



AIからアーカイブまでユニバーサルNAS ストレージのご紹介 VAST Data

2021年2月
ノックス株式会社
小幡 明広





全てのアプリケーション用のシングルデータプラットフォーム



Copyright © 2020 NOX Co., Ltd. All rights reserved.



新しいアーキテクチャの活用

新しいテクノロジーを組み合わせ、これまで考えられなかった方法でゼロから新しいタイプのストレージアーキテクチャを發明



NVME OVER FABRICS



QLC FLASH



ストレージクラスメモリ

オールフラッシュストレージ
わずか5Uにペタバイトデータを格納



物理容量
675TB~

購入単位
100TB単位

最小システムユニット
5U

VAST Dataとは

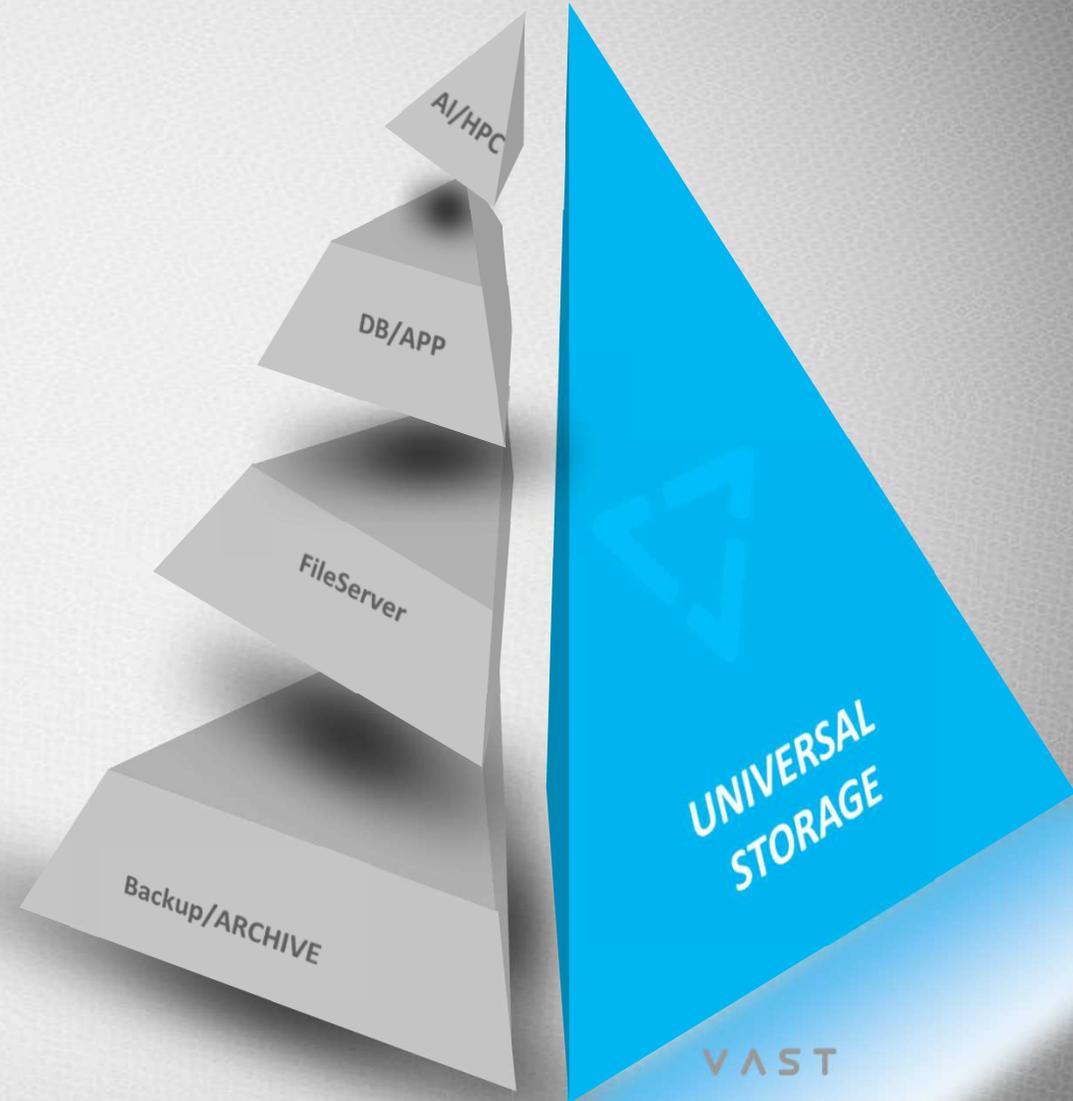
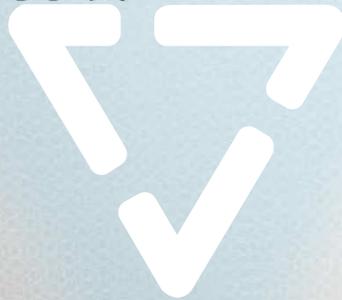
ミッション

ストレージ階層化の排除

NO MORE TIERS

何十年にも渡るストレージの複雑さとアプリケーションのボトルネックを終わらせること。

VASTは、一連のイノベーションを組み合わせ、すべてのデータとすべてのアプリケーションに対するフラッシュのイノベーションを起こします。



VAST Data 会社&チーム

2016年に設立。2019年2月製品出荷。本社はニューヨーク、開発拠点はイスラエル。



JEFF DENWORTH

VP, PRODUCTS AND
MARKETING, CO-FOUNDER



MIKE WING
PRESIDENT



SHACHAR FIENBLIT

VP, R&D AND CO-FOUNDER



RENEH HALLAK
CEO AND FOUNDER



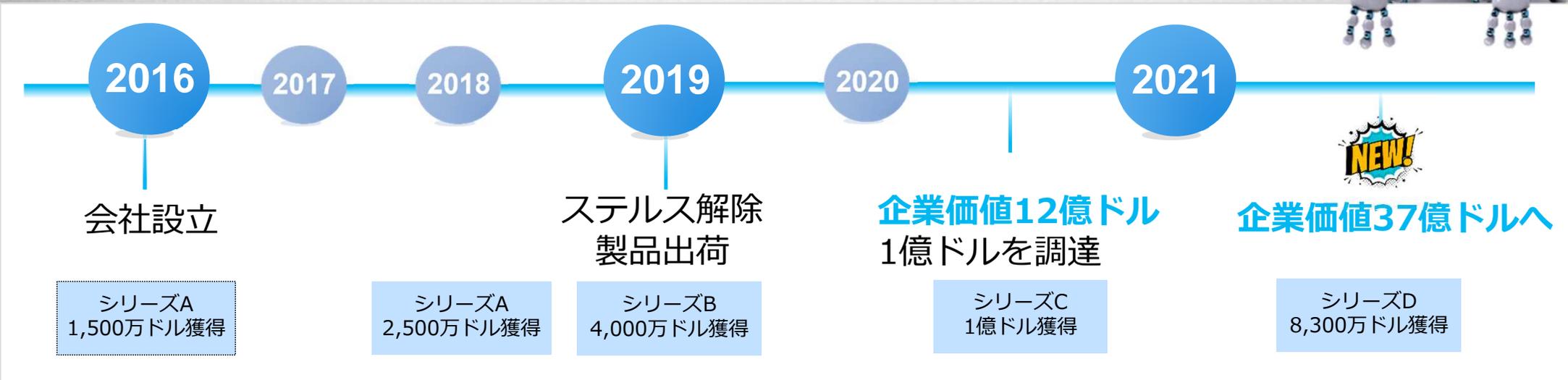
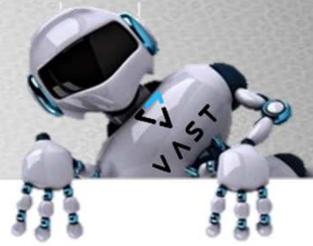
AVERY PHAM

VP OPERATIONS



タイムラインと企業価値

会社設立から5年、製品出荷からわずか2年で企業価値37億ドルの企業へ



Reuters - VAST valuation triples to \$3.7 bln after Tiger Global-led investmentSource
Captured: 04/05/2021 11:37:08

Bloomberg - Software Storage Firm Vast Data Valued at \$3.7 Billion in RoundSource
Captured: 04/05/2021 11:37:08

C-Tech - VAST Data raises \$83 million series D to triple valuation to \$3.7 billion in just one yearSource
Captured: 04/05/2021 11:37:08

