# Computing Infrastructure for the AI Era

Sachiko Oonishi

Executive Member of the Board, Executive Vice President, CCXO and Co-CAIO
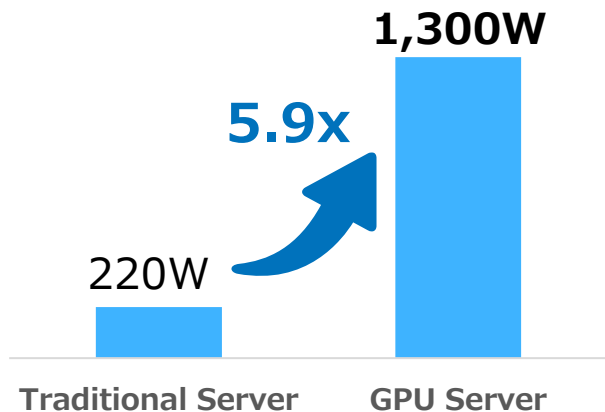Head of Research and Development Market Strategy
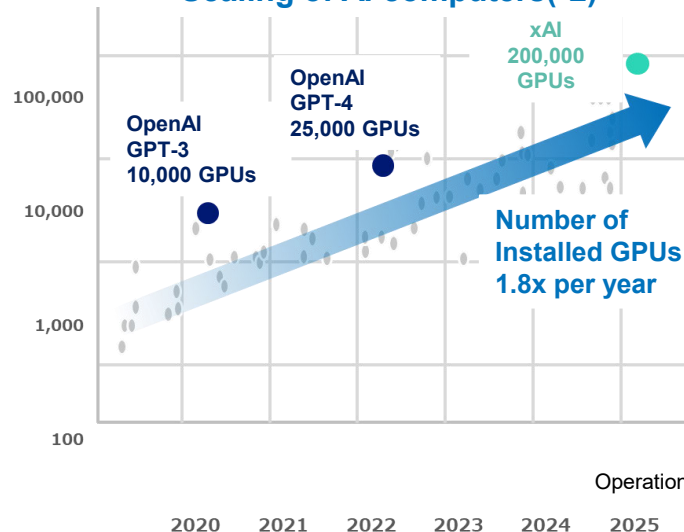NTT, Inc.

Oct. 6th, 2025

# The Rapid Surge in Power Consumption Due to the Wider Adoption of AI

- GPU servers consume 5.9 times the power of traditional servers.
- GPU installations are growing 1.8x annually, driving continuous power consumption increases.

**Power consumption of a server (＊1)**



**1,300W**

**5.9x**

**220W**

**Traditional Server**   **GPU Server**

**Scaling of AI computers(*2)**



**xAI 200,000 GPUs**

**OpenAI GPT-4 25,000 GPUs**

**OpenAI GPT-3 10,000 GPUs**

**Number of Installed GPUs 1.8x per year**

100,000

10,000

1,000

100

Operation start date
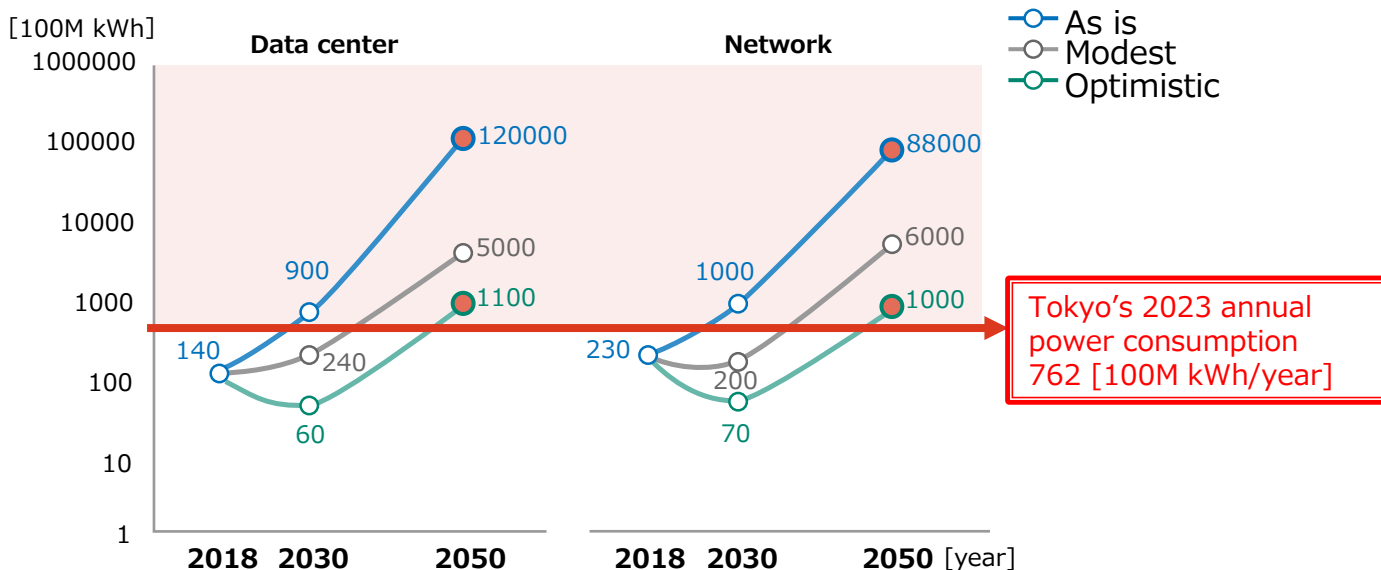
2020   2021   2022   2023   2024   2025

