

# DDNのHPCストレージへの貢献と 新しい課題への取り組み

---

DDN Storage

井原 修一





## アジェンダ

1. DDNについて
2. HPCストレージのトレンド
3. ExascalsierにおけるHPCストレージの課題に対する取り組み
4. 次世代のストレージシステム“DDN RED”について



# DDNは世界最大手の 非上場ストレージベンダー

世界規模のマーケットリーダー

- 20年以上に渡る業界リーダーシップ
- 10,000以上の顧客
- 世界10拠点にテクノロジーセンターを設置

## At Scale | Enterprise

AIストレージシステム

大規模HPCストレージ

スケーラブルデータ管理

仮想環境

ソフトウェアデファインド統合ストレージ

高性能統合ストレージ

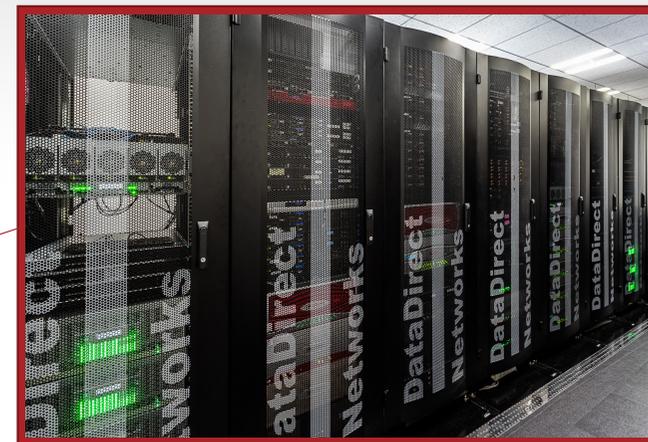




## 株式会社データダイレクトネットワークスジャパン (米国DDN社100%子会社)、2008年1月設立

- ストレージのフルスタックソリューションを提供する営業、エンジニアリング集団
- 高性能ストレージシステムの販売、設計、構築、保守
- 販売パートナー: 富士通、日本電気、日立製作所、SCSK、日本コムシスなど

© 2020 DDN



### 日本国内のストレージ導入実績

合計1.4エクサバイト(7年間)





# DDNにおける20年間の技術革新

技術と需要

分散ストレージ: EXAScaler NFS/SMB, POSIX & Object 超高速で高効率

データ完全性、デクラスタリング - データ保護、GEO分散、スケーリング

Flash、NVMe - 性能、耐久性、メモリクラス

仮想化、コンテナ - Openstack, Docker, Kubernetes...

