

LINE Yahoo! US Data Center Technology at the Forefront LLM (Large Scale Language Models) and the Challenge of Water Cooling Technology

Norifumi Matsuya

# Self-introduction



# **Self-introduction in JAPAN**

- Norifumi Matsuya
- From 2000 to 2014,
  I was in charge of infrastructure design, construction and operation at Yahoo!
- Main Achievements
  - BGP in production network
  - Obtained AS from JPNIC, APNIC, and ARIN (AS23816, AS24572, AS18140)
  - Backbone network virtualization and VPLS
  - Creating Japan's largest Openstack cluster
  - Construction of domestic data centers (Tokyo, Osaka, Chubu region, Kyushu region, Tohoku region)





US Data Center Tour '09

# **Self-introduction in the US**



Actapio Data Center

- Moved to the US in 2014 and seconded to Actapio
- Actapio is a US subsidiary of LINE Yahoo Japan Corporation, and a company that builds data centers and operates servers for AI and big data.
- Currently, while running a company, I am in charge of OCP (Open Compute Project) server operation and AI business development.

# AGENDA

- Today's announcements' overview
- Part 1: Water-cooling technology can solve the problem of insufficient airflow from air conditioners
  - Emergence of LLM and changes in server resources
  - Current state of data centers and predicted air conditioner air volume problems
  - Considering solutions to insufficient airflow from air conditioners
- Part 2: Server proximity constraints that can be solved with liquid cooling technology
  - The emergence of GB200 NVL72 (GPU server) and its background
  - Change power supply to 480V (under investigation)

#### summary

# Today's announcements' overview





- ChatGPT is a conversational LLM (large-scale language model)
- LLM uses a transformer model and is trained on large amounts of text data
- Transformer is a neural network architecture that demonstrates high performance in natural language processing



- This results in a significant increase in server workload.
- Check each resource in the data center and confirm that there is no problem
- Obtained information from a U.S. company about insufficient airflow in data center air conditioners



- water-cooling technology can allow
  - Improved cooling efficiency
  - Solves the problem of insufficient airflow
  - More servers can be installed



- In March 2024, NVIDIA will announce a new GPU server, the GB200 NVL72, at the GTC (GPU Technology Conference).
  - Water-cooling technology is essential for the GB200 NVL72
- Interview with a US company on why water cooling technology is necessary
  - The reason is that **high-bandwidth communication** between GPUs is required to shorten the model training time.





# Emergence of LLM and changes in server resources





# November 2022 ChatGPT 3.5 release



# November 2022 ChatGPT 3.5 release

How is it different from previous AI?

# Conventional "tabular data learning" and Differences in "Transformer Learning for LLM"

Features	Tabular format data	LLM (Large Scale Language Models)	
Data Type	Tabular data (e.g. Excel, CSV)	Text data	
Special skill	Numerical prediction, classification, data analysis	Text generation, translation, chatbots	
Training Data	Structured data (numerical, categorical)	Large amount of document data	
Model Scale	Relatively small scale (depends on features and number of data)	Large scale (billions of parameters)	
Flexibility	Specialized for specific datasets	Capable of handling a wide variety of tasks	

# Conventional "tabular data learning" and Differences in "Transformer Learning for LLM"

Features	es Tabular format data		LLM (Large Scale Language Models)
It requires more data compared to the previous AI.		CSV)	Text data
		sification, data	Text generation, translation, chatbots
		II, ca. (31)	Large amount of document data
Model Scale	Relatively small scale (depends on features and number of data)		Large scale (billions of parameters)
Flexibility	Specialized for specific datasets		Capable of handling a wide variety of tasks

LLM also available for Yahoo! JAPAN services



- Yahoo! JAPAN uses LLM to "understand language" in user posts
- Significant **improvement in accuracy** achieved by using LLM



However, there are also major challenges

LLM also available for Yahoo! JAPAN services



- Yahoo! JAPAN uses LLM to "understand language" in user posts
- Significant **improvement in accuracy** achieved by using LLM



However, there are also major challenges

# Changes in Server Workloads (with A100 GPU)

#### Tabular format data Workload 30%-80% for about a week



#### LLM Workload 80%-100% continue for 1-2 months







# Current status of data centers and predicted air conditioner air volume problems



Network: Comparison of server bandwidth between LLM and another system (Hadoop)

#### Hadoop server interface graph

	allaltet er dek and
ى الى الى الى الى الى الى الى الى الى ال	

LLM server interface graph



	Server port bandwidth capacity	Usage situation
Hadoop server interface	25Gbps	13.50 Gbps(54%)
LLM Server Interface	100 Gbps	23.5Gbps ( 24% )

- There is enough bandwidth and no errors or drops.
- No network issues have occurred
- \* Some ingenuity when creating the model also played a role.

