



DDN Lustre Seminar

大阪大学様 OCTOPUS用ストレージシステム

DataDirect Networks Japan, Inc.

2018/07/04 橋爪信明

Agenda

- DDN会社/製品紹介
- Lustreについて
- Lustreの特徴
- Lustreのアーキテクチャ
- Lustreの各コンポーネント
- OCTOPUSストレージ概要
- OCTOPUSファイルシステムの利用
- Q&A

Agenda

- DDN会社/製品紹介
- Lustreについて
- Lustreの特徴
- Lustreのアーキテクチャ
- Lustreの各コンポーネント
- OCTOPUSストレージ概要
- OCTOPUSファイルシステムの利用
- Q&A

DDN | DataDirect Networks 会社概要

会社名	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DataDirect Networks社100%子会社)	設立	日本法人 US本社	2008年1月 1998年
所在 本社製造拠点	日本法人 東京/大阪/名古屋 US カリフォルニア・チャッツワース	従業員	650名以上 / 20カ国 29名 / 日本法人	
事業内容	ストレージプラットフォーム・ファイルシステムソリューションの開発・販売・サポート 高性能ストレージシステムの販売・設計・構築・保守			
国内代理店	SCSK・CTCSP・日本コムシス・Thymos HPCソリューションズ など	販売パートナー	富士通・NEC・日立製作所・Crayジャパン・ HPE・新日鐵ソリューションズ など	



グローバル・マーケットで知られており、受賞実績も多数



DDN | DDNストレージの特長

● ハイパフォーマンスマニストレージ

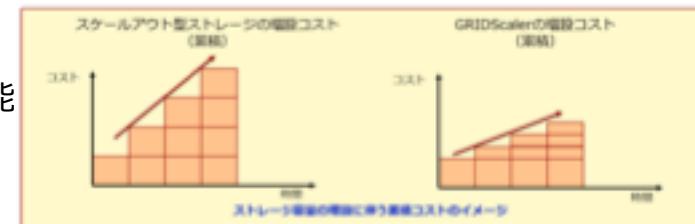
- ・最適に設計されたプラットフォームとユニークなOS機能による業界最速の並列ファイルシステムパフォーマンス（最大40GB/s：1ファイルシステムあたり）
- ・豊富なネットワークインターフェースをサポート（10/25/40/100Gb・16/32FC・Infiniband FDR/EDR・Intel Omnipath）

● 高密度設計

- ・4Uで90本の高密度設計が可能なディスクエンクロージャ
- ・1 Rack で、最大物理容量10.8 PBを搭載可能（12TB HDD使用時）
- ・40PBを1ファイルシステムで構築した実績

● ご要望に応じた拡張性とシステム可用性

- ・パフォーマンス（スループット/IOPS）と容量をフレキシブルに拡張可能
- ・すべてのコンポーネントレベルで冗長化

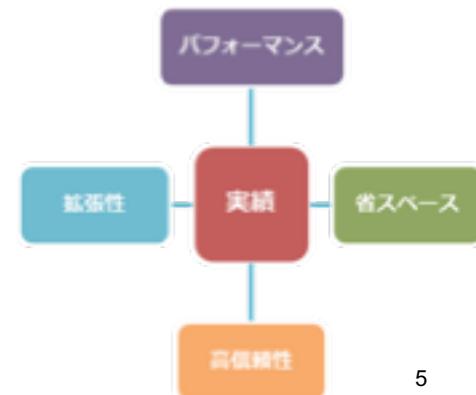


● 実績

- ・HPC/CAEの豊富な実績によりお客様が抱える様々なワークロードを解決
- ・ワールドワイドで1300サイト、日本で350サイトの導入実績
- ・HPC Top100ランクインの66%のお客様がDDNストレージを利用

● サポート

- ・日本全国24時間365日サポート対応可能な体制



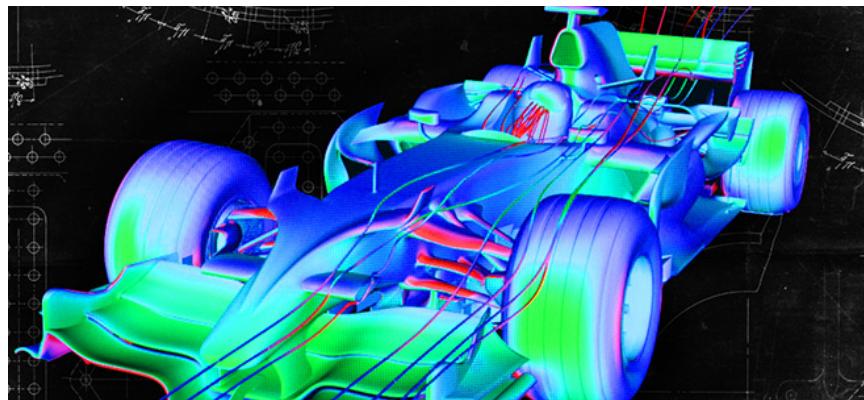
DDN | 導入実績① 學術研究機関

- 北海道大学
- 東北大学
- 筑波大学
- 東京大学
- 東京工業大学
- 名古屋大学
- 京都大学
- 大阪大学
- 九州大学
- JCAHPC
- 理研
- 産総研
- NICT
- JAXA
- JAMSTEC
- KAGRA
- NAOJ
- KEK
- 気象庁
- SACLAC
- 統数研
- 慶應大学
- JAIST
- NAIST
- OIST
- 物材研

DDN | 導入実績② ライフサイエンス・メディア・クラウド

- 東北メディカル・メガバンク機構
- 京都大学iPS細胞研究所
- 遺伝研
- かずさDNA研究所
- 東京大学Human Genome Center
- アステラス製薬
- 医薬基盤・健康・栄養研究所
- YAHOO Japan
- イマジカ
- 日本テレビ
- U-NEXT
- 白組

DDN | 導入実績③ エンジニアリング



- | | | |
|--------------------|----------------------|-------------------------------------------------------------------------------------------|
| ▶ マツダ株式会社 | 高速CAE用ファイルサーバ | https://ddn.co.jp/media/2016/11/29/24 |
| ▶ スズキ株式会社 | 高速CAE用ファイルサーバ | https://ddn.co.jp/media/2013/11/05/27 |
| ▶ 某大手自動車メーカー | 高速CAE（流体・構造）用ファイルサーバ | |
| ▶ 某大手自動車メーカー | 高速CAE・統合ファイルサーバ | |
| ▶ 某大手自動車メーカー | 高速CAE・統合ファイルサーバ | |
| ▶ 某大手自動車メーカー サプライヤ | 高速CAE用ファイルサーバ | |
| ▶ 某大手自動車メーカー サプライヤ | 高速CAE（流体・構造）用ファイルサーバ | |
| ▶ 某大手製薬企業 | 研究開発用大規模ファイルサーバ | |
| ▶ 某大手光学機器企業 | 高速CAEファイルサーバ | |
| ▶ 某鉄鋼メーカー | 高速大容量ファイルサーバ | |
| ▶ 某半導体企業 | 高速CAEファイルサーバ | |
| ▶ 某発動機・建機メーカー | 高速CAE・統合ファイルサーバ | |
| ▶ 某重工業メーカー | 高速CAE・統合ファイルサーバ | |
| ▶ 某総合電機メーカー | 高速CAE（流体）ファイルサーバ | |

TTDC
VMware Horizon Viewを利用したVDI環境を構築
<https://ddn.co.jp/media/2018/02/05/140>

DDN | 導入実績④ AI・マシンラーニング



産業技術総合研究所
人工知能研究センター

課題：研究競争力や産業界との連携の強化

解決策：GRIDScalerの採用で4PiBの実効容量、低遅延・高スループット、シングルスレッドI/O 3GB/秒を獲得

成果：Green500 Listで世界第3位を獲得、
AI研究開発用の共用計算プラットフォーム
として国内最大級の性能を獲得

https://ddn.co.jp/issue/machine_learning.html

- ▶ 国内 東京工業大学 学術国際情報センター <https://ddn.co.jp/media/2017/02/17/31>
- ▶ 国内 某自動車メーカー サプライヤ 自動運転研究開発
- ▶ 国内 某精密機械メーカー
- ▶ 海外 マイアミ大学 <https://ddn.co.jp/media/2017/06/05/69>
- ▶ 海外 某大手決済サービス会社
- ▶ 海外 フォーチュン100 SaaSプロバイダー
- ▶ 海外 某大手サービス会社

DDN | DDN®プロダクトポートフォリオ

Burst Buffer

Infinite Memory Engine (IME)

クライアントと並列ファイルシステムの中間レイヤーに位置づけられるBurst Bufferであり一時的なBurst IOに対応し、多数プロセスからShared Fileへ小さなランダムIOアクセスなど並列ファイルシステムが得意としているI/Oの高速化が可能。アプライアンスおよびソフトウェアによる提供。



IME240

File System

DDN EXAScaler®

- オープンソースのLustreファイルシステムをベースに構成
- 膨大な数のクライアントに対し高スループットIO環境を提供
- 多数のHPCサイトで多くの実績

DDN GRIDScaler®

- 高いIO性能に加えてマルチプロトコルアクセス環境
- Snapshot/Replication/DR/サイト間連携など多彩な機能を提供
- 低速ストレージ、テープ装置等を用いた階層化機能の実現

