



大規模AIクラウド計算システム「ABCI 3.0」と そのストレージサービス

産業技術総合研究所
谷村 勇輔



- 経産省「人工知能に関するグローバル研究拠点整備事業」(2016年)の一環として整備
 - 大規模なAIモデル開発の場を提供
 - 産総研が培ったスパコン構築・運用技術をAI計算基盤に応用
- 2018年8月1日：**ABC I運用開始** NVIDIA V100
 - 産学官によるA I 研究開発を加速するオープンイノベーションプラットフォーム
 - 高い計算能力を活用したAI技術の研究開発・実証、社会実装の推進、AI分野の最重要課題への挑戦が目的
- 2021年5月10日：**ABC I 2.0運用開始** NVIDIA A100
 - 経産省「人工知能に関する橋渡しインフラ拡張」によりアップグレードを実現
 - 2024年10月31日運用終了



Expert



LLM構築支援プログラム・ABCIグランドチャレンジ：
画期的な成果が見込まれる最重要課題への挑戦に
ABCIの計算資源を有償・無償で大規模に提供

Advanced & Intermediate



誰でも利用可能
すぐ使えるソフトウェア、データセット、
学習モデル等を提供

Beginner



初学者にも使いやすい統合開発環境を実現
大規模言語モデルハッカソン・ハンズオン：
LLM学習の実践的ノウハウを参加者と共有

「AIを試す場」
人工知能産業のための
オープンプラットフォーム形成
最先端のAI研究から
誰でも試して使えるAIまで

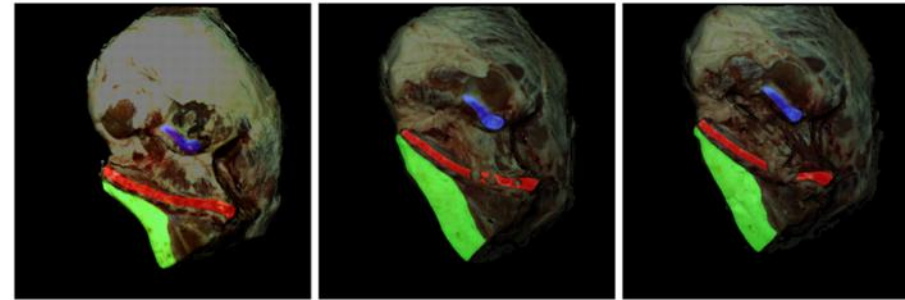
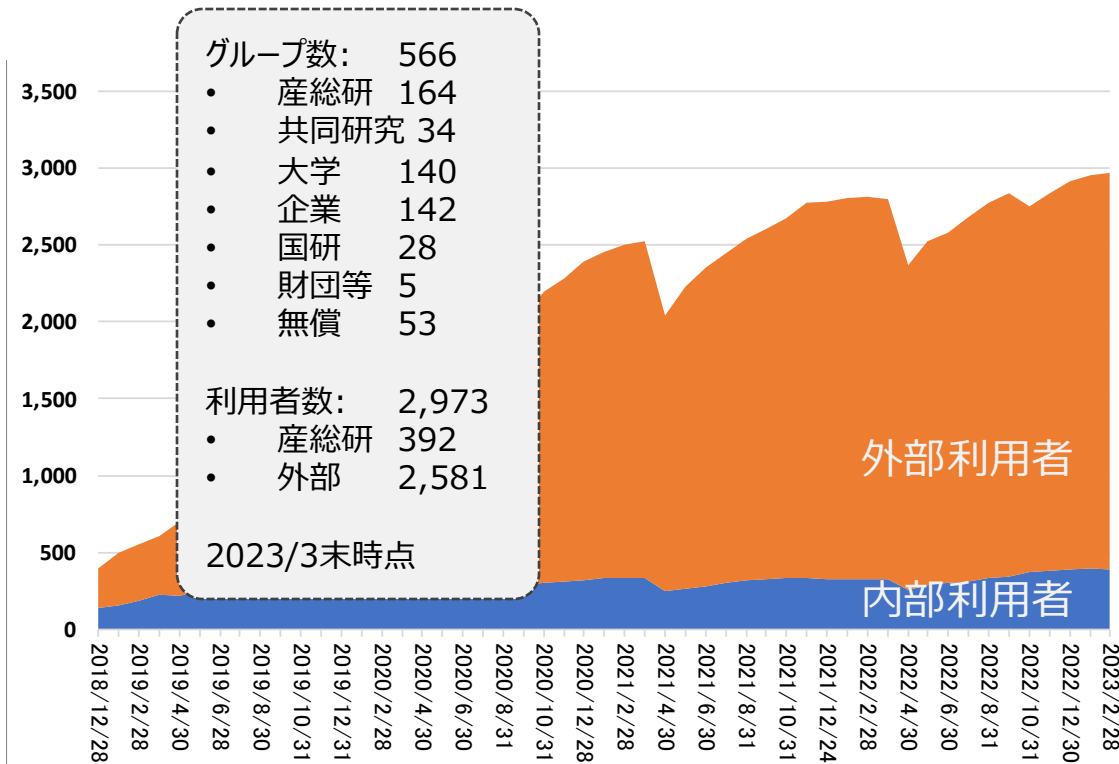


