

ACTAPIO



LINEヤフー米国データセンタ技術の最前線
LLM(大規模言語モデル)と水冷技術への挑戦

Norifumi Matsuya



自己紹介



自己紹介 in JAPAN

- 松谷 憲文(まつや のりふみ)
- 2000年から2014年まで
ヤフーにてインフラの設計、構築、運用を担当
- 主な実績
 - プロダクションネットワークの BGP化
 - JPNIC、APNIC、ARINでAS取得
(AS23816、AS24572、AS18140)
 - バックボーンネットワークを仮想化、VPLS化
 - 国内最大規模の Openstack クラスタ作成
 - 国内データセンターの建設
(東京、大阪、中部地区、九州地区、東北地区)



米国データセンター見学を経て、米国進出を決意



米国データセンターの見学 '09

自己紹介 in U.S.



Actapio データセンター

- 2014年に渡米し、Actapioに出向
- ActapioはLINEヤフー株式会社の米国法人でAI・ビッグデータ用のデータセンター建設、サーバ運用を行う会社
- 現在は会社を経営しながら、OCP(Open Compute Project)のサーバの運用やAIビジネス開発などを担当

AGENDA

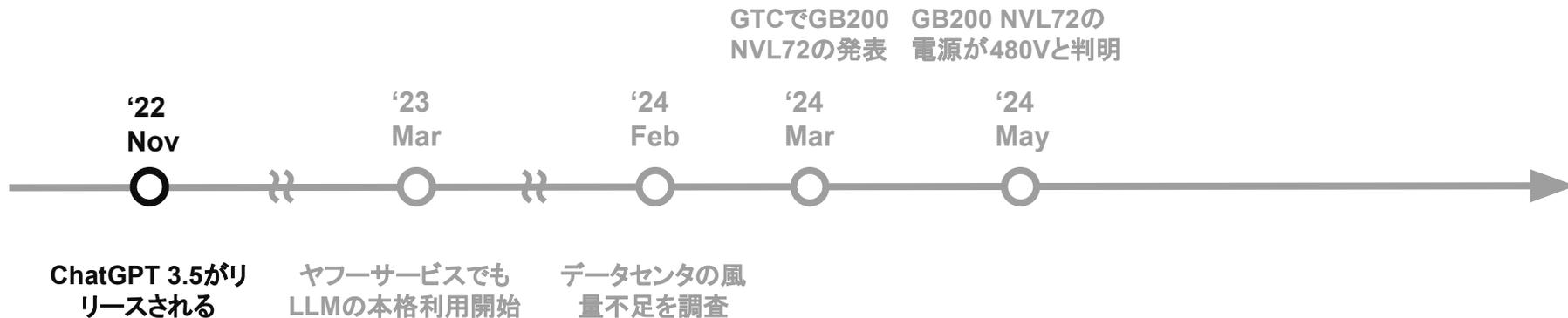
- **本日の発表ダイジェスト**
- **第一部：水冷技術で解決できる空調機の風量不足**
 - LLMの出現とサーバリソースの変化
 - データセンタの現状と予測される空調機の風量問題
 - 空調機の風量不足への解決策検討
- **第二部：水冷技術で解決できるサーバ近接性制約**
 - GB200 NVL72 (GPUサーバ)の登場とその背景
 - 電源の480V化(調査中)
- **まとめ**

本日の発表ダイジェスト



第一部

第二部



- ChatGPTは対話型の**LLM** (大規模言語モデル) である
- LLMは**トランスフォーマモデル**を使用し、**大量のテキストデータ**で訓練されている
- トランスフォーマは**自然言語処理で高い性能を発揮**するニューラルネットワークアーキテクチャ

第一部

第二部

GTCでGB200 NVL72の発表
GB200 NVL72の電源が480Vと判明

'22
Nov

'23
Mar

'24
Feb

'24
Mar

'24
May

ChatGPT 3.5がリリースされる

ヤフーサービスでもLLMの本格利用開始

データセンタの風量不足を調査

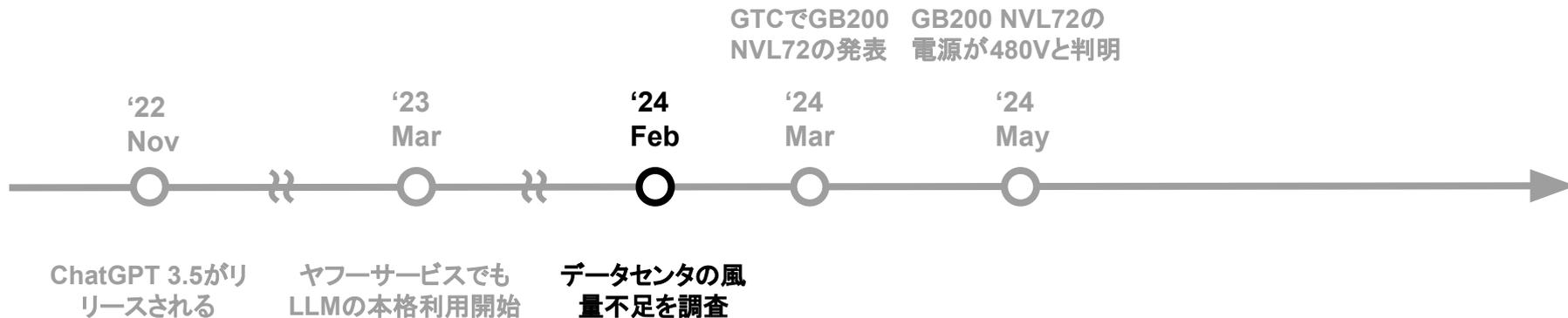
- LLMは**大量のデータ**と**膨大なコンピューティングリソース**が必要
- その結果、サーバの**ワークロードが大幅に増加**



- 何か問題がないか、データセンタの**各リソース**をチェックするも**問題ない**ことを確認
- 安心していたところに、データセンタ**空調機の風量が不足**する情報を米国企業から入手

第一部

第二部



- 空調機の風量不足問題を**水冷技術**で解決できるか検証
- 水冷技術の利用によって、
 - **冷却効率が上昇**
 - 空調機の**風量不足が解決**
 - **より多くのサーバを設置可能**

第一部

第二部

GTCでGB200 NVL72の発表
GB200 NVL72の電源が480Vと判明

'22
Nov

'23
Mar

'24
Feb

'24
Mar

'24
May

ChatGPT 3.5がリリースされる

ヤフーサービスでもLLMの本格利用開始

データセンタの風量不足を調査

- 2024年3月、GTC(GPU Technology Conference)でNVIDIAが新しいGPUサーバである**GB200 NVL72**を発表
 - GB200 NVL72には**水冷技術**が必須となっている
- 水冷技術が必要な理由を**米国企業からヒアリング**
 - 理由はモデル学習時間の短縮を図るためにGPU間で**高帯域の通信が必要**

第一部

第二部

GTCでGB200 NVL72の発表
GB200 NVL72の電源が480Vと判明

'22
Nov

'23
Mar

'24
Feb

'24
Mar

'24
May

ChatGPT 3.5がリリースされる

ヤフーサービスでもLLMの本格利用開始

データセンターの風量不足を調査

- ・GB200 NVL72の電源が**480V入力**
- ・現在、わかっている状況までを本編で説明



LLMの出現とサーバリソースの変化

イマココ！



第一部

第二部

GTCでGB200 NVL72の発表
GB200 NVL72の電源が480Vと判明

'22
Nov

'23
Mar

'24
Feb

'24
Mar

'24
May

ChatGPT 3.5がリリースされる

ヤフーサービスでもLLMの本格利用開始

データセンタの風量不足を調査



2022年11月 ChatGPT 3.5リリース



2022年11月 ChatGPT 3.5リリース

今までのAIと何が違うの？

従来型「タブラー形式データ学習」と 「LLMのトランスフォーマー学習」の違い

特徴	タブラー形式データ	LLM(大規模言語モデル)
データの種類	表形式データ(例: Excel, CSV)	テキストデータ
得意なこと	数値予測、分類、データ分析	文章生成、翻訳、チャットボット
トレーニングデータ	構造化データ(数値、カテゴリ)	大量の文書データ
モデルの規模	比較的小規模 (特徴量やデータ数に依存)	大規模(数十億のパラメータ)
柔軟性	特定のデータセットに特化	多様なタスクに対応可能

従来型「タブラー形式データ学習」と 「LLMのトランスフォーマー学習」の違い

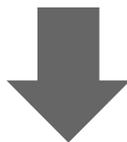
特徴	タブラー形式データ	LLM(大規模言語モデル)
	(Excel, CSV)	テキストデータ
	分析	文章生成、翻訳、チャットボット
		大量の文書データ
モデルの規模	比較的小規模 (特徴量やデータ数に依存)	大規模(数十億のパラメータ)
柔軟性	特定のデータセットに特化	多様なタスクに対応可能

今までのAIと比べて桁違いのデータ
が必要😱

ヤフーのサービスでもLLMを利用



- ヤフーではユーザの投稿など「言語を理解」するためにLLMを使用
- LLMの利用によって大きな精度向上を達成



ただし、大きな課題も



サーバワークロードの変化 (with A100 GPU)

タブラー形式データワークロード 30%-80%が1週間ほど



LLM ワークロード **80%-100%が1、2ヶ月継続**





データセンタの現状と予測される空調機の風量問題

イマココ！



第一部

第二部

GTCでGB200 NVL72の発表
GB200 NVL72の電源が480Vと判明

'22
Nov

'23
Mar

'24
Feb

'24
Mar

'24
May

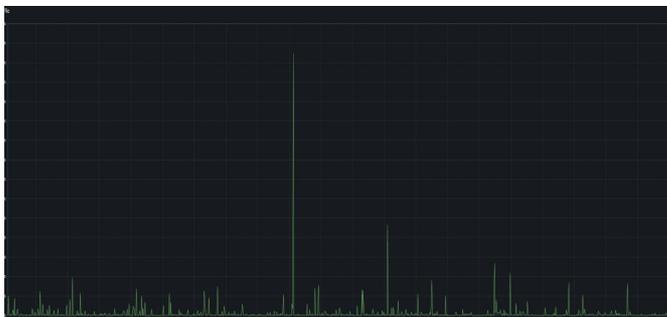
ChatGPT 3.5がリ
リースされる

ヤフーサービスでも
LLMの本格利用開始

データセンタの風
量不足を調査

ネットワーク: LLMと別システム(Hadoop)のサーバ帯域比較

Server for Hadoop Compute



Server for LLM



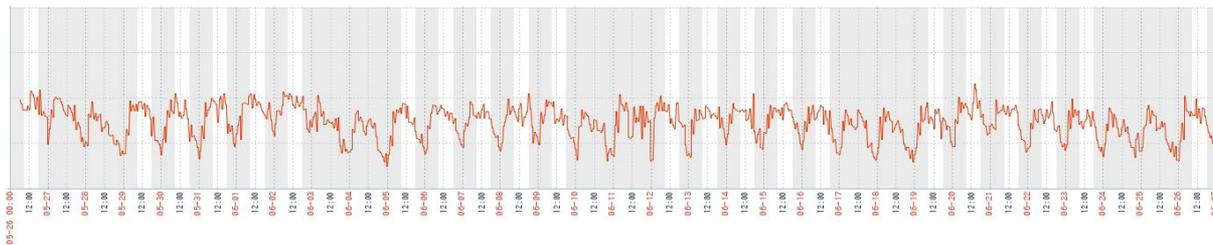
	サーバポートの帯域キャパシティ	利用状況
Server for Hadoop Compute	25 Gbps	13.50 Gbps (54%)
Server for LLM (with A100 GPU)	100 Gbps	23.5 Gbps (24%)

- 帯域も足りており、エラーやドロップも見当たらない
- ネットワークにおける課題は発生していない
- 波形パターンには大きな差異があり、LLMのワークロードは一度動き始めると長期間トラフィックを出し続ける傾向があることを確認
- モデル作成時に少し工夫していることも影響

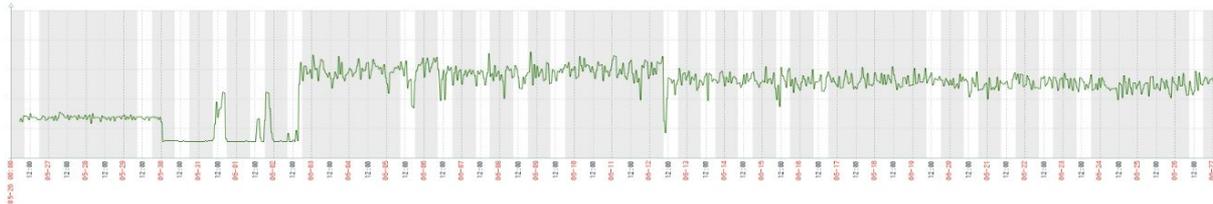


電力：LLMと別システム(Hadoop)のサーバ電力比較

Server for Hadoop Compute (最大稼働率：76.12%)



Server for LLM (最大稼働率：95.44%)



Hadoop Rack
(30 nodes / rack)