Hybrid Cloud - オークストレーション, AI, SW/Cloud

Enterprise & MultiCloud - 簡素化 & 管理

S2A - ASIC  
システム

EXAScaler - 並列  
ファイルシステム

SFA - スケーラブル  
ストレージエンジン

エンベデット  
システム

RED - エラスティック  
データサービス

2001

2003

2013

2014

2017

2018

2019

2020

DDN設立

最初の大規模  
並列ファイル  
システム導入

世界最速のシステムが  
DDNのシステム採用

世界最大のプライベート  
ストレージ会社に

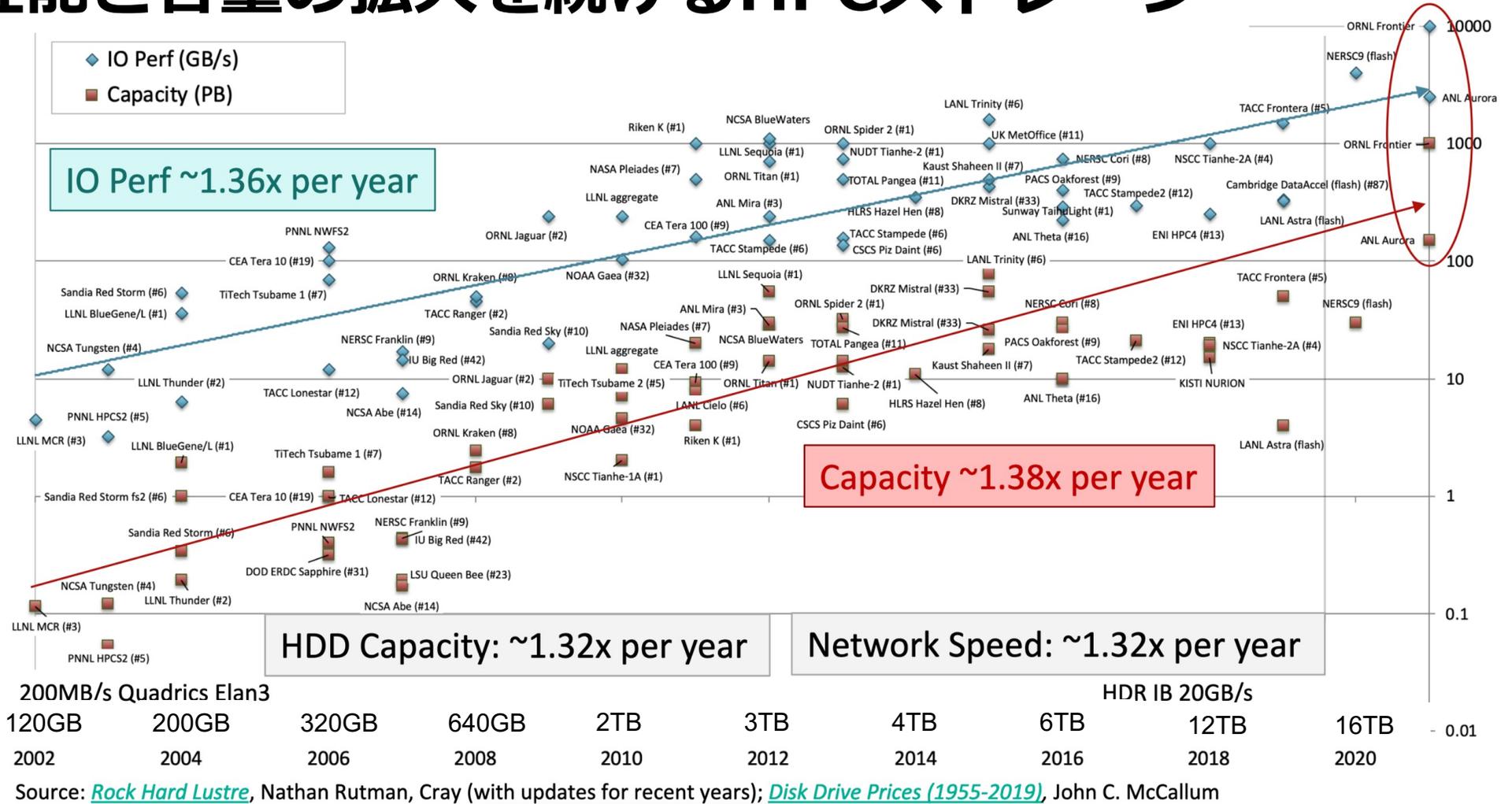
世界最速のNVMe  
システム導入

IntelからLustre  
部門を買収

At-ScaleとEnterprise  
部門を発足

世界最速のAIスパコンが  
DDNを採用

# 性能と容量の拡大を続けるHPCストレージ





## DDN AI400Xが支える NVIDIA SuperPODとDGX A100

DDNストレージがスケーラブルなソリューションであることが  
最大のSuperPODとDGX A100のプロダクション環境で証明

### Top500のNo7にランキングされるSelene

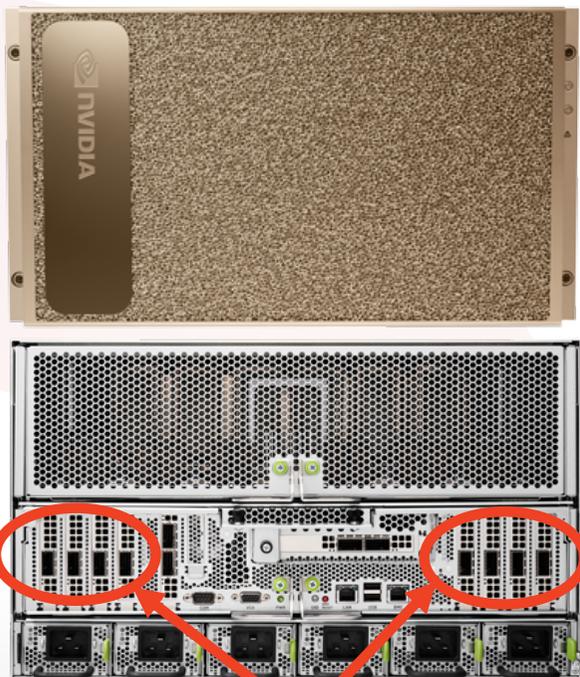
- 280 NVIDIA DGX A100 システム
- 170 Mellanox 200G HDR スイッチ
- 40 x DDN AI400X アプライアンス
  - 最大2TB/secのスループットと1.2億 IOPS
  - 最適かつ完全にDGX A100向けにインテグレーション
  - 最初の10システムを4時間で導入し、その後シームレスな拡張を提供

NVIDIA社によってテスト、検証されたAI400Xは 非常に簡素化されDGX A100, DGX PODまたはDGX SuperPODのいずれにおいても利用可能





