front of both GPU node and storage node.

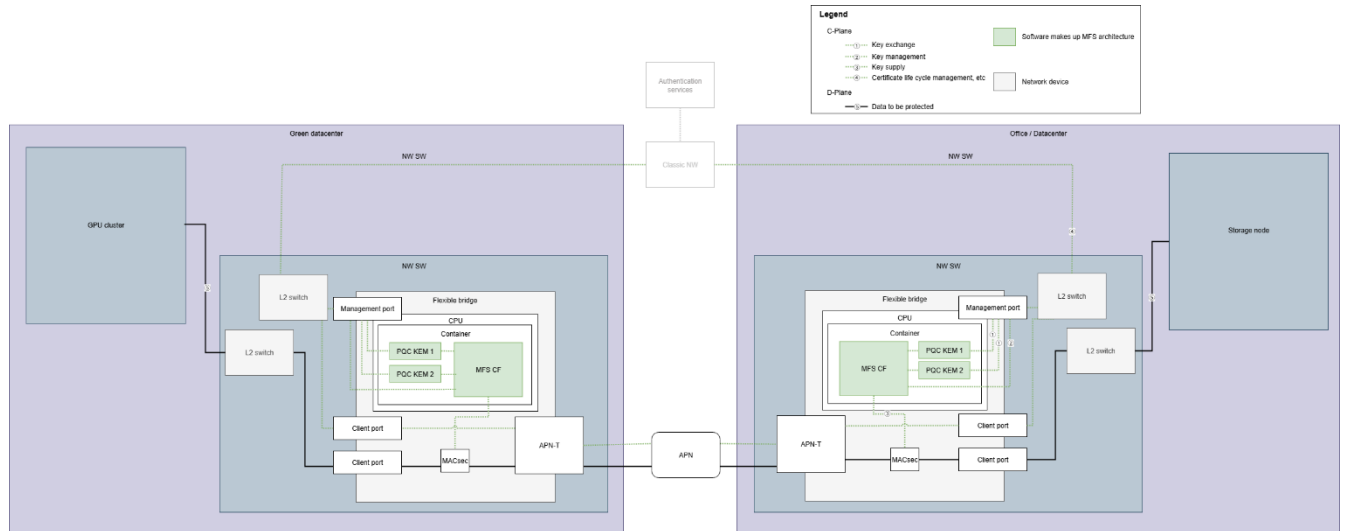


Figure 3.3-1: Structure of optional RIM 1 for protection of data in motion

3.3.2.1. Challenges of existing approaches

Conventional cryptographic systems in network equipment have vertically integrated functions such as key exchange, key management and encryption, which require significant costs for system migration.

3.3.2.2. Design overview

In this configuration, the MFS architecture is implemented in a Flexible Bridge which is a white-box transponder, and post-quantum encrypted communication using the MFS architecture is achieved between Flexible Bridges. Any two different types of post-quantum cryptography (PQC) key encapsulation mechanisms (KEMs) in the MFS Control Function (CF) are performed on the CPU of the Flexible Bridge, and the combined key is supplied to the MACsec hardware accelerator of the Flexible Bridge (or alternatively, the OTNsec can be used).

In Figure 3.3-1, the green dashed lines indicate C-Plane, which realizes the MFS architecture, while the solid black lines represent D-Plane, through which data is transmitted. Details on the MFS interface are provided in the MFS architecture [MFS FA]. In this structure, key exchange between PQC KEM1 and KEM2 is performed over the D-Plane.

The MFS concept is realized from combining PQC algorithms based on different mathematical problems. With this configuration, all data passing through a particular port is encrypted with MACsec using the key supplied by the MFS architecture.



Encryption of communication channels using the MFS architecture is not limited to implementation in a Flexible Bridge or an APN-T. In addition to implementation on a VM as shown in Figure 3.2-1, other implementations are possible, for example, such as implementation on a Smart NIC, or performing key exchange on the VM and offloading only the encryption via MACsec or IPsec to a Smart NIC. In these cases as well, the versatility of the configuration described in this section is not lost.

Overview for achieving a software security architecture

MFS CF

- Software, defined in the "MFS architecture ", that operates multiple key exchange method modules to derive encryption/decryption keys and provides the keys to a MACsec / IPsec device.

PQC KEM 1/2

- Software implementation of PQCbased Key Encapsulation Mechanism (KEM).