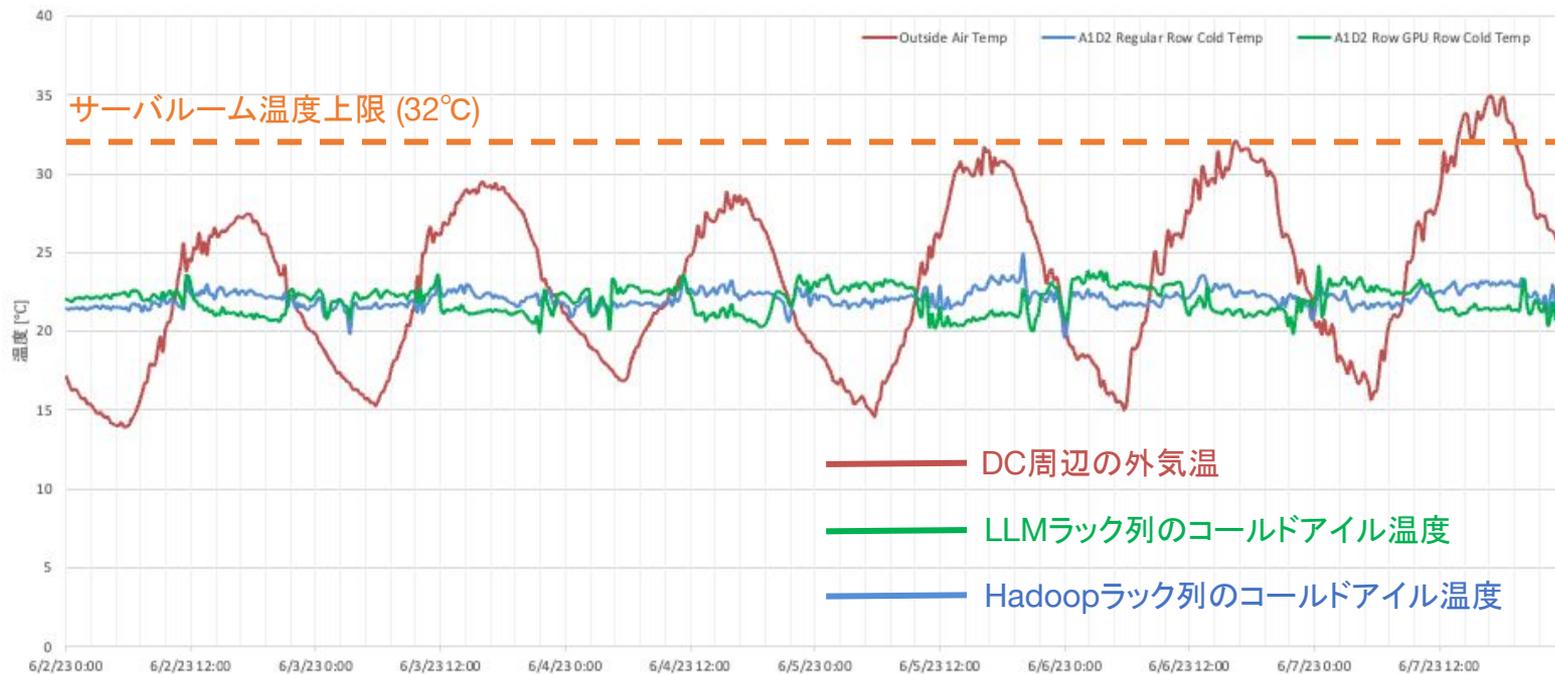


LLM GPU Rack
(2 nodes / rack)



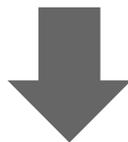
サーバワークロードが長期間にわたって高負荷状態にあるため、消費電力が上昇

冷却: LLMと別システム(Hadoop)の温度比較



温度の観点ではLLMとHadoopのラック列とで大きな差異はなく、
温度上の問題は発生していない

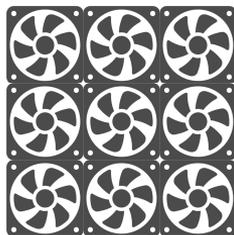
現時点ではネットワーク、電力、冷却に関して概ね問題なし？



- OCP Summit 2023での**米国企業と水冷技術**について**対話**
 - GPUの**TDP** (Thermal Design Power)が**増加傾向**
 - TDPの増加により**サーバ風量** (ファンスピード)が**最大に**
 - サーバ風量増加により**空調機の風量が増加する懸念**が明らかになった

データセンタの空調機が問題なく稼働するためには、
空調機の風量が**サーバの風量**を上回る必要があります

サーバの風量



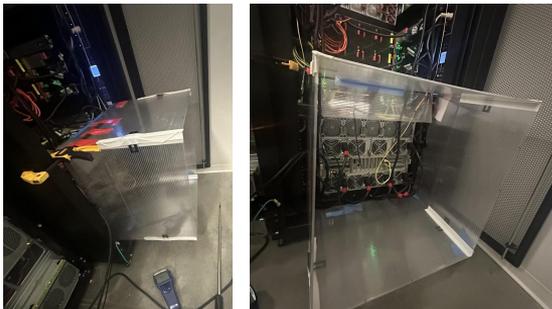
空調機の風量



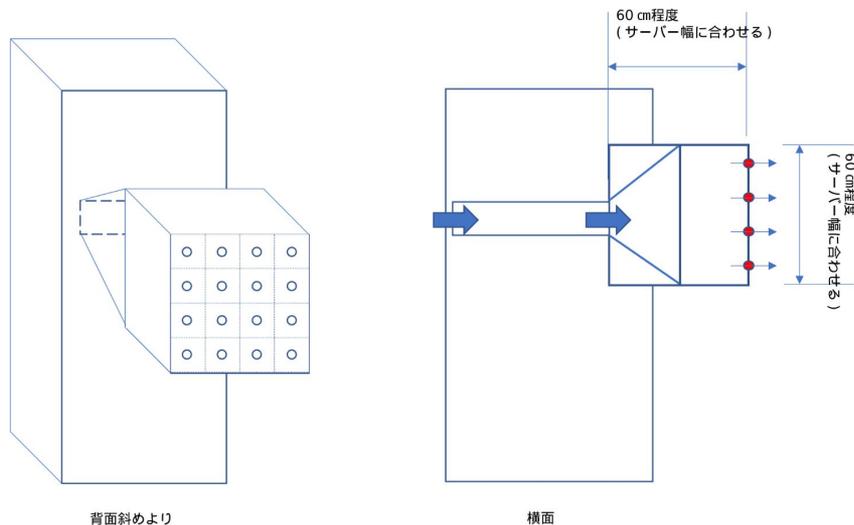
空冷GPUサーバの風量を測定

- サーバは**フル負荷**で測定
- サーバ背面をプラスチック製風防で囲う
- 断面を4x4の区間に糸で区切り、それぞれの区間を**風速計**で測定

実際の様子



測定方法



風速計

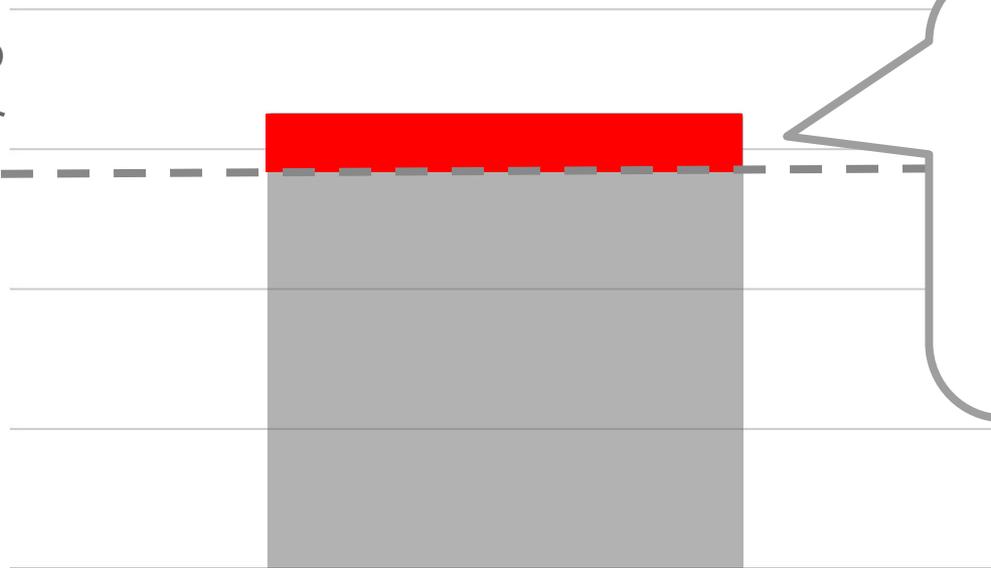


Actapioの2MWのサーバールームに LLM用の空冷GPUサーバをフル設置した場合のシナリオ

空冷GPUサーバの2MW分のトータル風量

空調機の
風量限界

サーバ風量



- 空冷GPUサーバ風量が空調機の**設備限界を上回る**
- サーバルームの**温度制御が出来なくなる**

空冷GPUサーバ



空調機の風量不足への解決策検討

イマココ！



第一部

第二部

GTCでGB200 NVL72の発表
GB200 NVL72の電源が480Vと判明

'22
Nov

'23
Mar

'24
Feb

'24
Mar

'24
May

ChatGPT 3.5がリリースされる

ヤフーサービスでもLLMの本格利用開始

データセンタの風量不足を調査

空調機の風量不足への解決案

案① 設置サーバ数の抑制

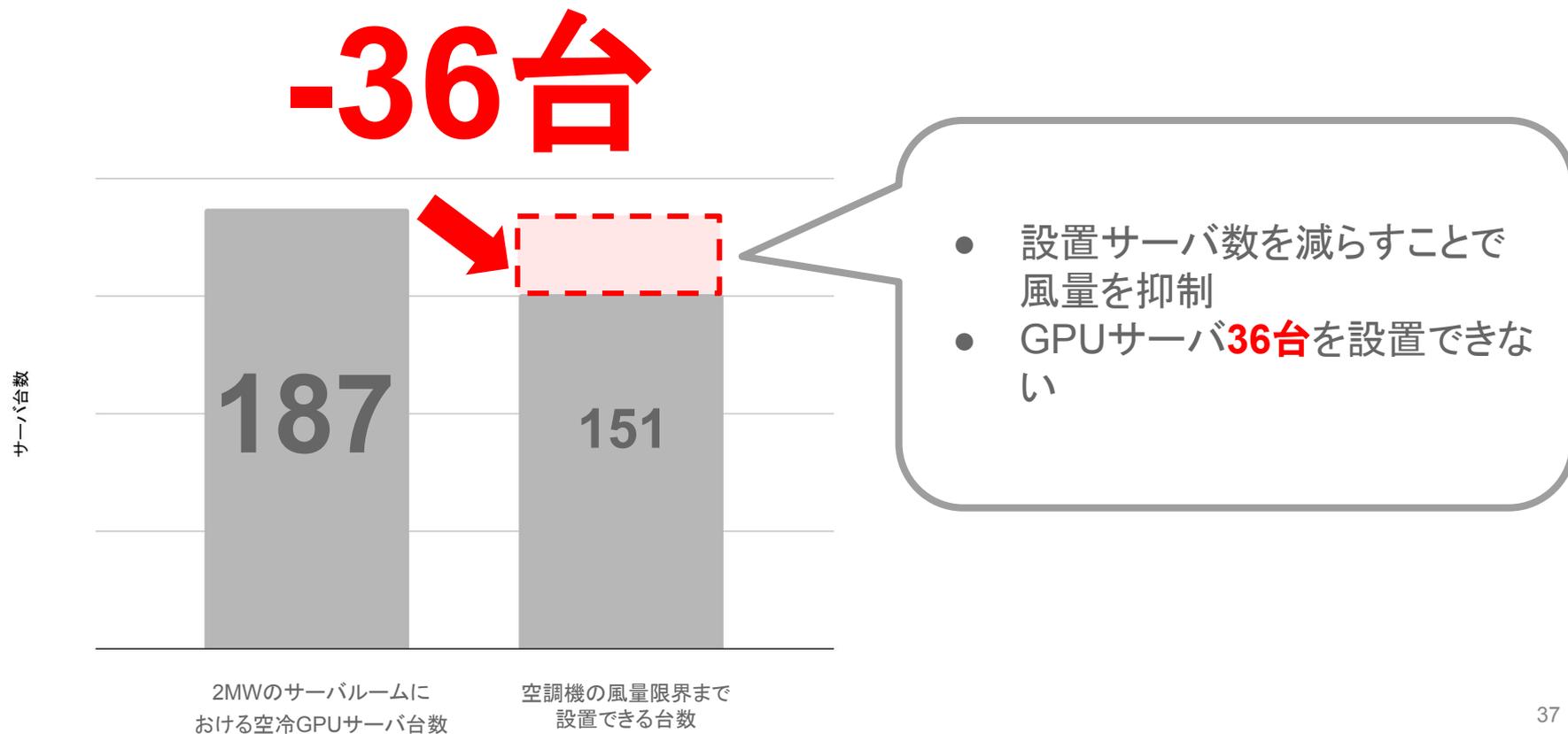
案② 空調機の増強

- a. 空調機の増設
- b. 空調機のスペック増強

案③ 水冷方式の検討: DLC (Direct Liquid Cooling)