Sources:
*1 Created by NTT from EY Japan "Requirements of data centers of GPU servers"
*2 Epoch AI

# The Rapid Adoption of AI Is Fueling a Power Crisis ⦿NTT

- At this rate, data center power will surpass Tokyo's 2023 annual consumption.

## Forecast of data center power consumption



[100M kWh]

Data center — As is / Modest / Optimistic

Data center: 2018: 140, 2030: 900 / 240 / 60, 2050: 120000 / 5000 / 1100

Network: 2018: 230, 2030: 1000 / 200 / 70, 2050: 88000 / 6000 / 1000

Tokyo's 2023 annual power consumption 762 [100M kWh/year]
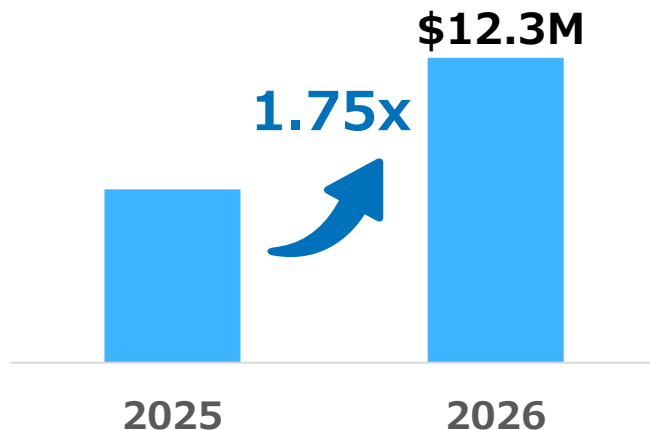
2018  2030  2050 [year]

Source: Subcommittee on Basic Policy, Advisory Committee for Natural Resources and Energy (56th Meeting)
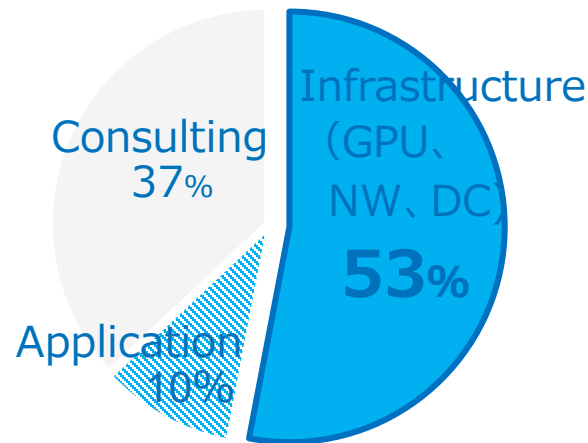https://www.enecho.meti.go.jp/committee/council/basic_policy_subcommittee/2024/056/056_005.pdf

# As AI Use Expands, User Costs Are Also Surging  ⵔ NTT

- As AI adoption grows, surging utilization costs are fueling concerns about declining ROI.
- The infrastructure domain will account for over half the AI market.