# Nから1へダウンスケール、性能はN倍



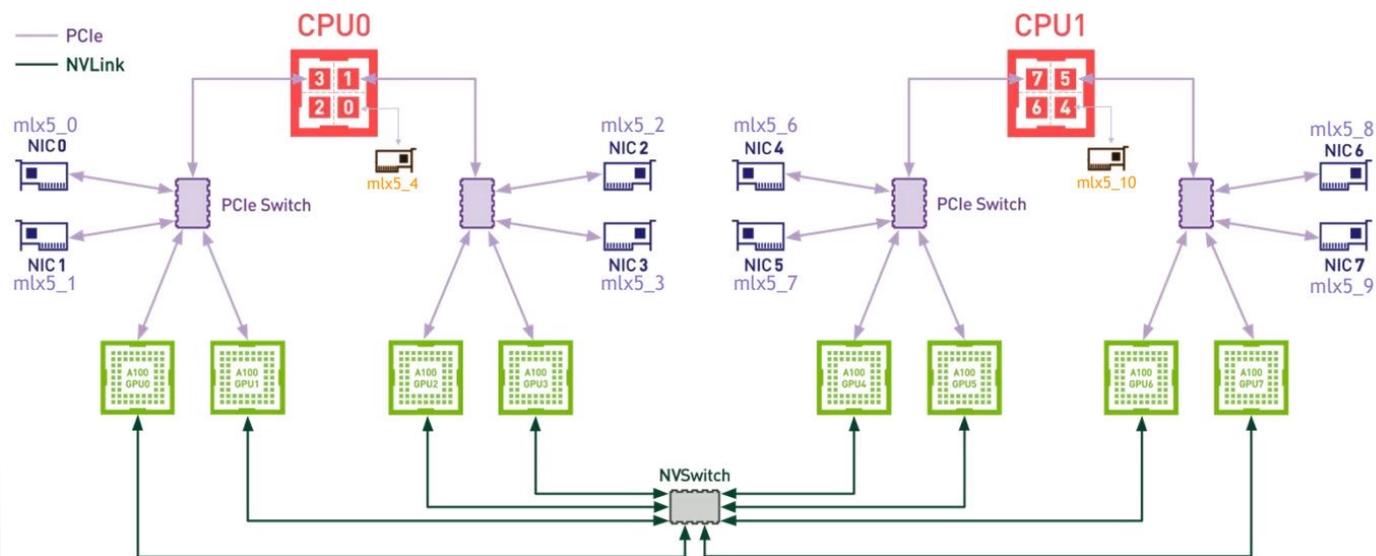
8 x IB-HDR200+  
2 x IB-HDR200

© 2020 DDN

## DGX A100

### High-level Topology Overview (with options)

Data plane (can be used as eth or IB)  
Compute plane (IB)



I/Oネットワーク帯域は一般的なクライアントの20倍以上\*

\*IB-EDRを比較した場合

# 自動運転の研究開発基盤にDDNを採用

## DDN EXAScalerを用いてハイパースケールR&Dを実用化(米国,2016)

- 全米規模のハイパースケール自動運転ソリューション
- 全体ストレージ容量は350PB以上、性能は1.6TB/s以上
- 各リージョンごとに非常に大きな容量を持つ単一ネームスペースに対する性能を、コスト効率の高いアプローチで実現
- 完全な大規模ストレージ管理環境を提供

## 大手自動車メーカー向けにモジュール型でスケラブルなAIストレージ(欧州, 2019-2020)

- 1つのモジュールは27PB、200GB/sec
- 複数サイトの複数モジュールにて構成
- システムの性能、状態だけでなく、ユーザとワークロードのIOアクティビティを詳細に把握可能





# 指数関数的に増加する大規模データ処理の需要

## AI & Analytics

自立型IT のマーケット規模は2025年に\$25B (USD)

75%のAIシステムを採用しているところは今後もAI環境のための予算が増加

## Government/Academia

\$150B (USD)規模がCovid19 関連の研究のために利用

データインテンシブのHPCでは 7% 平均成長率

データ解析17%の成長率

## Web/Cloud

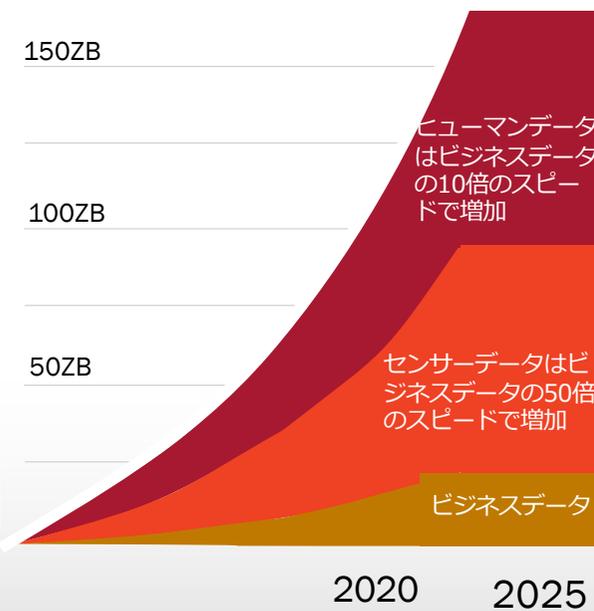
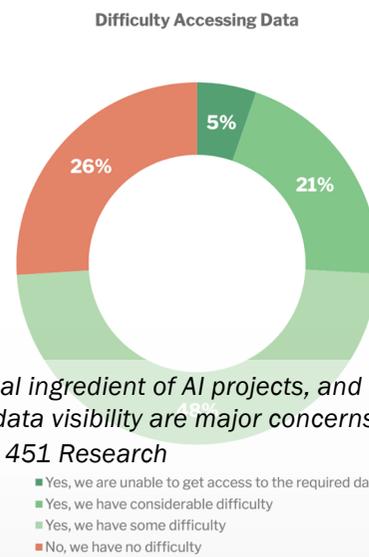
2025年には世界の約半分のデータがクラウドに