アワード

TechTarget エンタープライズ分野で2年連続Gold Winnerを受賞するなど高い業界評価



WINNER

2020 Storage
Trailblazers Winner



GOLD

Best Enterprise Storage
Array – 2019 & 2020

GIGAOM

LEADER

2020 Radar For File
And Object Storage



WINNER

2020 Firestarter Award



LEADER

2020 Research Map
For File Storage



Flash Memory Summit

BEST IN SHOW

2020 Most Innovative
Flash Memory Technology

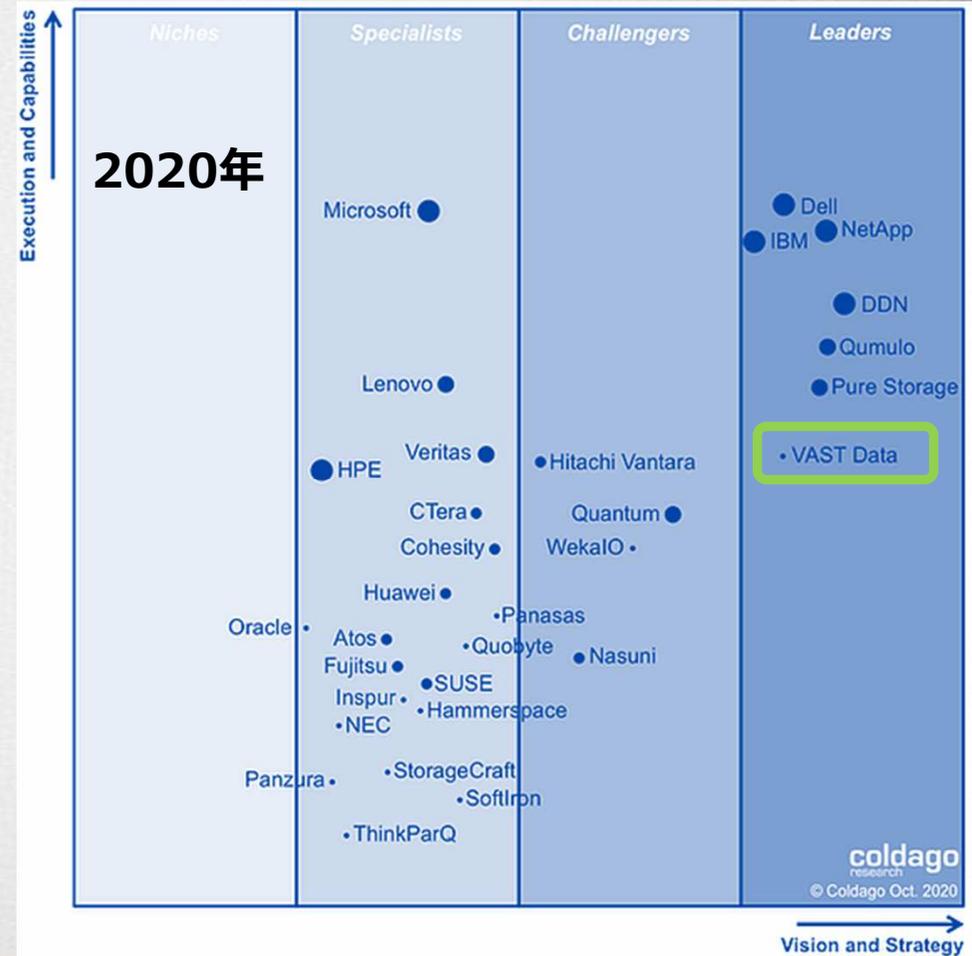
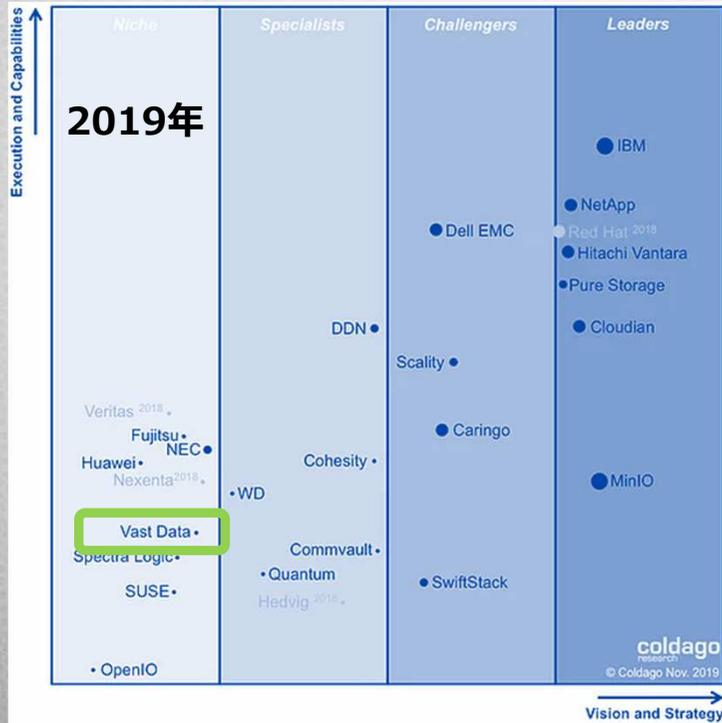
Gartner

COOL VENDOR

First 'Cool Vendor' In Data
Infrastructure Since 2018

2020年 リーダーポジションへ

2020年10月、HPC、エンタープライズ
ストレージベンダーマップ



導入先

業種・業界問わず幅広い分野で導入

PIXAR
ANIMATION STUDIOS

verizon[✓]

バイオテック企業



米国防衛大手企業



YAHOO!

Goldman
Sachs

メディカル企業





自動運転



地震探査会社



US政府系(軍関連)



ソフトウェア会社



学術大学



AIプロバイダ



PBテープリプレイス



最大手スポーツ放送局



SaaS会社



US政府系



最大手投資会社



アニメスタジオ



ジャズスポーツ
リーグ



ライフサイエンス



自動運転



地震探査会社



US政府系(軍関連)



ソフトウェア会社



学術大学



AIプロバイダ



PBテープリプレイス



最大手スポーツ放送局



SaaS会社



US政府系



最大手投資会社



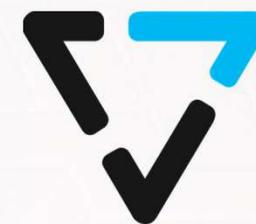
アニメスタジオ



ジャズスポーツ
リーグ



ライフサイエンス



カスタマー

US連邦政府機関への導入

国防総省、エネルギー省、国立衛生研究所、海洋大気庁、航空宇宙局、退役軍人省など、多くの連邦機関に導入され、**わずか2年で100PB以上の規模**でVASTが利用されています。



エネルギー省



航空宇宙局



海洋大気省



国防総省



退役軍人省



保険福祉省



インテリジェンス・コミュニティー

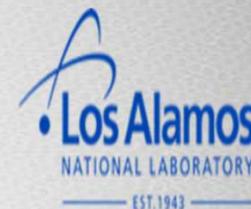
VAST DATA US GOVERNMENT AGENCY CUSTOMERS



OVER 100PB DEPLOYED
ONLY 2 YEARS

カスタマー

HPC分野



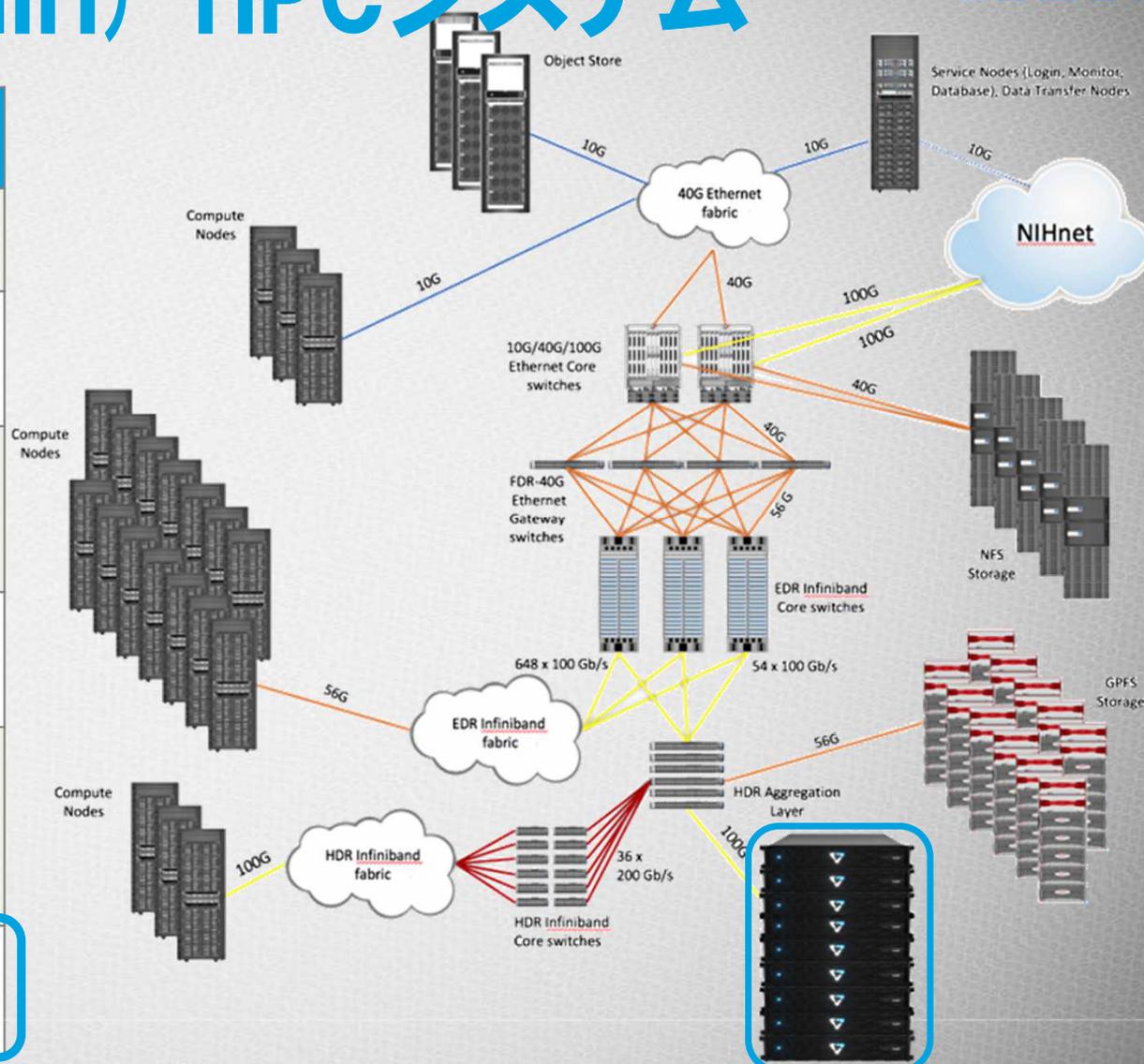
Success Story

- ・従来FileSystemと比べシンプルなNAS出来て良かった
- ・データ削減効果が75%データ削減された
- ・1日1枚だったユーザートラブル・チケットが不要になった
- ・HPCコードの実行が最大20倍高速化
- ・VASTの1クラスターで、home/scratch/trainingをサポート 22種類のクラスターに容易にマウント可能

国立衛生研究所 (NIH) HPCシステム



storage system	configuration	filesystem	network connectivity	usable storage (TB)
NetApp Cluster	10 x FAS8040 controllers SATA, SSD	NFS	16 x 10 Gb/s Ethernet	1100
DDN SFA12Ke	2 controllers 8 embedded file servers SATA SSD metadata	GPFS	16 x 56 Gb/s FDR Infiniband	1400
3X DDN SFA12KXe	2 controllers 8 embedded file servers NL-SAS SSD metadata	GPFS	16 x 56 Gb/s FDR Infiniband	6100
5X DDN SFA12KX-40	2 controllers 8 file servers NL-SAS SSD metadata	GPFS	16 x 56 Gb/s FDR Infiniband	21600
DDN SFA18K	2 controllers 8 file servers NVMe SSD NL-SAS SSD metadata	GPFS	8 x 56 Gb/s FDR Infiniband	4100
VAST	36 file servers NVMe Flash	NFS	72 x 100 Gb/s HDR-100 Infiniband	4900