Storage Platform + Operation System

SFA14KX



54GB/sec / 1.7M IOPS
Max 840 Drive

SFA7700X



12GB/sec / 450K IOPS
Max 396 Drive

SFA®OS

Drive Enclosure / Drive

SS8460/8462



6/12Gbps 84Drive
Enclosure

SS9012



12Gbps 90Drive
Enclosure

Drive



7.2Krpm NL-SAS
10/15Krpm SAS

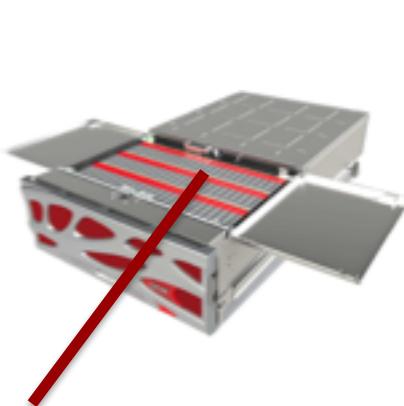
Server

DDN 1U Server



Intel Xeon skylakeベース
ファイルシステム管理用

DDN | ストレージプラットフォーム コントローラ SFA14KX 仕様



2.5インチSSD 72本搭載可能
もしくは48 個のNVMeを搭載可能



SFA14K+SS8462/10台接続構成

項目	製品仕様
フォームファクター	4U (Dual ControllerおよびBBU、2.5inch Drive slot x72を内蔵)
コントローラ	Active/Active Dual Controller
対応 RAIDレベル	RAID1(1+1)・RAID5 (8+1 or 4+1)・RAID6(8+2 or 4+2)
最大理論性能	シーケンシャル Read/Write 45GB/s
対応ドライブ数	ベースエンクロージャー 48 x PCIe NVMe SSD ・ 72 x 2.5inch SAS SSD 拡張エンクロージャー 1680 3.5inchドライブ ※ 20x SS8462構成時 1440 3.5inchドライブ ※ 16x SS9012構成時
キャッシュ	64GB / Controller (Battery-Backup)
ホストインターフェイス	6x EDR/FDR Infiniband 4 x Omni-path 4 x 100GbE イーサネット 24 x16GbpsFC
ソフトウェア機能	LUNマッピングおよびマスキング・Read QoS・ポートゾーニング・Direct Protect
管理機能	Web GUI管理・CUI管理・メールベースの障害発報・SNMP

DDN | ストレージプラットフォーム ディスク拡張エンクロージャ SS8462

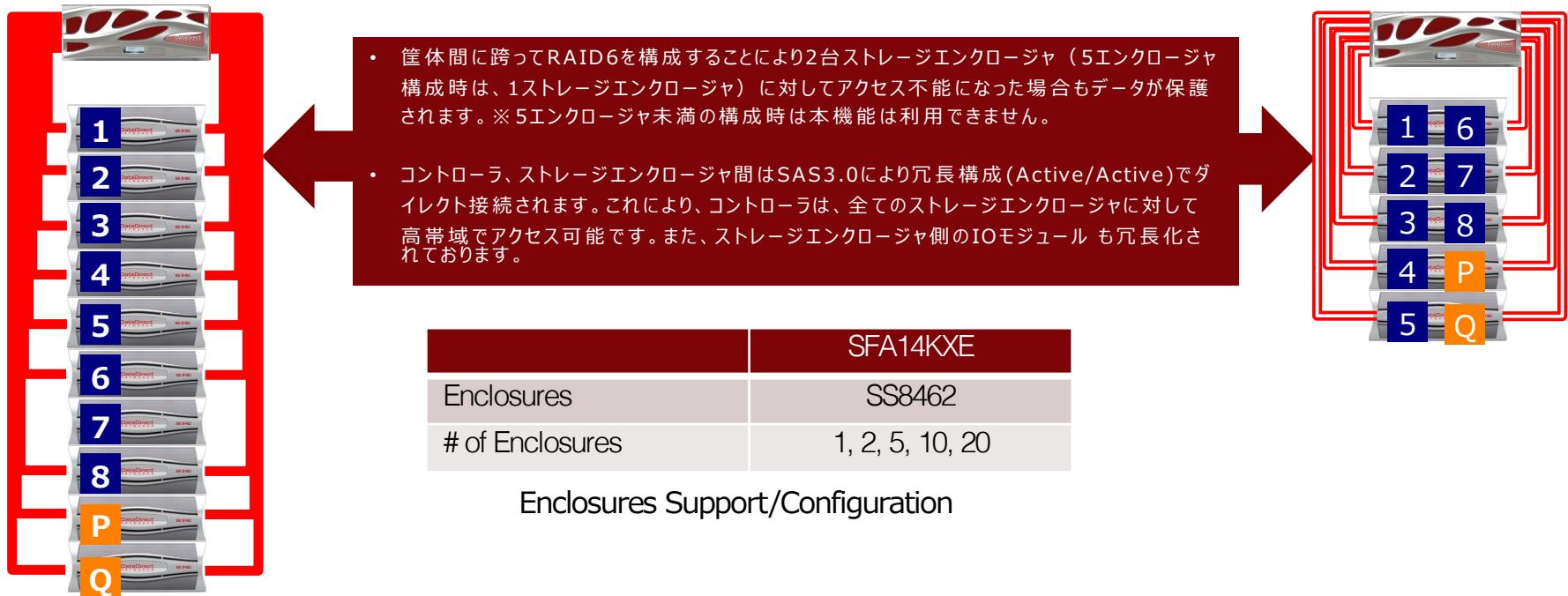
SS8462ディスクエンクロージャは、4Uで84本のドライブを収容可能な高密度ストレージエンクロージャです。SSD、SAS、SATAの各種ドライブを混在可能であり、I/Oバスも含めてコンポーネントが冗長化されております。



項目	製品仕様
フォームファクター	4U / 84 Drive
スロット数	84 (3.5inch/2.5inch)
対応ドライブ	SSD, SAS, NL-SAS (6Gb/s)
冷却ファン	2x 冷却ファン
電源	2 x 1200W 80PLUS Platinum
ホットスワップ対象部品	ドライブ、電源、冷却ファン、IOモジュール
LCD Display	Status, Power, Environmental Monitoring
LED	Power, Status, Monitoring, Drive Activity
冷却ファン	2x 冷却ファン

DDN | SFA14KX ダイレクト接続

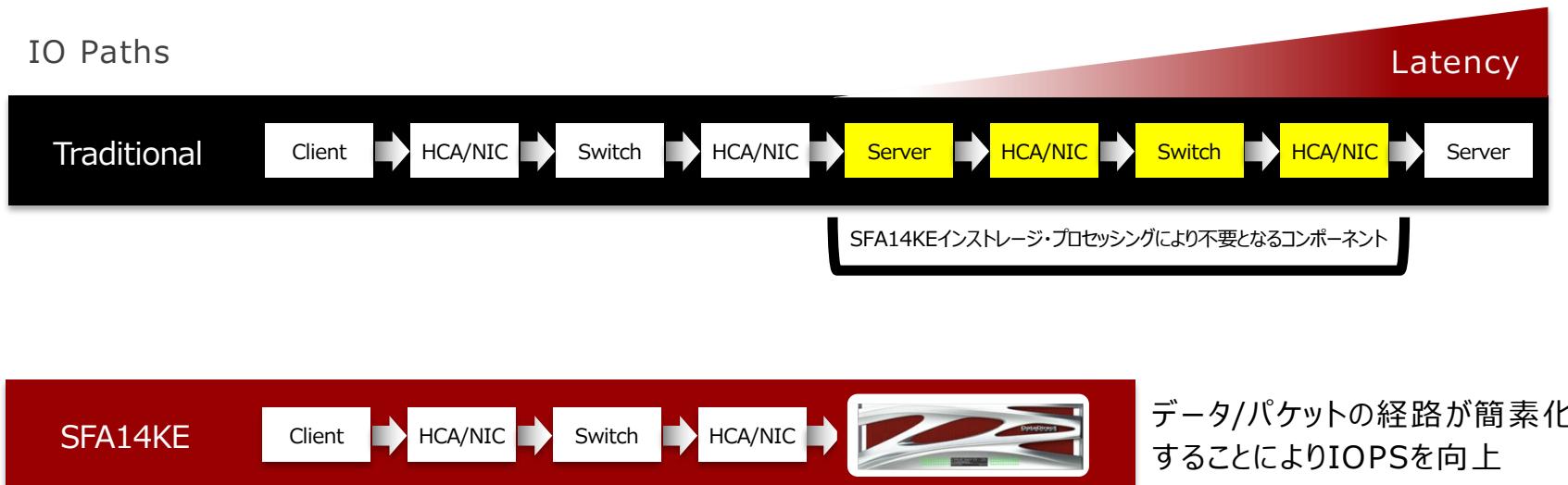
SFA14KXではコントローラ～ストレージエンクロージャ間がSAS 3.0インターフェイスにより直接接続されます。これにより、完全な冗長性と高内部帯域幅を提供します。



DDN | SFA14KXE インストレージ・プロセッシング

ストレージコントローラ内でファイルサービスを実装することにより、外部サーバ、ネットワークアダプタ、スイッチなどを必要としません。主なメリットは以下の通りです

- ・ 外部サーバ、ネットワークアダプタ、スイッチ購入コスト削減
- ・ 構成の簡素化による管理コストの削減
- ・ 低遅延でのストレージアクセスを実現（下記図参照）



Agenda

- DDN会社/製品紹介
- Lustreについて
- Lustreの特徴
- Lustreのアーキテクチャ
- Lustreの各コンポーネント
- OCTOPUSストレージ概要
- OCTOPUSファイルシステムの利用
- Q&A

Lustreについて

- コンピュートクラスタ(HPC)向け並列分散共有ファイルシステム
 - ▶ 多くのシステムから一つのファイルシステムを利用可能
 - ▶ 高IOスループットと大容量を実現
 - ▶ 10万台規模のクライアントを構成可能
- HPCおよびその他の分野で利用
- オープンソース、コミュニティ主導の開発
- 種々の高速ネットワークインターフェースをサポート
 - ▶ Infiniband, OmniPath等
- 高可用性運用構成も可能

Why Lustre?