2025年にはインターネットユーザーデータは300%の増加

## Enterprise at Scale

データ戦略における迅速な採用。ただ、74%がデータアクセスの困難さに直面。

爆発的に増加するヒューマンデータとセンサーデータ





## SSD Expands, but HDD Does Not Die - *Gartner*

- DDNは大規模HPCに多くのNVMeのストレージを提供
- ガートナーの予測では:
  - 2030年、エンタープライズにおけるPBの70%はHDD
    - 技術とサプライヤーの経済活動
    - HDDとSSDのギャップは供給とコスト/GB
  - これまでは~25%で推移していた 2020年のNANDベースSSDの生産コストの削減は10-15%
  - 2030年までにすべてのエンタープライズストレージがNANDベースになるのは懐疑的
- 今後もHDDとNVMe/SSDとのハイブリッドは非常に重要



## 全てのデータを1つの場所へ

- 従来の“ビッグデータ” HDFSやデータレイクは限界
  - HDFSはバッチ処理、データレイクはコールドデータを想定
  - これまで全体の1%が、今後は30%以上のデータ解析が必要
  - データと計算リソースの分離は加速。データ領域は共有され、多くの他種多様の計算解析エンジンで処理。
  - データ保護、コンプライアンスは計算リソースと独立
  - 計算リソースの導入、拡張サイクルはストレージと異なる
- 新しいパラダイムは、単にデータをインGESTして共有するだけでなくどのようにデータから価値を生み出すか。リアルタイム性も要求 (Kafka, Spark等)
- ファイルやオブジェクトに対する高いIOスループットとIOPS
- スケールアウトデザイン、多次元的な性能、大規模な並列アーキテクチャ



# DDN EXAScaler

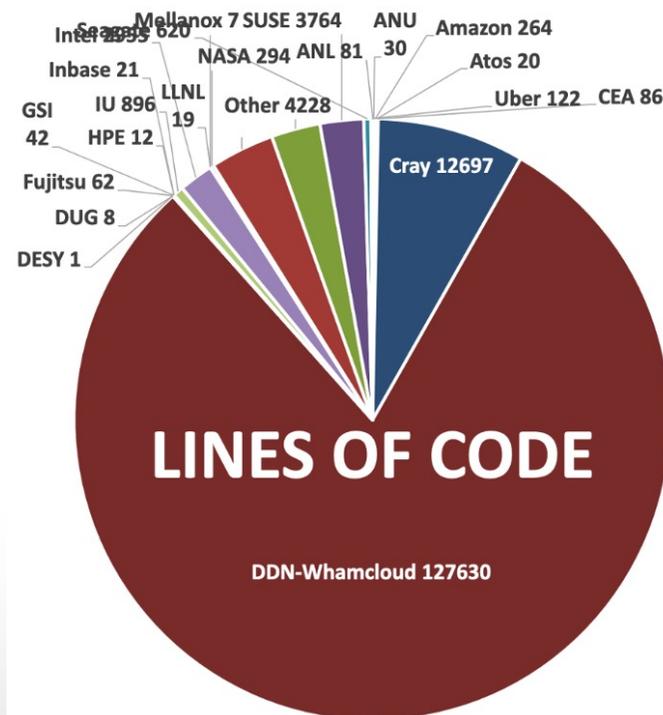
HPC/HPDAストレージへの課題への取り組み



## HPC/AIの大規模なストレージ環境においても データ管理、共有が要求される時代

- 高速なデータ解析のためのデータアクセラレータ
- データ管理とデータ相互交換
- 様々なプロトコルでのデータサービス
- クラウドシステムとの連携
- マルチテナントとセキュリティ
- データアクセスにおける性能の見える化