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

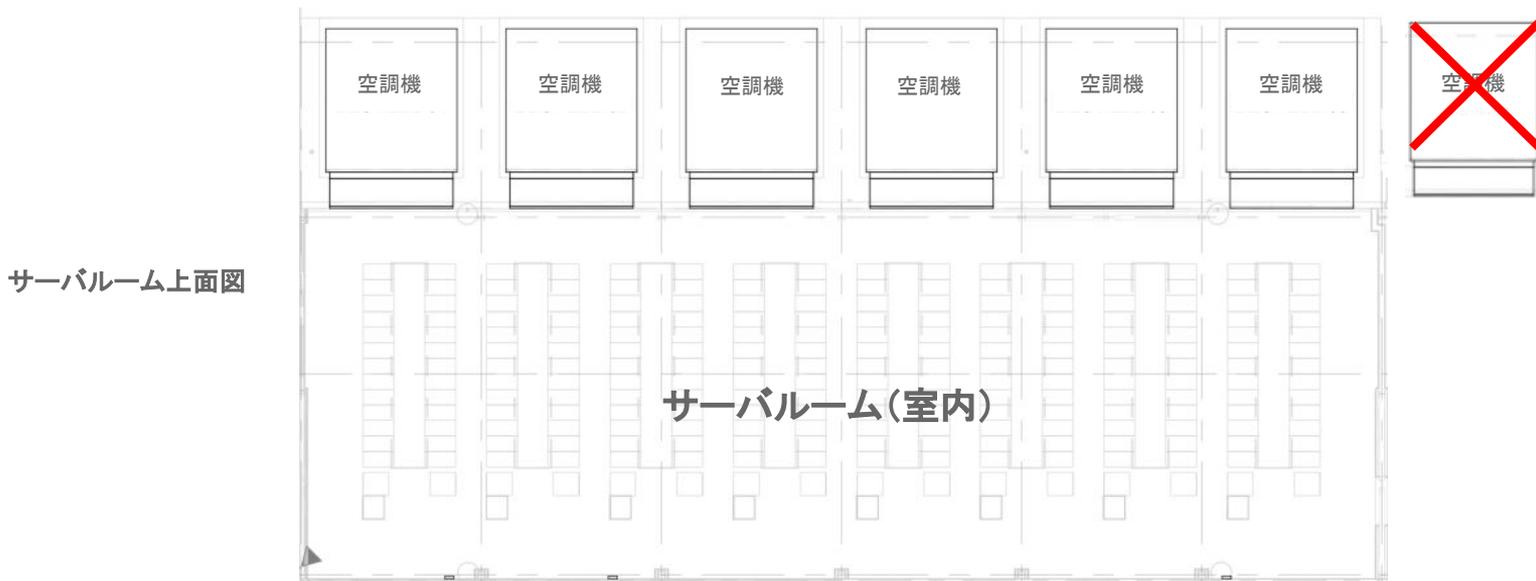
案① 設置サーバ数の抑制



案② 空調機の増強

a. 空調機の増設

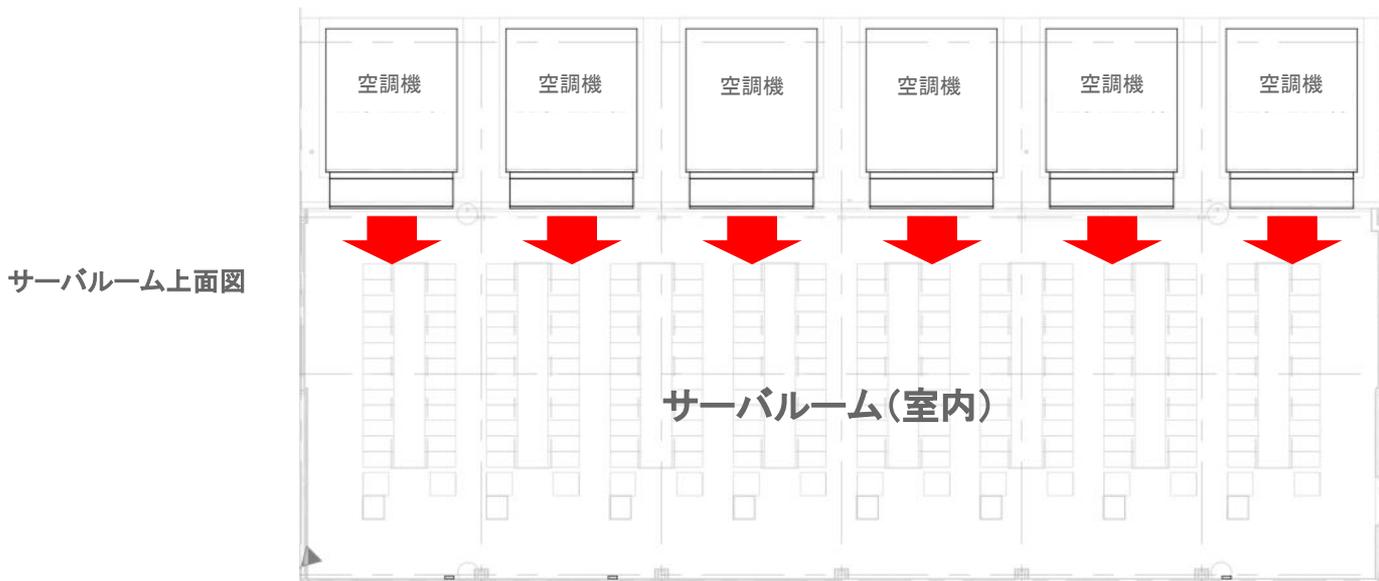
データホールのレイアウト(下図)のため物理的に**増設負荷**



案② 空調機の増強

b.空調機のスペック増強

既存空調機のスペックアップはコスト、工事レベル含め**増強は困難**



案③ 水冷方式の検討

水冷方式(DLC)は2種類に分類できます

DLC(Direct Liquid Cooling)