**Enterprise AI adoption costs (＊1)**
**(Average Annual LLM Budget per Company)**

**$12.3M**

**1.75x**

2025          2026

**Forecast of**
**AI market segmentation, 2027(*2)**

Consulting
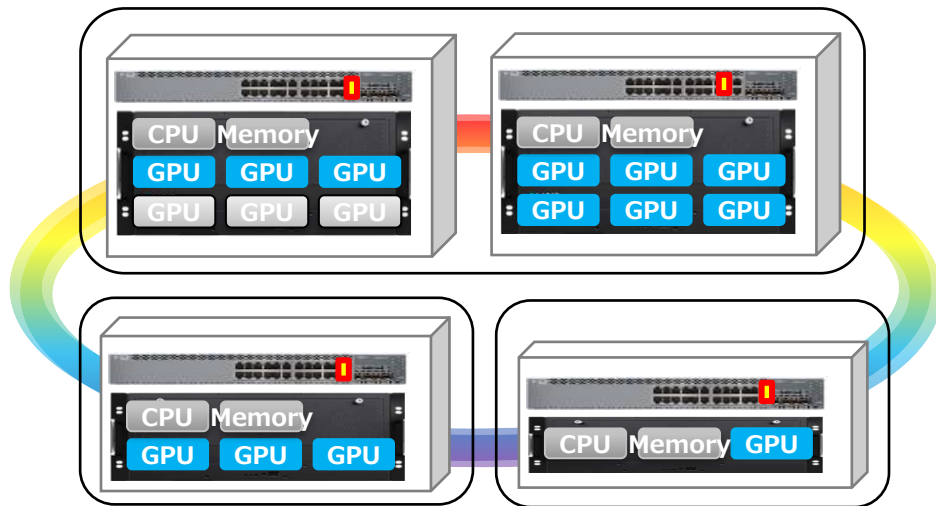37%

Infrastructure
(GPU、NW、DC)
**53%**

Application
10%

Sources:
*1 Andreessen Horowitz "How 100 Enterprise CIOs Are  Building and Buying Gen AI in 2025"
*2 Created by NTT based on multiple reports from various sources, including Fuji Chimera Research Institute,
    Yano Research Institute, and Deloitte Touche Tohmatsu MIC Research Institute.

3

# Requirements for Infrastructure in the AI Era

- Efficient operation and reduced energy consumption are essential for AI computing infrastructure.
- IOWN is a crucial component for balancing AI benefits with efficient operation and reduced power consumption.
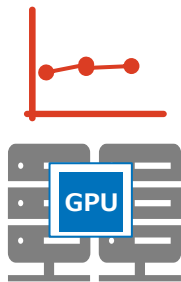


Efficient Infrastructure Operations

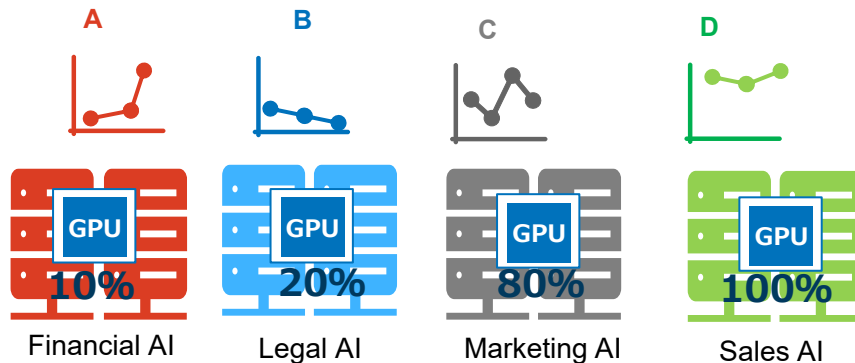Reducing Power Consumption of Infrastructure Equipment

# Wider AI Adoption is Reshaping Compute Utilization

- As AI use expands across business processes, fluctuating compute utilization drives up power and costs.

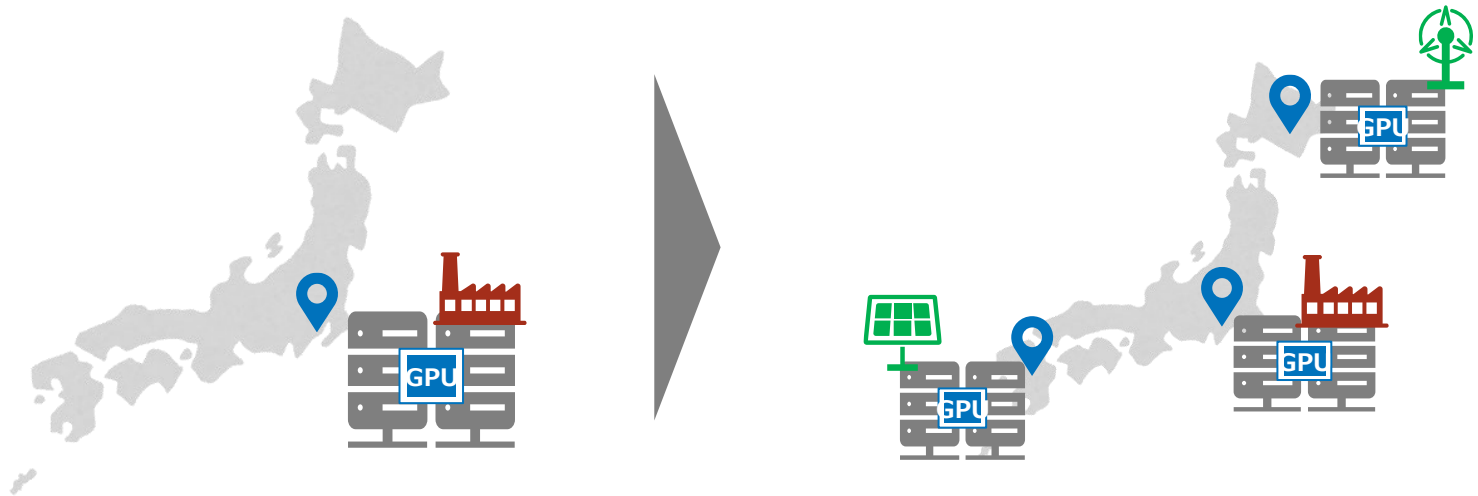Previously: Limited AI applications led to high GPU utilization for specific tasks