製品コンセプト

コンセプト

シングルティアの全く新しいアーキテクチャ
独自のデータ保護、データ削減技術、Optane™Technologyの最大限の活用から、
安価で高速なQLCフラッシュを活用したプライマリ領域からアーカイブ領域まで活用可能

オールフラッシュ パフォーマンス

Optane™ Write性能
QLCフラッシュ Read性能

TIER-5 コスト効率

独自のデータ削減技術より
安価なフラッシュを活用した
HDDレベルのコスト効率

シンプル

シングルティアのNASアプリケーション
エクサバイトクラスまで容易にスケール可能
NFSv4・NFS・SMB・S3のサービス提供可能

全てのデータのためのOne Data Platform

あらゆるワークロードに対応

オールフラッシュNASストレージ

シングルティア&シングルネームスペース

ペタバイト&エクサバイトスケール

業種別ソリューション

-  ライフサイエンス
-  フィナンシャル
-  コンテンツ
-  HPC

セカンダリストレージ

-  rubrik
-  COMMVAULT
-  VEEAM

ビッグデータ&AI

-  splunk
-  APACHE Spark
-  TensorFlow

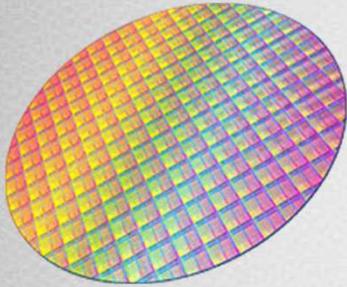
エンタープライズ

-  vmware
-  kubernetes
-  OPENSIFT



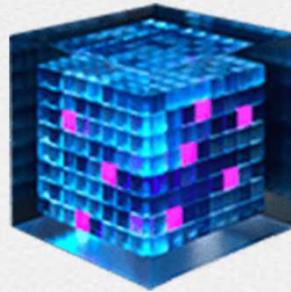
フラッシュ イノベーション

ストレージを根本的に再考



低コストQLCフラッシュ 85%節約

QLCフラッシュ専用開発された
File System
従来ストレージシステムで利用する
QLCフラッシュよりも、VASTクラ
スタのQLCフラッシュは20倍の耐久
性実現



次世代イレイジャーコーディング 最大66%節約

VASTローカルデコード可能なイ
レイジャーコーディングは、
6000万年の復元を提供。
必要なオーバヘッドはわずか9%。

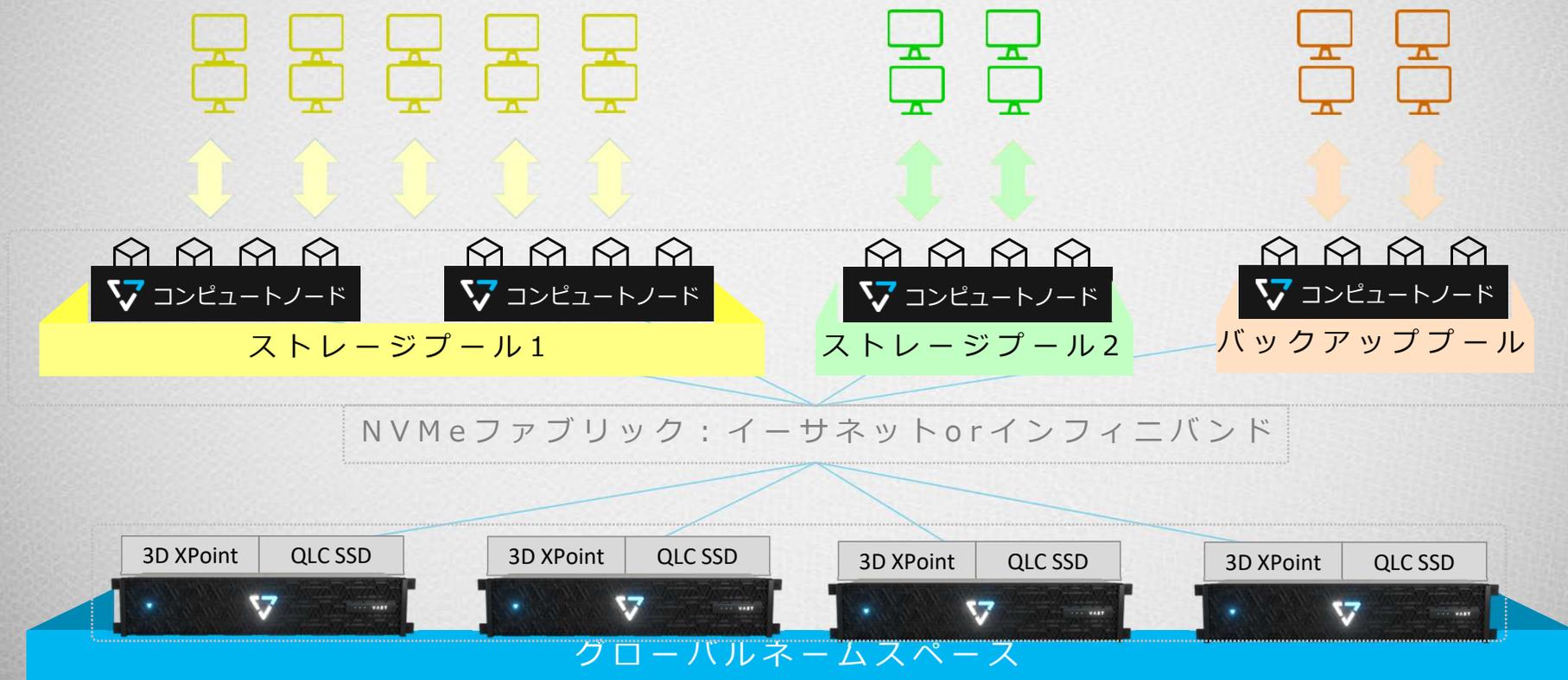


新世代データ削減 さらに50~90%節約

類似性ベースのデータ削減は、
グローバル重複排除と組み合わせ、
前例のないストレージ効率を実現。

DASE サーバプール

マルチプロトコルアクセス： NFS, NFS+RDMA、NFS+RDMA+GPUDirect™, SMB, S3, K8S CSI



新しいデータ削除効果

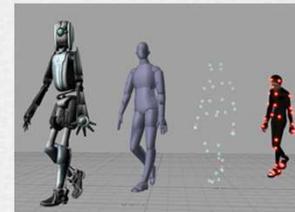
独自の圧縮アルゴリズムにより、従来の重複排除よりも**4,000~128,000倍の粒度**によるデータ削減を実現。非構造データ、バックアップデータ、圧縮済みデータにも有効。