# EXA5

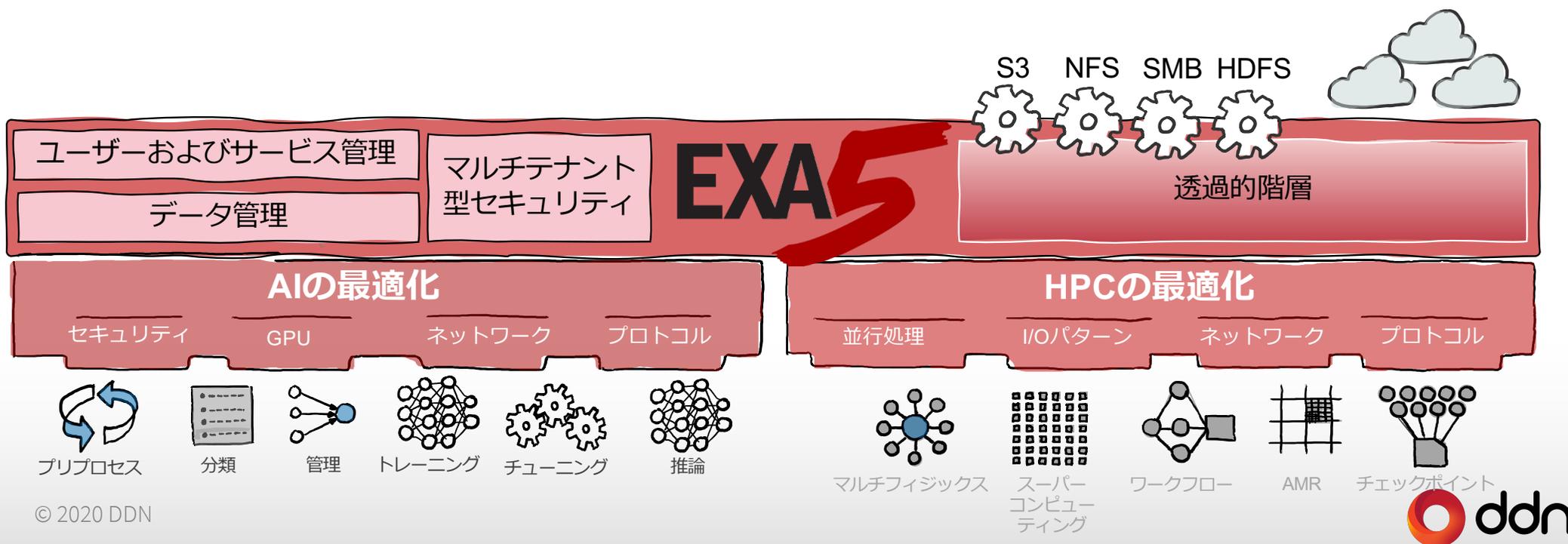


Lustre-2.13への貢献

# EXA5 (ExaScaler5)

## AIおよびHPC向けに最適化されたストレージソフトウェア環境

- ▶ AIとHPCの双方に対する詳細な最適化により最高の効率性と適正な機能を実現
- ▶ データを必要な場所に必要なタイミングで提供

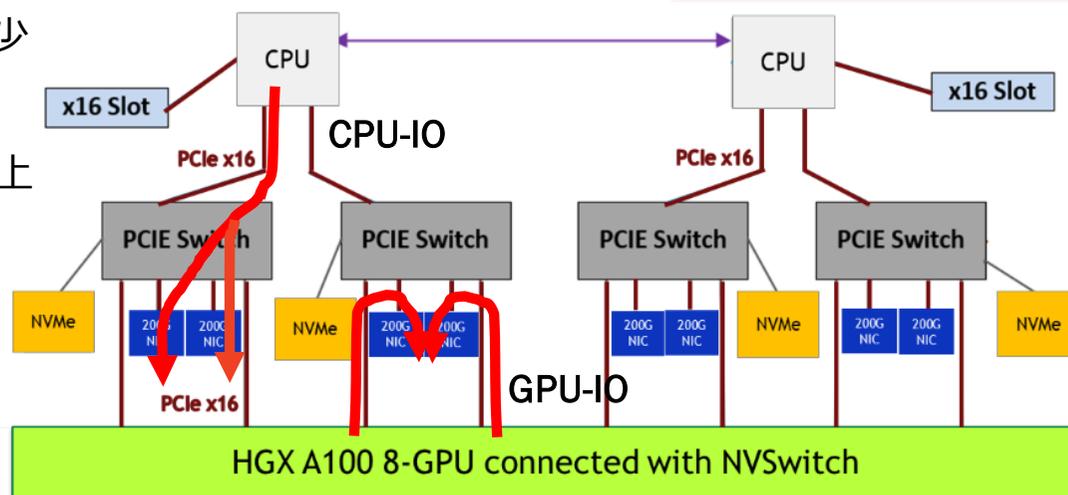
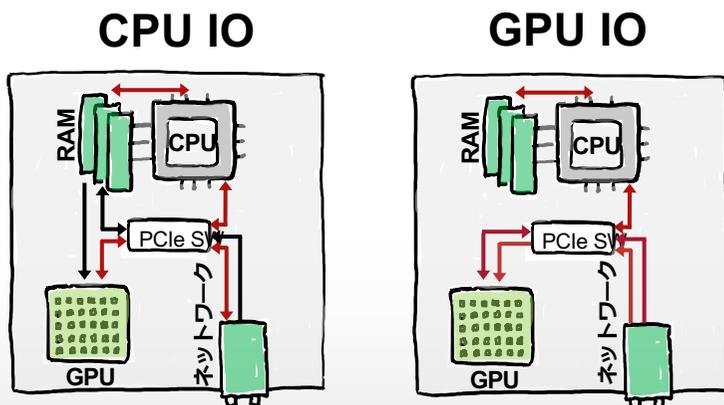




# GPU DIRECT STORAGE

GPUから直接ストレージアクセスを可能に

- 不必要なメモリコピーを排除し低レイテンシーを実現
- ホストCPUやメモリサブシステムの消費量が減少
- ストレージとGPU間のBWデータ転送が高速化
- AI、深層学習、HPCアプリケーションの性能向上