The expansion of AI and task diversity is causing fluctuating GPU utilization



| A | B | C | D |
|---|---|---|---|
| GPU | GPU | GPU | GPU |
| 10% | 20% | 80% | 100% |
| Financial AI | Legal AI | Marketing AI | Sales AI |

✓ **Only 50% of the total resources are being utilized.**
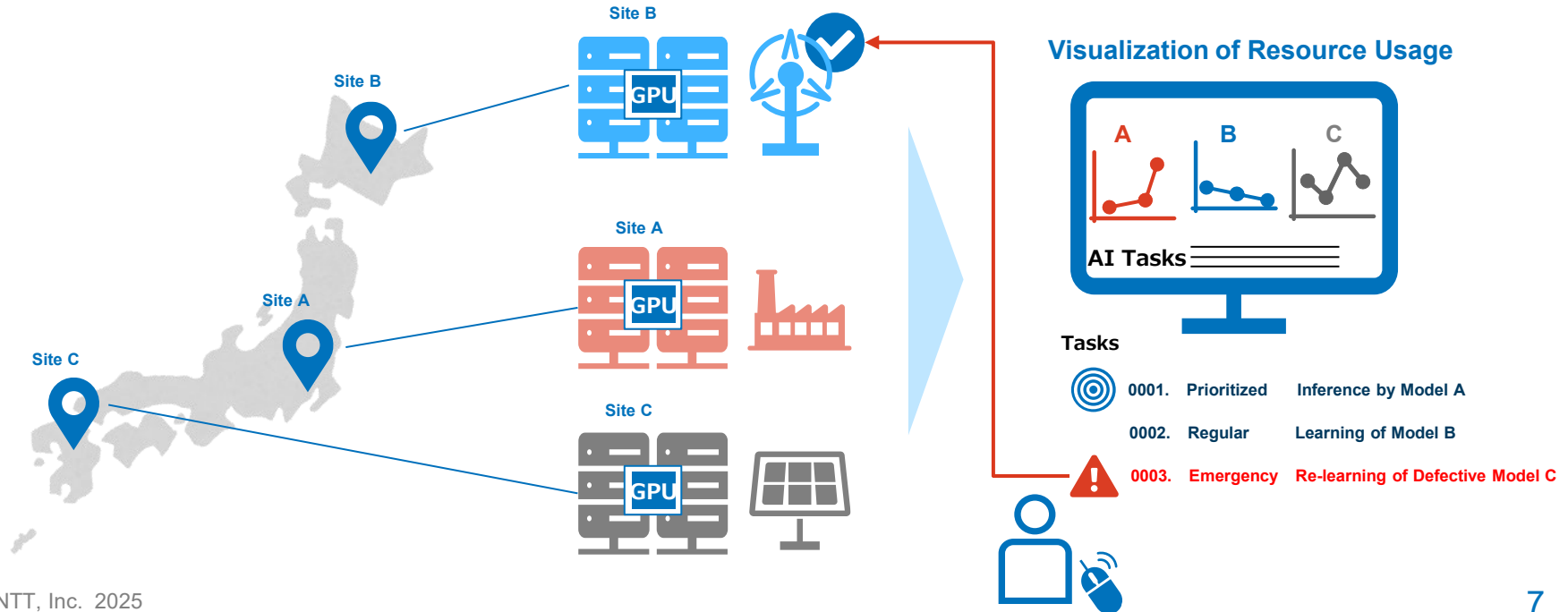**⇒ Optimizing scattered resources via GPUaaS can halve hardware investment.**

# The Reality of Compute Infrastructure in the AI Era  ⏻NTT

- Growing infrastructure demand, led by AI, is maxing out power supply in the Tokyo Metropolitan Area.
- Data centers and GPUs are increasingly being distributed to regions with power surpluses.