- ネットワークファイルシステム：CIFS, NFS,,,
- 通常はボトルネックがあり、低パフォーマンス
- 理由は少ないサーバ台数、低いネットワーク帯域等々、、、
- HPCおよびIOインтенシブアプリケーションにとって大きな問題
- Lustreは、多くのクライアントからなる大規模システム向けに高いパフォーマンスを提供可能

Lustreと言えば、 、、、

- 古典的ファイルシステムでは、 、、、
 - ▶ RAIDストレージ、 非RAIDストレージ
 - ▶ ノードから直接ファイルシステムにアクセス可能
 - ▶ スケーラビリティ・ボトルネック
- 並列分散ファイルシステム
 - ▶ ネットワークベースの分散データシステム
 - ▶ サーバ数に比例してパフォーマンスが向上
 - ▶ データの一貫性確保のための機能実装
- **Lustreは代表的な並列ファイルシステム**

Lustreの歴史

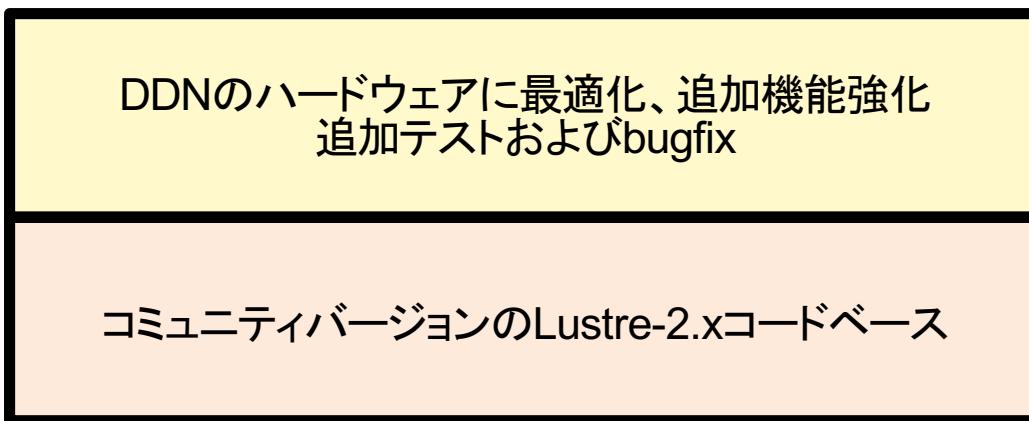
- 2000年代半ば、米国エネルギー省(DOE)科学部および国家核安全保障局(NNSA)のファンドによって開発開始
- ClusterFS(CFS)がファンドを受ける
- 2007年、サンマイクロシステムズがCFSを買収
- 2010年、オラクルがサンを買収したが、オラクルはHPCを知らなかつた。。。
- Lustreエンジニアはオラクルを離れWhamCloud/Xyratex/DDNへ
- 2011年、インテルがWhamCloudを買収
- XyratexがLustre商標、ロゴ、関連資産(サポート)をオラクルから引継ぎ(その後、Communityに返還)
- 2014年、SeagateがXyratexを買収
- 2017年、CrayがSeagateからClusterStor製品ライン(Xyratex)を買収
- **2018年、DDNがインテルのLustre Business Unitを買収**

DDNによるIntel Lustre Business Unit買収

- 買収後も100% Lustre CommunityにCommit
 - ▶ 独自Branchの作成はしません
- Intel Lustreの顧客に対するサポートを継続
 - ▶ 他社製HW上のLustreもサポートします
- DDN ExaScaler
 - ▶ 従来通り、Lustre Community EditionをCoreとし付加価値を追加しパッケージング
- Lustre開発者の殆どがDDNに移籍
 - ▶ 元々CFS/Whamcloudの人間が殆どで顔見知りであり、Communityで協業していた

DDN ExaScalerとLustre

- DDN Lustre(EXAScaler)は、Lustre Community Editionをベースに性能および機能強化したDDNが提供するLustreディストリビューションです。全世界で300を超えるシステムで動作しています
- DDNのハードウェアに最適化され、お客様に提供するシステムと同等構成にてテストおよび検証されています。また、DDNのレベル1, 2, 3サポートにてストレージハードウェアとLustreのシームレスなサポートをDDNよりワンストップにて提供します



DDN Lustre
(EXAScaler)

Lustre Community
Edition

Agenda

- DDN会社/製品紹介
- Lustreについて
- **Lustreの特徴**
- Lustreのアーキテクチャ
- Lustreの各コンポーネント
- OCTOPUSストレージ概要
- OCTOPUSファイルシステムの利用
- Q&A

Lustreの特徴

パフォーマンス	1TB/secを越えるスループットを実現可能
アーキテクチャ	オブジェクトベースファイルシステム
適応性	様々なネットワークインターフェースとストレージ製品をサポート
スケーラブル	分散ファイルオブジェクト処理により10万クライアント規模をサポート
安定性	製品自体の安定性とフェイルオーバ機能
モジュラー型	機能拡張が容易なインターフェース
高可用性	HA構成によりSPOF無し
拡張性	オンライン容量拡張
OSサポート	ほとんどのLinuxディストリビューションで動作
互換性	POSIX互換
アカウンティング	ユーザ、グループおよびプロジェクトクォータ
オープン	オープンソース
コミュニティ	LAD(Lustre Admins Developers) & LUG(Lustre User Group)
Get Involved!	OpenSFS and/or EOFS

Lustre Numbers

- TOP100システム中の75以上のシステムで使用
- 導入システムでの最大値

Rawバンド幅と比較したFile IO%	~90%
実効シングルOSS IO性能	> 8 GB/s
実効シングルクライアントIO性能	> 4 GB/s
実効総IO性能	> 1 TB/s
平均メタデータ操作回数	50,000+ ops/s ???
最大ファイルサイズ	64PB

Lustre Numbers

- 実装上の制限値

最大ファイルシステムサイズ	1 EiB
最大OST数	8,150
最大MDT数	256 / ファイルシステム
最大ファイル数 / ディレクトリ	1000万(ldiskfs)、 2^{48} (ZFS)
最大ファイル数 / ファイルシステム	40億(ldiskfs)、256兆(ZFS)
最大ストライプ数	2,000
最大ストライプサイズ	< 4 GiB
最小ストライプサイズ	64 KiB
最大オブジェクトサイズ	16 TiB(ldiskfs)、256 TiB(ZFS)
最大ファイルサイズ	31.25 PiB(ldiskfs)、8 EiB(ZFS)

Lustre Limitations

- データとメタデータが厳密に分離
 - ▶ 2.11よりDoM(Data on MDT)が実装されたため、2.11以降は異なる
- シングルMDTのメタデータ性能
 - ▶ DNE(Distributed Namespace Environment)の利用により改善
- RPC Size : 1MB, 16MB(Large Bulk IO)
- 最大OSTサイズはLustreのバージョンによる
 - ▶ 実システムでは256TiB程度が上限
- ext4(ldiskfs) バックエンドファイルシステムによる制限
 - ▶ ZFSも利用可能
- ディレクトリクォータ
 - ▶ 以前はユーザ、グループクォータのみ
 - ▶ プロジェクトクォータにてディレクトリクォータをサポート

Lustreの利用

- Lustreに適するのは、 、 、
 - ▶ ラージファイル
 - ▶ シーケンシャルアクセス
 - ▶ プロセス毎に異なるファイルに書き込む並列アプリケーション
- Lustreが苦手とするのは、 、 、
 - ▶ 大量の小さいファイル
 - ▶ 莫大なメタデータ操作、 ただしバージョンアップ毎に改善
 - ▶ OST容量の平準化

今後リリース予定の機能

- 障害復旧処理の迅速化
- ChangeLogsを利用したHSMインテグレーション
- ZFS、BtrFSをバックエンドファイルシステムとして利用
- ただし、Roadmapはかなり流動的です、、、
- OpenSFSとEOFSが新機能の実装について決定します

Agenda

- DDN会社/製品紹介
- Lustreについて
- Lustreの特徴
- Lustreのアーキテクチャ
- Lustreの各コンポーネント
- OCTOPUSストレージ概要
- OCTOPUSファイルシステムの利用
- Q&A

Lustre略語集

DNE	Distributed Namespace Environment	単一ネームスペースを複数MDTで管理する機能
FID	File IDentifier	1ファイルシステム内のオブジェクトを示す128bitの一意な識別子
LMV	Logical Metadata Volume	複数MDTを扱う為のクライアントソフトウェアレイヤ
LNET	Lustre NETtwork	Lustre Networkの総称
LOD	Logical Object Device	複数のMDT, OSTへのアクセスを管理するMDSソフトウェアレイヤ
LOV	Logical Object Volume	複数OSTを扱う為のクライアントソフトウェアレイヤ
MDC	Metadata Client	MDSへのインターフェースとなるクライアントソフトウェアレイヤ
MDD	Metadata Device Driver	POSIXファイルアクセスを管理するMDSソフトウェアレイヤ
MDS	MetaData Server	クライアントからファイルシステムネームスペース(inode, path, permission等)へのアクセスを管理するソフトウェアサービス
MDT	MetaData Target	ファイルシステムメタデータ(属性, inode, ディレクトリ, xattr等)を保持するストレージデバイス
MGS	ManaGement Server	クライアント・サーバの構成情報を管理するサービス
MGT	ManaGement Target	コンフィグレーション情報を保持するストレージデバイス
NID	Network IDentifier	ノードIDとネットワークタイプによってLustre Networkのエンドポイントを一意に識別するために使用

Lustre略語集

OFD	Object Filter Device	ファイルIOを扱うOSSソフトウェアレイヤ
OSC	Object Storage Client	OSSへのインターフェースとなるクライアントソフトウェアレイヤ
OSD	Object Storage Device	ldiskfs(ext4)やZFSなどのバックエンドディスクファイルシステムへアクセスするMDDとOFDを抽象化するサーバソフトウェアレイヤ
OSP	Object Storage Proxy	あるMDSから別のMDSもしくはOSS上のOSDにアクセスするインターフェースとなるサーバソフトウェアレイヤ
OSS	Object Storage Server	ファイルシステムデータへのアクセス(read, write, truncate等)を管理するソフトウェアサービス
OST	Object Storage Target	ファイルシステムデータを保持するストレージデバイス

Lustreを構成するノードタイプ

- Management Servers (MGS)
 - ▶ システムに一つ、 MDSと共に存可能
- Meta-Data Servers (MDS)
 - ▶ ファイルシステム毎に最低一つ
 - ▶ 1 MDSにつき一つ以上のMeta Data Target(MDT)
- Object Storage Servers(OSS)
 - ▶ 通常ファイルシステム毎に複数存在
 - ▶ 各OSSが実ファイルデータを保持する一つ以上のObject Storage Target(OST)を提供
- クライアント
 - ▶ Lustreクライアントソフトウェアが Lustre ファイルシステムへのアクセスを提供

Lustre MGS

- ManaGement Server
- Lustreシステムの構成DBを管理するサーバ
- クライアント、OSSおよびMDSはMGSにアクセスできなければならぬ
- MDSと共存可能
 - ▶ DDNが国内で提供するLustreはほぼ100%共存構成

Lustre MDS

- Meta Data Server
- ネームスペースを管理し、ファイルパーティションやOST上のデータオブジェクトへのポインタなどをMDT上に保持する
- 通常、MDSはActive/Standby HA構成とする
- 一つのMDTは、一つのMDSで管理する必要がある
- メタデータのサイズは2KiB～16KiB