Power: Comparison of server power consumption between LLM and another system (Hadoop)

Hadoop Server (Maximum operating rate: 76.12%)



#### LLM Server (Maximum operating rate: 95.44%)





Server workloads remain high for long periods, causing **increased power consumption** 

### Cooling: Temperature comparison between LLM and another system (Hadoop)



In terms of temperature, there is no significant difference between the LLM and the Hadoop rack rows. There are no temperature issues

# Currently no issues with network, power and cooling



Discussion with US company about water cooling technology

- GPU TDP (Thermal Design Power) and other factors are on the rise
- Increased TDP maximizes server fan speed
- There is a concern that the increased airflow from servers will lead to an increase in the airflow from air conditioners.

# Currently no issues with network, power or cooling



Discussion with US company about water cooling technology

- GPU TDP (Thermal Design Power) and other factors are on the rise
- Increased TDP maximizes server fan speed
- There is a concern that the **increased airflow** from servers will lead to an increase in the airflow from air conditioners.

Air conditioning system and cross-section of Actapio data center



- Increased airflow in server rack in ①
- An increase in pressure in the hot aisle in ② and a decrease in pressure in the cold aisle in ③ occured
- Sufficient airflow from the air conditioner is required to exhaust and discharge the air to areas 2 and 3.
- If the air conditioner does not have enough airflow, the pressure balance inside the data center cannot be maintained and temperature control becomes impossible.



For data center air conditioners to operate without problems, airflow **of the air conditioner** must be more than the server airflow

Server airflow



Air conditioner air volume





Measuring the airflow of an air-cooled GPU server

- Servers are measured at full load
- Surround the rear of the server with a plastic wind shield
- The cross section was divided into 4x4 sections with string, and each section was measured with an anemometer



\* Measurement method: JIS A 1431 (https://kikakurui.com/a1/A1431-1994-01.html)



# Actapio's 2MW server room Scenario for fully installed air-cooled GPU servers for LLM

Total airflow for 2MW of air-cooled GPU servers



# Actapio's 2MW server room Scenario for fully installed air-cooled GPU servers for LLM







# Considering solutions to insufficient airflow from air conditioners



### Solutions to insufficient airflow from air conditioners

**Option 1: Limiting the number of installed servers** 

### **Option 2: Increase air conditioning capacity**

- a. Addition of air conditioners
- b. Increasing air conditioner specifications

### **Option 3: Consideration of water cooling method: DLC (Direct Liquid Cooling)**

- a. L to L: Liquid to Liquid
- b. AALC: Air Assisted Liquid Cooling

# Option 1: Limiting the number of installed servers



For a 2MW server roomAir conditioner's maximum airflowNumber of air-cooled GPU serversNumber of units that can be installed



Option 2: Increase air conditioning capacity

### a.Additional air conditioners

The layout of the data hall (see below) requires additional physical capacity.





Option 2: Increase air conditioning capacity

b. Increasing the specifications of air conditioners Upgrading the specifications of existing air conditioners is costly and difficult, including the construction required



Option 3: Consider water cooling method

Water cooling methods (DLC) can be classified into two types

DLC (Direct Liquid Cooling)

- 1. L to L: Liquid to Liquid
- 2. AALC: Air Assisted Liquid Cooling

Among the water cooling methods, we are considering AALC and we will explain why


## Heat absorption

- **①The cold coolant** is sent to the **②Cooling plate**
- **2Cooling plate 3absorbs heat** from the processor
- **(1)** The warm coolant is sent to the heat exchanger



## Heat absorption

- **1**The cold coolant is sent to the **2**Cooling plate
- **2Cooling plate 3absorbs heat** from the processor
- **(1)** The warm coolant is sent to the heat exchanger



#### Heat absorption

39

- **①The cold coolant** is sent to the **②Cooling plate**
- Ocooling plate Obsorbs heat from the processor
- **(1)** The warm coolant is sent to the heat exchanger



## Heat absorption

40

- **①The cold coolant** is sent to the **②Cooling plate**
- **2Cooling plate 3absorbs heat** from the processor
- **(1)** The warm coolant is sent to the heat exchanger



- **①**Chillers and cooling towers are used to make **②**chilled water
- 27The cold water exchanges heat with the cooling liquid in the server room at 32LtoL
- **4** The cold liquid is sent to the server's cooling plate
- **5**The warm coolant is returned to the **4**CDU, where it exchanges heat with the cold water.
- **6**The hot water is returned to **1**the chiller or cooling tower to produce **2**cold water again.