3.3.2.3. Example implementation choices

Table 3.3-1: Example implementation choices for security RIM 1 for protection of data in motion

#	Hardware diagram block	Example implementation choice
1	Flexible Bridge	<ul style="list-style-type: none"> • 100G client port • 100G/200G/300G/400G line port • 10G management port • Hardware accelerator for MACsec/OTNsec

3.3.3. Optional RIM 2 structure

This sub section outlines another structure that utilizes quantum key distribution (QKD), a distinct type of post-quantum key exchange method.

Figure 3.3-2 shows post-quantum encrypted communication with MFS architecture between Flexible Bridges installed in front of both the GPU node and the storage node.

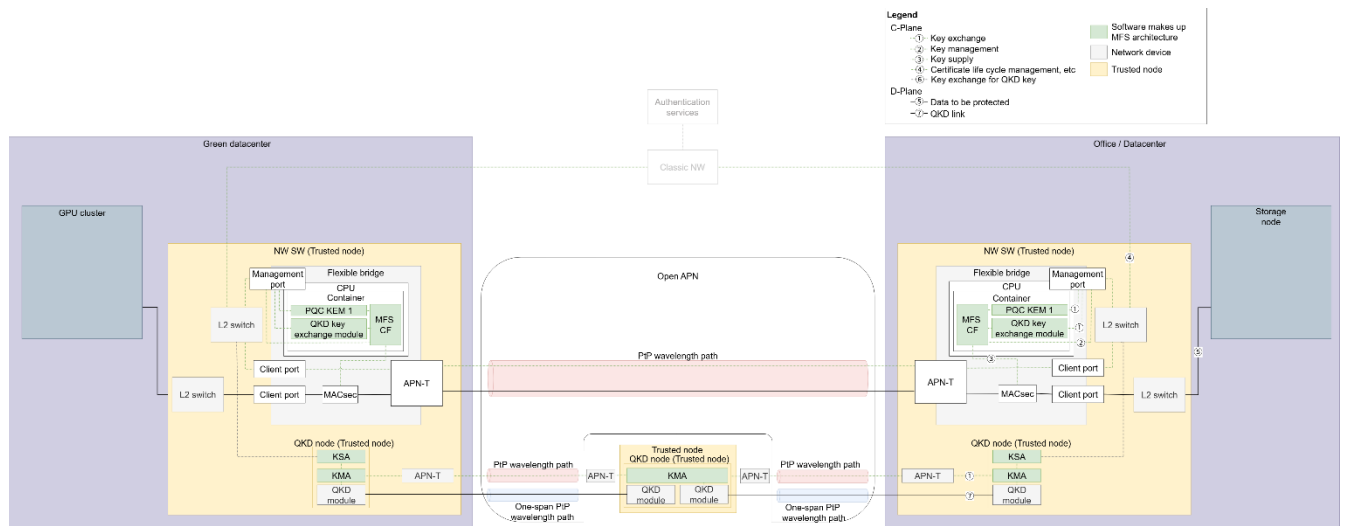


Figure 3.3-2: Structure of security RIM 2 for protection of data in motion

3.3.3.1. Challenges of existing approaches

The security of the PQC-based key exchange algorithm relies on computational difficulty, making it inherently susceptible to sudden compromise. QKD offers another approach to key exchange, with its security guaranteed by information-theoretic principles. However, QKD networks require trusted relays, which inherently introduce additional security risks such as insider threats. To address the limitations of both methods, a cryptographic system with high crypto-agility is desirable through the MFS architecture that combines PQC and QKD.

3.3.3.2. Design overview

In this configuration, the MFS architecture is implemented in a Flexible Bridge which is a white-box transponder Post-quantum encrypted communication using the MFS architecture between Flexible Bridges is achieved via the PtP wavelength path over Open APN [Open APN FA]. Hybrid of key provided by post-quantum cryptography (PQC) key encapsulation mechanisms (KEMs) and quantum key distribution (QKD) in the MFS Control Function (CF) are performed on the CPU of the Flexible Bridge, and the combined key is supplied to the MACsec hardware accelerator of the Flexible Bridge (or alternatively, the OTNsec can be used). With this configuration, all data passing through a particular port is encrypted with MACsec using the key supplied by the MFS architecture.