Lustre OSS

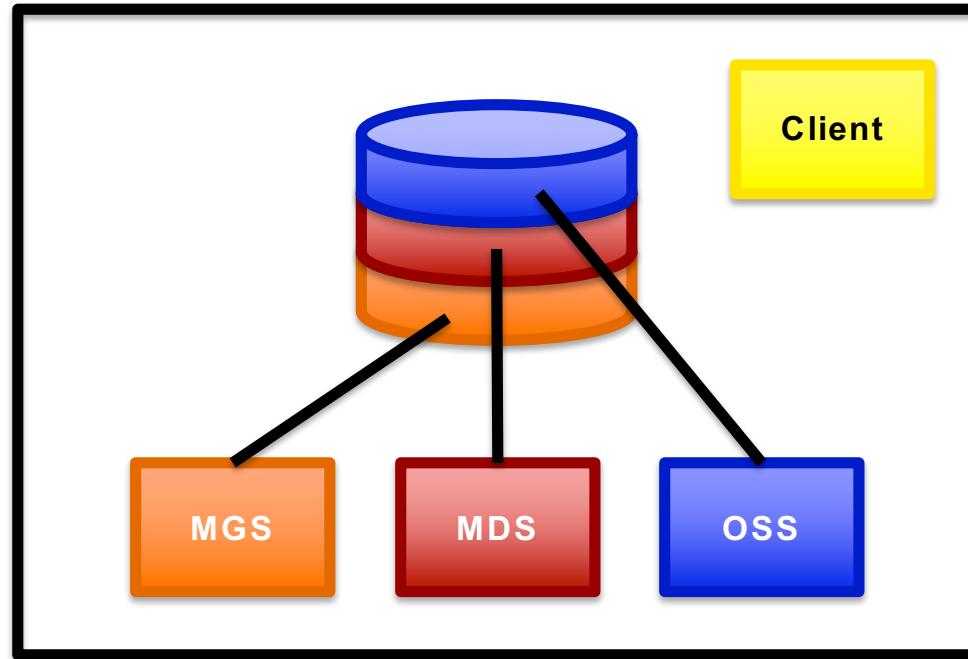
- Object Storage Server
- 実際のファイルIOを提供し、配下のOSTへのリクエストを処理する
- 通常、OSSはActive/Active HA構成とする
- OSSが複数のOSTを管理する構成が一般的

Lustreクライアント

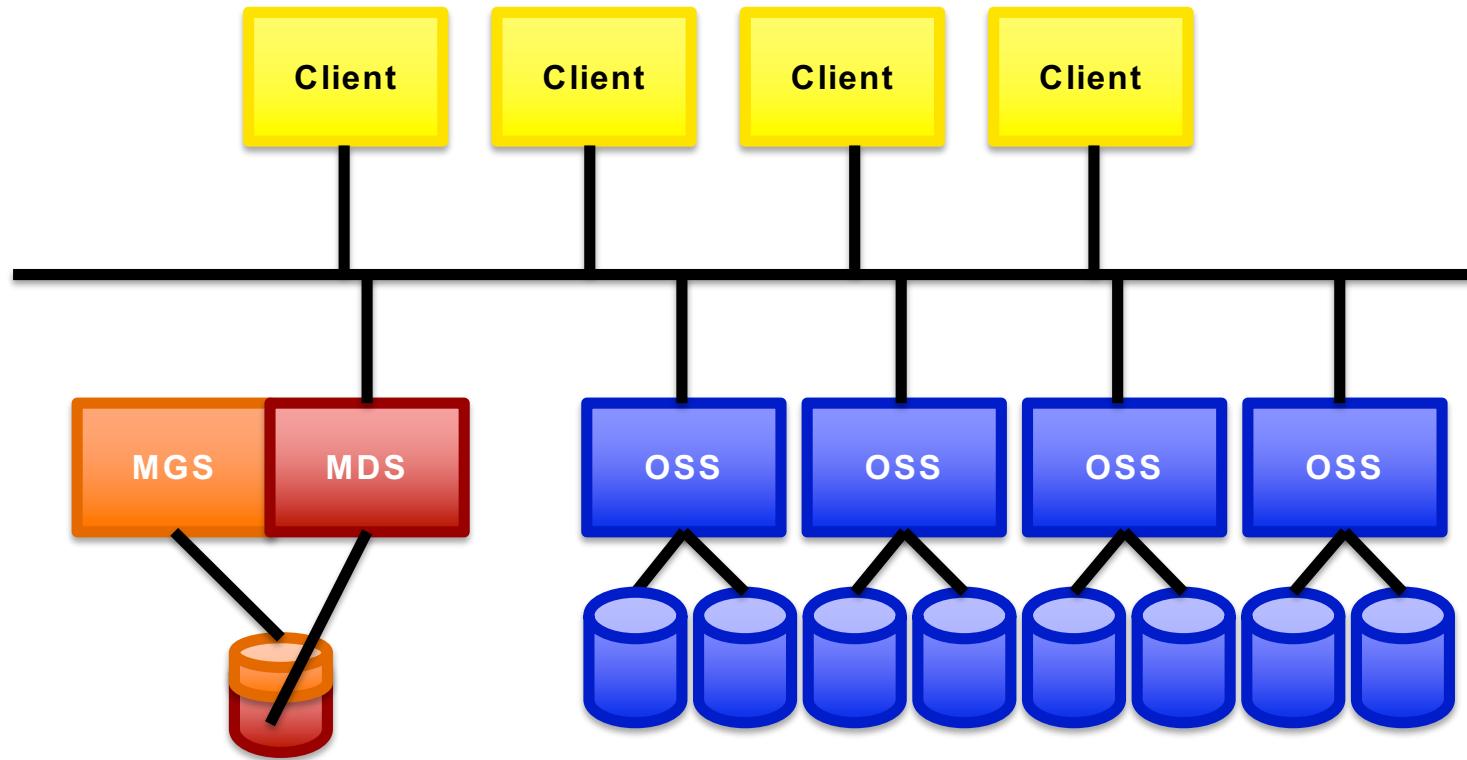
- POSIX互換ファイルアクセスレイヤ
- DNE環境下での複数MDTへのアクセスを管理するLogical Metadata Volume(LMV)レイヤを持つ
- 複数のOSTへのファイルストライピングを管理するLogical Object Volume(LOV)レイヤを持つ
- 各サーバに対するクライアントレイヤを持つ
 - ▶ Object Storage Client(OSC) : OSSと協調
 - ▶ Metadata Client(MDC) : MDSと協調

Lustre最小構成例

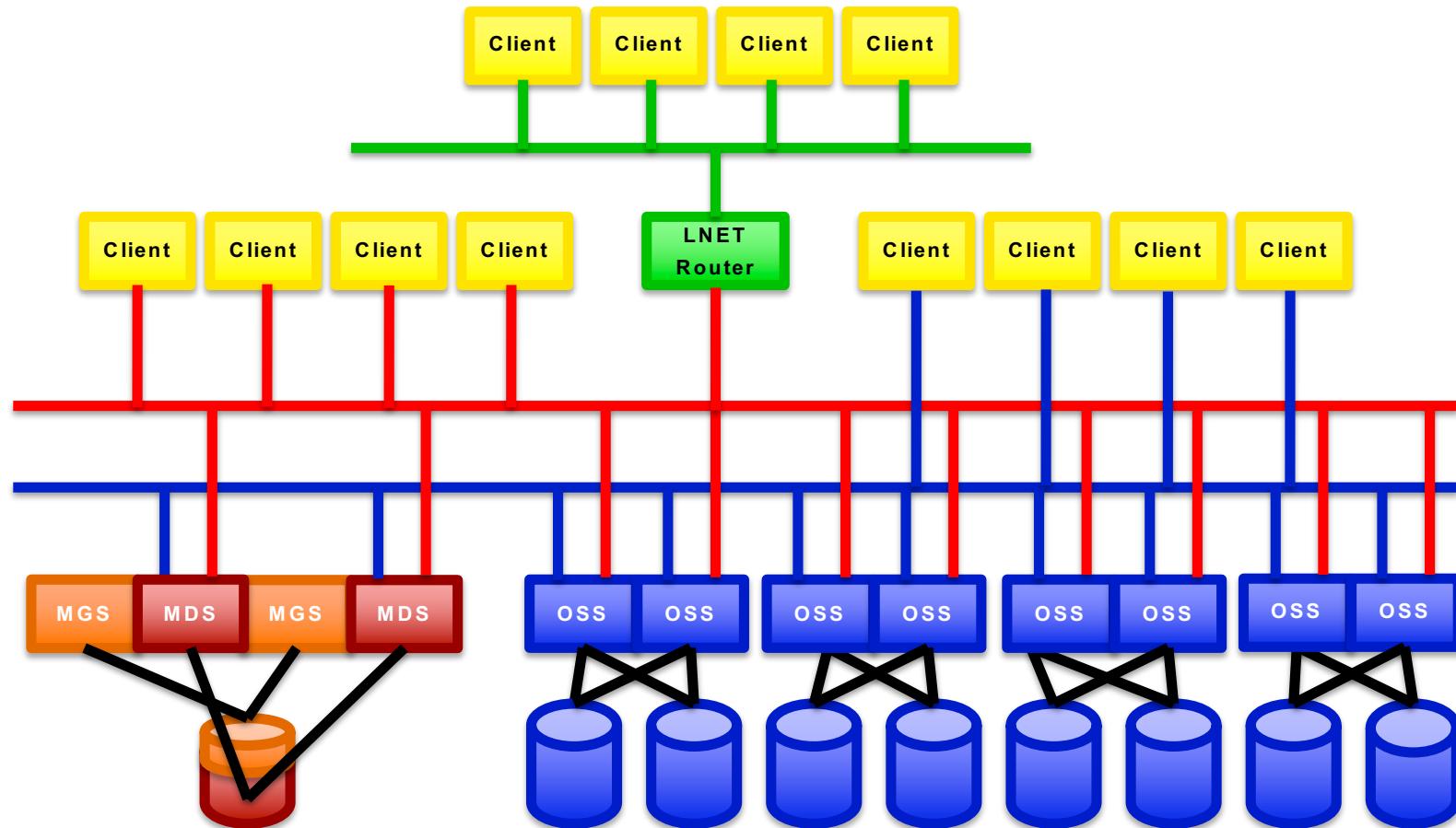
- 単一ノード構成
- 全機能を1サーバーノードに実装



Lustre基本構成



複雑な構成



Agenda

- DDN会社/製品紹介
- Lustreについて
- Lustreの特徴
- Lustreのアーキテクチャ
- Lustreの各コンポーネント
- OCTOPUSストレージ概要
- OCTOPUSファイルシステムの利用
- Q&A