# Infrastructure Optimization Initiatives at Leading Companies
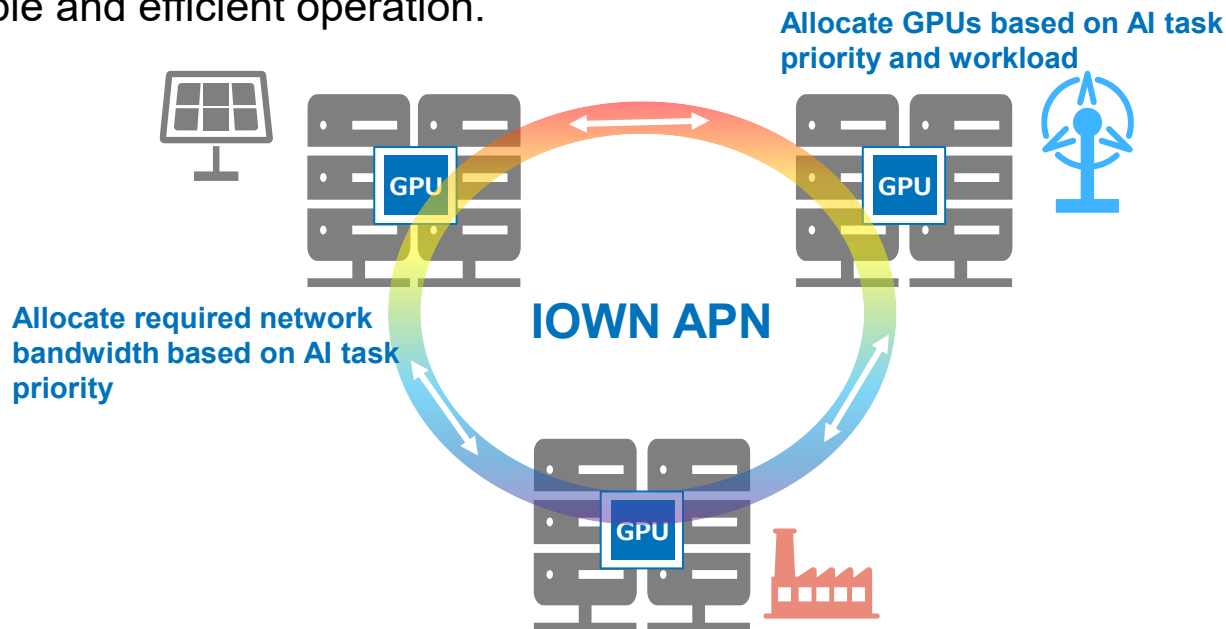
- Leading companies are starting to efficiently operate infrastructure by visualizing resource status (e.g., GPU usage and power consumption) and prioritizing AI workloads in resource-rich locations.

# Infrastructure Requirements in the AI Era:
# 1. Efficient Infrastructure Operations

- High-speed APN links to connect multi-tiered GPUs and distributed GPUs.
- Flexible allocation of compute resources, such as GPUs, based on AI application and utilization.
- Dynamically assign AI tasks to sites with sufficient power, based on application and workload, ensuring stable and efficient operation.



**Allocate GPUs based on AI task priority and workload**

**IOWN APN**

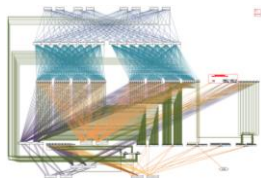**Allocate required network bandwidth based on AI task priority**

# As AI Utilization and Data Volume Increase, Electrical Processing Is Reaching Its Limit

- Increased AI processing → More GPUs → Greater intra-computer communication
- Increased power consumption and heat generation present a limit for electrical communication → Transition to optical communication.
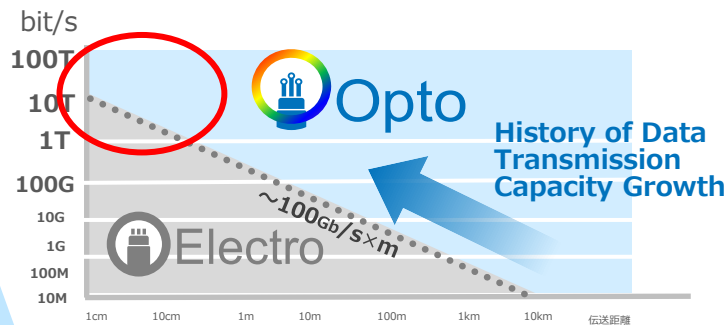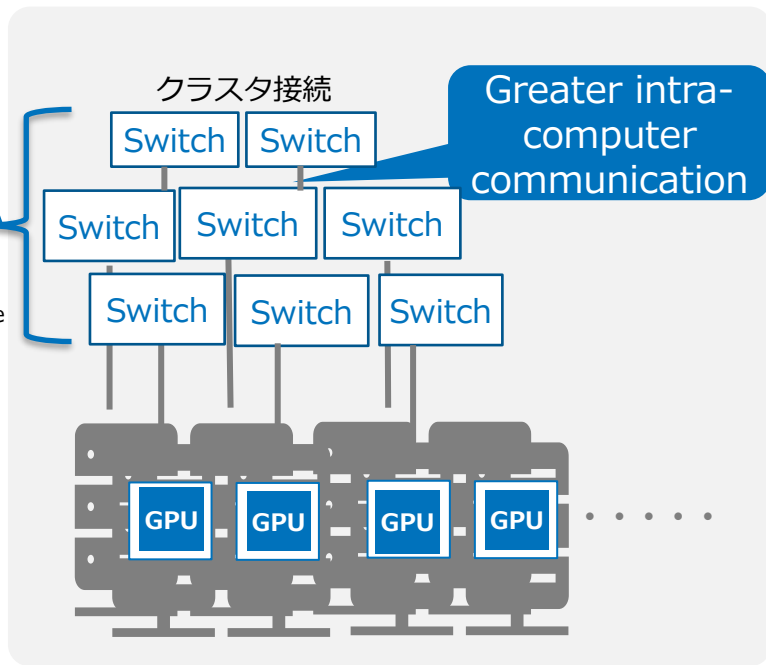
Increased GPU processing drives higher data transfer volumes.