QKD-based key exchange is enabled through a QKD network deployed over the Open APN. Since QKD signals cannot pass through optical amplifiers, and advanced long-distance techniques like quantum repeaters remain immature, the trusted key relaying scheme is the only practical solution for network-wide sharing of end-to-end keys. End-to-end key sharing across the network is achieved by key generation between a pair of neighboring QKD modules via QKD links over the one-span PtP wavelength path [Open APN FA], and by trusted key relaying between end-to-end key management agents (KMAs) over the APN communication channel provided by the PtP wavelength path. Once the keys are shared, key supply agents (KSAs) at both ends provide them to the respective Flexible Bridges. QKD modules, KMAs, and KSAs must be located

within a QKD node (trusted node), with implementation examples described in Table 3.3-2. This network configuration highlights the potential for low-cost, wide-area QKD deployment without requiring dedicated QKD infrastructure.

Additionally, the diagram assumes a scenario in which the Flexible Bridge is located within the trusted node along with the QKD node. In this setup, no additional security measures are required for the link between the KSA and the Flexible Bridge, demonstrating the feasibility of achieving end-to-end, information-theoretically secure QKD-based key exchange between MACsec encryptors.

This structure can also be implemented on a VM or smart NIC. Even in these cases, the versatility of the configuration described in this section remains intact.

Overview for achieving a software security architecture

QKD key exchange module

- A functional element installed in the Flexible Bridge container that requests keys from the KSA in the QKD network and receives the supplied keys.

Key management agent (KMA)

- A functional element to manage keys generated by one or multiple QKD modules in a QKD node. It consists of key storage, key relay, and key life cycle management functions [ITU-T Y.3802].

Key supply agent (KSA)

- A functional element to supply keys to a cryptographic application, being located between a KMA and the cryptographic application. It consists of a key supply function and optionally, a key combination function [ITU-T Y.3802]

To avoid redundancy, MFS CF and PQC KEM, as described in the RIM structure 1 are not repeated here.

3.3.3.3. Example implementation choices

Table 3.3-2: Example implementation choices for security RIM 2 for protection of data in motion

#	Hardware diagram block	Example implementation choice
1	QKD node (Trusted node)	To protect the information assets and the QKDN physical entity assets inside a QKD node against the security threats such as unauthorized physical access and unauthorized network access, some implementations examples of security measures [ITU-T X.1713] are listed up as follows.

#	Hardware diagram block	Example implementation choice
		<ul style="list-style-type: none"> • Access control: physical access to trusted node facilities (e.g. racks, server rooms) is strictly controlled through multi-factor authentication using an IC card and biometric authentication. • Surveillance camera: unauthorized physical access can be prevented through remote surveillance. Surveillance images are stored for a certain period of time and used for forensic investigations. • Firewall: a firewall establishes a barrier between a trusted, secure internal network and another external network that is assumed not to be secure and trusted. It relays all communications between trusted nodes, and the trusted node and the Flexible Bridge, and monitors and controls the incoming and outgoing network traffic on the communications. • Intrusion detection and prevention systems: unusual behaviors caused by malicious activities in the QKD network or the trusted node system can be monitored, identified and reported.
2	QKD module	<ul style="list-style-type: none"> • QKD protocol: BB84, Continuous-variable (CV) etc. • QKD link: Quantum channel, Clock channel, Classical channel (10GbE etc.)



4. Performance Considerations

This section discusses the impact of the RIM described in the previous section on training time and storage access speed, which are key benchmarks for the Remote GPU use case.

In general, when protecting data with encryption technology, performance degradation due to encryption processing overhead is unavoidable. In the RIM described in this document, performance degradation occurs due to two major sources of overhead: communication channel encryption and TEE memory encryption.

Regarding the encryption of communication channels between CVMs, to maintain confidentiality in E2E, encryption processing cannot be offloaded to the NIC in the current product; so therefore, encryption processing is performed on the CPU. Although the impact on throughput depends on network conditions and CPU performance, means of encrypted communication (i.e., MACsec, IPsec or TLS) should be selected not only based on performance but also according to operational requirements. On the other hand, regarding the communication channel encryption between Flexible bridges, described as a subsystem option for the communication part in section 3.3, since a hardware accelerator performs the encryption process at L2, the impact of this RIM on latency and throughput is expected to be negligible for many services that use Remote GPU.