## GPU ダイレクトストレージ

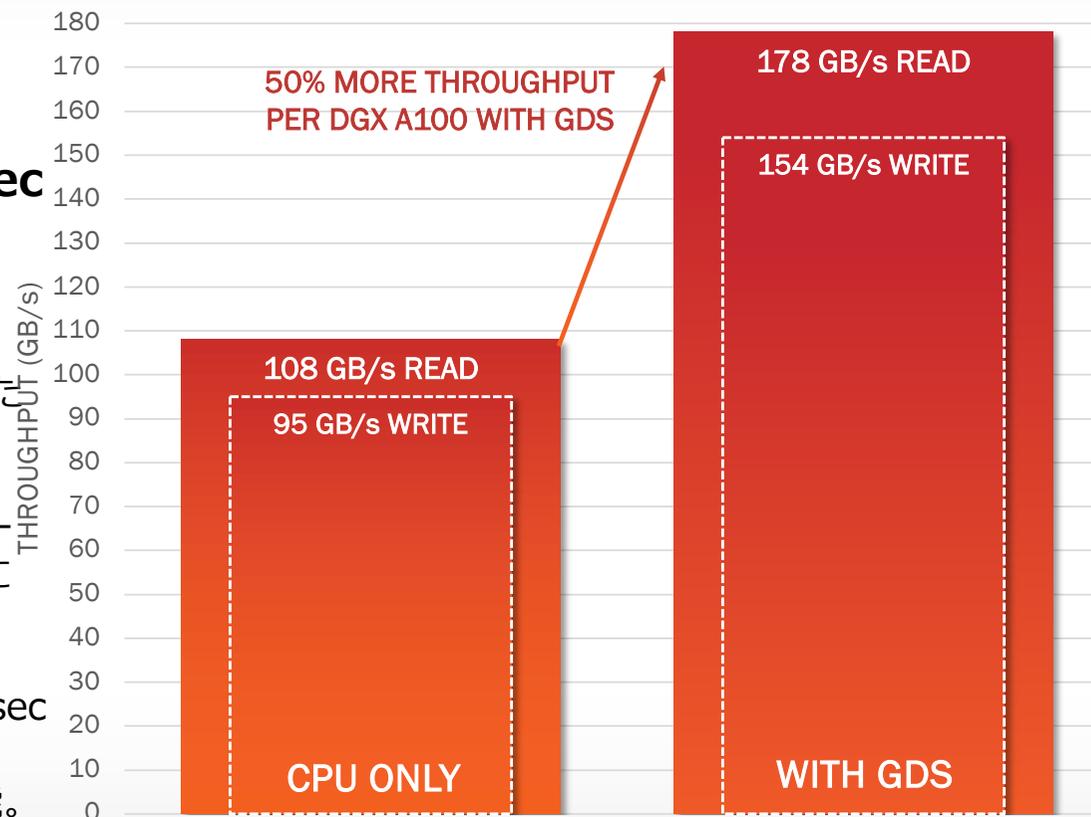
ハードウェアの制限をバイパスし178GB/sec  
性能を達成

DDNはNVIDIA GPUDirect Storageをシステムに統合した  
最初のストレージベンダー。

GPUメモリへの直接データ転送を可能とし、システムアー  
キテクチャにある制限をバイパス、レイテンシを最小限に  
抑えることによりGPUがもつ最大IO性能を達成。

GPU Direct Storageにおいて、DDN AI400Xは178GB/sec  
Readと154GB/sec Write性能を達成。  
同システムにおけるCPUだけのIOと比較して1.7倍の性能。

DDN AI400X PERFORMANCE SCALING  
WITH SINGLE DGX A100 SYSTEM AND GDS



four DDN AI400X, one DGX A100 client

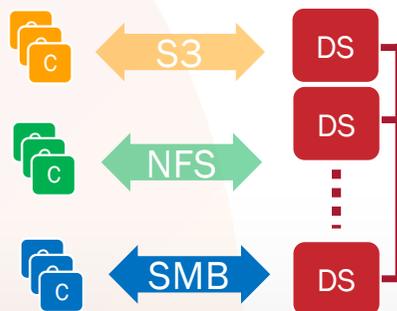


# Exascaler データサービスとデータ相互交換



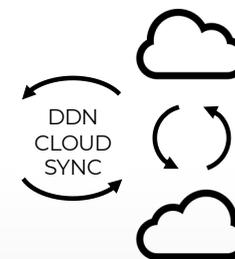
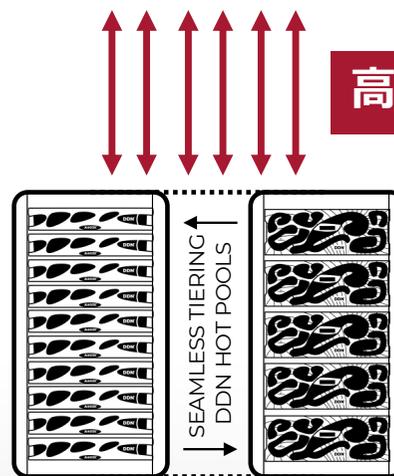
データサービス

様々なプロトコルで  
同じ名前空間を共有



高速データアクセス

高速並列アクセス  
高帯域、低レイテンシネットワーク  
Infiniband, Ethernet(RoCE)