数百の研究機関・大学・企業による利用・協業、
数千の研究者・エンジニアによる利用を促進

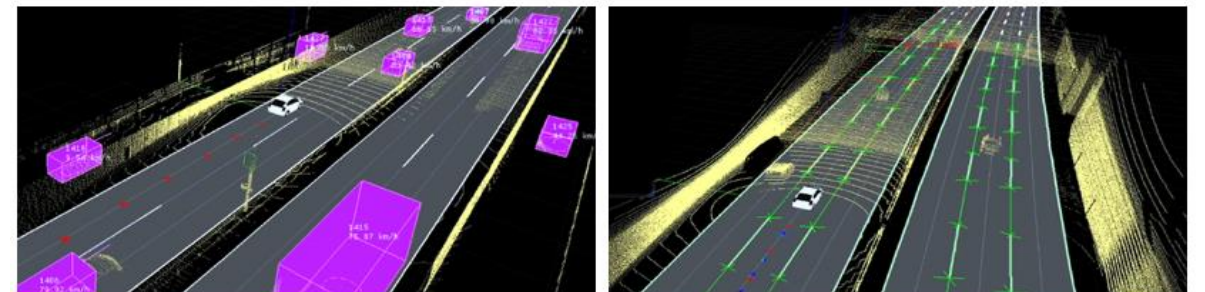
ABCIの利用者数・利用事例

- 利用者数は約3000人（うち外部利用が約87%）
- AIスタートアップから総合電機メーカーまで利用が拡大し、世界に伍する研究が可能になるとともに、我が国のAI研究全体を支える重要基盤に成長

<利用事例（https://abci.ai/ja/use_case/から）>



画像認識で食肉から骨を見つけ出す、食肉加工機械の進化
(株式会社前川製作所)



AI・自動運転技術で新しい物流インフラを構築する
(株式会社T2)

現在（生成AI時代）のABCI

- **大規模言語モデル構築支援プログラム(FY2023)、大規模生成AI研究開発支援プログラム(FY2024)**
 - 生成AI・LLM開発の提案を公募し、有識者による審査の結果採択された提案に、計算資源を優先的に割り当て
 - FY2023: Preferred NetworksのPLaMoやELYZAの日本語LLM、東工大-産総研のSwallowなどの成果を創出
 - FY2024: 産総研、NICT、理研、ELYZAの計6提案を採択



- **国内生成AI人材交流・育成の取り組み (LLMハッカソン, FY2023)**

- 産総研研究者やGPUベンダ(NVIDIA)技術者等がチューターとなり、LLM学習やファインチューニングに関する実践的ノウハウを共有
- 20チーム (民間16、大学3、国研1)が参加

LLM構築支援プログラム採択者 第1回



第2回



ELYZA



現在（生成AI時代）のABCI

- 生成AI開発への機運の高まりを踏まえ、昨年8月より大規模言語モデル構築支援プログラムを開始
- その結果利用率100%が恒常化し、当該プログラム以外の開発利用に大きな障害を来している

2024年度

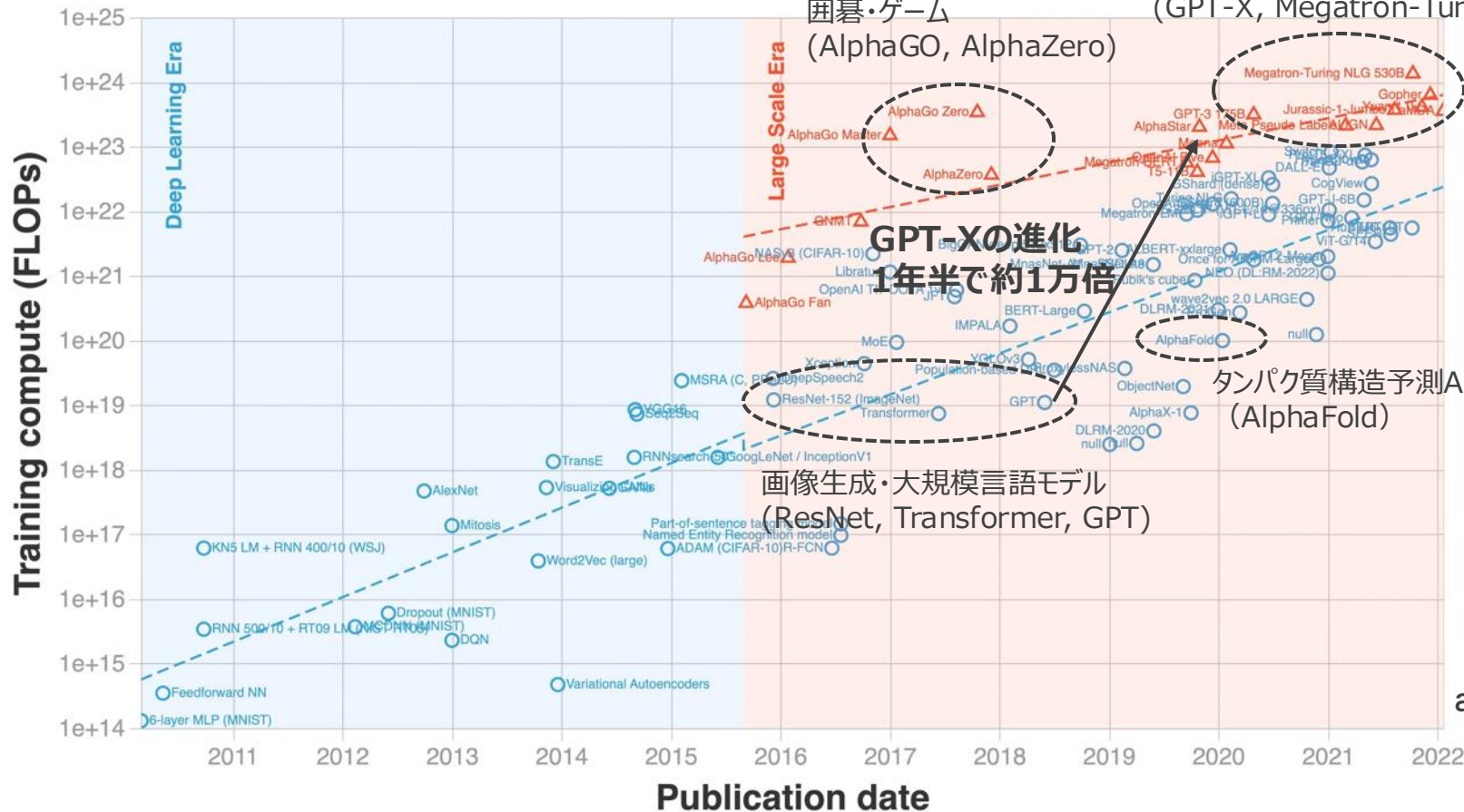


基盤モデル構築に必要な計算量とABCIの性能のギャップの広がり

- 世界的に競争力のある生成AIの開発・活用には計算資源の抜本的な拡充が必要
- 世界最先端の基盤モデル開発には、ABCI 2.0の10倍以上計算能力が利用されたとされる