Lustreコンポーネント配置図

Lustre Client

Manage stripes

Logical Object Volumes

OSC

Lustre Client File System

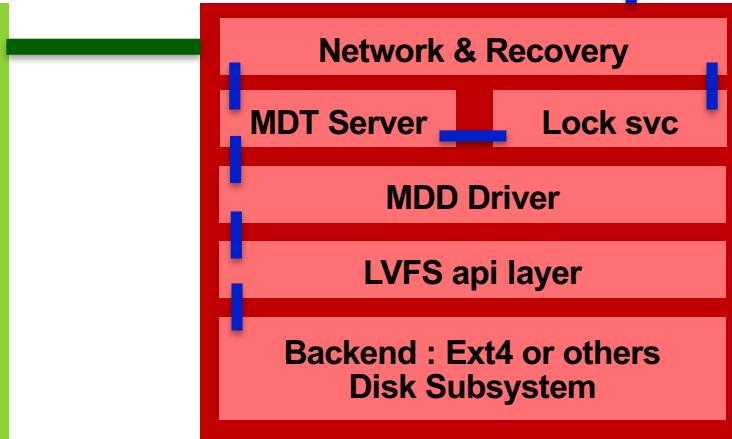
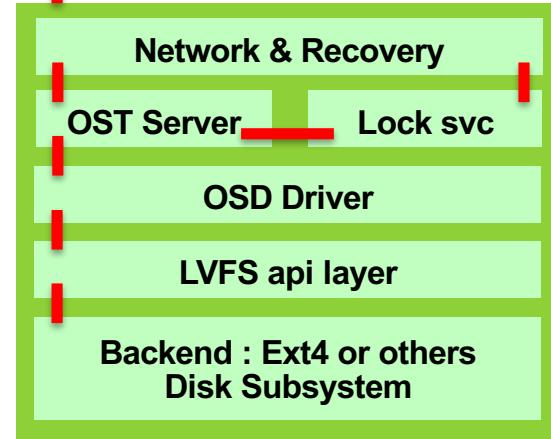
Logical Metadata Volumes

MDC Lock svc MDC

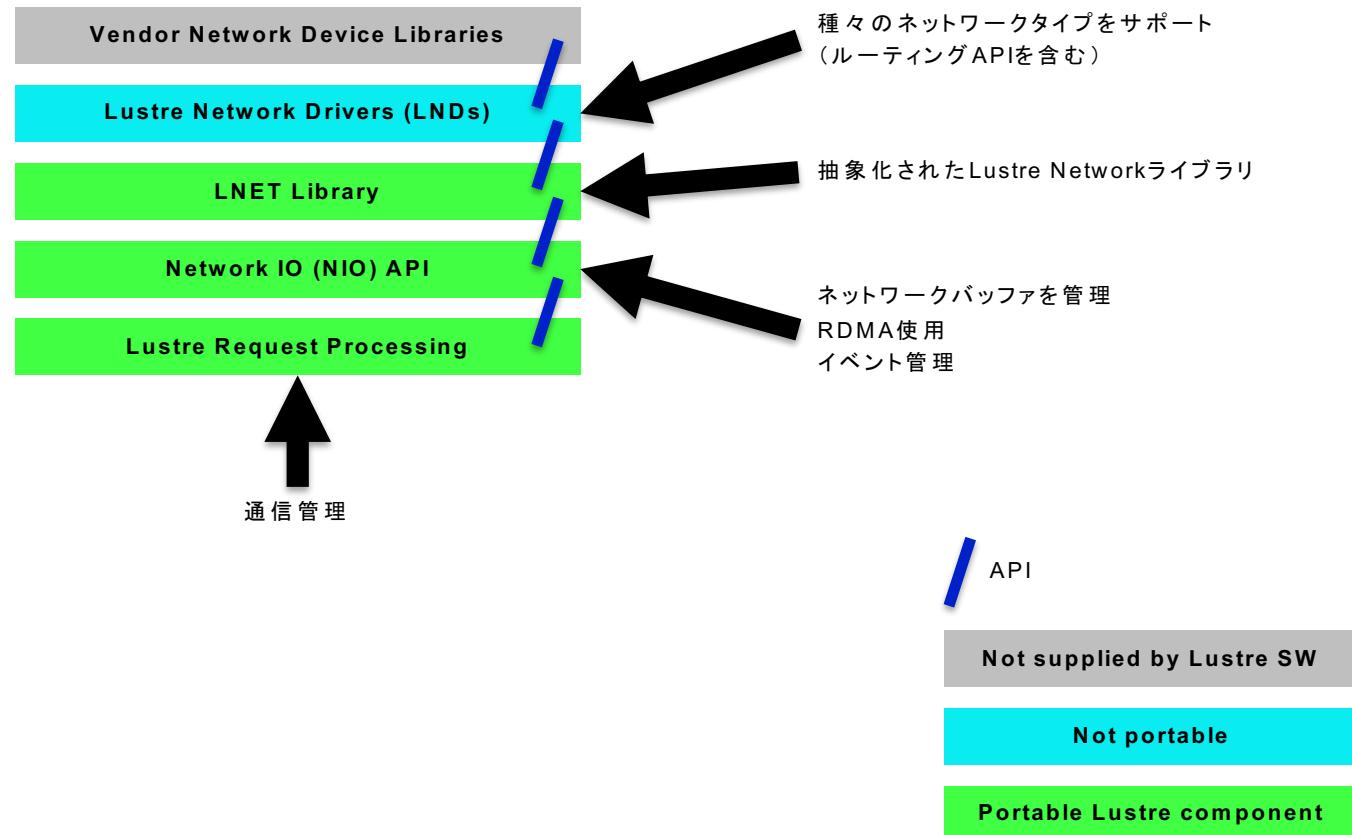
Network & Recovery

Data Object &
Lock protocol

Metadata &
Lock protocol



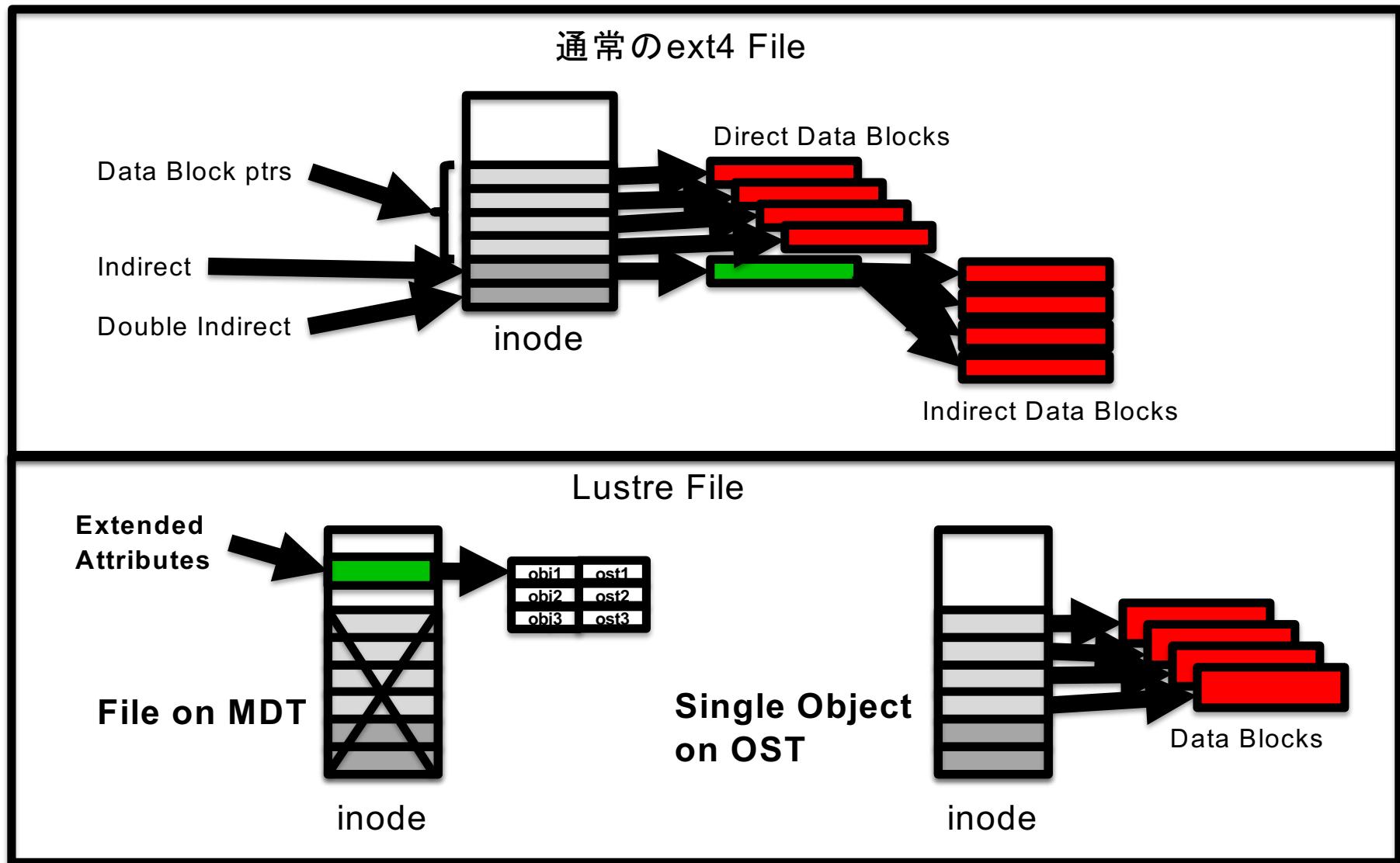
階層型ソフトウェアモジュール設計



Lustre FS Storage & IO

- FID (Lustre File IDentifiers)
 - ▶ 128bit
 - 64bit sequence number (1 FS内でユニーク)
 - 32bit object id
 - 32bit version number
- OST上のファイルデータの位置
 - ▶ MDTオブジェクトのExtend Attributeに保持
- クライアントからのIO(read/write)操作要求
 - ▶ MDS/MDTからファイルデータの位置情報を取得
 - ▶ 位置情報を元にOSS/OSTにファイルIOを要求