クラウド連携

AWS, Azure, GogoleCloud  
上でのExaScaler  
オンプレミスなデータとの  
データコピー

ALL-NVME

NVME/HDD混在

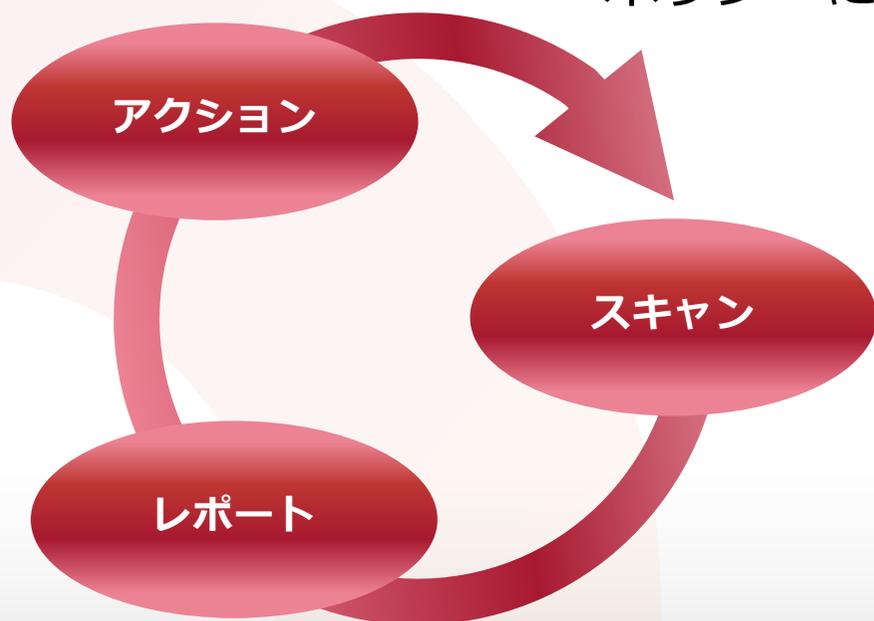
低レイテンシ高いIOPS性能

大容量、高いIO帯域性能

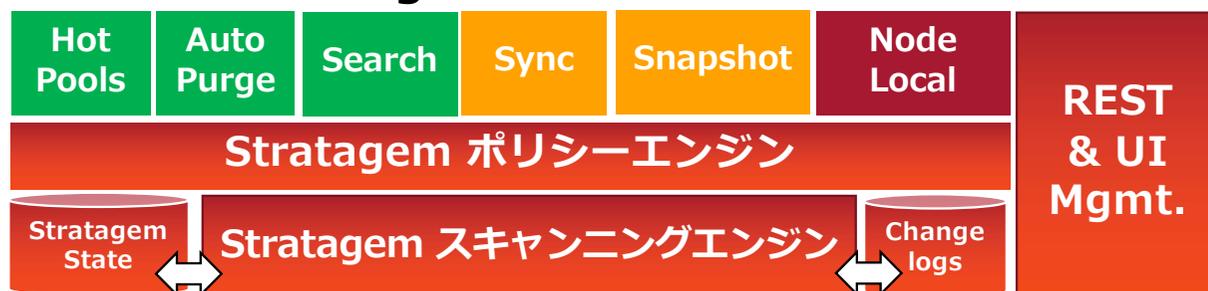


## Exascaler データ管理フレームワーク - Stratagem

ファイルシステムを高速にスキャンし、レポートし  
ポリシーに基づきアクションを実施



### Stratagem サービス



Job履歴

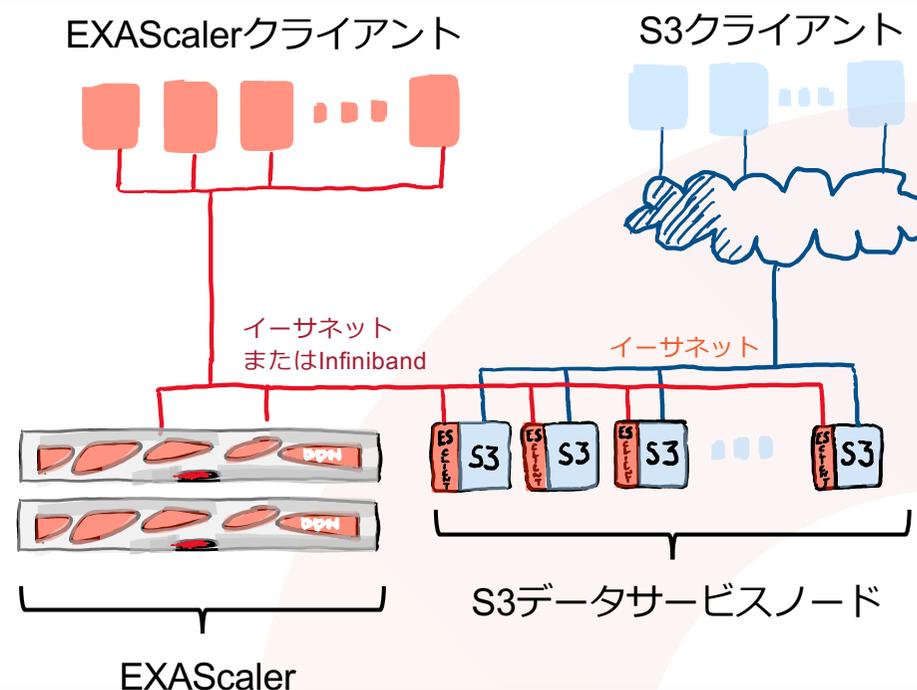
高速で並列化されたスキャンニング



# EXAScaler S3データサービス

## 統合されたスケールアウト型S3データサービス

- ▶ ExaScaler ストレージ上にてS3 APIを提供
- ▶ S3とPOSIXのデータアクセスの統合
- ▶ S3経由でデータの書き込み、S3もしくはExaScalerからデータの読み込み
  - 格納される:  
"/s3\_content/<バケット>/<ファイル名>..."
- ▶ ExaScalerのデータをS3にて読み出し
- ▶ S3経由にてデータをインジェストし、追加のデータコピーなしに深層学習の入力データとして直接利用可能