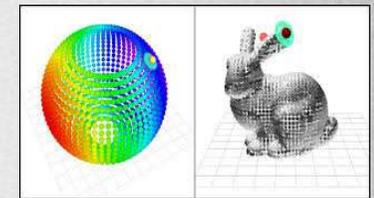
ログデータ
4:1



HPC
3:1



CGI/アニメ
2:1



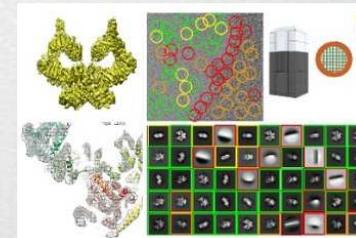
レンダリングデータ
~10:1



気象データ
2.5:1



AI/機械学習データ
3:1



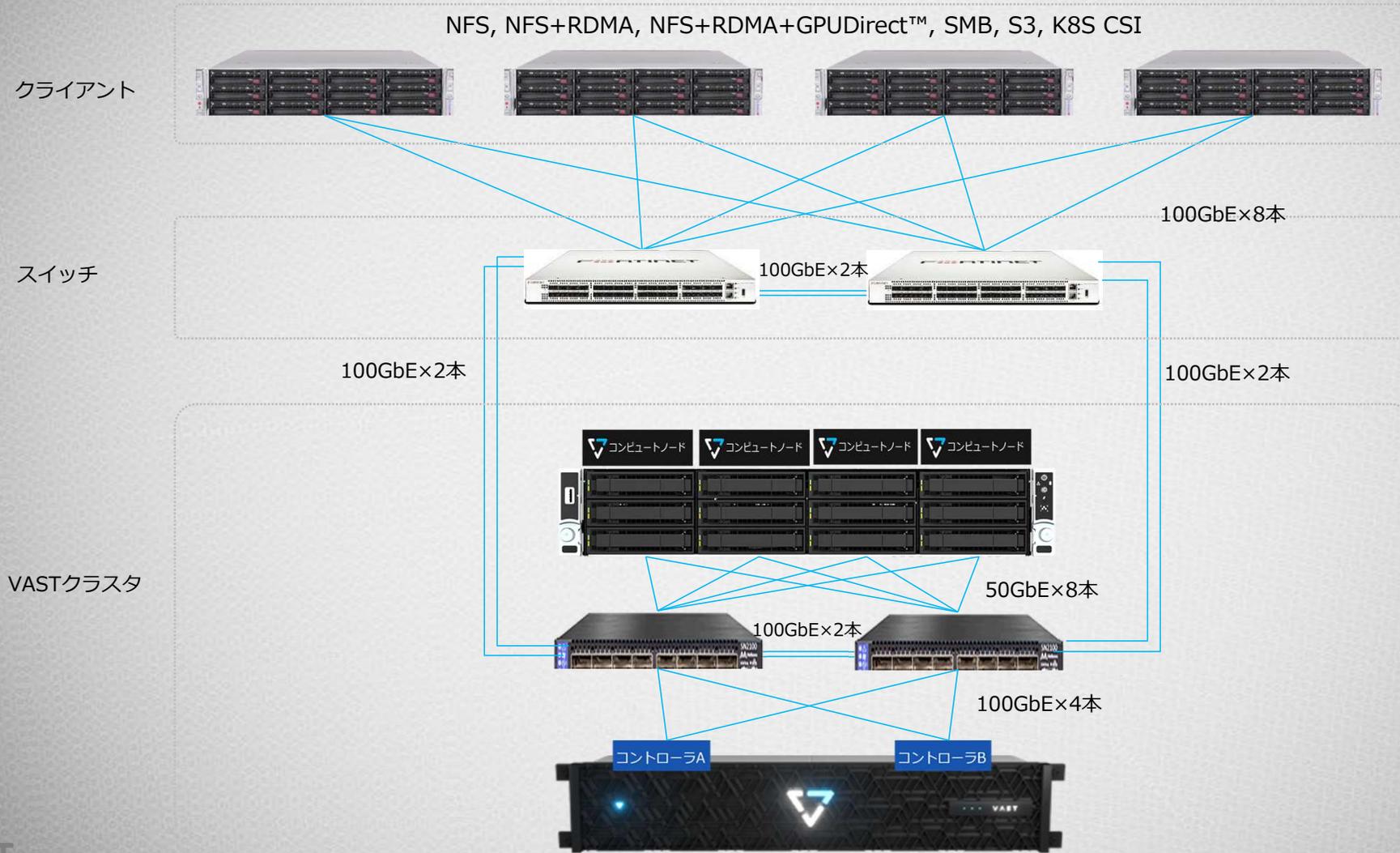
クライオデータ
4:1



ゲノムデータ
2:1

機器構成

構成イメージ



テクノロジー &機能

特許取得済みの独自テクノロジー

グローバルデータ圧縮やドライブ寿命を延ばす方法など各種特許を取得

- SYSTEM AND METHOD FOR GLOBAL DATA COMPRESSION (特許番号 : 20190379394)
- SYSTEM AND METHOD FOR USING FREE SPACE TO IMPROVE ERASURE CODE LOCALITY (特許番号 : 20200348855)
- TECHNIQUES FOR PROLONGING LIFESPAN OF STORAGE DRIVES (特許番号 : 20200174678)
- DISTRIBUTED SCALABLE STORAGE (特許番号 : 20190377490)
- STORAGE SYSTEM INDEXED USING PERSISTENT METADATA STRUCT
- METHOD AND SYSTEM FOR PROVIDING IMPROVED EFFICIENCY SNAF
- RESILIENCY SCHEMES FOR DISTRIBUTED STORAGE SYSTEMS (特許番
- DISTRIBUTED SCALABLE STORAGE (特許番号 : 20190377490)
- TRANSACTION MANAGER (特許番号 : 20200073964)



DASEクラスタ メリット

VASTのDASEクラスタアーキテクチャが提供するグローバルアルゴリズムは、以下のメリットを得られます。



キャッシュコピーレテンシーの課題無し

共有リソースに対するキャッシュの一貫性を保ちながら、スケーリングできます。



バッテリー不要

VASTクラスタは100%不揮発性。UPSやバッテリーの心配不要。



サーバ障害時のリビルド無し

ステートレスのため、コンピュートノード障害時のリビルドを排除します。



Dockerベースの自動スケーリング

コンテナ化されたアーキテクチャが、自律的で適応性のあるクラスタのスケーリングを行います。



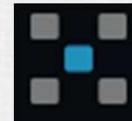
グローバルネームスペース

VASTクラスタは、柔軟に拡張できるメタデータとデータを1つのグローバルネームスペースを実現します。



画期的なフラッシュ効率

VASTの新しいグローバルフラッシュトランスレート(FTL)により、10年以上のQLCの耐久性を実現します。



優れたシステム使用率

VASTのデータ保護に対するグローバルレイジャコーディングは、わずか3%のオーバーヘッドの効率性を実現します。

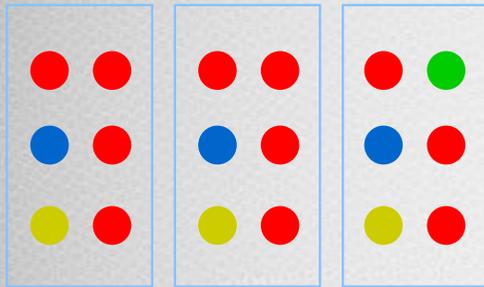


革新的なデータ削減

新しいグローバルデータ削減アルゴリズムは、すでに圧縮されたデータに対しても優れたデータ削減効率を提供します。

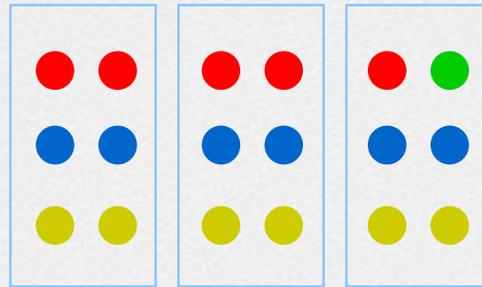
類似性データ削減テクノロジー

圧縮



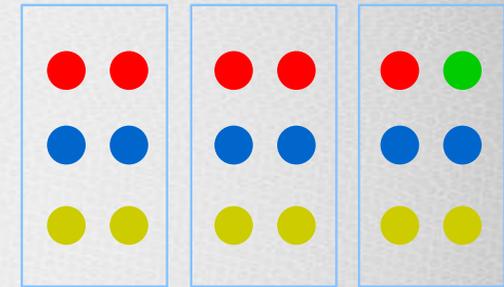
粒度=細かい、ローカル

重複排除



グローバル、粒度=粗い

VAST DATA 類似性



グローバル&細かい粒度

類似性効果

3:1 データ削減済み
バックアップデータ

3:1 圧縮済み
ログファイル

2:1 ライフサイエ
ンスデータ

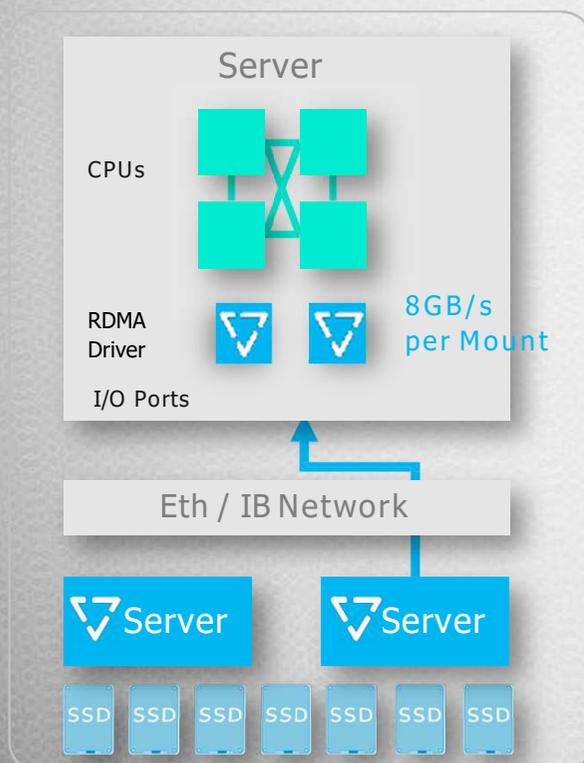
3:1 HPC
データ

3:1 アニメーション
データ

8:1 時系列データ

NFS over RDMA 提供

VASTはNFSプロトコルによるRDMAサーバ拡張機能を備えています。VASTが提供する「NFSoverRDMA」ドライバをLinuxサーバにインストールすることで、NASのシンプルさとNFSプロトコルの欠点となる性能を向上できます。

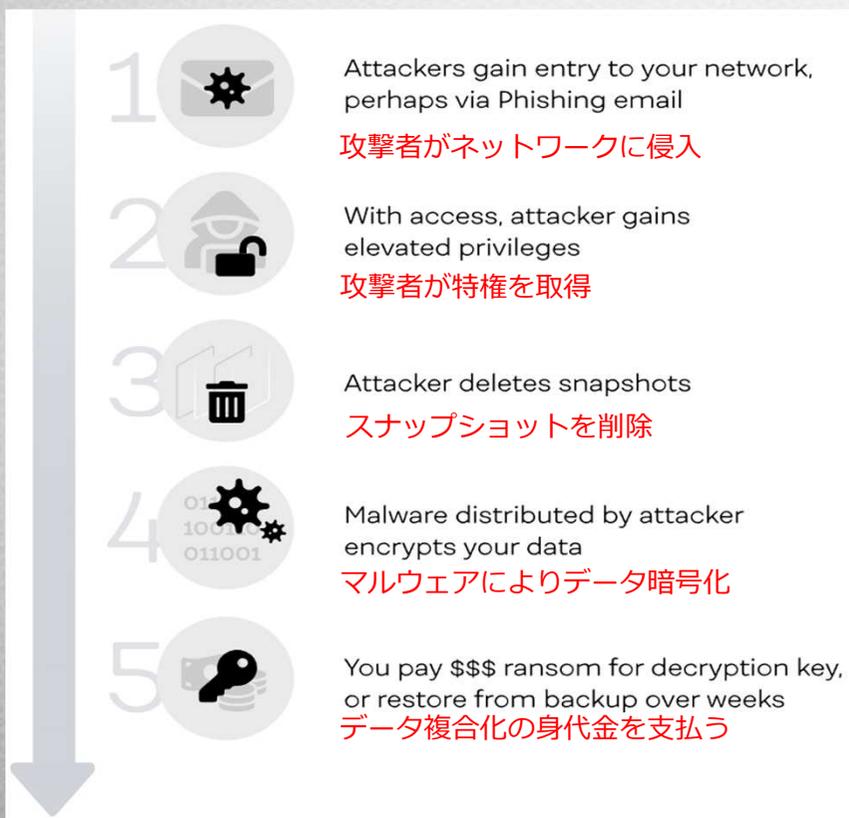


	オールフラッシュ NASアプライアンス	フラッシュベース HPC ファイルシステム	VAST
シングルマウント GB/Sec	2GB/s (TCP制限)	10GB/s+ (RDMA-enabled)	8.7GB/s+ RDMA
接続方式	TCP Ethernet	Infiniband, RoCE Ethernet	Infiniband, RoCE Ethernet
シンプルさ	シンプル	複雑なインテグレーション	シンプル
スケール	ペタバイトスケール	ペタバイト～エクサバイト	エクサバイトスケール