NTT DATA's GPUaaS

The computing infrastructure for developing NTT's LLM, "tsuzumi"

High-volume intra-computer communication

クラスタ接続

Switch Switch
Switch Switch Switch
Switch Switch Switch

GPU GPU GPU GPU

Greater intra-computer communication



History of Data Transmission Capacity Growth

$\sim 100\text{Gb/s}\times\text{m}$

Opto

Electro

bit/s
100T
10T
1T
100G
10G
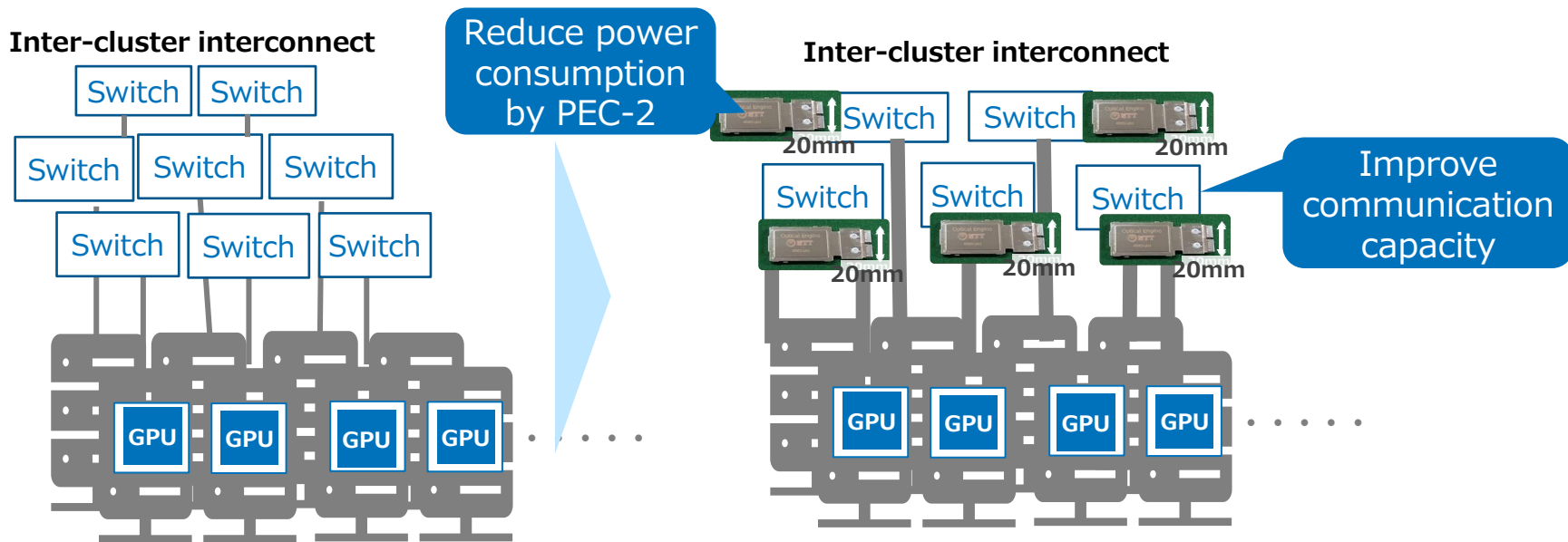1G
100M
10M

1cm  10cm  1m  10m  100m  1km  10km  伝送距離

IOWN's Photonics-Electronics Convergence devices can suppress heat generation and power consumption even with high-capacity communication.