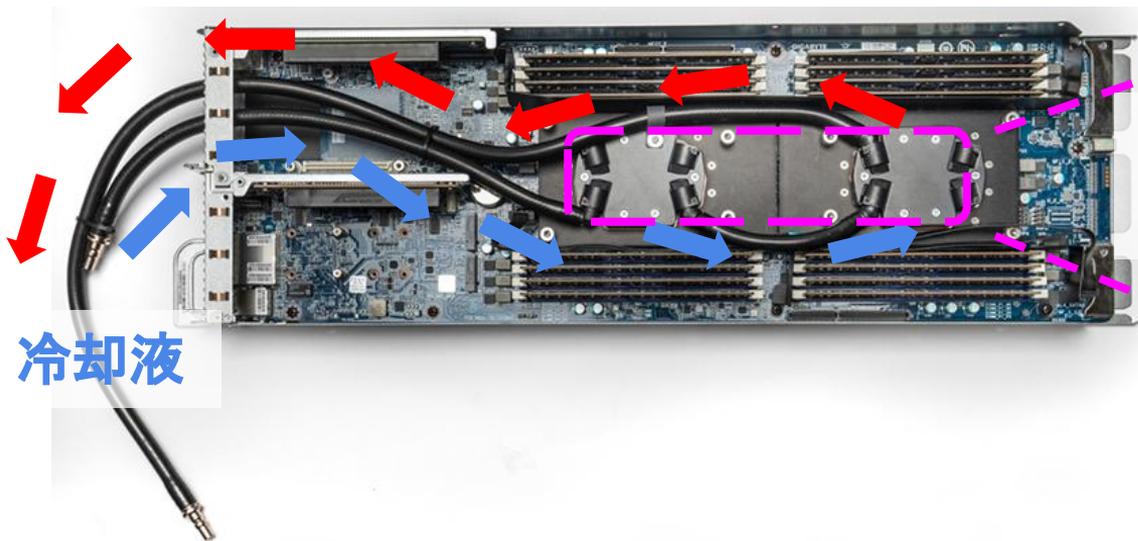
1. L to L : Liquid to Liquid

2. **AALC : Air Assisted Liquid Cooling**

水冷方式の中でも**AALC**を検討しています、
その選択理由を説明します

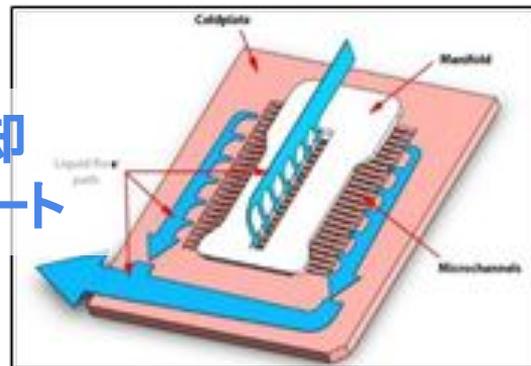


DLC (Direct Liquid Cooling) について



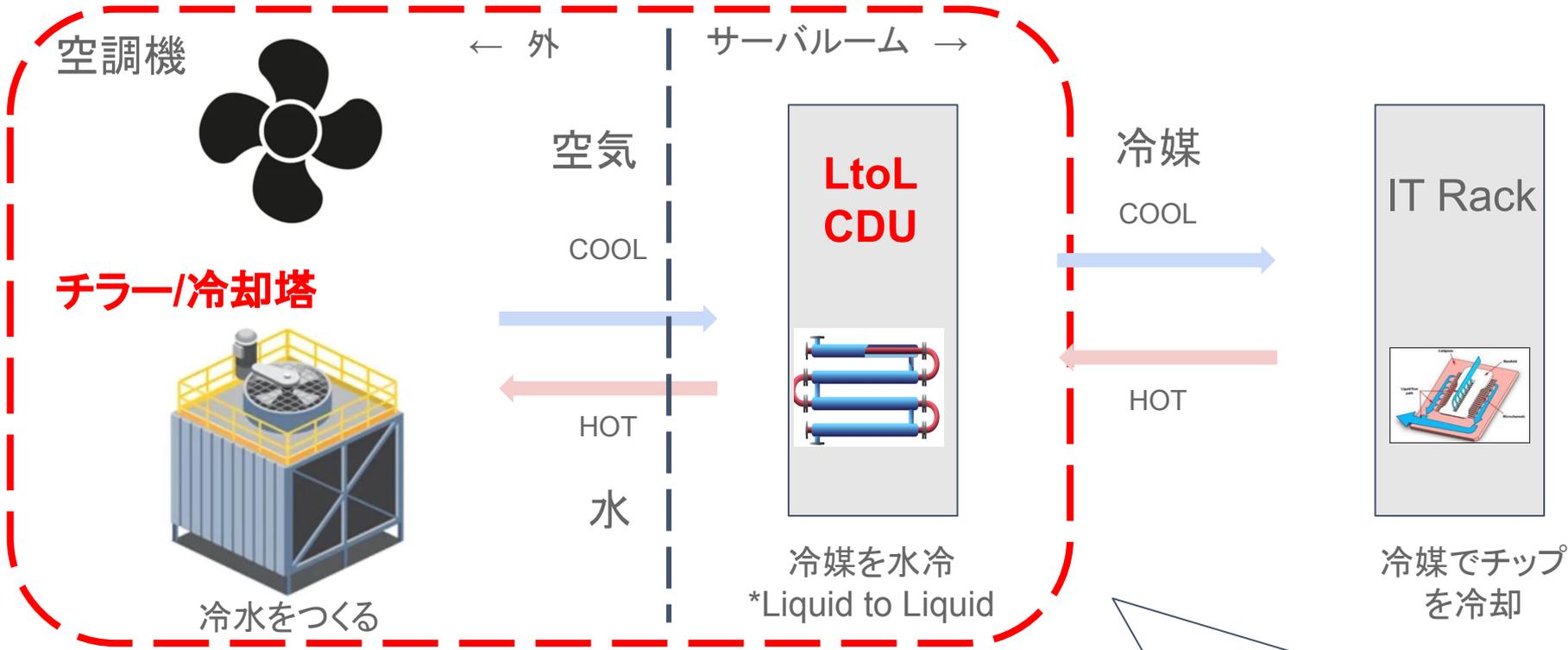
冷却液

冷却
プレート



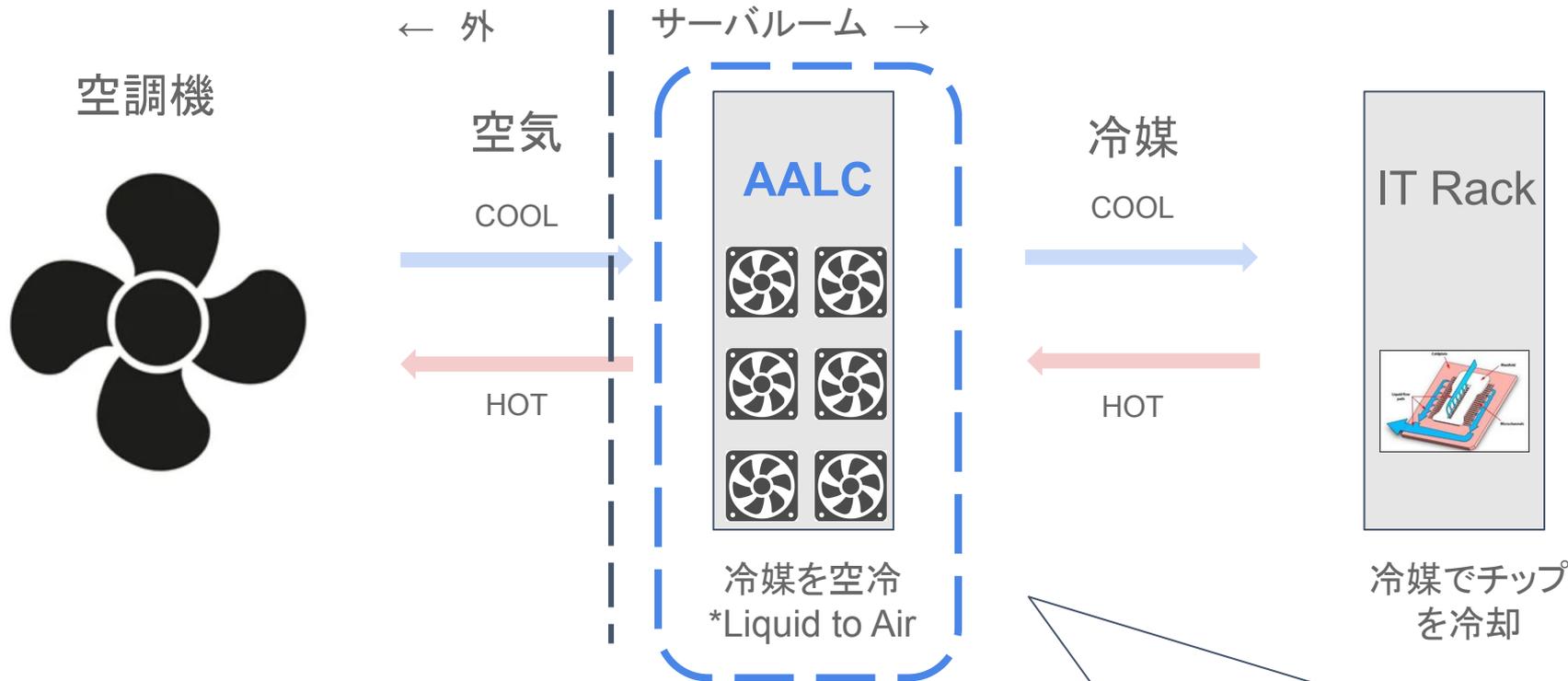
- 冷却液を冷却プレートに送る
- 冷却プレートがプロセッサに直接接触して冷却

LtoLによる冷却スキーム



サーバ風量は確実に減るが、冷水設備とサーバールーム内への配水管が必要で**コストや納期がかかる**

AALCによる冷却スキーム



- サーバルームに設置するだけなので **低コスト、低納期**
- 風量が削減されるか検証

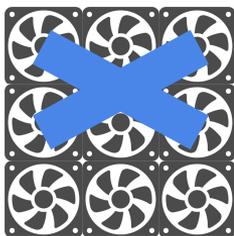
AALCでサーバールーム風量が減少するか検証を行った



AALCを利用することでサーバールーム風量は減少した

サーバ風量とAALCの風量が空調機の風量を下回ればOK

サーバ風量



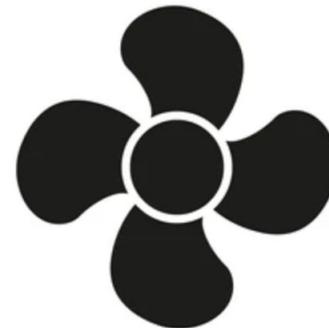
AALC風量



ラック型AALC

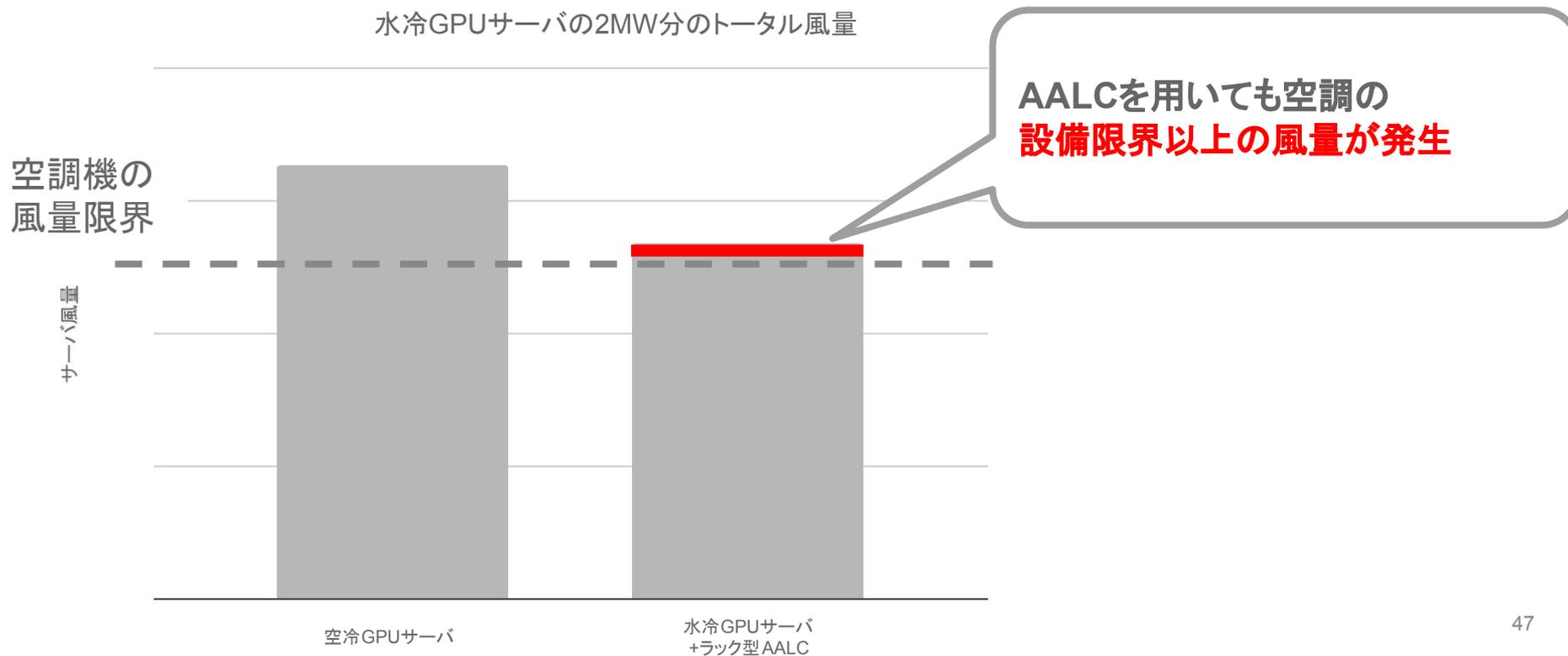


空調機の風量

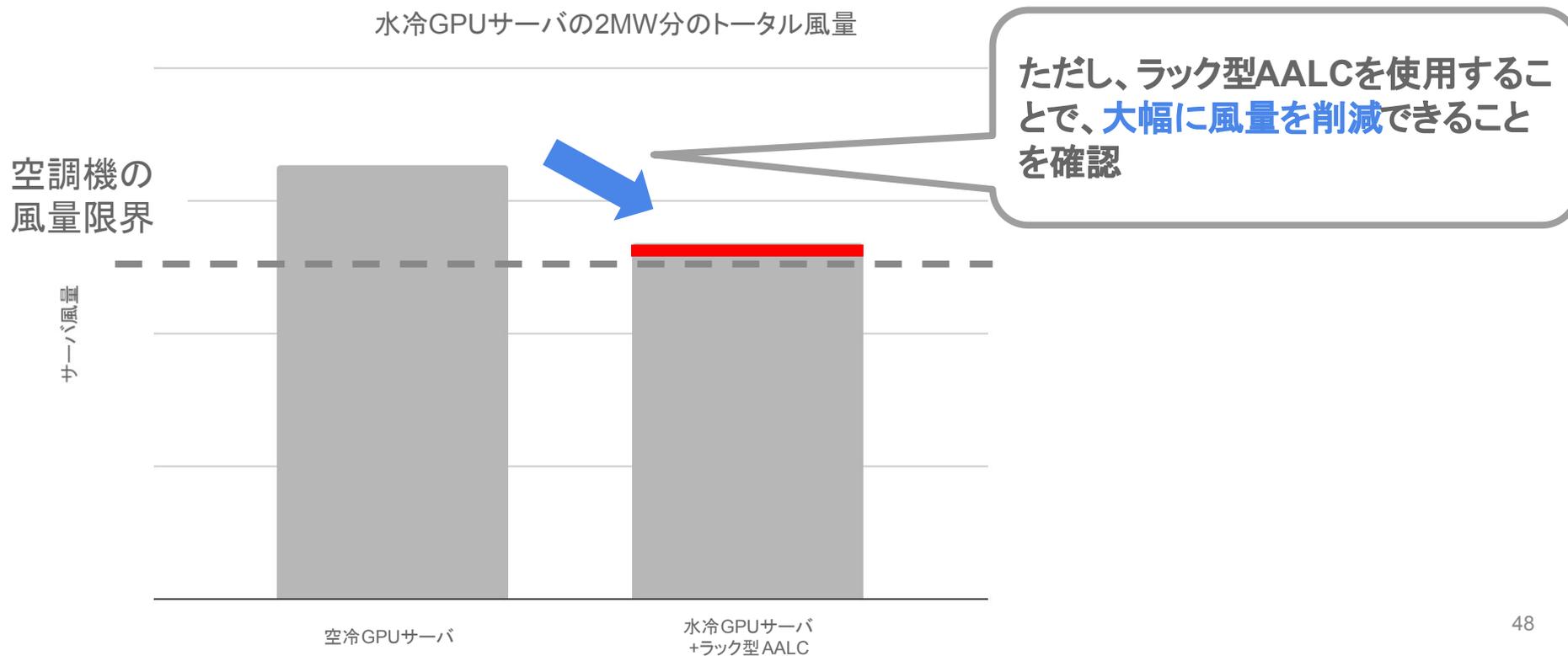


ポイント：GPUとCPU分の風量が減少した分、AALCの風量は増加する

Actapioの2MWのサーバールームにLLM用の水冷GPUサーバとAALCをフル設置した場合のシナリオ



Actapioの2MWのサーバールームにLLM用の水冷GPUサーバとAALCをフル設置した場合のシナリオ

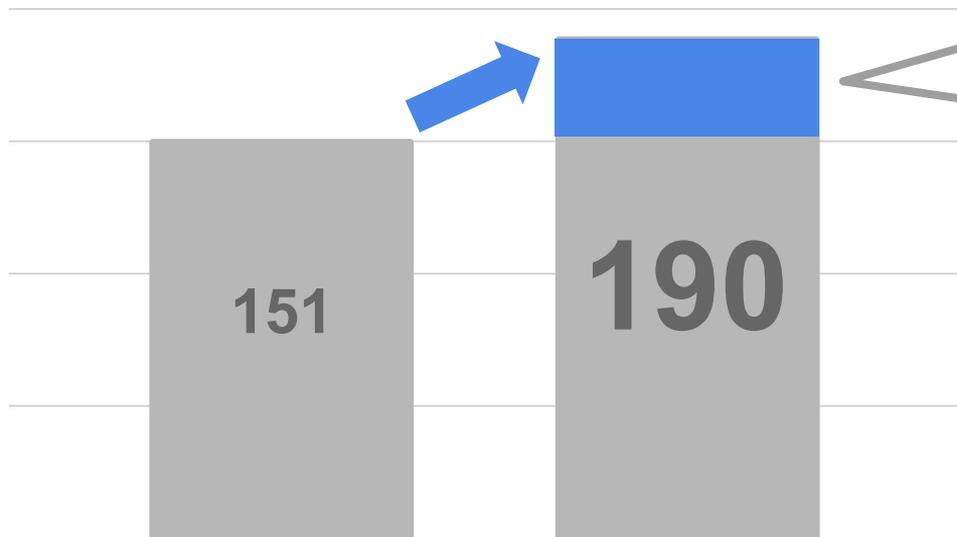




Actapioの2MWのサーバールームにLLM用の水冷GPUサーバとAALCをフル設置した場合のシナリオ

+39台

サーバ台数



AALCを利用することで、
最大で**39台**多く設置可能

空冷GPUサーバ

水冷GPUサーバ
+ラック型AALC

第一部サマリ

第一部では、

- LLMは**大量のデータ**と**膨大なコンピューティングリソース**を必要とする
- その結果、**サーバワークロードが増加し、データセンタリソースを逼迫**させる
- 米国企業から空調システムの**風量不足の可能性**を教えてもらった
- 空調システムの風量不足を**AALCでの解決**