Training compute (FLOPs) of milestone Machine Learning systems over time

n = 99



計算量の増大：3年で10~100倍

(出典) Sevilla et al., Compute trends across three eras of machine learning, 2022

- **生成AIモデルをはじめとした最先端AI技術の研究開発能力の強化**を目的として、公的計算基盤として国内最大規模の計算能力（6.22EFLOPS、従来比約7倍）を整備、国内の産学官に提供
 - 経済産業省「生成AIの基盤的な開発力強化に資する計算資源の整備」（令和5年度補正）の一環として整備
- スケジュール
 - 2024年11月中旬：**一部システム（半精度演算性能0.89 EFLOPS相当分）**の試験運用を開始
 - 2024年12月末：全系導入完了
 - 2025年1月中旬：**一般提供を開始**予定

ABCI 3.0 ハードウェア構成

Compute Node(H)

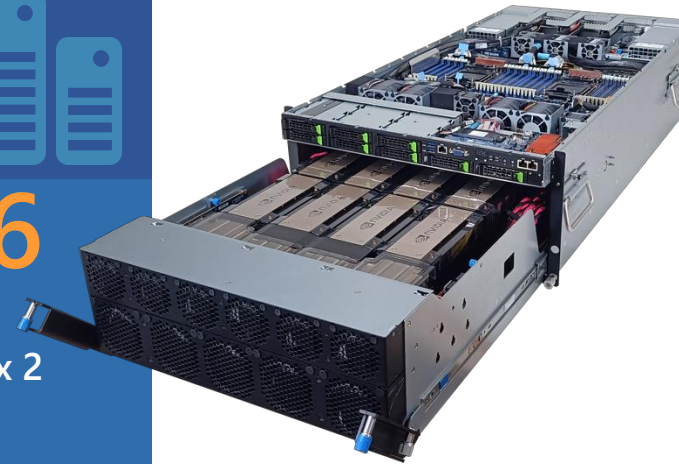
6,128 GPUs , 73,536 CPU cores
1.53 PiB Memory , 11.77 PB NVMe SSD



Node cfg.

× 766

- GPU NVIDIA H200 SXM5 (141GB) x 8
- CPU Intel Xeon Platinum 8558 (2.1GHz/48cores) x 2
- Memory DDR5 2 TiB
- Local Storage NVMe SSD 7.68TB x 2
- Interconnect InfiniBand NDR (200 Gbps) x 8



Compute Network (InfiniBand NDR)

Storage Network (InfiniBand HDR)

Service Network (10/100G Ethernet)

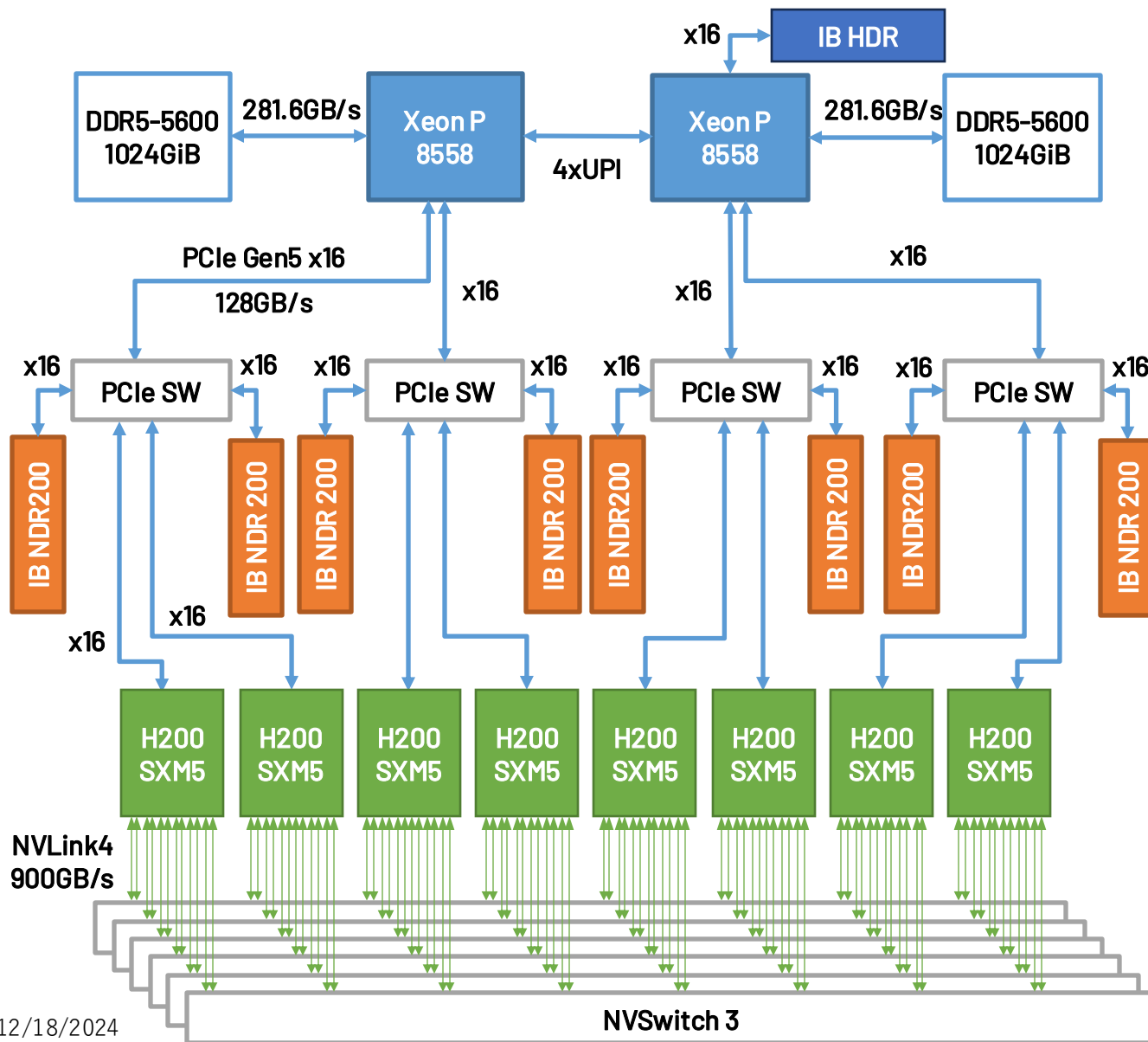
All Flash Storage System
75 PB (Lustre FS, S3)