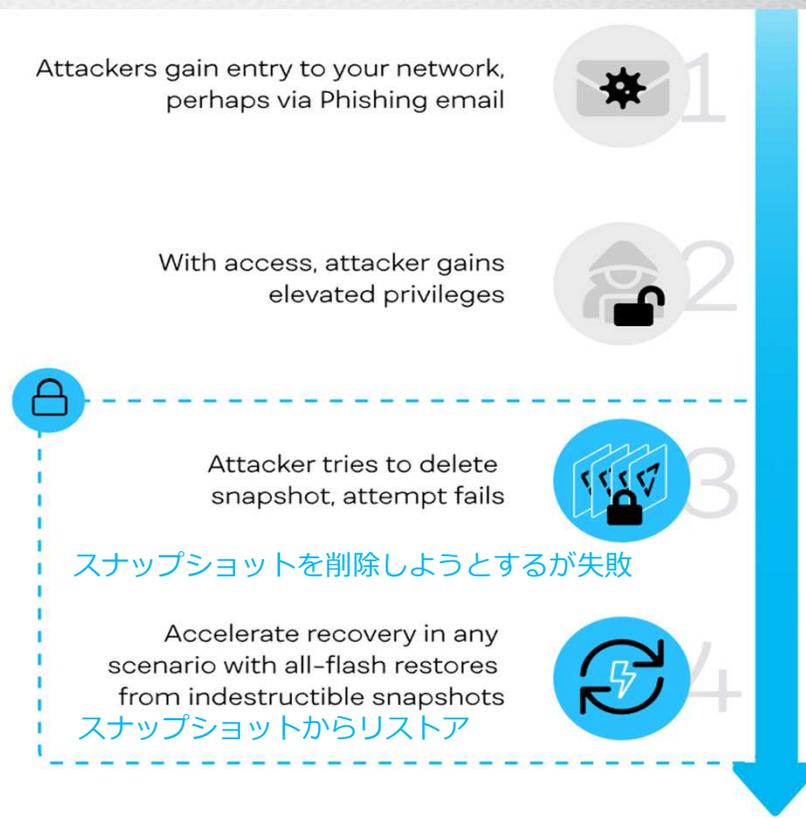
ランサムウェア対策

インディストラクティブル(破壊できない)・スナップショット機能により、個人が有効期限が切れる前にスナップショットが削除できない仕組みによる保護レイヤーを提供。

Indestructibleスナップショットがない場合



Indestructibleスナップショットがある場合



バックアップ製品連携

バックアップデータの格納先として利用することで、従来よりも高いデータ削減率を実現できます。



	Rubrik Compression	+	Additional VAST Data Reduction	=	Total Combined Data Reduction Ratio
VMware バックアップ	4:1		1.3:1		5.2:1
非構造化データ	3.8:1		1.48:1		5.6:1
VMware SQL サーバ	5:1		1.68:1		8:1



COMMVAULT

	Commvault Compression: GZIP	Commvault Compression: LZO	Commvault Compression: Disabled
VMware バックアップ	3.4:1	6:1	8:1
非構造化データ	3.5:1	4.5:1	4.6:1
データベースサーバ	3.4:1	5.3:1	9:1

VAST ソフトウェア・スタック

クラスタ マネージメント

GUI
CLI
REST (SWAGGER)

コールホーム VASTリモート監視、サポート

VASTマルチプロトコルアクセス NFS v3、NFSv3 over RDMA、SMB2、S3

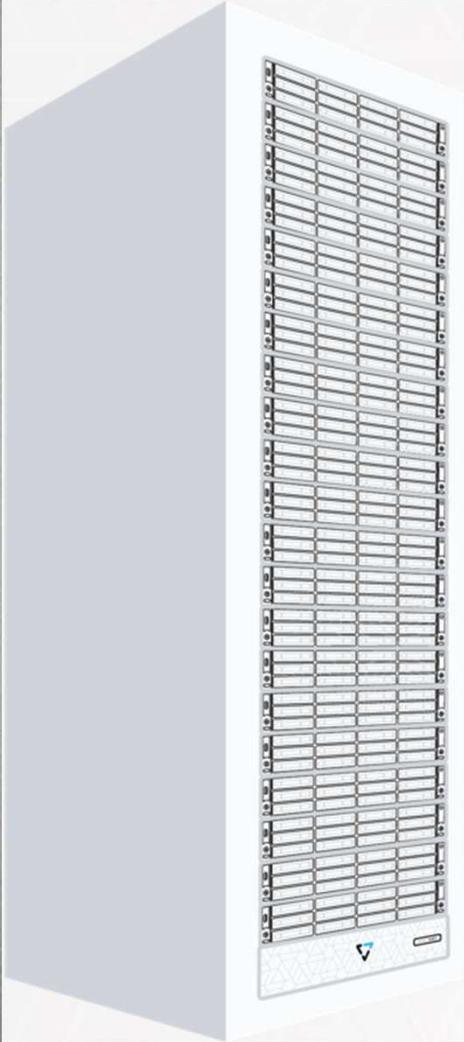
VAST データサービス

- エクサバイト スケールネームスペース
- 分散ファイルロックング
- LDAP、AD、
- グローバルデータ削減
- POXIX ACLs
- スナップショット
- クォータ機能
- S3 出力
- 暗号化
(256-bit AES-XTS)

VAST インフラサービス

- トランザクションストレージシステム
- 分散永続メタデータ
- グローバル QLCフラッシュトランスレーション
- リビルドエラー修正
- 活性交換アーキテクチャ

VAST クラウドバックアップ



クラウド



オンプレ



製品インテグレーション

アイデンティティ プロバイダー



バックアップアプリケーション



アラート、モニタリング&通知



ビッグデータ&ディープラーニング



自動化&仮想化



GPUダイレクトストレージ

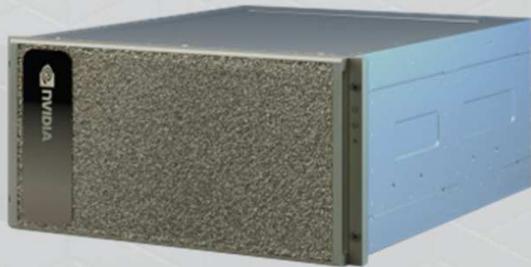
LIGHTSPEED

AI用ユニバーサルストレージ

SIMPLICITY, PERFORMANCE, AND SCALE TO POWER

これからの10年のためのストレージプラットフォーム

VAST + NVIDIA GPUダイレクトストレージ



DGX-A100



DGX-2

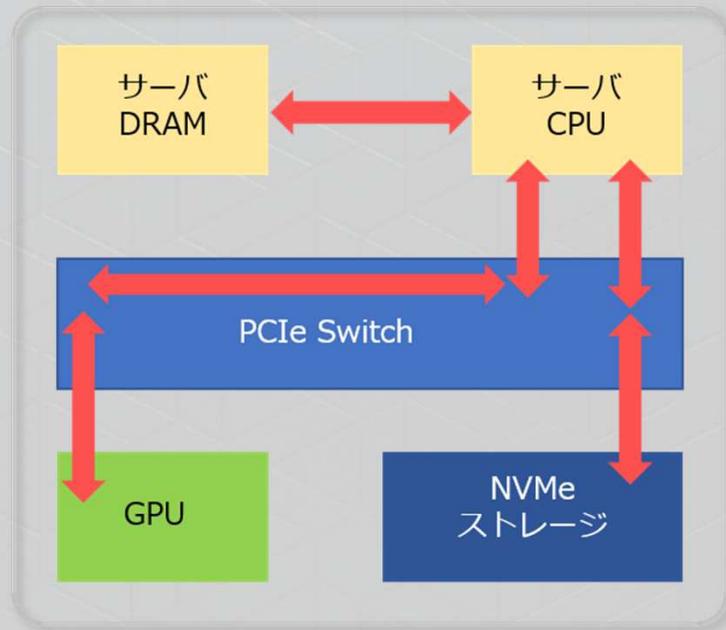


GPU Direct Storageとは？

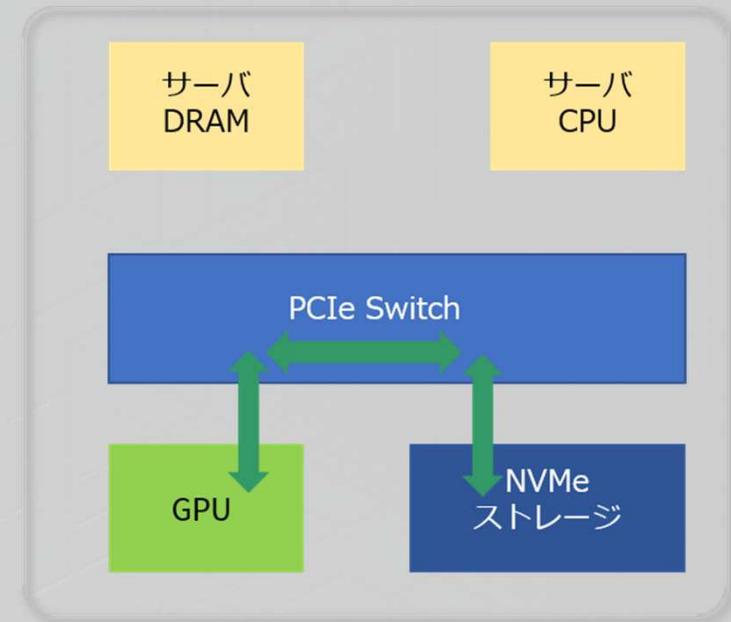
GPU Direct Storage (GDS) ソフトウェアにより、GPUメモリとNVMeストレージがGPUと直接通信できるようにします。

NVMe-over-Fabricsからアクセスもできるため、ホストサーバーのCPUとDRAMは不要になり、ストレージとGPU間のIOパスはより高速になります。

従来のデータフロー



GPUダイレクトストレージ



マスター タイトルの書式設定

VASTは、GDS認定センター3社のうちの1社になります。



ACCELERATED COMPUTING

MOFED and Filesystem Requirements

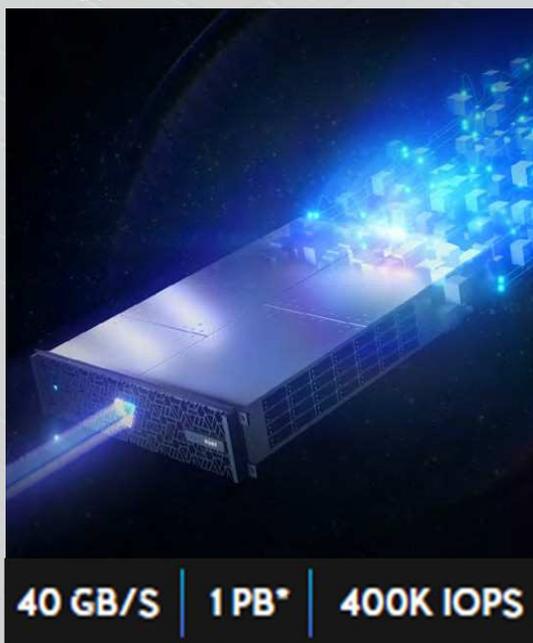
Here are the requirements:

- Ubuntu 18.04 and 20.04
- MOFED 5.1-0.6.6.0 and later, which supports NVMe NVMeoF, NFSoRDMA (VAST) on Linux kernel 4.15.x and 5.4.X
- The following distributed filesystems:
 - WekaFS 3.8.0
 - DDN Exascaler 5.2
 - VAST 3.4

GPUクラスタ アーキテクチャ

LIGHTSPEED

最大16GPUクライアント



LightSpeed エンクロージャー x1
VASTコンピュートノード x1

PENTAGON

最大80GPUクライアント



LightSpeed エンクロージャー x5
VASTコンピュートノード x5

DECAGON

最大160GPUクライアント



LightSpeed エンクロージャー x10
VASTコンピュートノード x10

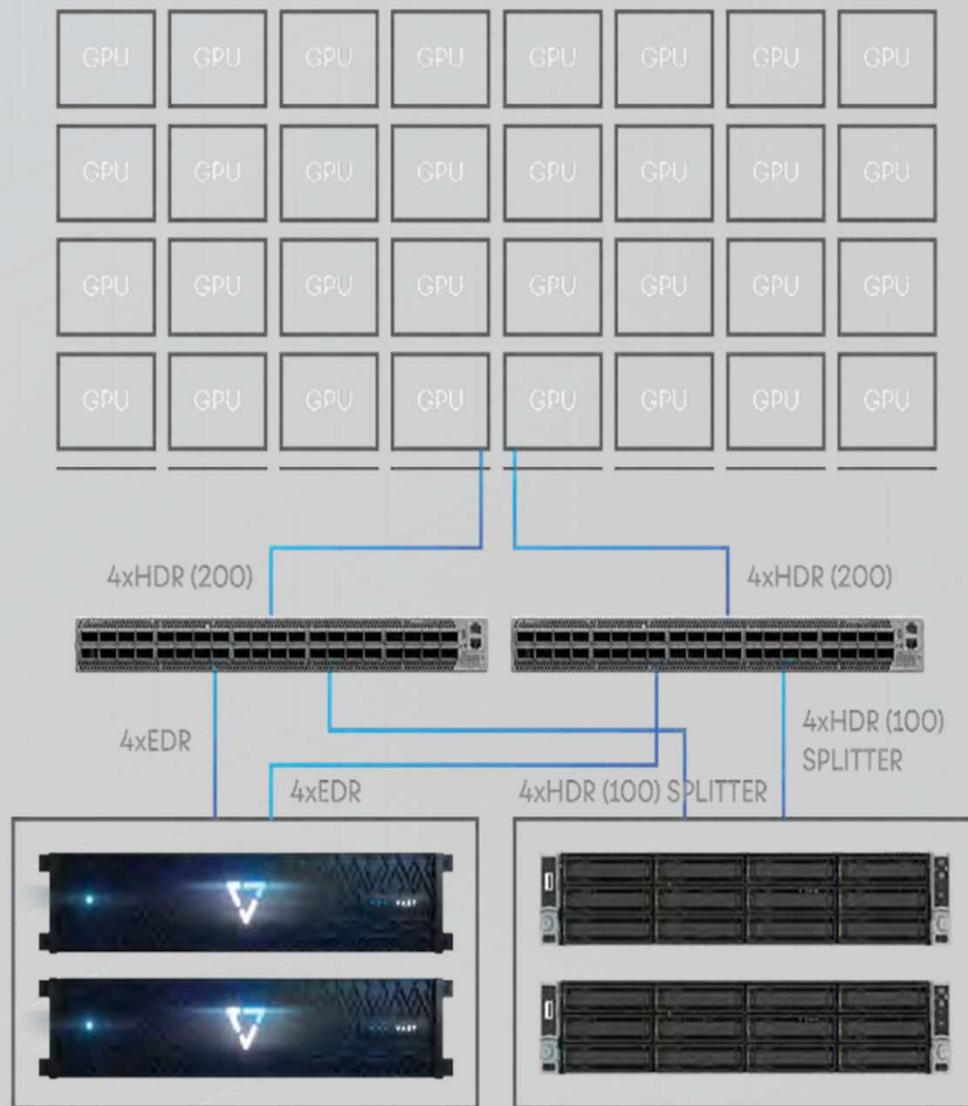
LIGHTSPEED

構成:

- 32 GPU クライアント
- VAST 2x2 (2 LightSpeed エンクロージャー)
 - 80GB/s Read
 - 10GB/s Write
- 26RU (~31RU with gaps)

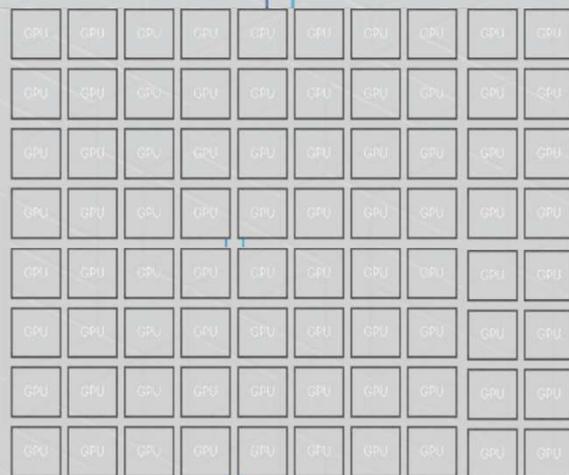
ネットワーク:

- GPUs: 8xHDR (200)
- ISL: 4xHDR (200)
- エンクロージャー: 8xEDR
- コンピュートノード: 8xHDR (100) splitters

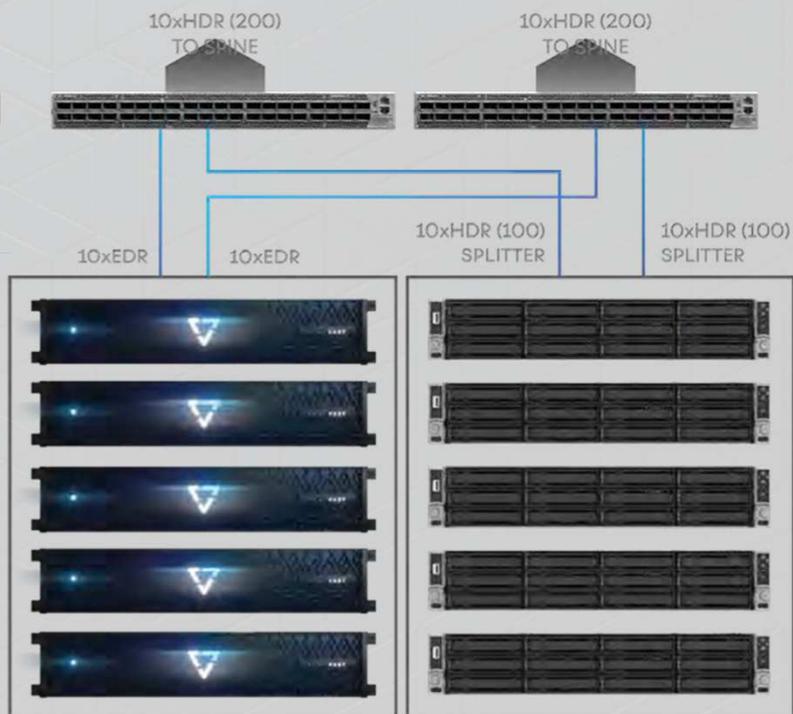


PENTAGON

COMPUTE RACKS



STORAGE RACKS



200 GB/S | **5 PB*** | **2M IOPS**

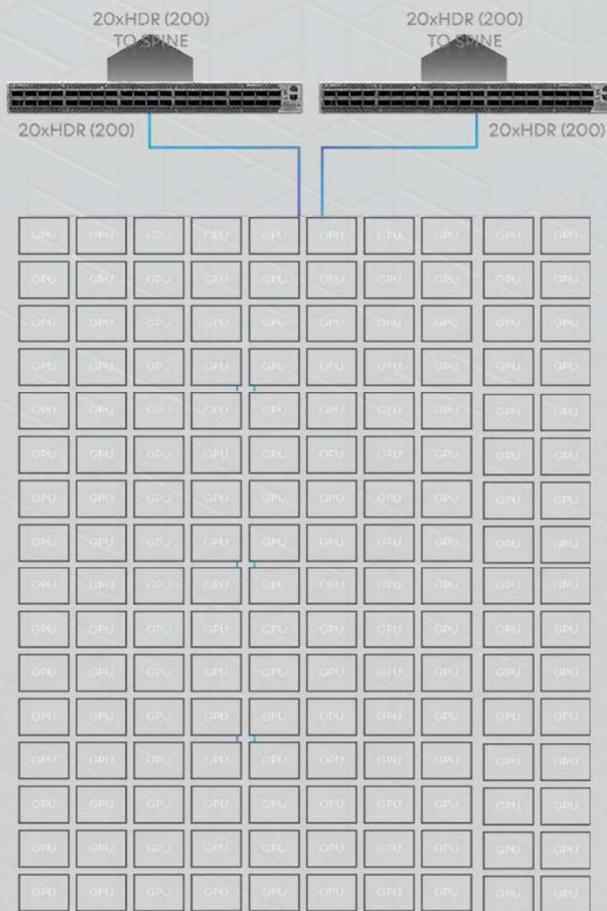
構成:

- 80 GPU クライアント
- VAST 5x5 (5 LightSpeed Enclosures)
 - 200GB/s Read
 - 25GB/s Write
- 60RU (~31RU with gaps)