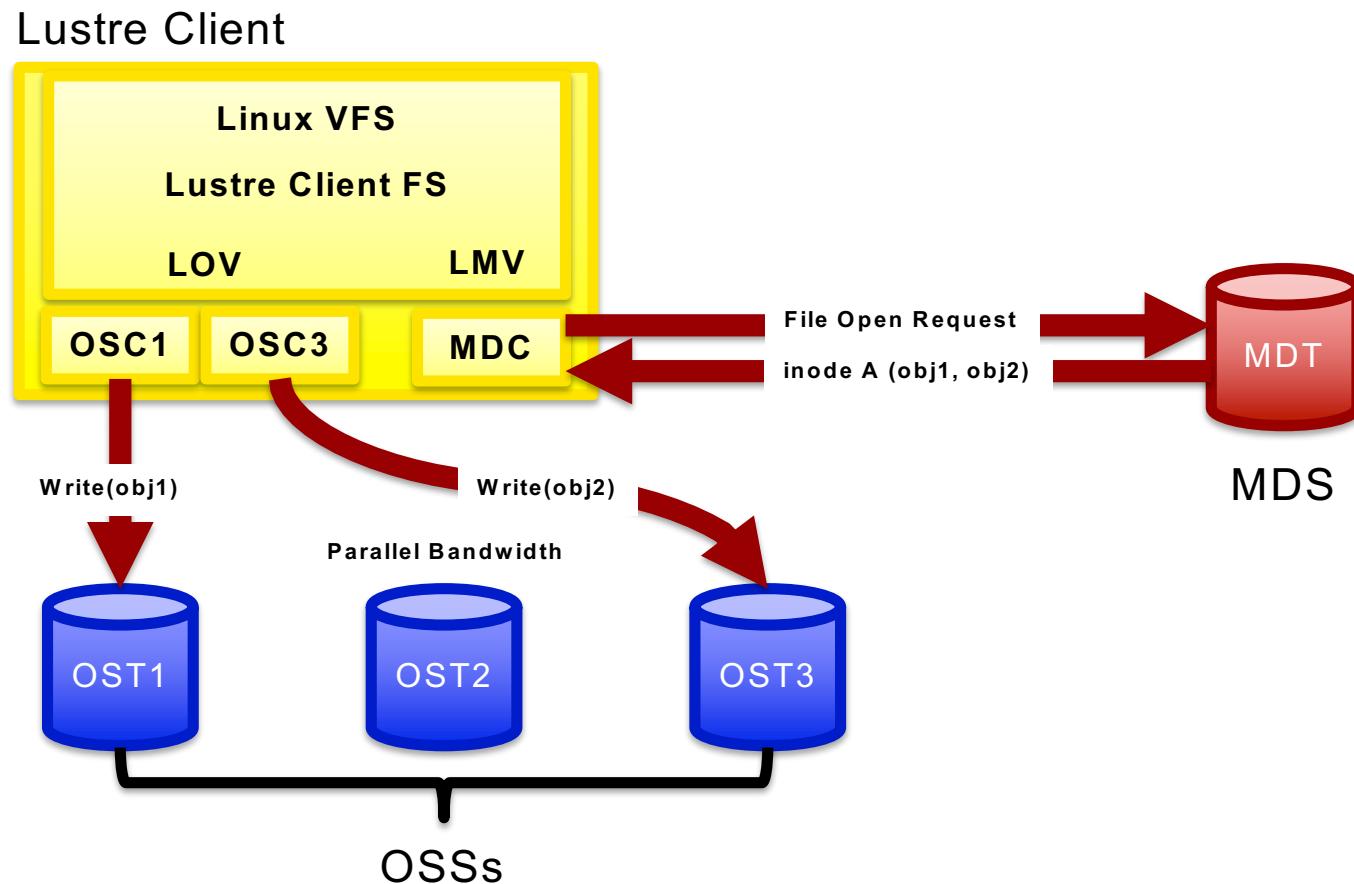
Disk Layout



Lustre File IO

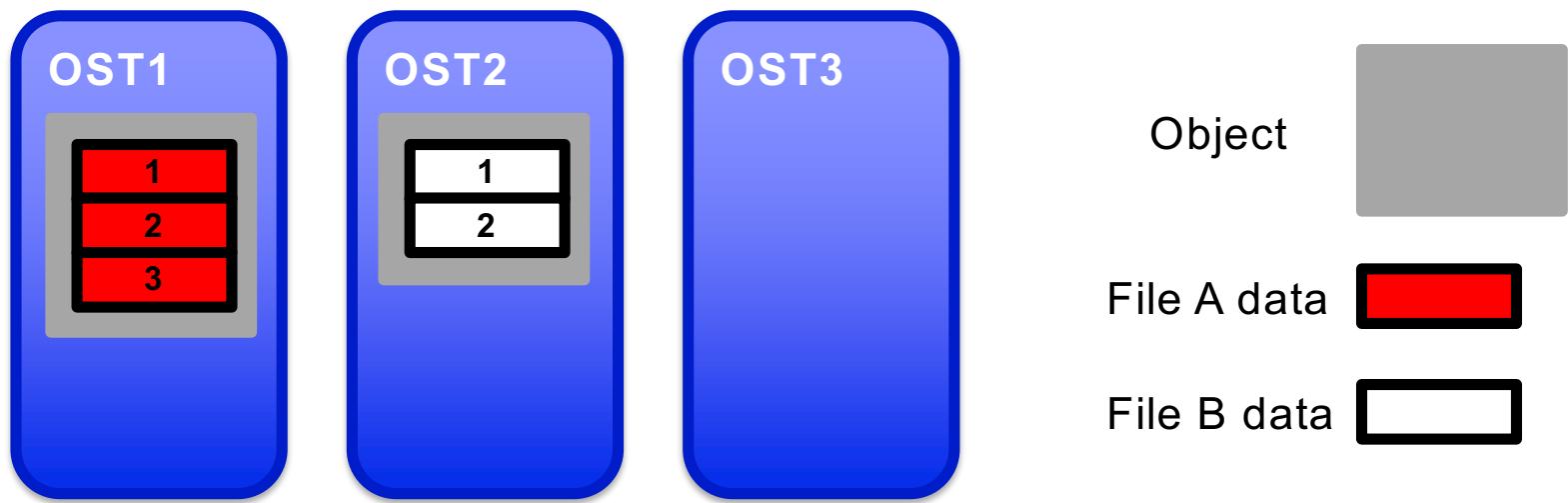
- クライアントがファイルOpenを要求
 - ▶ LMVがMDSへリクエストを送信
 - ▶ MDSはdirectory treeをロックし、ファイルをCreate/Openした後、directory treeのロックを解除
 - ▶ MDSはOpenしたファイルのinodeをLMVに返信
- LMVがクライアントアプリケーションにファイルディスクリプタを返す
 - ▶ inodeはOST位置情報を含むextended attributesを持つ
 - ▶ LOVはストライプを含むOSTの情報を元にファイルIOを行う
- Open後のファイルIOは主にクライアントとOSS間で行われる

Lustre File IO



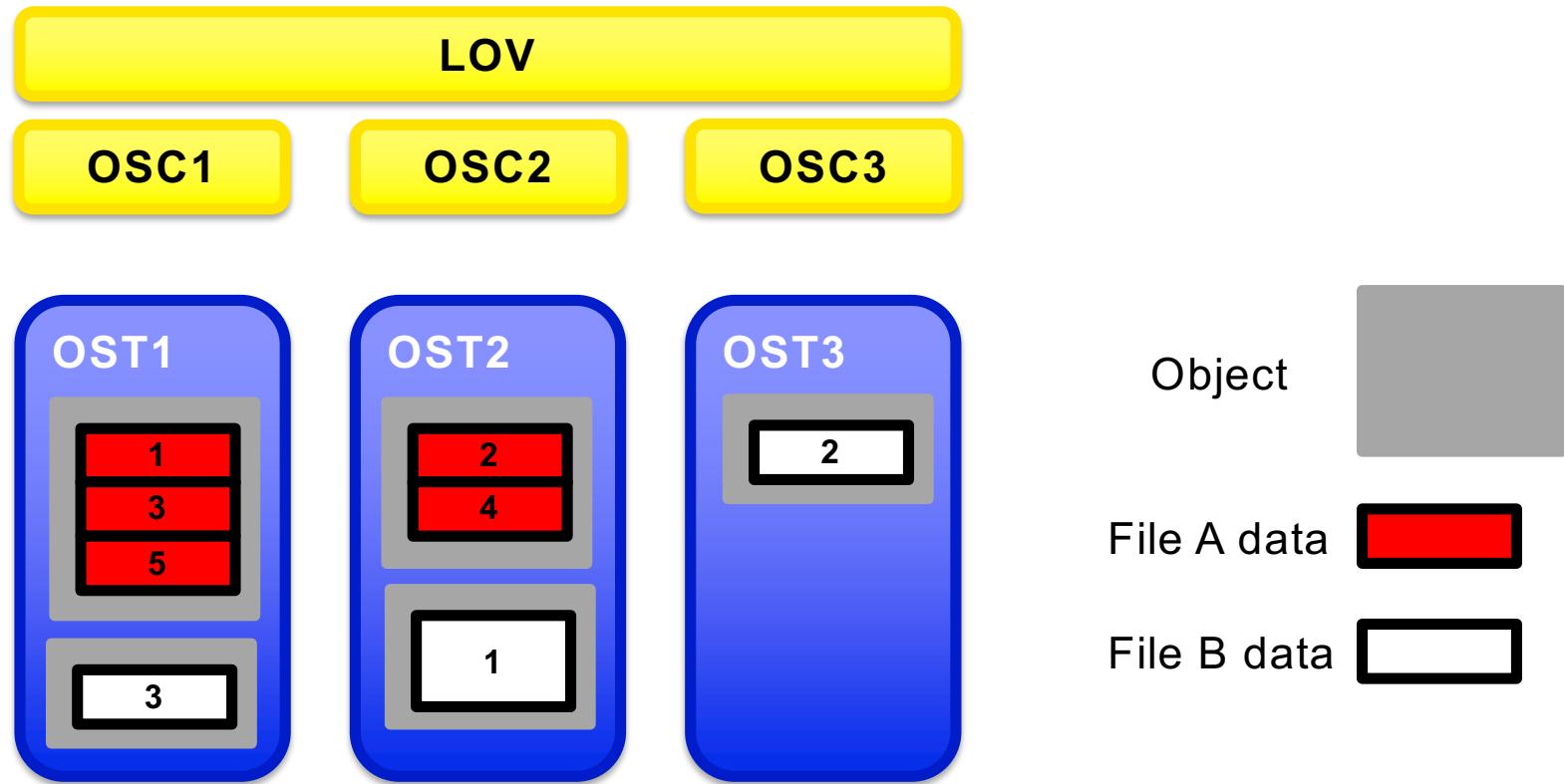
Lustre File Striping

- File Stripingを利用することによりファイルを複数のOSS/OSTに分散して書き込むことが可能
 - ▶ Stripeをしない場合、Objectの最大サイズ(16TiB) == 最大ファイルサイズ
- Single Stripe(Stripe無し)例