Interactive Node

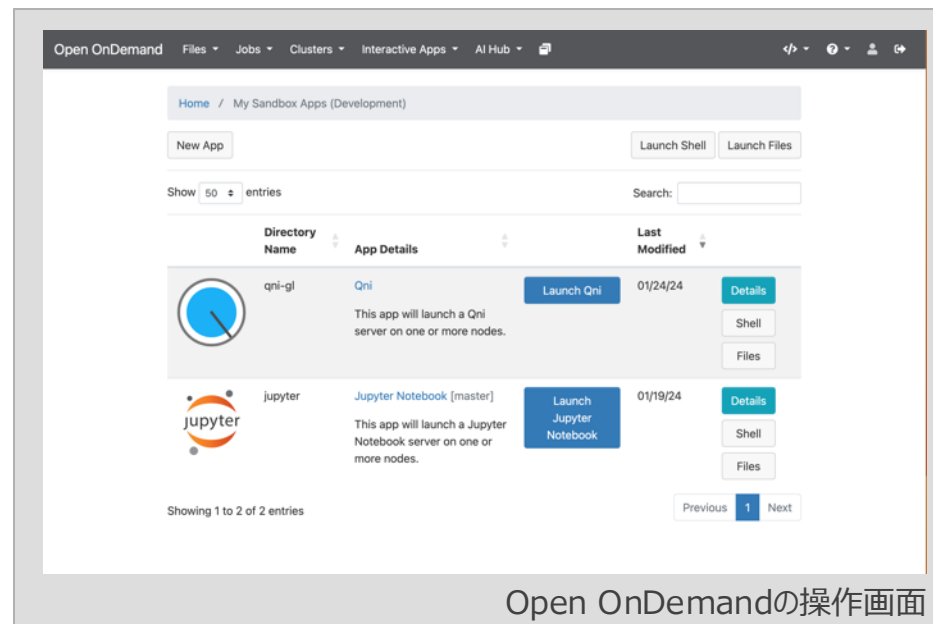


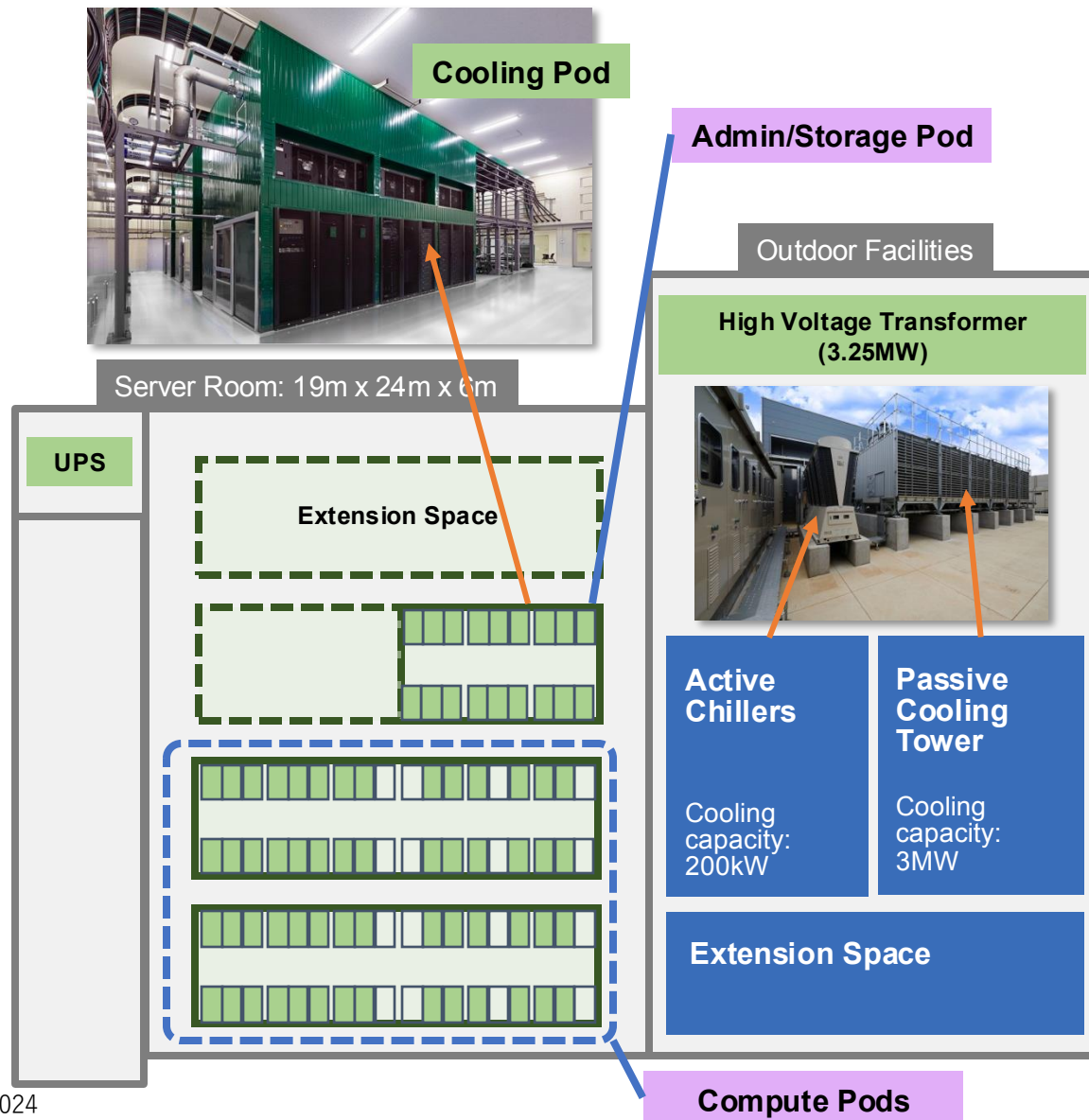
ABCI 3.0 計算ノード詳細



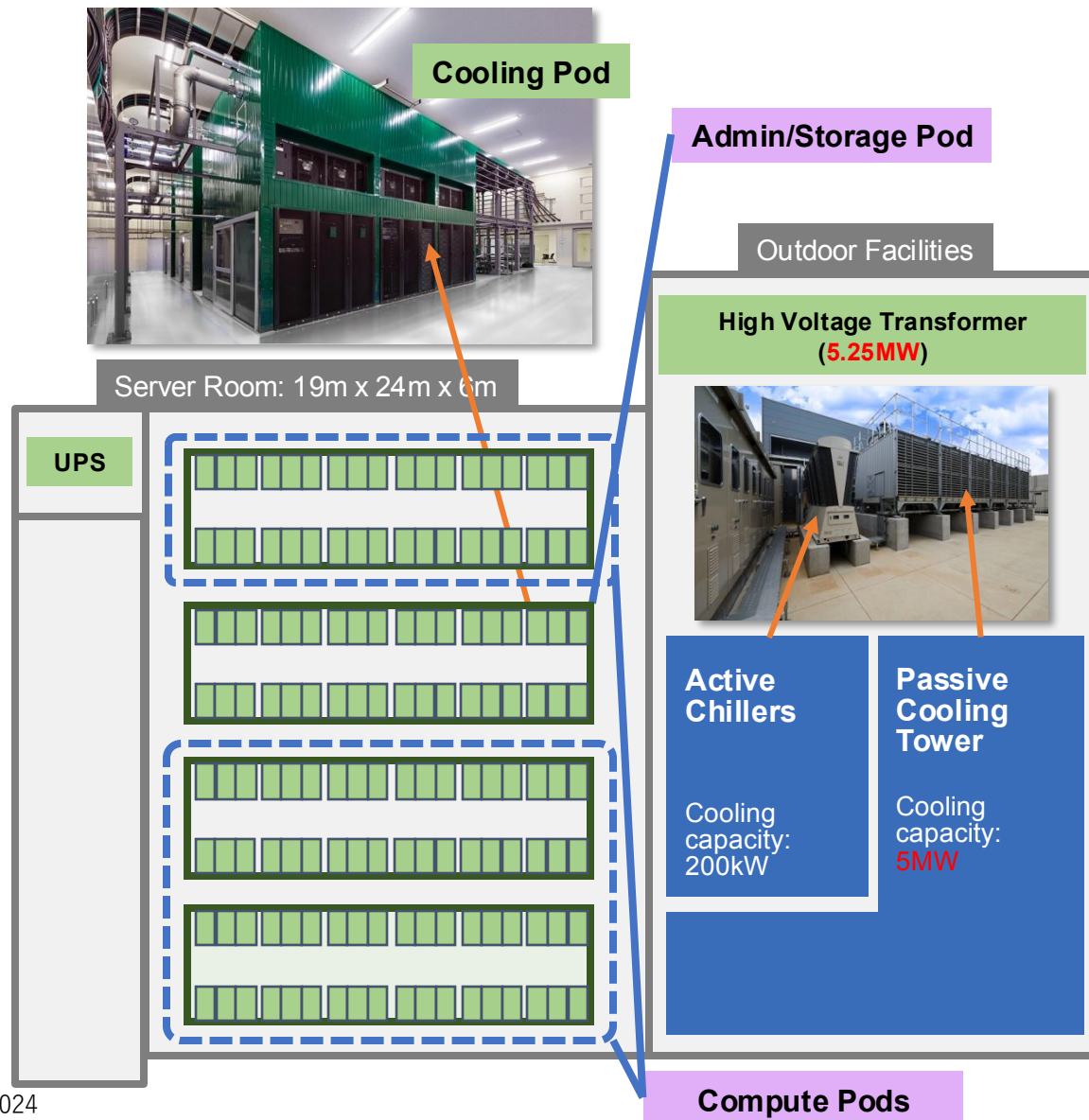
HPE Cray XD670 Server (5U)	
CPU	Intel Xeon Platinum 8558 (Emerald Rapids, 2.1GHz, 48core) x2
GPU	NVIDIA H200 SXM5 141GB x8
Memory	2048GiB DDR5-5600 DIMM
SSD	U.3 NVMe SSD 7.68TB x2
Compute Interconnect	InfiniBand NDR200 (200Gbps) x8
Storage Interconnect	InfiniBand HDR (200Gbps) x1
Cooling	Air cooling

- ABCI 2.0と同等の、利用者にとって使いやすい利用サービスを提供
 - ジョブスケジューラは PBS Pro (≠ Altair Grid Engine)
 - バッチ/インタラクティブジョブ、ノード予約、仮想的なノード分割を提供
 - コンテナランタイムとしてSingularity CE (or PRO)を導入
 - AI開発に必要な基盤ソフトウェア (CUDA, NCCL, MPI等) をEnvironment Modulesで提供
 - 複数バージョン提供、必要なものを選択し、切り替えて利用可能
- Open OnDemandも正式サービスとして提供