- **Orbillers and cooling towers are used** to make
- 27 The cold water exchanges heat with the cooling liquid in the server room at 30 LtoL
- **4** The cold liquid is sent to the server's cooling plate
- **5**The warm coolant is returned to the **4**CDU, where it exchanges heat with the cold water.
- **6**The hot water is returned to **1**the chiller or cooling tower to produce **2**cold water again.



- **OChillers and cooling towers are used** to make
- **2The cold water** exchanges heat with the cooling liquid
- **(1)** The cold liquid is sent to the server's cooling plate
- **6**The warm coolant is returned to the **4**CDU, where it exchanges heat with the cold water.
- **6**The hot water is returned to **1**the chiller or cooling tower to produce



- **OChillers and cooling towers are used** to make
- **2**The cold water exchanges heat with the cooling liquid
- **4** The cold liquid is sent to the server's cooling plate
- **5**The warm coolant is returned to **3**LtoL and exchanges heat with the cold water.
- **6**The hot water is returned to **1**the chiller or cooling tower to produce



- **OChillers and cooling towers are used** to make
- **2**The cold water exchanges heat with the cooling liquid
- **(1)** The cold liquid is sent to the server's cooling plate
- **5**The warm cooling liquid is returned to **3**LtoL and exchanges heat with the cold water.
- **6**The hot water is returned to **1**the chiller or cooling tower to produce



- **OChillers and cooling towers are used** to make
- **2**The cold water exchanges heat with the cooling liquid
- **(1)** The cold liquid is sent to the server's cooling plate
- **5**The warm cooling liquid is returned to **8**LtoL and exchanges heat with the cold water.
- **6**The hot water is returned to **1**the chiller or cooling tower to produce

Source: Dell Technologies Direct Liquid Cooling Support for New PowerEdge Servers | Dell Technologies Info







- Air conditioner 2 creates cool air
- **2**The cold air exchanges heat with the cooling liquid in
- **The cold liquid** is sent to the server's cooling plate
- **5**The warm cooling liquid is returned to **8**AALC and exchanges heat with the cold air.
- **6**The warm air is returned to **1**the air conditioner, where it again creates



- Air **conditioner 2** creates cool air
- **2**The cold air exchanges heat with the cooling liquid in
- **(1)** The cold liquid is sent to the server's cooling plate
- **5**The warm cooling liquid is returned to **3**AALC and exchanges heat with the cold air.
- **6**The warm air is returned to **1**the air conditioner, where it again creates



- Air conditioner 2 creates cool air
- 2The cold air exchanges heat with the cooling liquid in
- **(1)** The cold liquid is sent to the server's cooling plate
- **5**The warm cooling liquid is returned to **4**AALC and exchanges heat with the cold air.
- **6**The warm air is returned to **1**the air conditioner, where it again creates

Source: Dell Technologies Direct Liquid Cooling Support for New PowerEdge Servers | Dell Technologies Info



- Air conditioner 2 creates cool air
- **2**The cold air exchanges heat with the cooling liquid in
- **4** The cold liquid is sent to the server's cooling plate
- **5**The warm cooling liquid is returned to **4**AALC and exchanges heat with the cold air.
- **6**The warm air is returned to **1**the air conditioner, where it again creates

Source: Dell Technologies Direct Liquid Cooling Support for New PowerEdge Servers | Dell Technologies Info



- Air **conditioner 2** creates cool air
- **2**The cold air exchanges heat with the cooling liquid in
- **(1)** The cold liquid is sent to the server's cooling plate
- **5**The warm cooling liquid is returned to **4**AALC and exchanges heat with the cold air.
- **6**The warm air is returned to **1**the air conditioner, where it again creates



- Air **conditioner 2** creates cool air
- **2**The cold air exchanges heat with the cooling liquid in
- **(1)** The cold liquid is sent to the server's cooling plate
- **5**The warm cooling liquid is returned to **4**AALC and exchanges heat with the cold air.
- **6**The warm air is returned to **1**the air conditioner, where it again creates

Source: Dell Technologies Direct Liquid Cooling Support for New PowerEdge Servers | Dell Technologies Info



Source: Dell Technologies Direct Liquid Cooling Support for New PowerEdge Servers | Dell Technologies Info Hub



## We tested the reduction in airflow in server rooms with AALC.

### By using AALC, the amount of air flow in the server room was reduced.



## We tested the reduction in airflow in server rooms with AALC.

#### By using AALC, the amount of air flow in the server room was reduced.

the server airflow and AALC airflow are lower than the air conditioner airflow.