# EXAScaler マルチテナント型セキュリティ

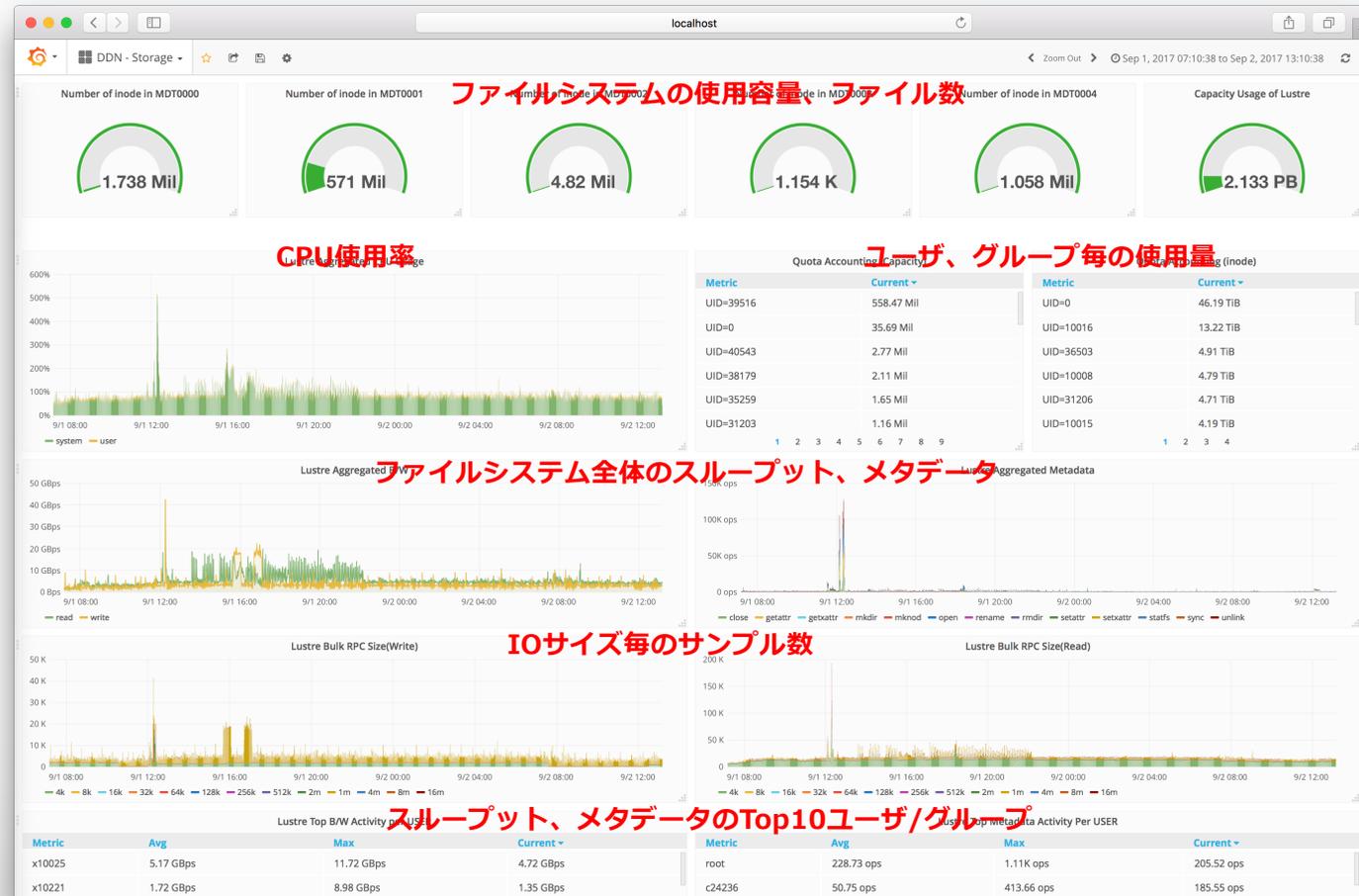
大規模な分散環境に対応する高度なセキュリティとデータの安全性



- ▶ 不正なクライアント、VM、およびコンテナはファイルシステムのroot権限を取得不可
- ▶ クライアント/コンテナのセキュアな共有キーまたはケルベロス認証
- ▶ マルチテナントで特定のデータセットのみをユーザーにエクスポート
- ▶ セキュアな監査ログ
- ▶ ドライブおよびNVMe\*における保存データ(Data at Rest)の暗号化

© 2020 DDN

# EXAScaler I/O モニタリング



© 2020 DDN

詳細なIO性能をリアルタイムにモニタリングし、ユーザ毎のIO統計情報の取得やアプリケーションのI/Oパターンまたは特徴などの解析に有用





ddn