第一部

第二部

GTCでGB200 NVL72の発表
GB200 NVL72の電源が480Vと判明

'22
Nov



ChatGPT 3.5がリリースされる

'23
Mar



ヤフーサービスでもLLMの本格利用開始

'24
Feb



データセンタの風量不足を調査

'24
Mar



'24
May





GB200 NVL72の登場とその背景

イマココ！



第一部

第二部

GTCでGB200 NVL72の発表
GB200 NVL72の電源が480Vと判明

'22
Nov

'23
Mar

'24
Feb

'24
Mar

'24
May

ChatGPT 3.5がリリースされる

ヤフーサービスでもLLMの本格利用開始

データセンタの風量不足を調査

GB200 NVL72とは？

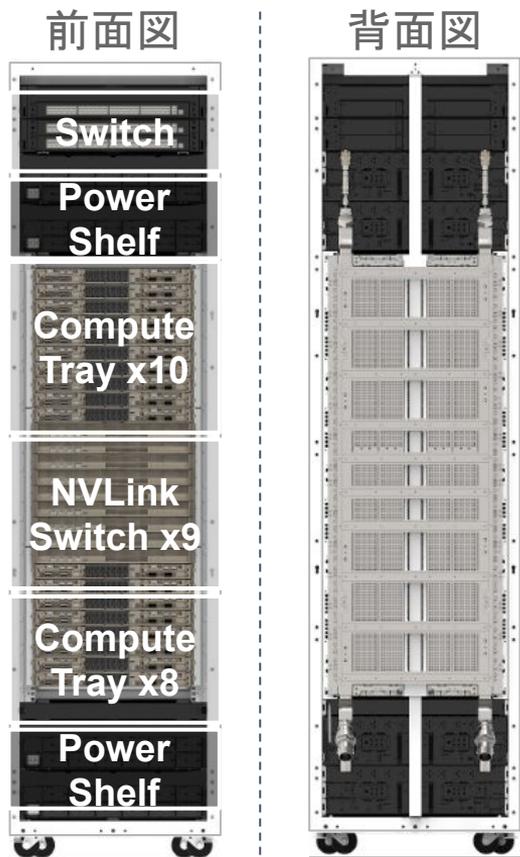
- 最新のGPU72基が1つのラックに搭載されたLLM向けのGPUシステム
- サーバ単体ではなく、ラック一式で提供
- ラック内のGPU全てがNVLinkにより相互接続され、高速なGPU間通信を実現
- 前モデルのチップであるH100比較にて、トレーニング速度が4倍、電力効率は1/25(GPT-MoE-1.8Tモデルの場合)



※出典：<https://www.nvidia.com/en-us/data-center/gb200-nvl72/>



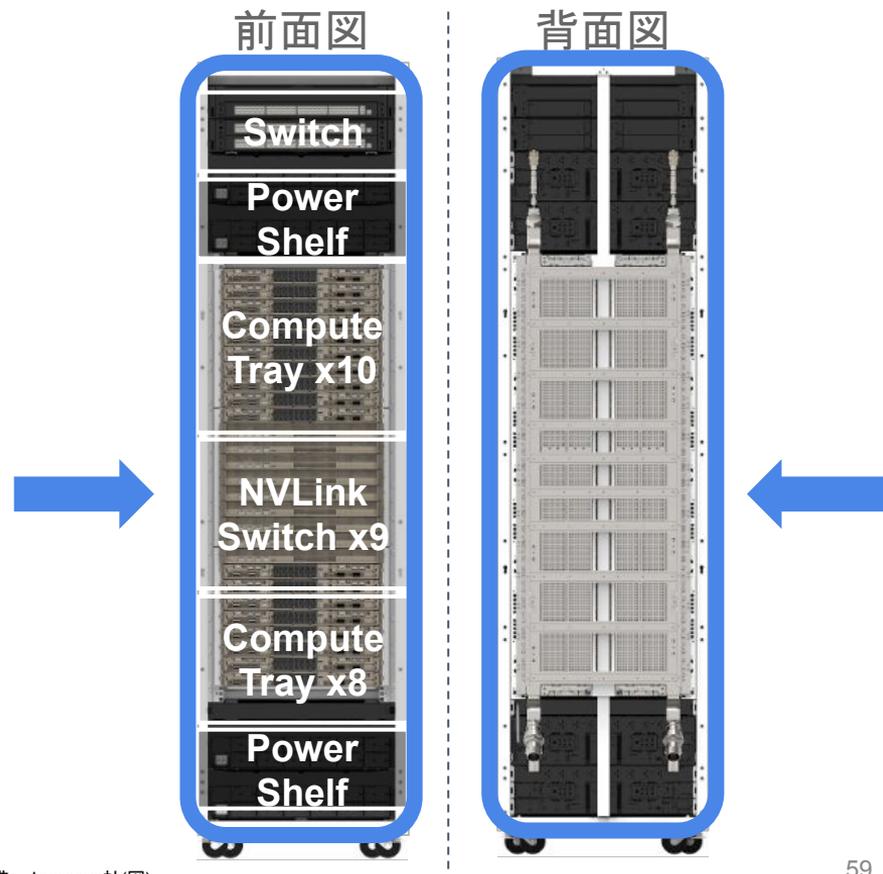
GB200 NVL72の特徴



GB200 NVL72の特徴

ラック：OCP仕様のORv3ラック

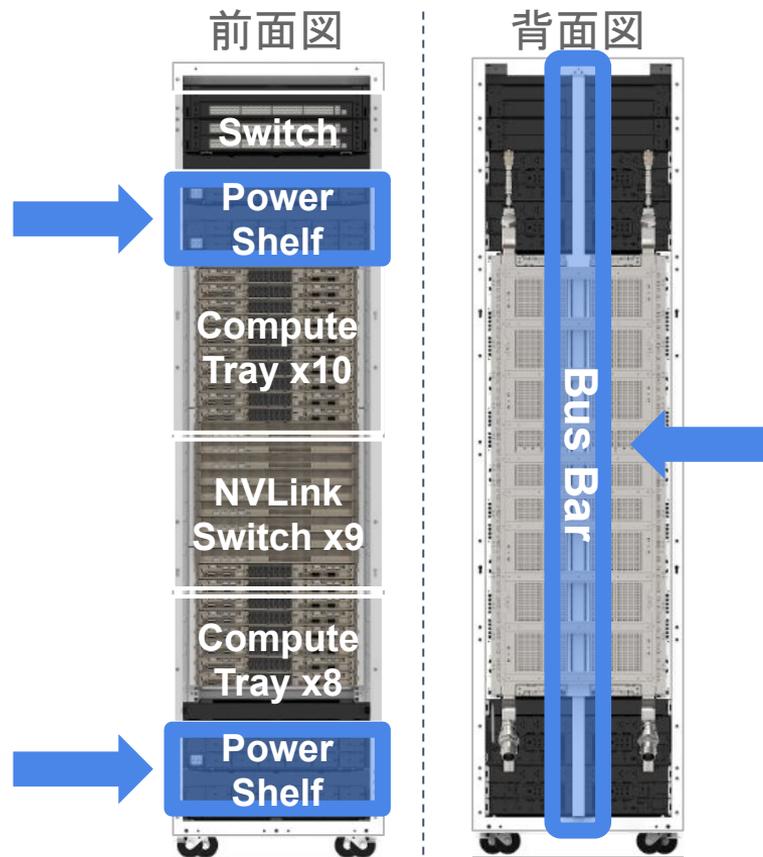
- ORv3ラックはラック幅が**21インチ**
- **19インチ、21インチ両方**の機器を設置可能
- GB200 NVL72ではPower Shelf以外は19インチ



GB200 NVL72の特徴

電源：OCP仕様のBus BarとPower Shelf

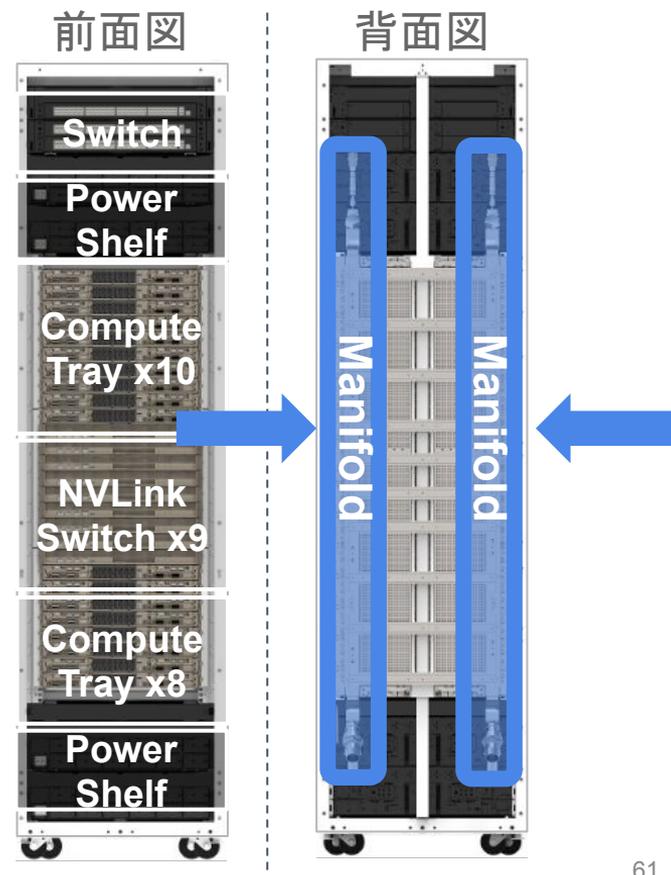
- Bus Barは48V直流電圧
- Power Shelfは必要に応じて電力やバックアップシステムを追加可能



GB200 NVL72の特徴

Manifold: OCP仕様の水冷用の集合配管

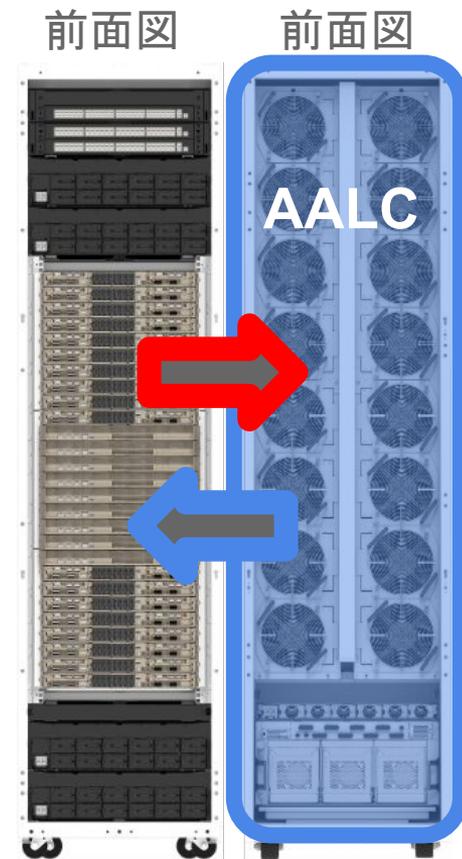
- ラック内の各IT機器に**冷却液を分配**
- ユニバーサルクイックディスコネクトを採用し、**迅速で安全な接続・切断**が可能
- OCPで標準化することによって、**互換性が確保**



GB200 NVL72の特徴

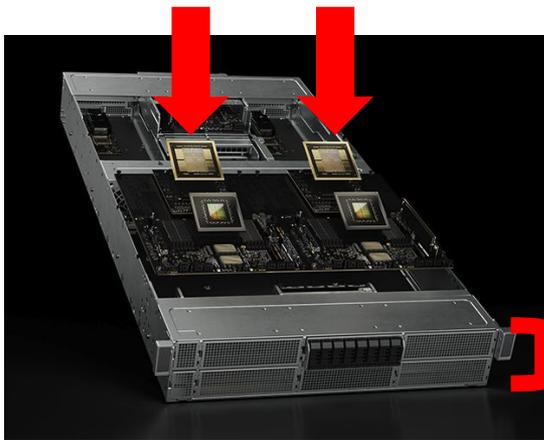
冷却システム：**水冷方式(DLC)**が必須

- 各社**ラック型AALC**で繋ぐモデルで考えている模様



なぜ？ 水冷方式が必須

- NVIDIA Blackwellの **GPUチップが2つ**配置されているモデルで、**空冷ヒートシンクが2U**サイズ



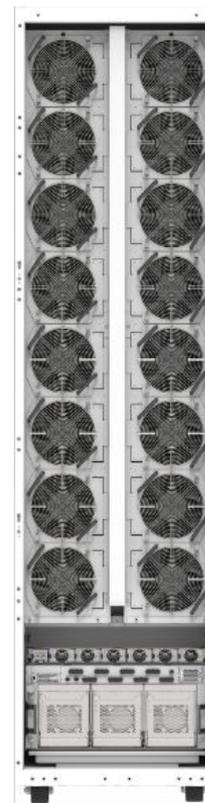
2U

GB200 NVL2

前面図

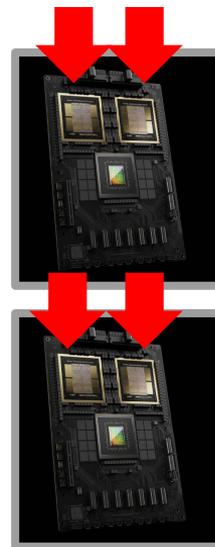


前面図

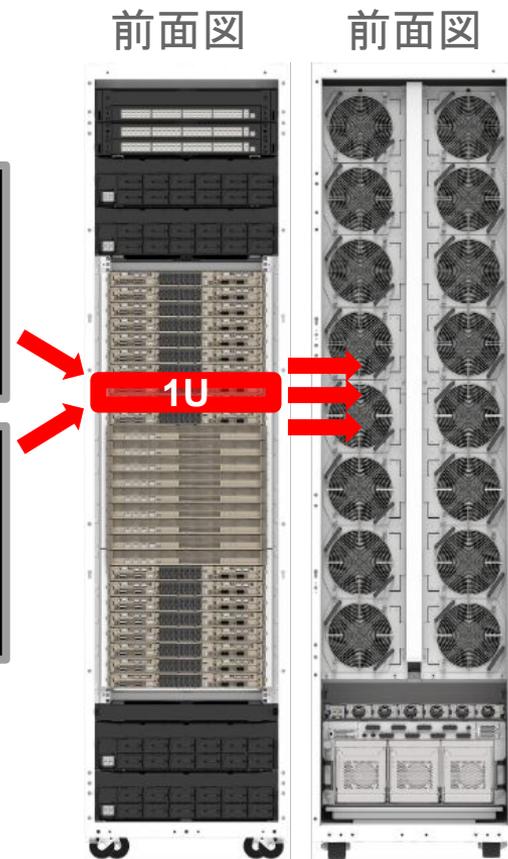


なぜ？ 水冷方式が必須

- GB200 NVL72は、NVIDIA Blackwell GPUチップが4つ、1Uに配置
- 1Uに空冷ヒートシンクを配置するのはスペース的に難しい
- そのため水冷技術を使用し、熱を運ぶ必要があった