Lustre File Striping

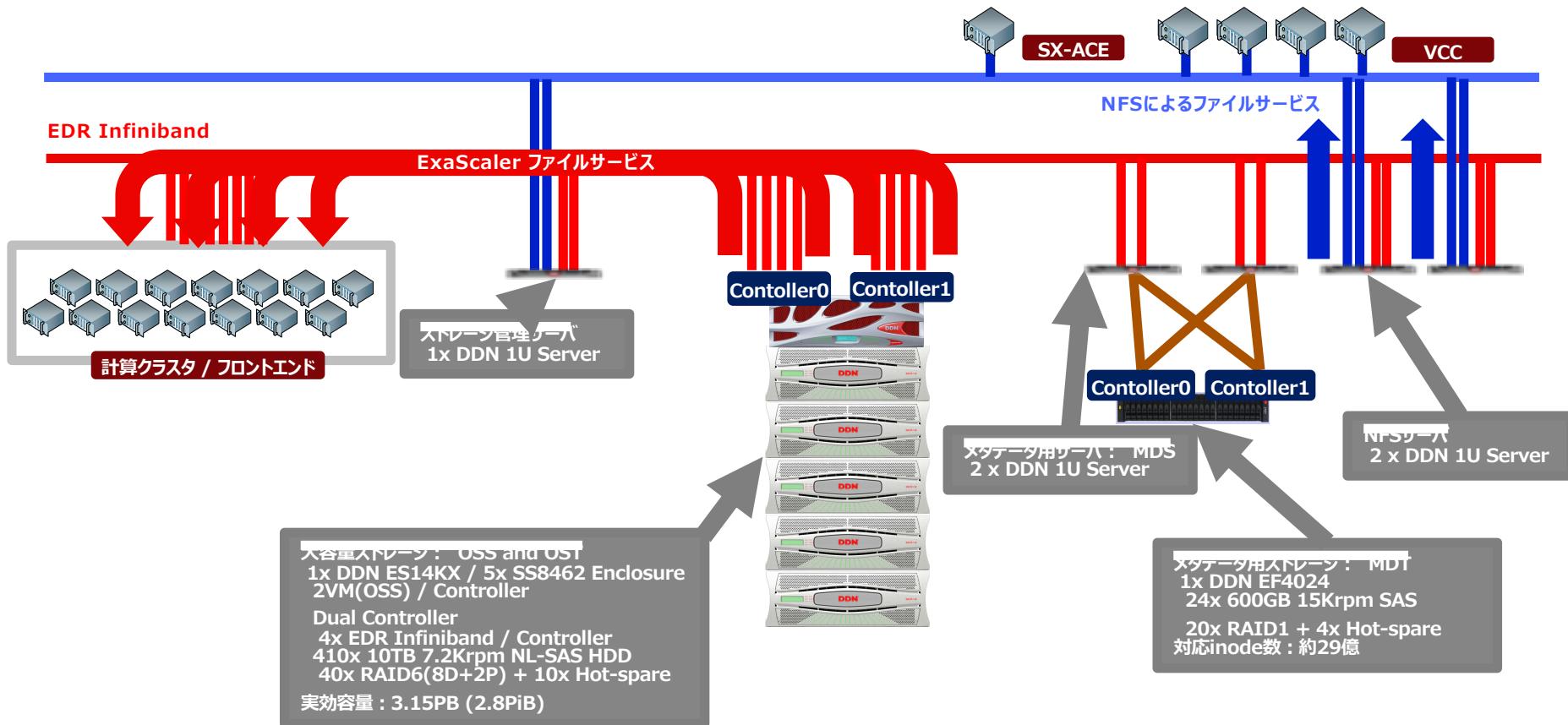
- Multiple Stripes 例



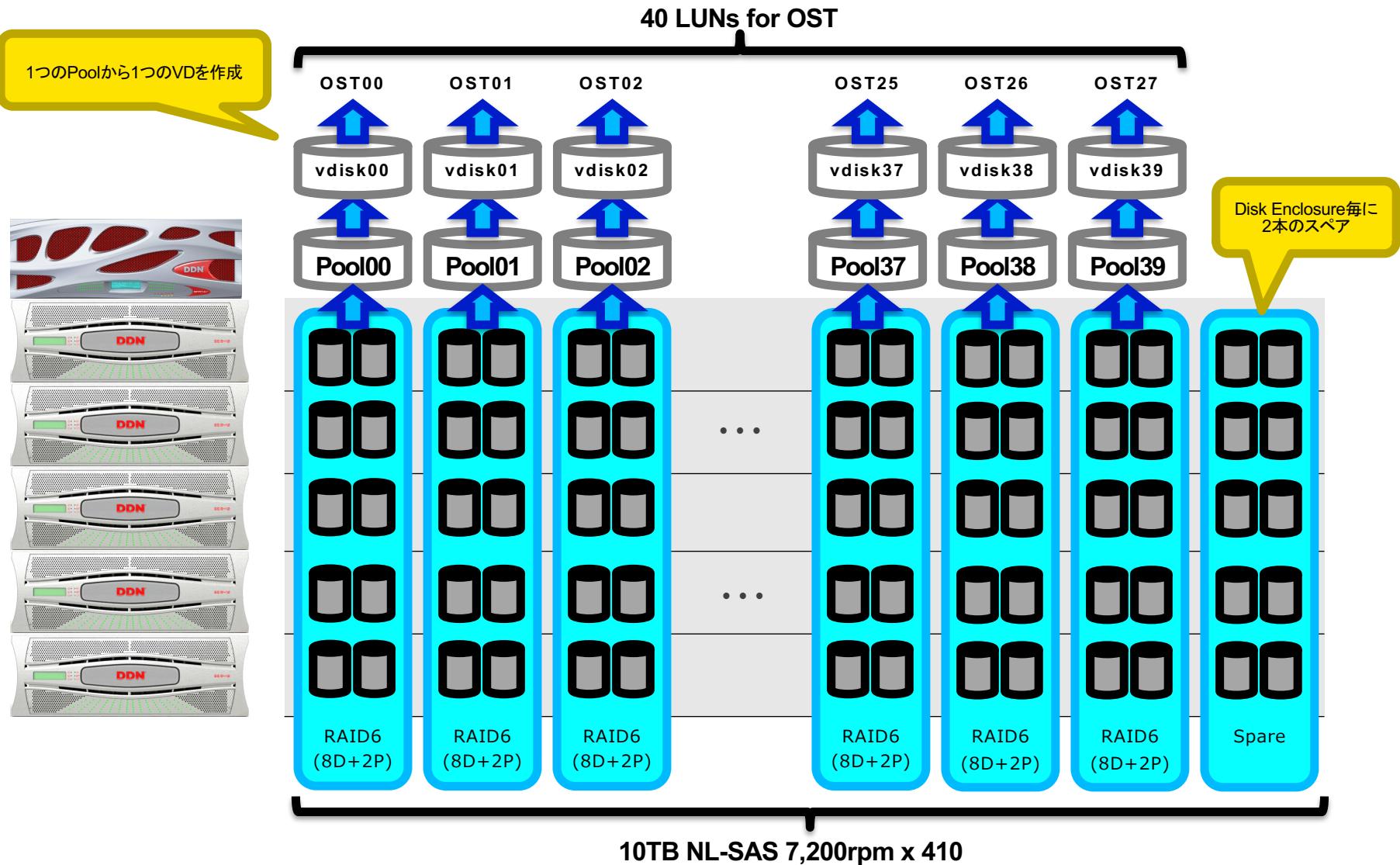
Agenda

- DDN会社/製品紹介
- Lustreについて
- Lustreの特徴
- Lustreのアーキテクチャ
- Lustreの各コンポーネント
- OCTOPUSストレージ概要
- OCTOPUSファイルシステムの利用
- Q&A