- ラック数
 - 計算ノード用：72 (うち 56 を使用)
 - 管理・ストレージ用：18
- 電力容量：3.25 MW
- 冷却容量：3.2 MW
 - 冷却塔：3MW
 - チラー：200kW

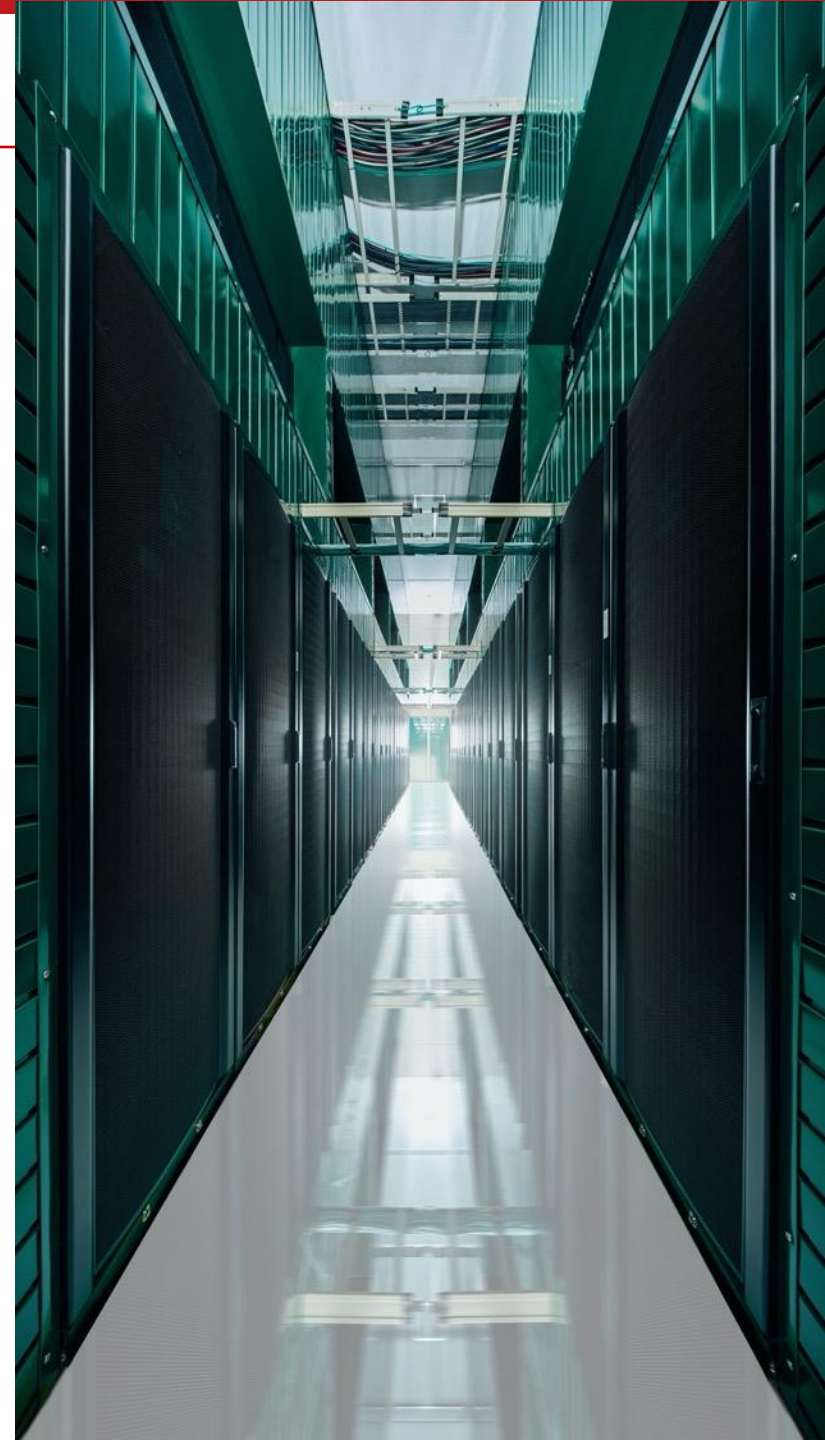


- ラック数
 - 計算ノード用 : 108
 - 管理・ストレージ用 : 36
- 電力容量 : 5.25 MW
- 冷却容量 : 5.2 MW
 - 冷却塔 : 5MW
 - チラー : 200kW
- 12月上旬に電力・冷却容量増強完了

ABCI 3.0のストレージサービス

ストレージの特徴

- 計算ノードローカルストレージ
 - 7.68 TB NVMe SSD x2
 - 利用形態（検討中）
 - 単純なローカルスクラッチ
 - Lustre キャッシュ（Hot Nodes） - Read-only
 - BeeONDによる、複数ノードのローカルストレージを束ねた、一時的な共有ファイルシステム
- 共有ストレージ
 - 75PBオールフラッシュのLustreファイルシステム
 - QLC SSDを採用
 - S3互換オブジェクトストレージ機能を提供
- NVMe SSD、Lustre の両方において GPU Direct Storage の有効化（検討中）