9

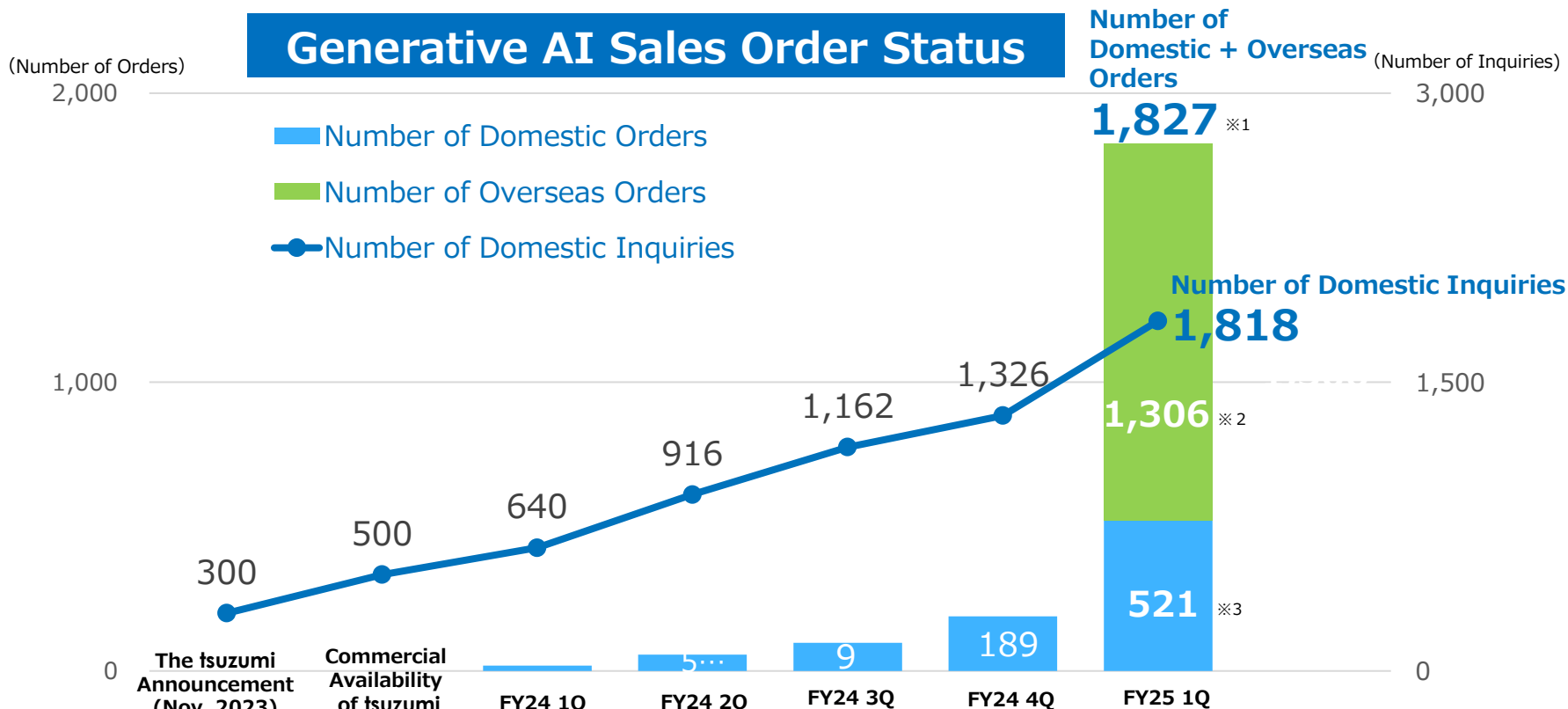# Infrastructure Requirements in the AI Era: 2. Reducing Power Consumption

- By implementing Photonics-Electronics Convergence devices (PEC-2) in network switches, we can achieve both increased communication capacity and reduced power consumption.
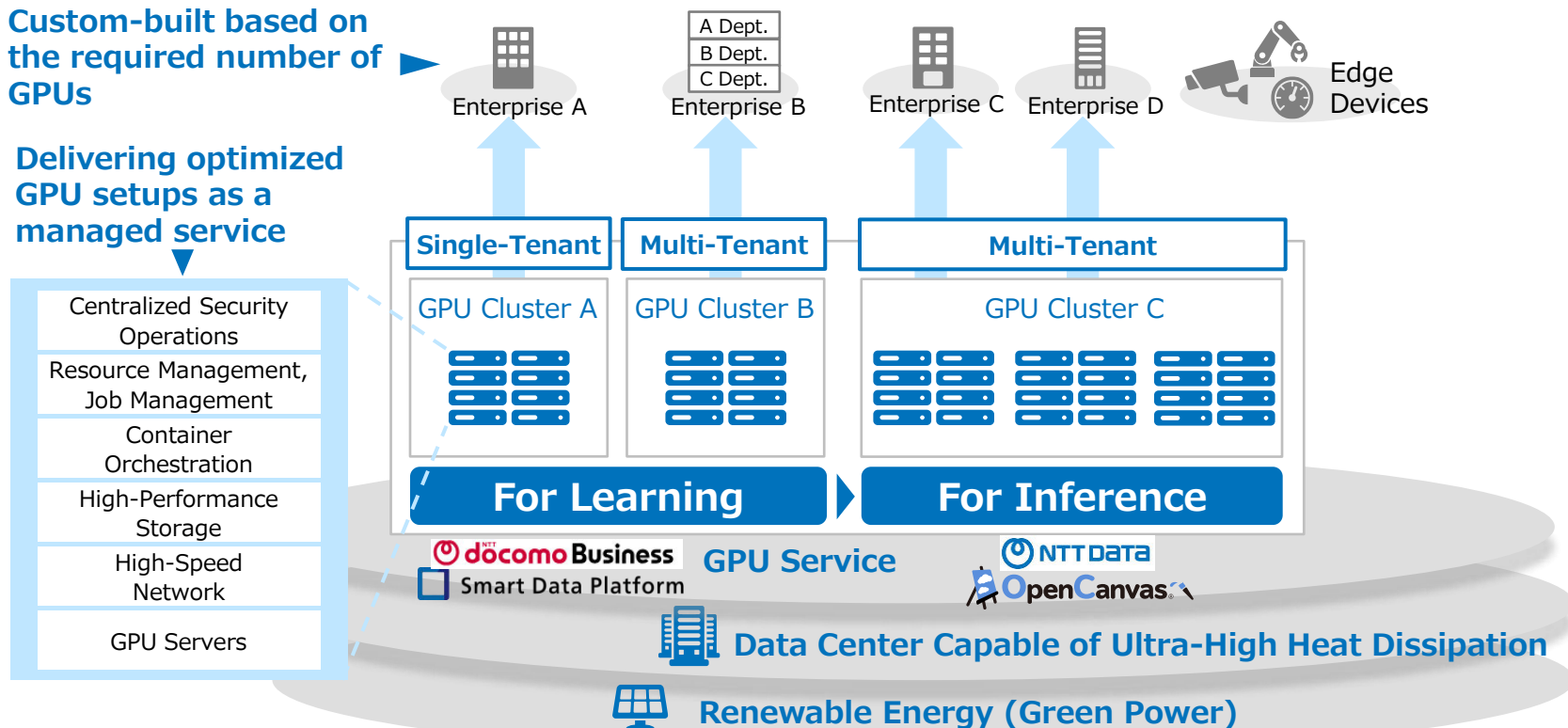
# AI Adoption Is Progressing Steadily