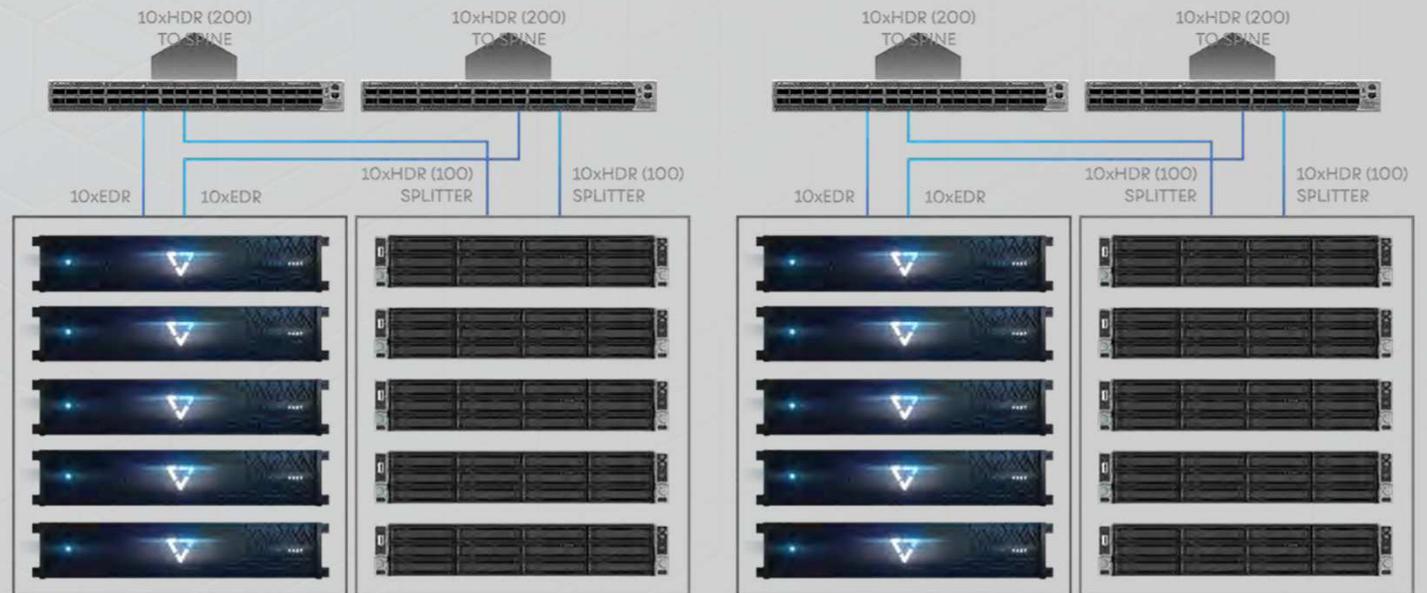
ネットワーク:

- GPUs: 20xHDR (200)
- UPLINKS: 40xHDR (200)
- エンクロージャー: 20xEDR
- コンピュートノード: 20xHDR (100) splitters

COMPUTE RACKS



STORAGE RACKS



構成:

- 160 GPU Clients
- VAST 10x10 (10 LightSpeed Enclosures)
 - 400GB/s Read
 - 50GB/s Write

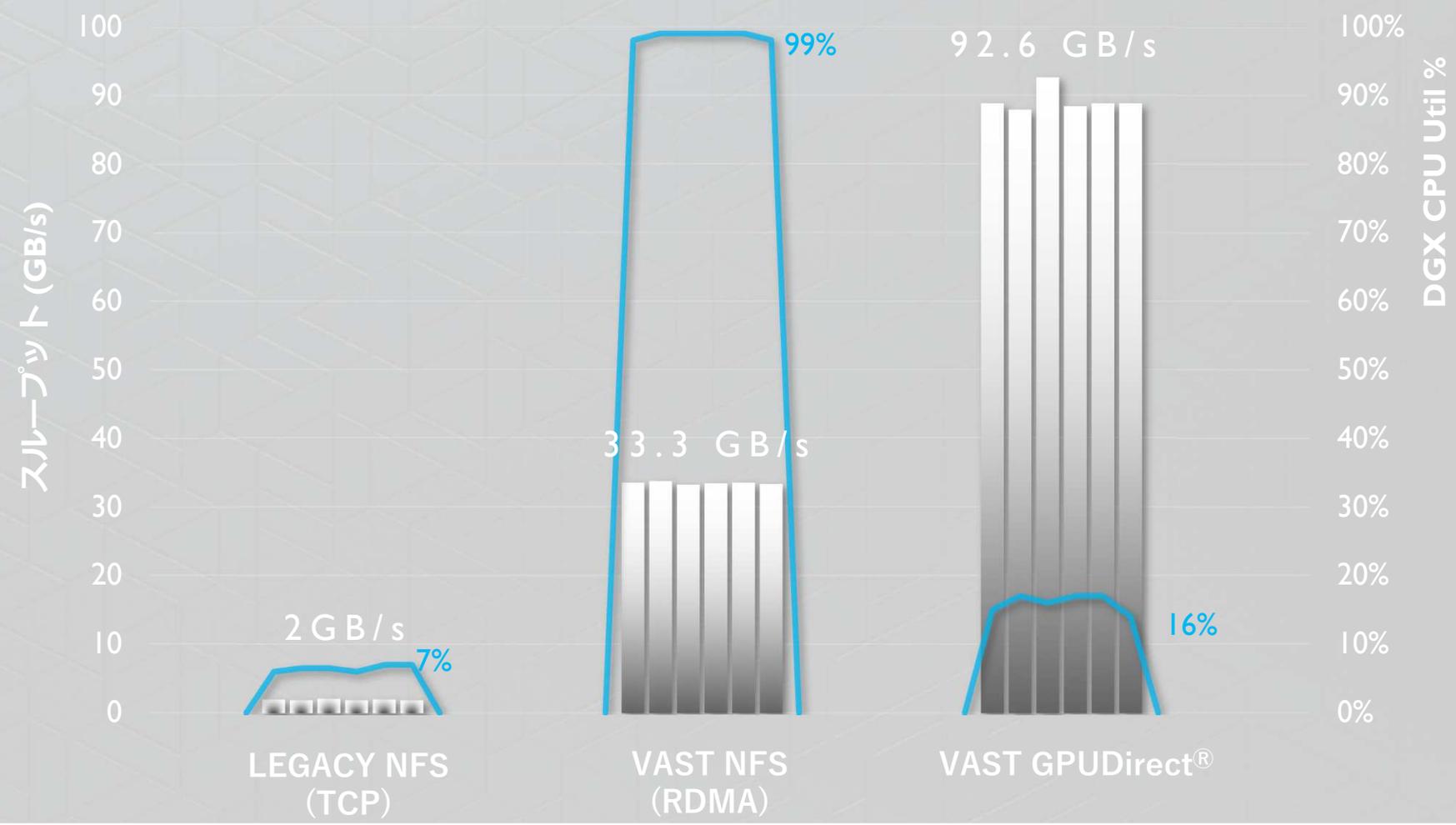
ネットワーク:

- GPUs: 40xHDR (200)
- UPLINKS: 60xHDR (200)
- DBOX: 40xEDR
- CBOX: 40xHDR (100) splitters

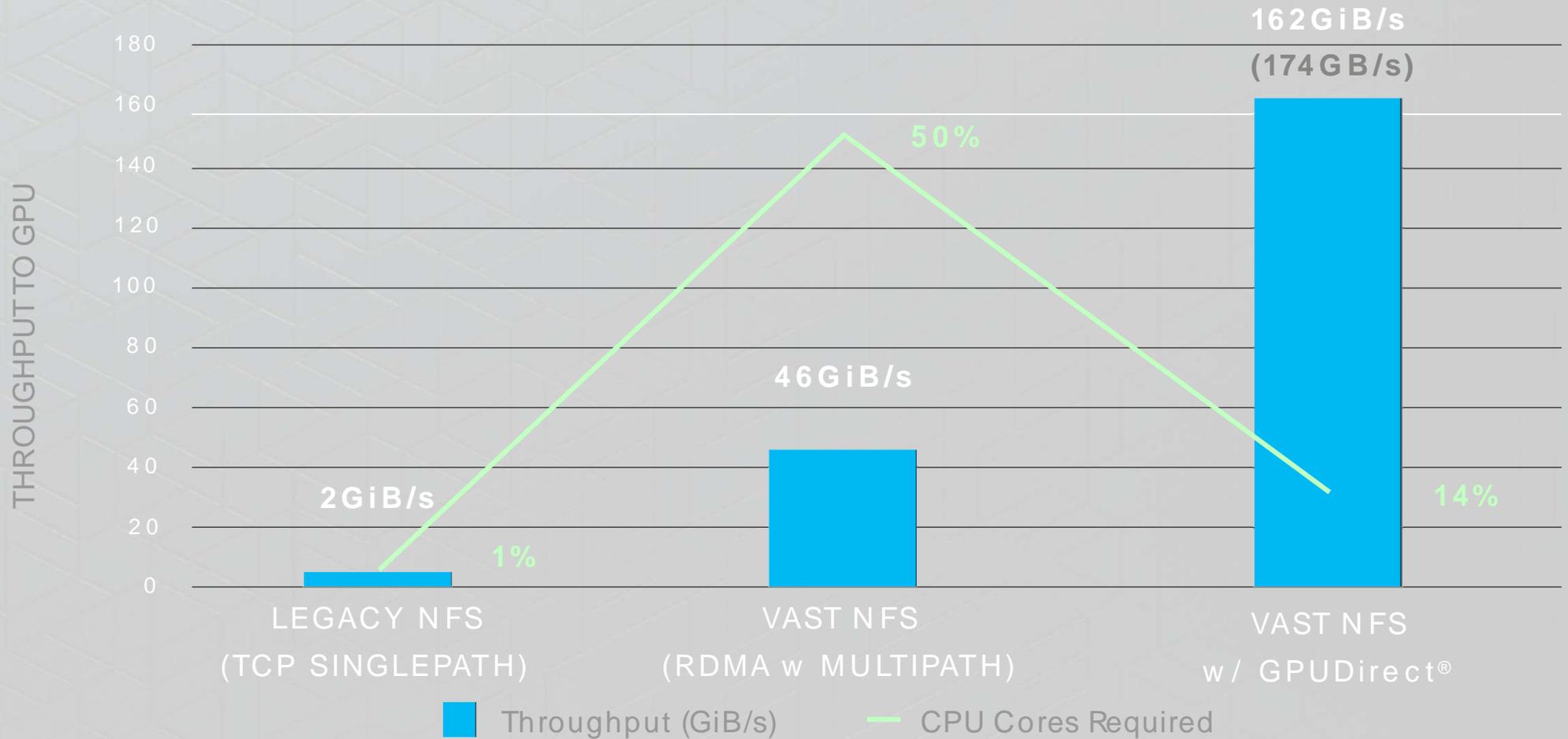
NVIDIA DGX-2 ベンチマーク

シングルマウントポイント

Throughput (GB/s) DGX-2 CPU Util %



DGX-A100 ベンチマーク



TESTED WITH GDSIO • 1 x DGX-A100 • 5 x VAST CBOX SERVER CHASSIS & 5 x VAST LIGHTSPEED ENCLOSURES • 4MB I/O SIZE • 4GB FILE SIZE • 96 THREADS x 8GPUs