Server air volume

#### AALC Air Volume

#### Air conditioner air volume



## Actapio's 2MW server room with water-cooled GPU servers for LLM and a scenario with a full AALC installation



### Actapio's 2MW server room with water-cooled GPU servers for LLM and a scenario with a full AALC installation



Server air volume

Actapio's 2MW server room with water-cooled GPU servers for LLM and a scenario with a full AALC installation





# Part 1 Summary

#### In part 1 we covered

- LLMs require large amounts of data and enormous computing resources
- As a result, server workloads increase, straining data center resources
- I was told by a U.S. company that there may be a shortage of airflow in air conditioners
- air conditioning system airflow shortages with AALC





- A GPU system for LLM with 72 of the latest GPUs installed in one rack
- Provided as a complete rack, not as individual servers
- All GPUs in the rack are interconnected via NVLink, enabling high-speed inter-GPU communication.
- Compared to the previous model chip, H100, the training speed is 4 times faster and the power efficiency is 1/25 (for GPT-MoE-1.8T model)



- A GPU system for LLM with 72 of the latest GPUs installed in one rack
- Provided as a complete rack, not as individual servers
- All GPUs in the rack are interconnected via NVLink, enabling high-speed inter-GPU communication.
- Compared to the previous model chip, H100, the training speed is 4 times faster and the power efficiency is 1/25 (for GPT-MoE-1.8T model)



- A GPU system for LLM with 72 of the latest GPUs installed in one rack
- Provided as a complete rack, not as individual servers
- All GPUs in the rack are interconnected via NVLink, enabling high-speed inter-GPU communication.
- Compared to the previous model chip, H100, the training speed is 4 times faster and the power efficiency is 1/25 (for GPT-MoE-1.8T model)



- A GPU system for LLM with 72 of the latest GPUs installed in one rack
- Provided as a complete rack, not as individual servers
- All GPUs in the rack are interconnected via NVLink, enabling high-speed inter-GPU communication.
- Compared to the previous model chip, H100, the training speed is 4 times faster and the power efficiency is 1/25 (for GPT-MoE-1.8T model)



- A GPU system for LLM with 72 of the latest GPUs installed in one rack
- Provided as a complete rack, not as individual servers
- All GPUs in the rack are interconnected via NVLink, enabling high-speed inter-GPU communication.
- Compared to the previous model chip, H100, the training speed is 4 times faster and the power efficiency is 1/25 (for GPT-MoE-1.8T model)



- A GPU system for LLM with 72 of the latest GPUs installed in one rack
- Provided as a complete rack, not as individual servers
- All GPUs in the rack are interconnected via NVLink, enabling high-speed inter-GPU communication.
- Compared to the previous model chip, H100, the training speed is 4 times faster and the power efficiency is 1/25

(for GPT-MoE-1.8T model)



Part 2: Server proximity constraints that can be solved with liquid cooling technology

#### Features of GB200 NVL72

 $\dot{\sim}$ 





Sources: - Ingrasys (Diagram) - https://www.youtube.com/watch?v=u9yB67nau7Q&ab channel=InventecDataCenterSolutions (Specifications)

#### Features of GB200 NVL72

Rack: OCP ORv3 rack

- The ORv3 rack has a rack width of **21 inches**.
- Both 19-inch and 21-inch devices can be installed
- GB200 NVL72 is 19 inches except for the Power Shelf



Sources: - Ingrasys (Diagram)

- https://www.youtube.com/watch?v=u9yB67nau7Q&ab\_channel=InventecDataCenterSolutions (Specifications)

Part 2: Server proximity constraints that can be solved with liquid cooling technology

#### Features of GB200 NVL72

Power: OCP bus bar and power shelf

- Bus Bar is 48V DC voltage
- allows you to add additional backup units and electrical power units if needed



Sources: - Ingrasys (Diagram) - https://www.youtube.com/watch?v=u9yB67nau7Q&ab\_channel=InventecDataCenterSolutions (Specifications)

Part 2: Server proximity constraints that can be solved with liquid cooling technology

### Features of GB200 NVL72

# Manifold: Collective piping for water cooling

- Distributes coolant to each IT equipment in the rack
- Universal quick disconnect for quick and secure connection and disconnection
- Standardization with OCP ensures compatibility



Sources: - Ingrasys (Diagram) - <u>https://www.youtube.com/watch?v=u9yB67nau7Q&ab\_channel=InventecDataCenterSolutions</u> (Specifications)
#### Features of GB200 NVL72

Cooling system: Liquid cooling (DLC) is required

• It seems that each company is thinking about a model that connects with a rack-type AALC



Why? Water cooling is essential

 The GB200 NVL2 is equipped with two NVIDIA Blackwell GPU chips and has an air-cooled heat sink in a 2U size.