**Generative AI Sales Order Status**

(Number of Orders)

**Number of Domestic + Overseas Orders**
(Number of Inquiries)

- Number of Domestic Orders
- Number of Overseas Orders
- Number of Domestic Inquiries

**1,827** ※1

**Number of Domestic Inquiries**
**1,818**

**1,306** ※2

1,326

1,162

916

640

500

300

**521** ※3

189

9

5···

The tsuzumi Announcement (Nov. 2023) | Commercial Availability of tsuzumi (Mar. 2024) | FY24 1Q | FY24 2Q | FY24 3Q | FY24 4Q | FY25 1Q

※1 Total Generative AI Inquiries and Orders: Compiled from NTT Group companies (Docomo Business, East, West, and DATA).
※2 Overseas order count includes only the cumulative total for Q1 2025.
※3 Including pre-approved projects.

11

# NTT's GPU as a Service: Delivering Compute Resources When and Where Needed, Based on AI Utilization

**Custom-built based on the required number of GPUs** ▶

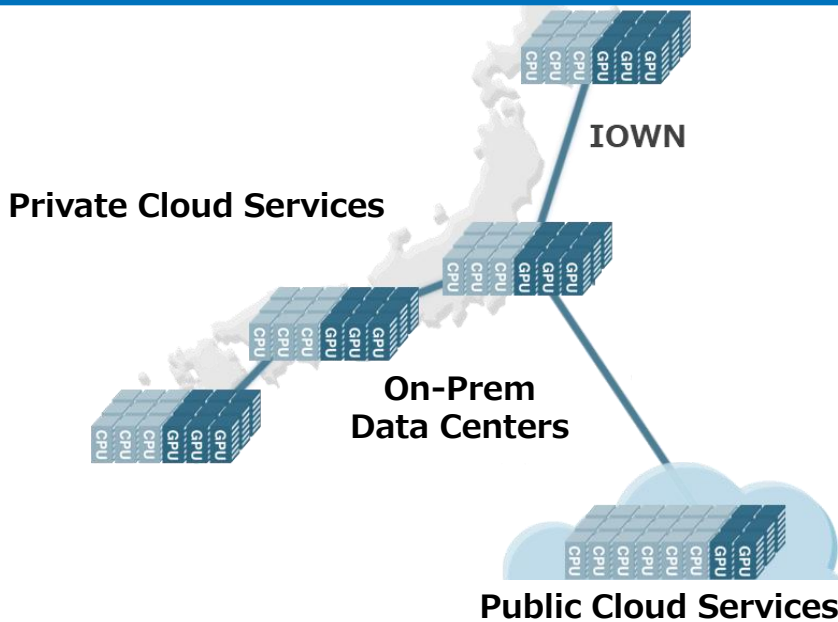**Delivering optimized GPU setups as a managed service** ▼

| Enterprise A | Enterprise B | Enterprise C | Enterprise D | Edge Devices |
|---|---|---|---|---|
| | A Dept. / B Dept. / C Dept. | | | |

| Centralized Security Operations |
|---|
| Resource Management, Job Management |
| Container Orchestration |
| High-Performance Storage |
| High-Speed Network |
| GPU Servers |

**Single-Tenant** | **Multi-Tenant** | **Multi-Tenant**

GPU Cluster A | GPU Cluster B | GPU Cluster C

**For Learning** ▶ **For Inference**

docomo Business — Smart Data Platform

**GPU Service**

NTT DATA — OpenCanvas

**Data Center Capable of Ultra-High Heat Dissipation**

**Renewable Energy (Green Power)**

# Computing Infrastructure for the AI Era

- Efficiently managing compute resources to boost GPU utilization and reduce overall power consumption.
- Introducing Photonics-Electronics Convergence devices enables further reduced power consumption by running the same infrastructure with optimized GPUs and equipment.

## Future Computing Infrastructure



## Infrastructure Power Consumption



Efficient operations

PEC

2020　　　　2030　　　　2040 [year]

14

# Innovating a Sustainable Future for People and Planet