提供モデルと諸元

販売モデル (Gemini)

ハードウェア



- VAST クアッドサーバエンクロージャ
- NVMeファブリック イーサネットスイッチ
- VAST NVMeエンクロージャ

買い取り

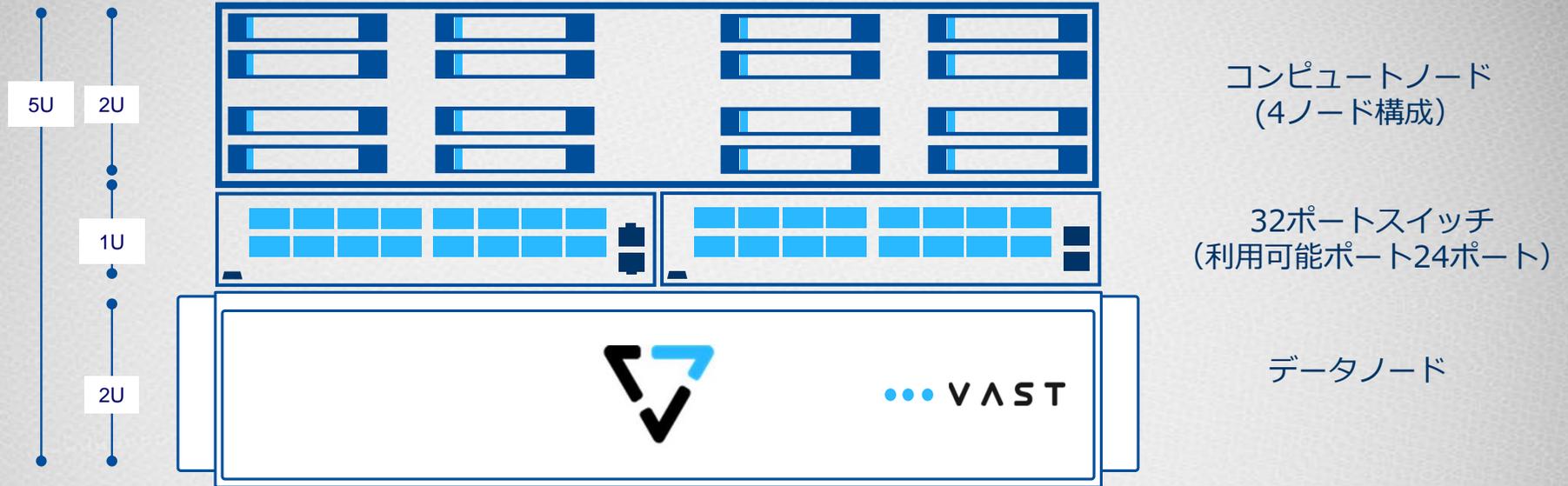
ソフトウェア



- 100TB単位
- 1~10年契約

サブスクリプション

VAST クラスタ構成



ストレージ容量 (物理容量)	
NVMe Flash容量(TB)	675.8
3DXpoint(TB)	18.0
ストレージ容量 (実効容量)	
NVMe Flash容量(TB)	552

ネットワークポート数	
データノード	4
コンピュータノード	8
スイッチポート(利用可能)	24

合計平均電力(kW)		2.392
データノード(kW)	1.4	
コンピュータノード(kW)	0.8	
スイッチ(kW)	0.2	

ピークストレージ性能(データノード)	
シーケンシャル Write	3.6GB/Sec
ランダム Write	5.1GB/Sec
ランダムRead	22GB/Sec
ランダム Write IOPS	102K
ランダム Read IOPS	182K

電源ケーブル数(C14/C15)		8
データノード	4	
コンピュータノード	2	
スイッチ	2	

ラックスペース	
ラックDepth(mm)	1193.8
ラックユニット数(最小)	5U

性能と制限事項

性能	4×コンピュータノード	8×コンピュータノード
シーケンシャル Write	5.1GB/Sec	5.3GB/Sec
ランダム Write	5.3GB/Sec	5.4GB/Sec
ランダムRead	19.2GB/Sec	22GB/Sec
ランダム Write IOPS	102,000	127,000
ランダム Read IOPS	182,000	322,000

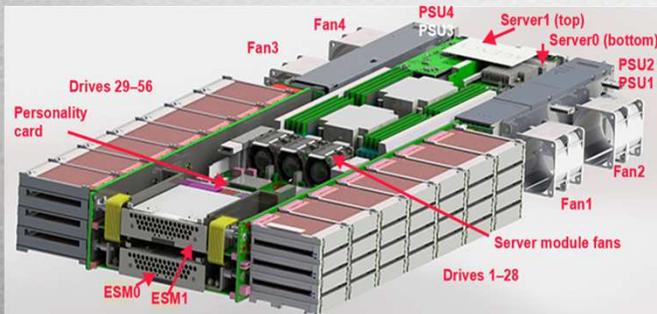
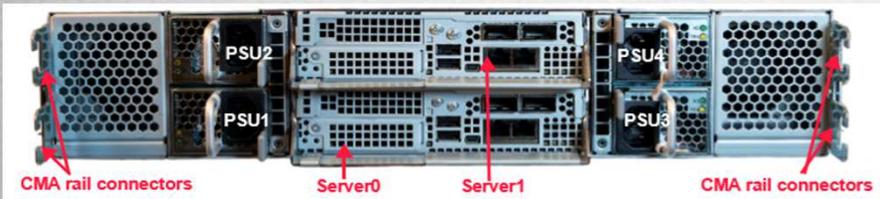
制限事項 (最大)	4×コンピュータノード
ファイル/フォルダ数(データノードあたり)	15億
ファイル/サブディレクトリ数(フォルダあたり)	100万
アクティブコネクション数(コンピュータノードあたり)	1,500
エクスポート数	200
ファイル/ディレクトリ文字長	255
ファイルサイズ	16TB
NFSファイルロック数(ファイルあたり)	100



エンクロージャー 諸元



VAST DF-5615

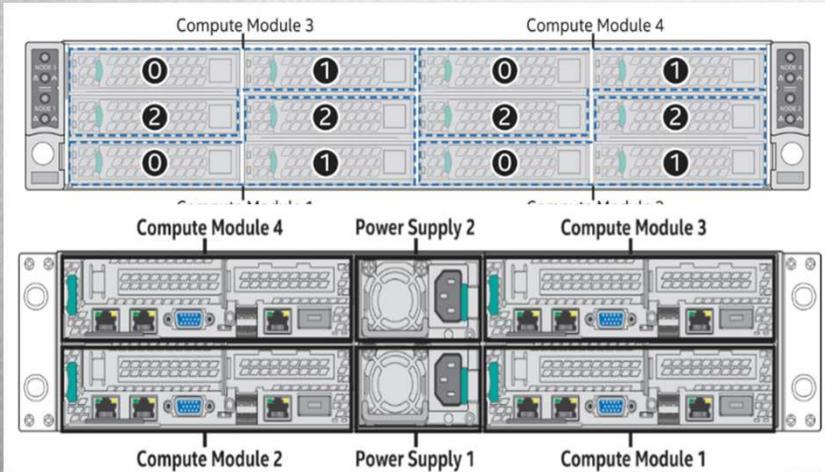


VAST DF-5615	
I/Oモジュール	Active/Active モジュール
ネットワーク	4×100GbEまたは4×100GbE InfiniBand
NVMe フラッシュ	38×15.36TB QLC
3Dxpoint	18×960GB
寸法	2 U、 H:81.28mm W:447.04 D:949.96
重量	38.55kg
電源	1500W×4
消費電力	1200W(平均) 、 1450W(最大)
最大拡張数	1,000(エンクロージャー)
冗長性	コントローラ 電源 NIC QLC SSD(イレージャーコーディング) 3DXPoint(ミラーリング) ※各部位活性交換可能

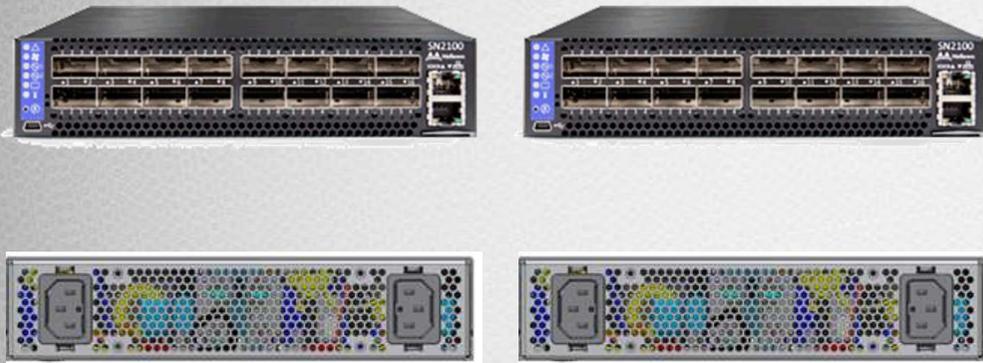
コンピュータノード 諸元



コンピュータノード	
サーバ	4 × ステートレスサーバ(4ノード)
ネットワーク	8 × 50GbE または 4 × 100GbE InfiniBand 4 × 1GbE (オプション)
CPUコア	80 × 2.4GHz (ノードあたり20コア)
メモリ	32 × 32GB (ノードあたり256GB)
寸法	2U、H:86.87mm W:437.90 D:733.04
重量	35.38kg
電源	1600W × 2
消費電力	750W (平均)、900W (最大)
最大拡張数	10,000 (サーバ)
冗長性	ノード 電源 DISK(SSD) NIC ※各部位活性交換可能



NVMeスイッチ 諸元



MSN-2100

Mellanox スイッチ1U 16ポート

NVMeスイッチ	
型番	Mellanox MSN-2100
構成、ポート数	16ポートスイッチ×2台
ユニットサイズ	1U ※2台のスイッチを1Uに集約
速度	10/25/40/50/56/100GBR/ポート
寸法	1U、H:43.8mm W:200 D:508
重量	4.54kg
消費電力	94.3W(平均) 、 248.6W(最大)
冗長性	スイッチ 電源 ※各部位活性交換可能

ご清聴ありがとうございました

ノックス株式会社
営業本部
小幡 明広

obata@nox.co.jp
03-5731-5551