**GB200 NVL2** 



Sources: - Ingrasys (rack diagram)

#### Why? Water cooling is essential

- GB200 NVL72 has four NVIDIA Blackwell GPU chips in 1U
- It is difficult to place an air-cooled heat sink in 1U due to space limitations.
- Therefore, it was necessary to use water cooling technology to carry heat



Why? Water cooling is essential Why 1U?

- To connect 72 GPUs with NVLink , you need to connect them with copper wires.
- **The length** of the copper is limited.
- Therefore, it is necessary to place GPUs (servers) close to each other



Why? Water cooling is essential

Why 1U?

Why? Do I have to connect via NVLink?

- The challenge of LLM is to shorten the model learning time
- To shorten the model learning time, **high-speed communication between GPUs** is essential.
- High-speed communication is possible by connecting GPUs with NVLink.

Speed companson (bi-directionally)		
	GB/s	Gbps
NVLink 5th	1,800	14,400
PCle-7 x16	512	4,096
800G Ether	200	1,600
400G Ether	100	800

Speed comparison (bi directionally)



77

## Why? "Water cooling is essential" Why 1U?

Why? Do I have to connect via NVLink?

#### " To shorten the model learning time, which is an issue for LLM."



**Extra point:** Processor to Processor Networks : The Future of Connectivity?

#### **NVlink**



# VS

#### **UALink (Ultra Accelerator Link)**



Developer: NVIDIA

Main use: Data transfer between processors Features: High bandwidth, low latency, connects 576 GPUs

Scalability: Connect up to 576 GPUs (GH200 shown)

Developers: Intel, AMD, Broadcom, Cisco, HP, Google, Microsoft, Meta

Main use: Data transfer between processors Features: High bandwidth, low latency, open standard Scalable: Connect up to 1,024 AI accelerators



#### GB200 NVL72 is a 120kw monster machine with 1 rack





High power of 120 kW is required because it is packed into 1 Rack.

#### GB200 NVL72 is a 120kw monster machine with 1 rack





#### Power Supply Scheme For 480V







#### Only about 10 sets can fit in a 2MW server room.



## Part 2 Summary

#### In part 2 we covered

- Reducing model learning time is a major issue
- The GB200 NVL72 technology announced at GTC is a solution to this problem.
- Proximity makes it difficult to cope with heat sinks, and water cooling technology is required for efficient cooling
- The servers packed into a single rack **created a new challenge: the power supply.**



### Summary

- The popularity of LLM has brought new challenges to our data center operations.
- LLM creates increased server workloads and associated cooling challenges
- Specifically, it surfaced as a problem of insufficient airflow in the air conditioning system and server proximity constraints

• We have spoken with U.S. companies and hyperscalers at GTC, OCP Summits, and elsewhere, and carefully considered that water cooling technology is the key to solving these "essential challenges"



#### Summary



#### Summary















How to avoid mistaking the solutions to future problems for those who have not faced any challenges yet

- Gather information on what the essential challenges are for those who face challenges
- In terms of LLM, hyperscale companies such as US companies
- It is important to immediately put the information gathered into practice in the near future (in the United States)



- Actapio's problem-solving approach
  - We collect, verify, and implement information from people who have faced challenges (U.S. companies), so we are able to select the best technology
- Vision for the future
  - The "users" are left behind
  - In order to lead the industry, we must become a "pioneer in problem solving"
  - As we look towards the future, will we aim to be on the "user side" that quickly adopts technology, or will we aim to be "pioneers in solving problems" and become industry leaders?
  - Whichever path you choose, let's start thinking together today about what we need to do to pave the way for the future.



# **THANKS**

## **Our Team**

Aaron Kirkpatrick Andrew Arnold Anthony Skwiat Atsuko Ishigaki Eiji Kawauchi Jason Van Winkle J.D. Salling Kai Fukazawa Ken Tairabune Koyo Uemizu Kyoya Torikai Hideki Mikami Hisatomo Tanaka Marie Rudolph

Masahiko Matsui Masayuki Ashida Osamu Kurino Shinichiro Okamoto Susan McKee Takashi Watanabe Takuya Kitano Takumi Uematsu Tatsumi Yuusuke Travis Mather Yoshimura Tsubasa Yuji Kohata Yukihito Imai

# **Our Datacenter**