Regarding the use of TEE and TEE-enabled GPU, overhead is incurred due to memory encryption and encrypted communication between the CVM and the GPU. There is a report that disk I/O and network I/O performance may degrade by up to 30–50% in VM-based TEE due to memory encryption [Performance of CVM]. Additionally, when using TEE-enabled GPUs, although the performance of encryption processing depends on each CPU and GPU product, if the ratio of GPU computation to input data volume is low, using a TEE-enabled GPU is said to degrade machine learning training time by approximately 40% [CC on GPU].



5. Conclusion

This document describes RIM, which implements IOWN PETs and MFS architectures for the Remote GPU use case where sensitive data is handled.

The RIM demonstrates that a seamless data protection mechanism with post-quantum security, extending from storage to GPU can be realized using the IOWN PETs architecture and MFS architecture with existing hardware. As future work, the detailed performance of this RIM needs to be clarified through PoC. In addition, since the optimal configuration will change with advances in hardware, it is necessary to continue studying configurations that are expected to further improve performance.

While the RIM described as an option of a subsystem for the communication part is not specific to the Remote GPU use case, it can be widely referenced as an implementation example of applying MFS for communication channel encryption between DCs.



6. References

[IOWN PETs FA]Functional Architecture for Protection of Data in Use: IOWN Privacy Enhancing Technologies, https://iowngf.org/wp-content/uploads/2025/02/IOWN-GF-RD-PETs_Functional_Architecture-1.0.pdf

[MFS FA]Functional Architecture for Protection of Data in Motion: Multi Factor Security Key Exchange and Management, https://iowngf.org/wp-content/uploads/2025/02/IOWN-GF-RD-MFS_Functional_Architecture-1.0.pdf

[Open APN FA]Open All-Photonic Network Functional Architecture, <https://iowngf.org/wp-content/uploads/2025/06/IOWN-GF-RD-Open-APN-Functional-Architecture-3.1.pdf>

[ITU-T Y.3802]ITU-T, Quantum key distribution networks - Functional architecture, 7 December, 2020, <https://www.itu.int/rec/T-REC-Y.3802-202012-l/en>

[ITU-T X.1713]ITU-T, Security requirements for the protection of quantum key distribution nodes, 29, April, 2024, <https://www.itu.int/rec/T-REC-X.1713-202404-l/en>

[Remote GPU UC RIM PoC Ref]Reference Implementation Model and Proof-of-Concept Reference of Green Computing with Remote GPU, https://iowngf.org/wp-content/uploads/2025/03/IOWN-GF-RD-Remote_GPU_Use_Case_RIM_PoC_Ref-1.0.pdf

[GC with Remote GPU UC]Green Computing with Remote GPU Service for Generative AI / LLM Use Case - Light Speed Data Transfer for AI Training -, https://iowngf.org/wp-content/uploads/2025/02/IOWN-GF-RD-GC_with_Remote_GPU_Use_Case-1.0.pdf

[GPU TEE]Survey of research on confidential computing - Feng - 2CC on GPU024 - IET Communications - Wiley Online Library, <https://ietresearch.onlinelibrary.wiley.com/doi/full/10.1049/cmu2.12759>

[X86 system-level TEE]Survey of research on confidential computing - Feng - 2024 - IET Communications - Wiley Online Library, <https://ietresearch.onlinelibrary.wiley.com/doi/full/10.1049/cmu2.12759>

[IETF Hybrid key exchange in TLS 1.3 2024]IETF, Hybrid key exchange in TLS 1.3, draft-ietf-tls-hybrid-design-11, 7 October 2024, <https://datatracker.ietf.org/doc/draft-ietf-tls-hybrid-design/11>

[NIST Getting Ready for Post-Quantum Cryptography 2021]NIST Cybersecurity White Paper, Getting Ready for Post Quantum Cryptography: Exploring Challenges Associated with Adopting and Using Post-Quantum Cryptographic Algorithms, April 28, 2021 <https://nvlpubs.nist.gov/nistpubs/CSWP/NIST.CSWP.04282021.pdf>

[Performance of CVM]Performance Overheads of Confidential Virtual Machines, <https://kartikgopalan.github.io/publications/mascots23.pdf>



[CC on GPU]NVIDIA, Confidential Compute on NVIDIA Hopper H100,

<https://images.nvidia.com/aem-dam/en-zz/Solutions/data-center/HCC-Whitepaper-v1.0.pdf>



History

Revision	Release Date	Summary of Changes
1	October 2025	Initial Release