- 3つの領域

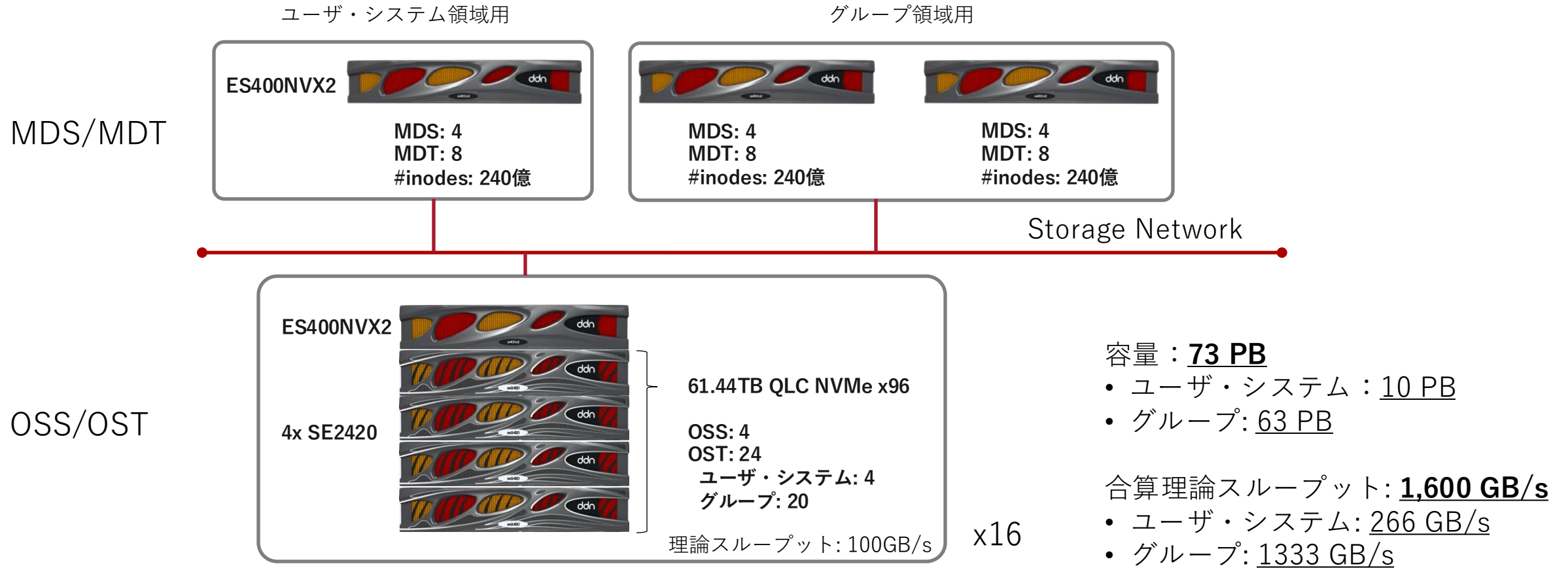
名称	# of inode	Size (PB)	用途
ユーザ・システム領域	240億	10	ユーザのHOME (1TB/ユーザ) システム領域 (管理ソフトウェア、アプリ等格納)
グループ領域	480億	63	グループ領域 (1TB単位で利用できる、グループ内共有領域)
オブジェクト領域	240億	1	S3アクセス領域