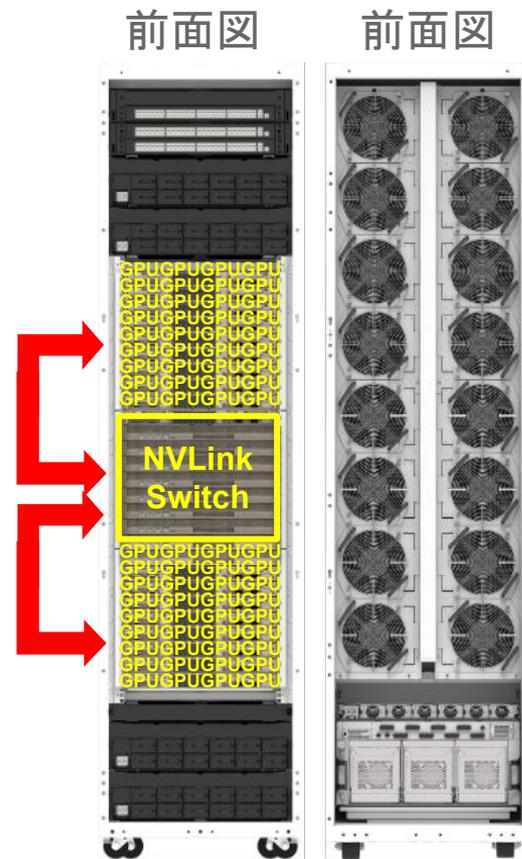
GB200 Superchip



なぜ？ 水冷方式が必須

なぜ？ 1Uにしている？

- 72個のGPUを**NVLink**で接続するには**Copper(銅線)**で接続する必要がある
- そのCopperの**長さ**に**制限**がある
- そのためGPU(サーバ)間を**近くに配置**する必要がある



なぜ？ 水冷方式が必須

なぜ？ 1Uにしている？

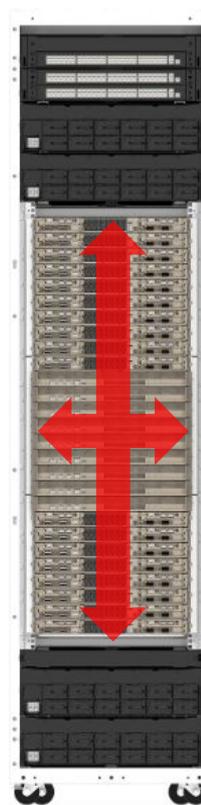
なぜ？ NVLinkで接続しないといけない？

- LLMの課題は**モデル作成時間の短縮**
- そのモデル作成時間の短縮には**GPU間的高速通信が必須**となっている
- **NVLink**でGPU間を接続することにより高速通信が可能になる

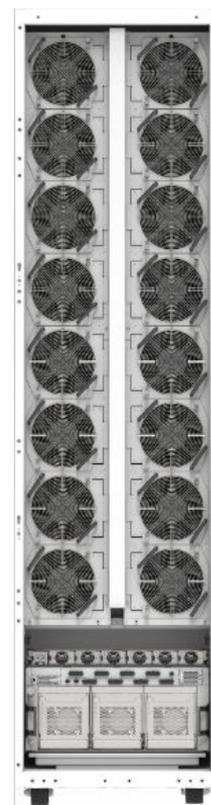
速度比較 (bi-directionally)

	GB/s	Gbps
NVLink 5th	1,800	14,400
PCIe-7 x16	512	4,096
800G Ether	200	1,600
400G Ether	100	800

前面図



前面図



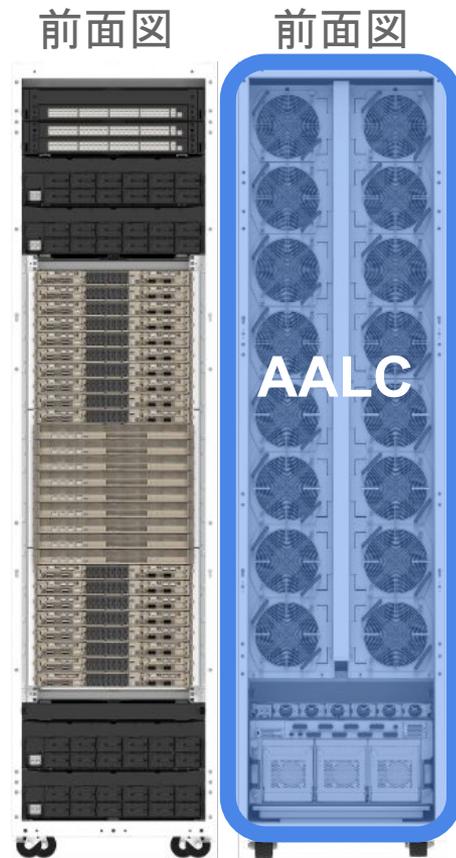
なぜ？「水冷方式が必須」

なぜ？1Uにしている？

なぜ？NVLinkで接続しないといけない？



なぜならば、
「LLMの課題であるモデル時間の短縮のため」

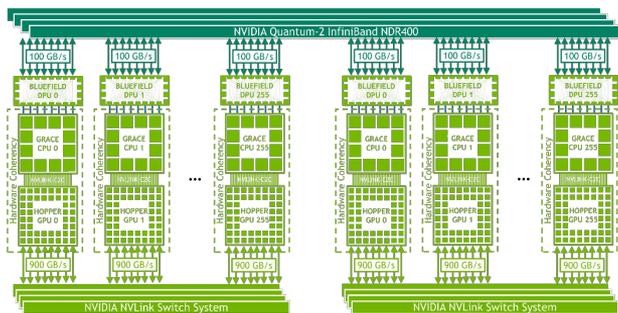




少し脱線

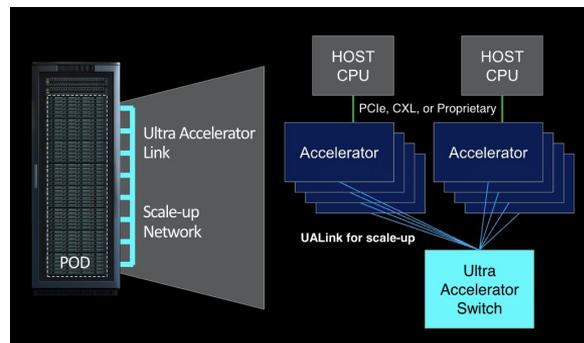
Processor to Processor networks: Future of Connectivity?

NVlink



VS

UALink(Ultra Accelerator Link)



開発者：NVIDIA

主な用途：プロセッサ間のデータ転送

特徴：高帯域、低レイテンシー、576のGPUを接続

拡張性：最大576のGPUを接続(図はGH200)

開発者：Intel, AMD, Broadcom, Cisco, HP, Google, Microsoft, Meta

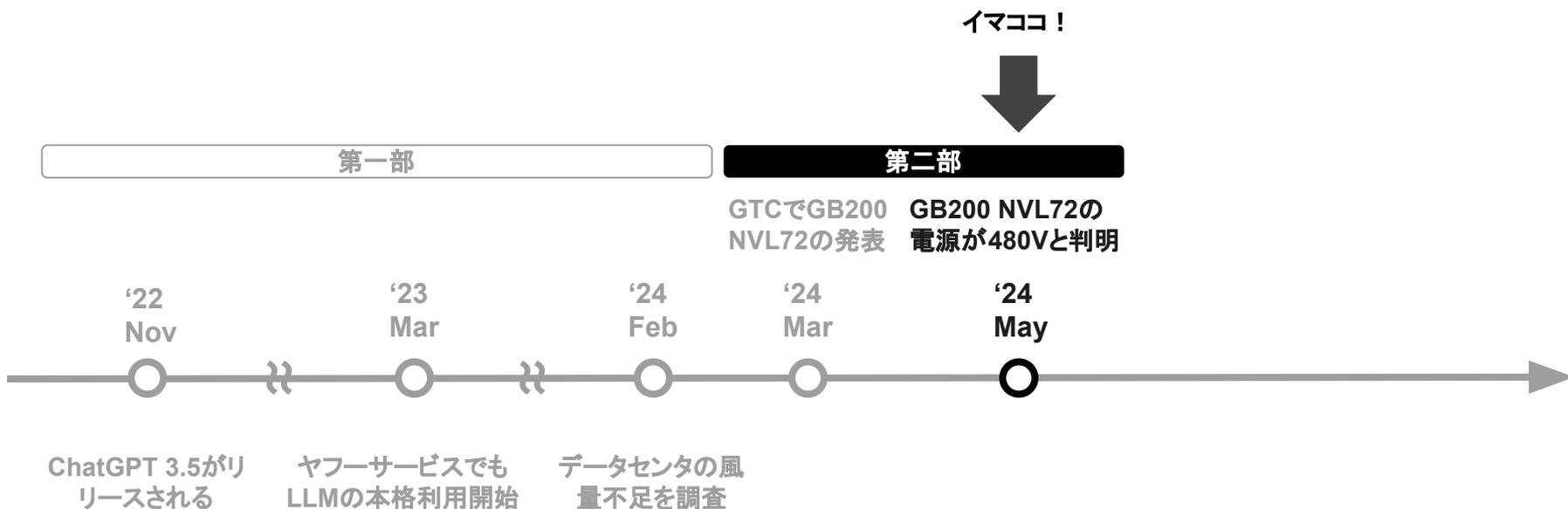
主な用途：プロセッサ間のデータ転送

特徴：高帯域、低レイテンシー、オープンスタンダード

拡張性：最大1,024個のAIアクセラレータを接続

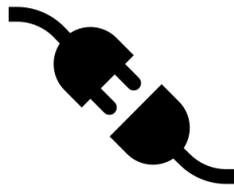
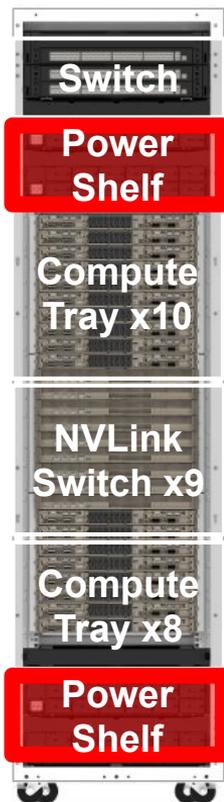


電源の480V化(新しい課題)