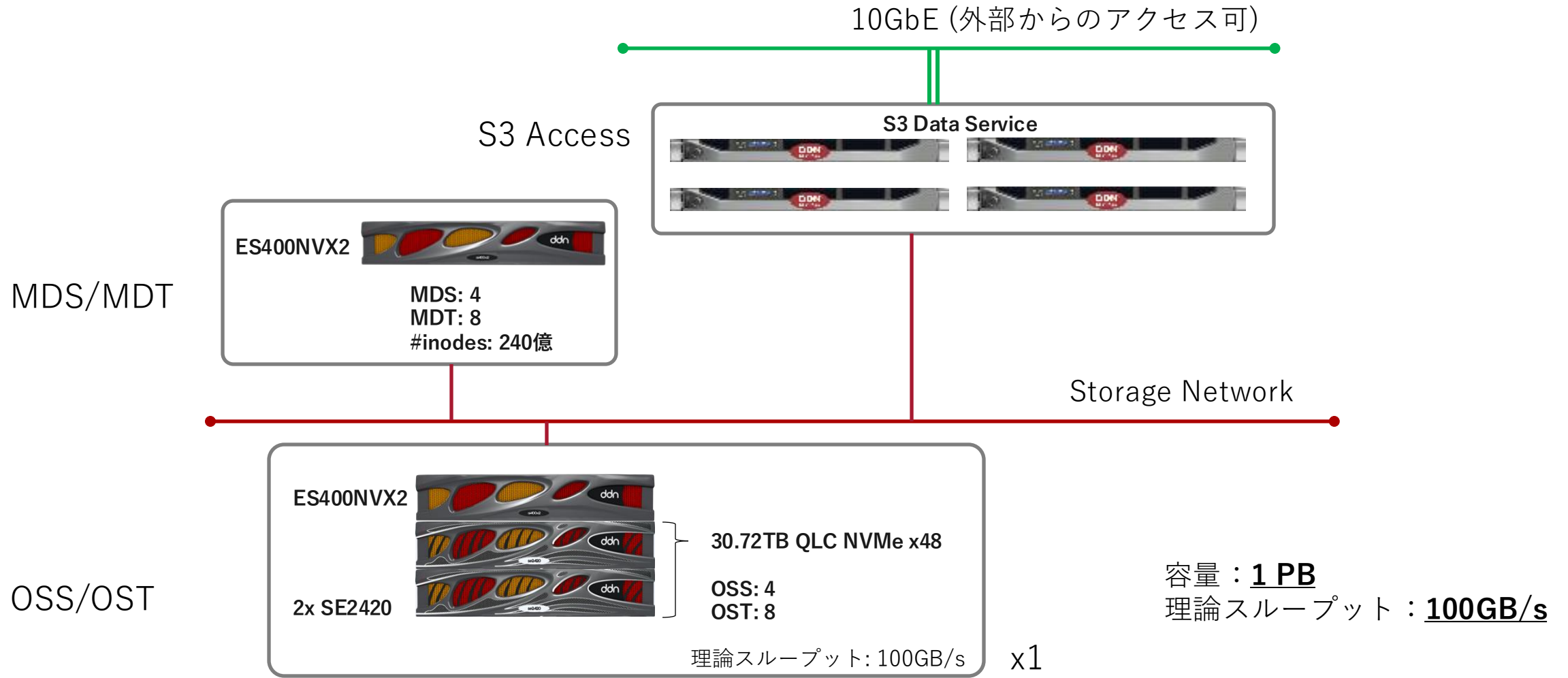
- ユーザ・システム領域、グループ領域の運用形態はABCI 2.0同様

- DDN社アプライアンスを導入

- Lustreアプライアンス: ES400NVX2
- ディスクエンクロージャ: SE2420
- S3アクセス: S3 Data Service

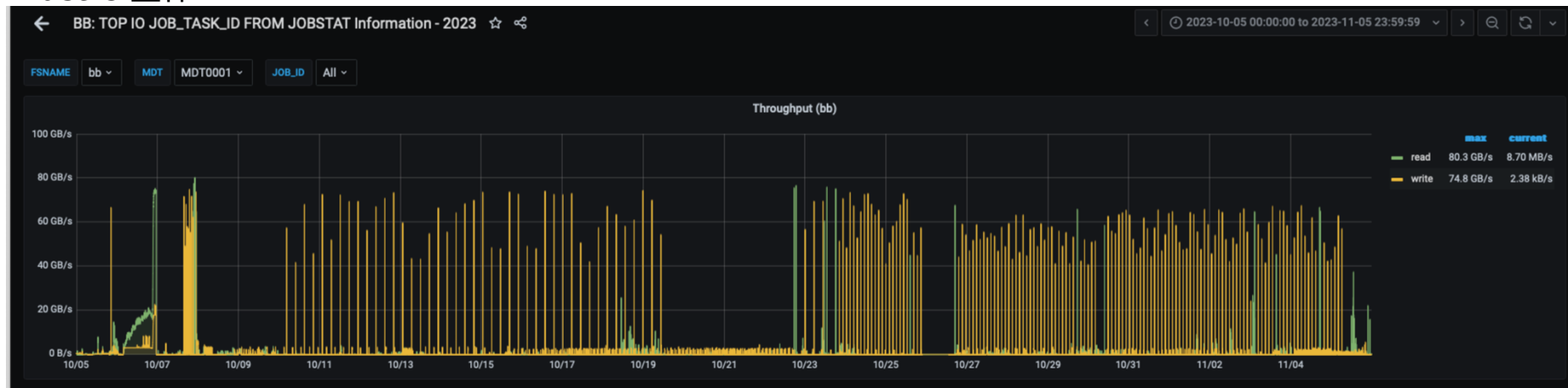
ユーザ・システム領域 と グループ領域



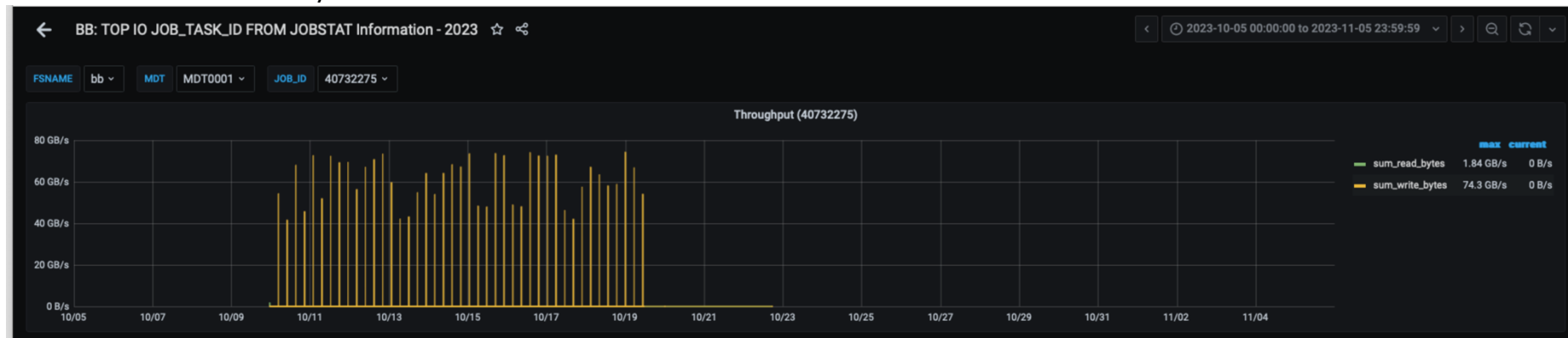


- ファイルシステム全体の利用統計をとるだけでなく、ジョブ毎のI/O性能を取得して活用
 - ジョブ実行履歴、計算ノードやネットワーク等のモニタリング情報と照会し、生成AI、基盤モデル構築による GPU 資源の利用効率化を図る

Lustre 全体



特定1ジョブのLustre I/O



ストレージ要件の背景

- オールフラッシュ & 単一ストレージ
 - SSD/HDD 混在は使いにくい
 - 全利用者に対して安定的に高い性能を提供したい
 - 省スペース化を図りたい
 - 複数領域に対する容量割当（inode 割当含む）の手間を削減したい
- ローカルディスクの活用
 - 大規模 LLM 開発でもローカルディスクの利用率は低い
 - 利用者自身でのステージングはなかなか浸透しない
- オブジェクトストレージの運用コストと公平な利用者負担
 - 性能や安定性に大きな問題はないが、保守運用費が安いわけではない
 - 有償利用があまり増えなかった

これらがどれくらい改善されたかは、1月中旬以降のサービス開始後に検証していく予定

- ABCI 3.0
 - 11月から一部システムの試験運用を開始し、12月までにアップグレード完了予定
 - 基本的な利用サービスと半精度演算性能0.89 EFLOPS相当の計算ノードを提供
 - 2025年1月中旬までに「ABCI 3.0」として一般提供を開始
 - 6000基強のNVIDIA H200 GPUを搭載
- Lustreに関して、75 PBのオールフラッシュのファイルシステムを導入
 - 容量、性能は従前のABCI 2.0の倍以上としつつ、同等の使い勝手を提供
- 多数の最新のGPU、高速・大容量ストレージにより、生成AI、基盤モデルの研究開発を支援



<https://abci.ai/>