GB200 NVL72は1ラック120kwのお化けマシン

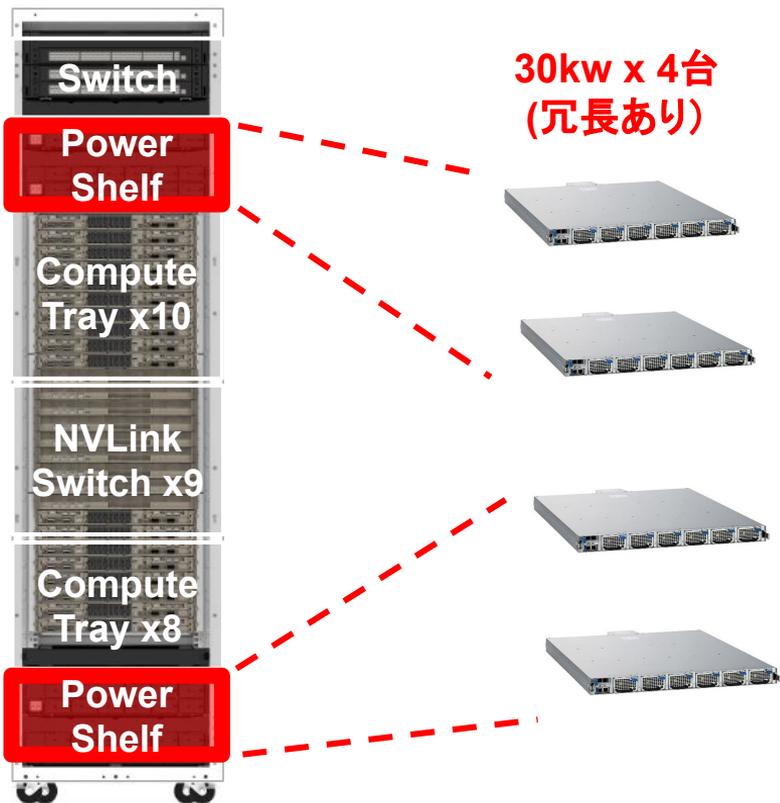
前面図



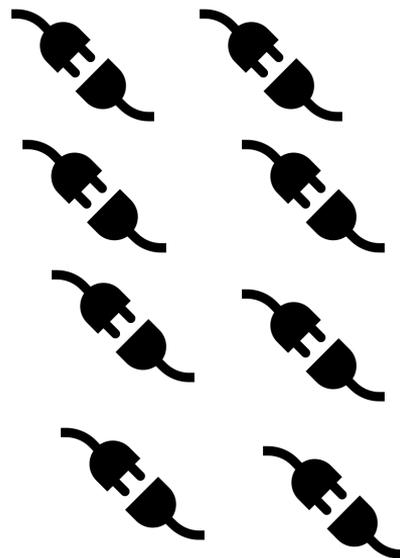
1Rackに詰め込んでいるため
120kw分の高電力が必要

GB200 NVL72は1ラック120kwのお化けマシン

前面図

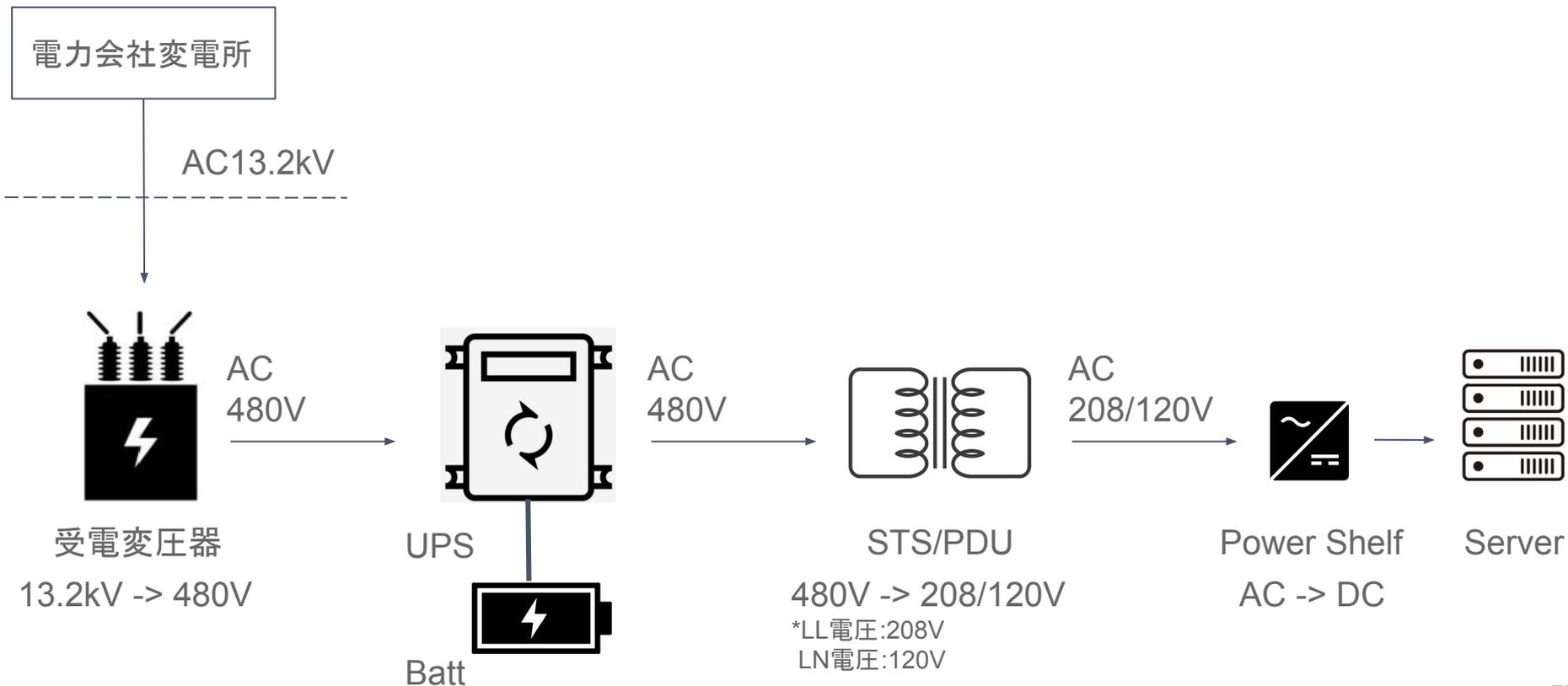


480V x 8本
19.0A
(故障時 38.0A)



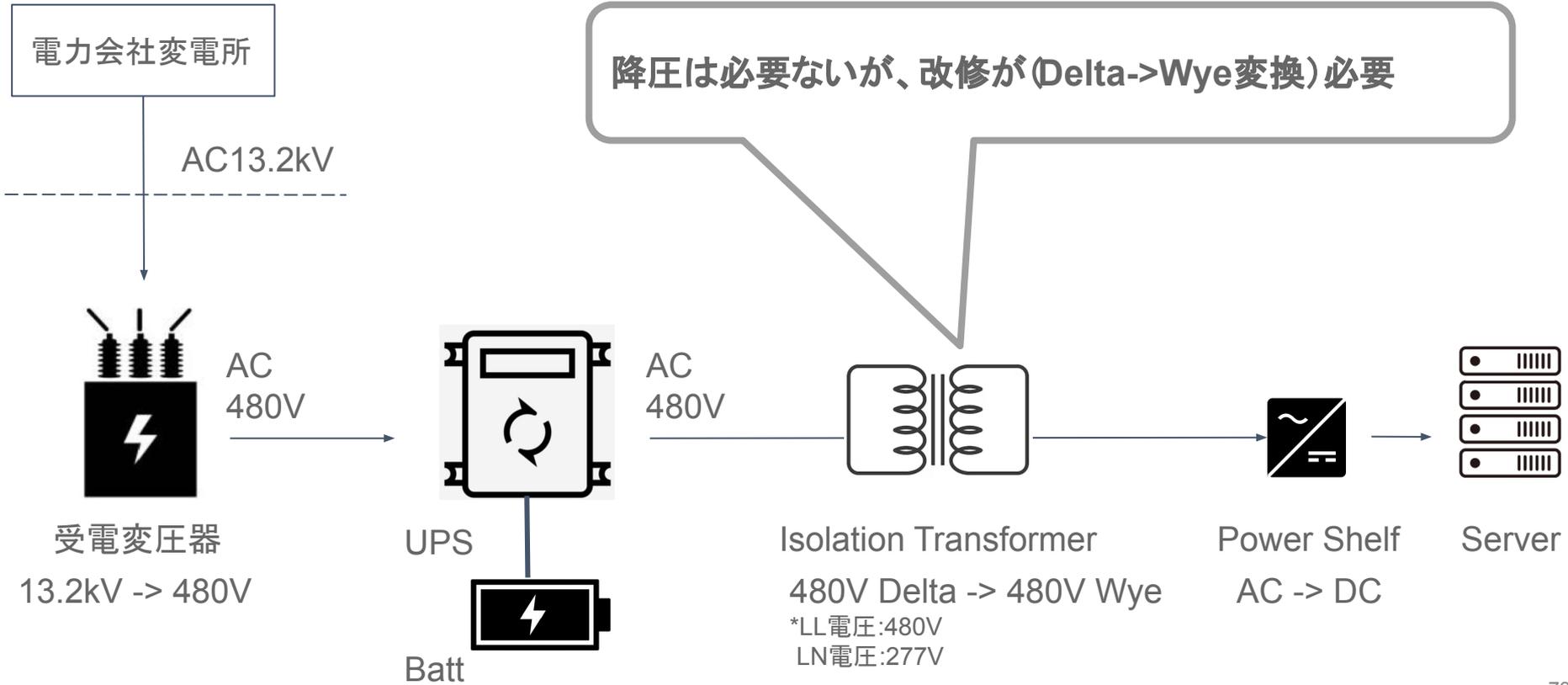
サーバラックまでの接続に高電圧が使用されるため、安全性に関する懸念

電力供給スキーム





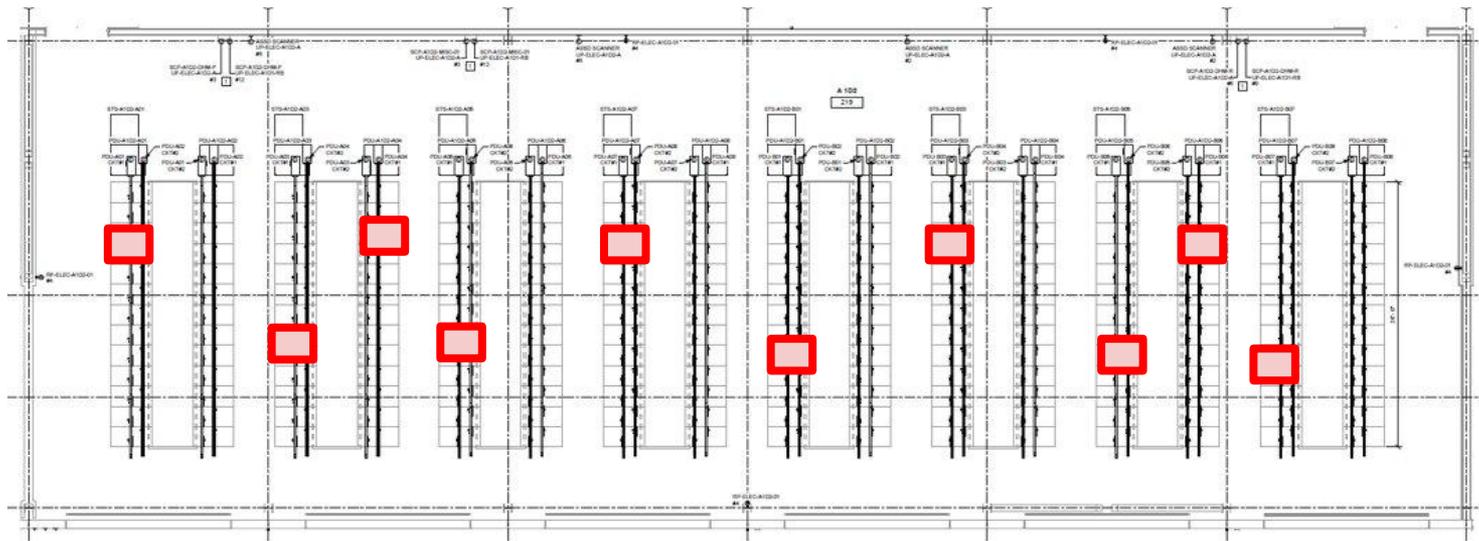
電力供給スキーム 480Vの場合





To be continued 🥲

2MWのサーバルームに10セットくらいしか置けない



スペース効率は悪いが、データセンタの**土地、建物のOpex**は**とても小さい**ので気にしない👩🏻

第二部サマリ

第二部では、

- LLMの作成時に、**モデル時間の短縮が大きな課題**になっている
- GTC(GPU Technology Conference)で発表されたGB200 NVL72は、**この問題を解決できる技術**
- 近接性のためにはヒートシンクでの対応が難しく、効率的な冷却のため**水冷技術(DLC)**が必要
- 1ラックに集積されたサーバは、新たな課題として**電力供給の問題を生み出した**

第一部

第二部

GTCでGB200 NVL72の発表
GB200 NVL72の電源が480Vと判明

'22
Nov

'23
Mar

'24
Feb

'24
Mar

'24
May

ChatGPT 3.5がリリースされる

ヤフーサービスでもLLMの本格利用開始

データセンタの風量不足を調査

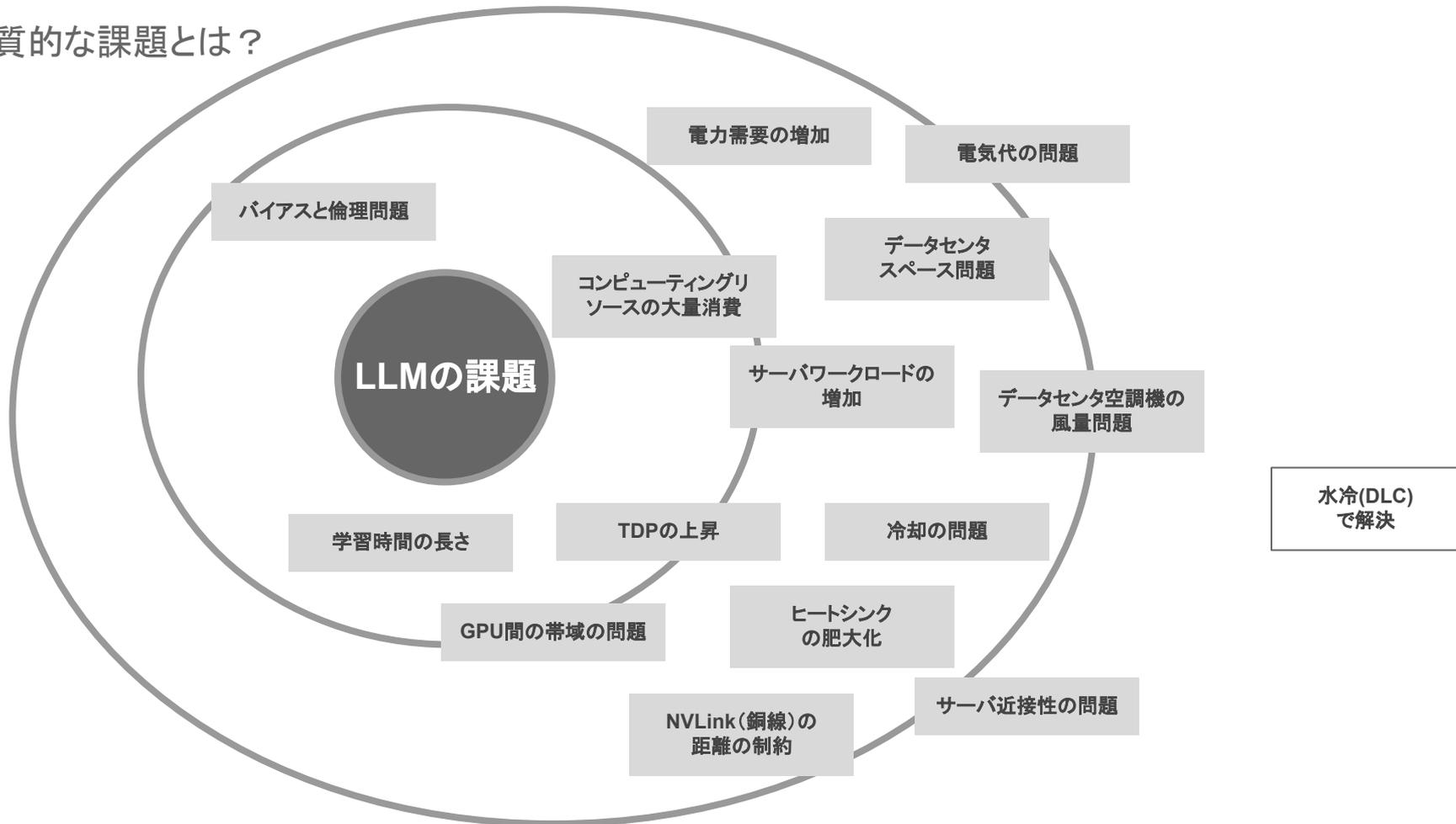
まとめ

- **LLMの普及**は、私たちのデータセンター運用に新たな挑戦をもたらした
- LLMによってサーバワークロードの増加とそれに伴う冷却の課題が発生
- 具体的には**空調システムの風量不足**や**サーバ近接性制約**という問題として表面化

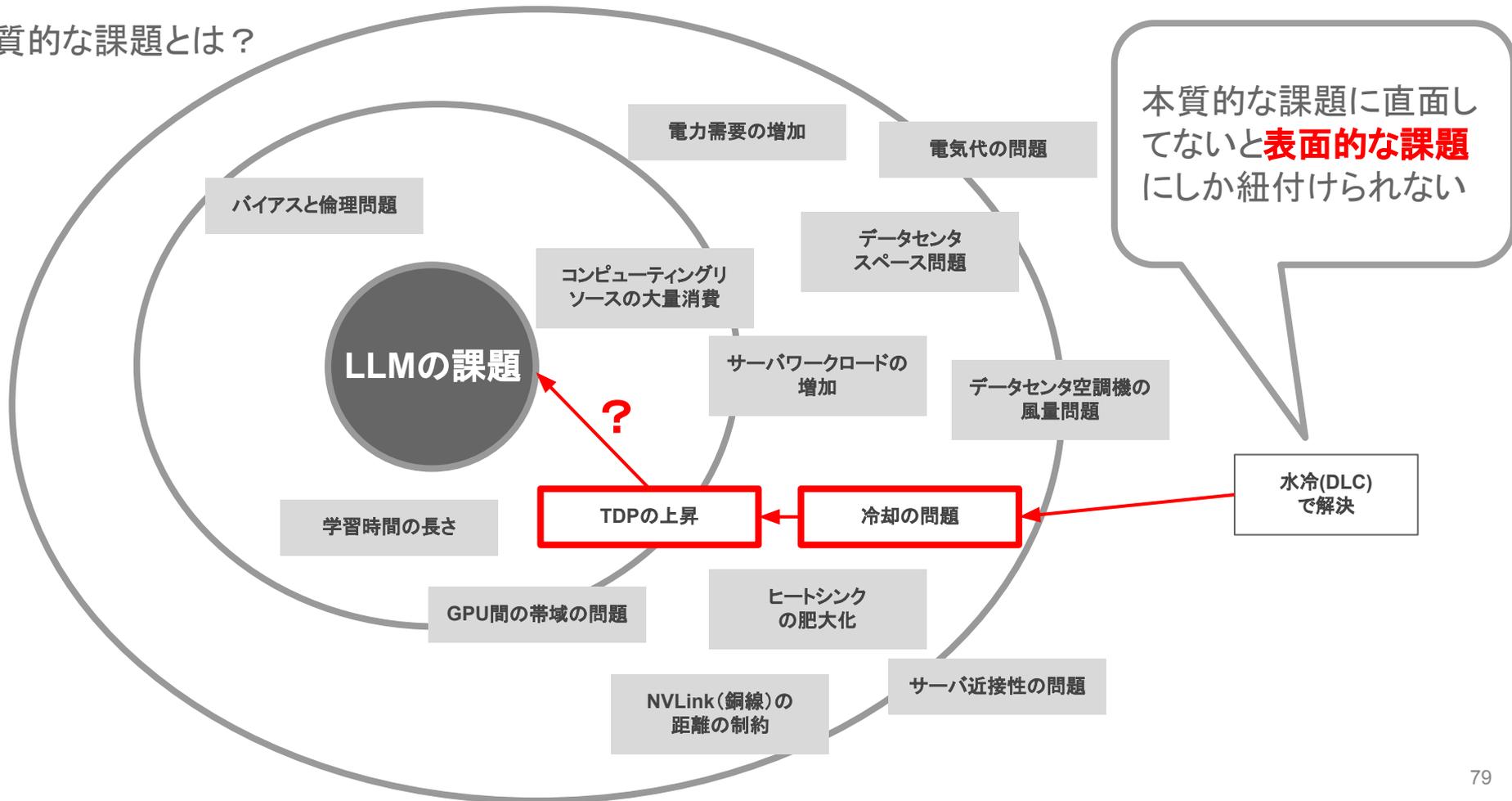


- 私たちは、GTCやOCP Summitなどで米国企業やハイパースケールとの対話を通じて、**水冷技術(DLC)**がこれらの**「本質的な課題」**を解決する鍵であることを慎重に検討しました

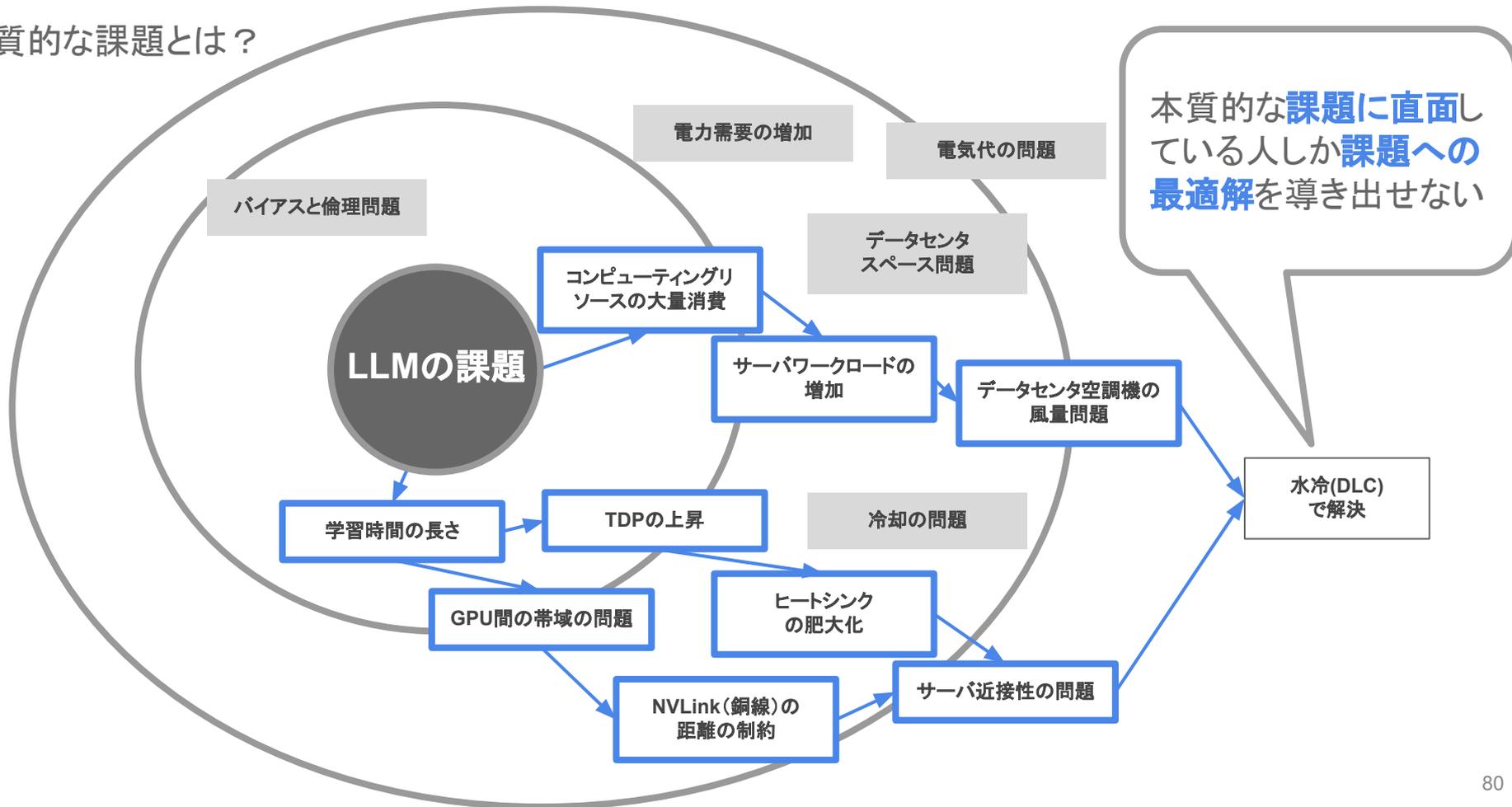
本質的な課題とは？



本質的な課題とは？



本質的な課題とは？



第一部

第二部

GTCでGB200 NVL72の発表
GB200 NVL72の電源が480Vと判明

'22
Nov

'23
Mar

'24
Feb

'24
Mar

'24
May

ChatGPT 3.5がリリースされる

ヤフーサービスでもLLMの本格利用開始

データセンターの風量不足を調査



米国企業(ハイパースケール)などはこの時点で**直面した課題**の解決を進めている

第一部

第二部

GTCでGB200 NVL72の発表
GB200 NVL72の電源が480Vと判明

'22
Nov

'23
Mar

'24
Feb

'24
Mar

'24
May

ChatGPT 3.5がリリースされる

ヤフーサービスでもLLMの本格利用開始

データセンタの風量不足を調査



- ・我々はやっとこの時点でいろいろな課題に**遭遇し始める**
- ・そもそも課題に遭遇するのが**米国企業より遅い**

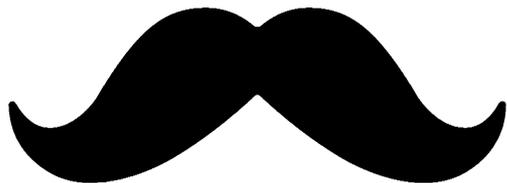
課題に直面してない人が、
これから直面するであろう課題の解決方法を見誤らない方法



- **課題に直面**した人に、**本質的な課題**は何かを情報収集
- LLMで言うと**米国企業などハイパースケーラ**
- 集めた情報を**近く(米国)ですぐに実践**することが大事

- Actapioの課題解決アプローチ
 - 課題に直面した人(米国企業)から情報を収集・検証・実施しているので**最善の技術を選択**できている
- 未来へのビジョン
 - 「**利用する側**」は**後手**にまわる
 - より業界をリードするには「**課題解決の先駆者**」にならないと
 - 私たちが未来に向けて目指していくのは、迅速に技術を取り入れる「利用する側」に徹するのか、「課題解決の先駆者」を目指し業界のリーダを目指すのか。どちらなのでしょう？
 - どちらの道を選ぶにしても、**今日から**必要な取り組みを共に考え、未来を切り拓きませんか？





THANKS

Our Team

Aaron Kirkpatrick
Andrew Arnold
Anthony Skwiat
Atsuko Ishigaki
Eiji Kawauchi
Jason Van Winkle
JD 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
Tsubasa Yoshimura
Yuji Kohata
Yukihito Imai

Our Datacenter