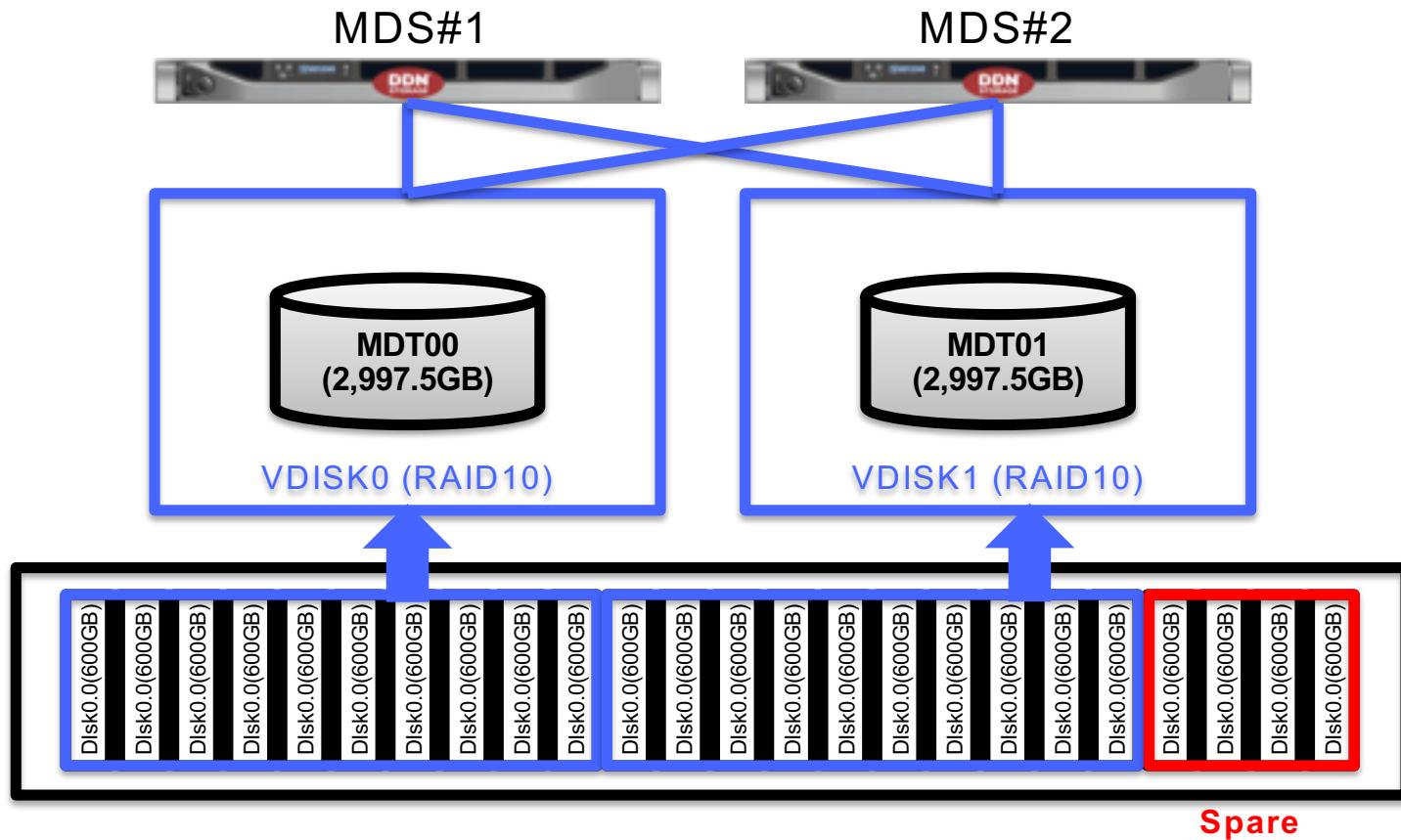
OCTOPUSストレージ構成図



OCTOPUS OST構成



OCTPUS MDT構成



OCTOPUS Lustre FS構成

octlfs : /octlfs
Capacity : 2.8PiB/3.15PB (FS format後)
of inodes : 2,927,242,752 (約29億個)

MDS#1 oct-mds01

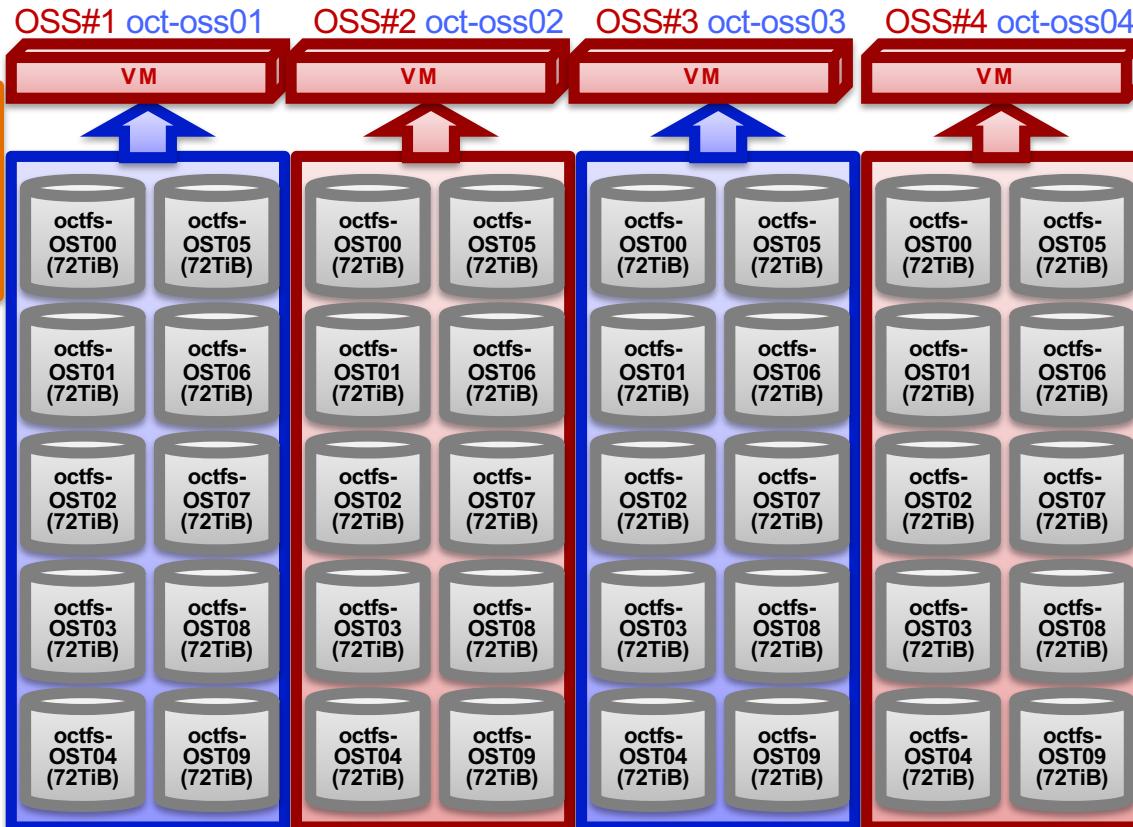
MDS#2 oct-mds02

octlfs-MDT
(5.45TiB)
Active

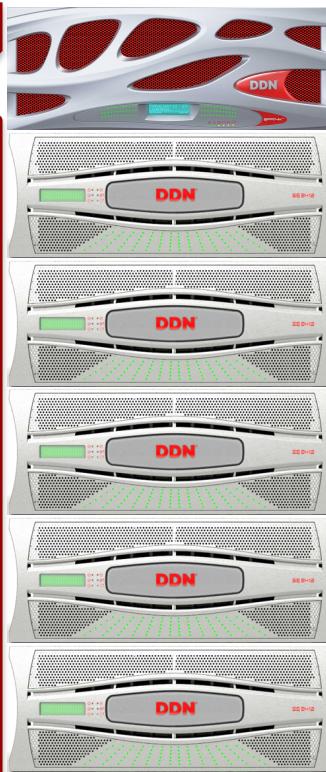
MDT00
(2.7TiB)

MDT01
(2.7TiB)

DDN EF4024 oct-mdt
600GB 15Krpm SAS x 24
(including Spare x 4)



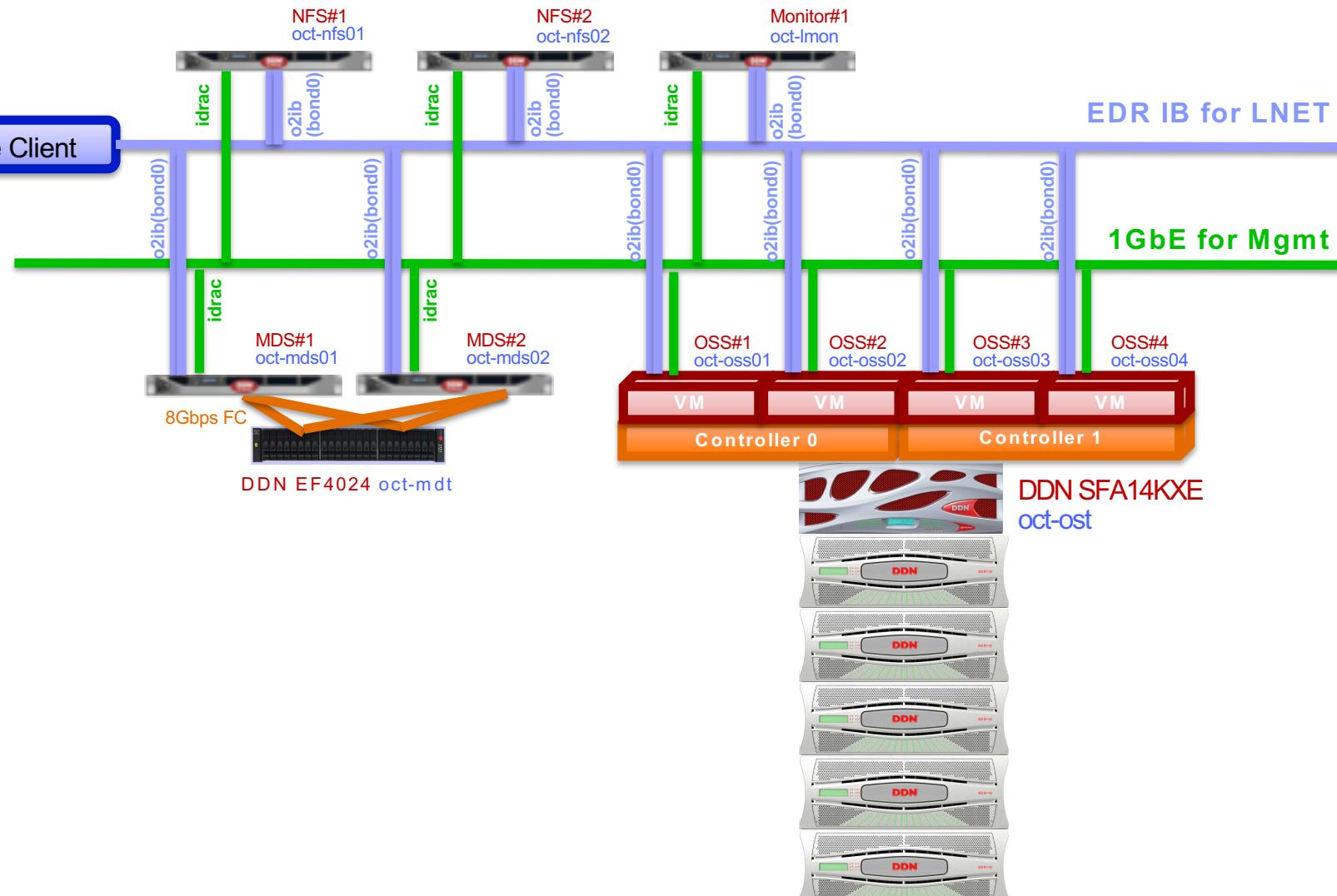
DDN SFA14KXE oct-ost
SS8462 x 5
10TB 4kn 7.2Krpm
NL-SAS x 410
(including Spare x 10)



OCTOPUS Lustre FS概要

- MDSは DDN 1U Server 2台でActive/Standby構成
- MDTはDDN EF4024(600GB 15Krpm SAS x 24) 1基で構成
 - ▶ 10本のDiskを用いたRAID1+0(5D+5D) x 2で構成し、OS上でLVMを用いて1つのMDTとして構成
- OSSはDDN SFA14KXE上のVM 4台で2台ずつペアのActive/Active構成
- OSTはDDN SFA14KXE(4TB 4Kn 7.2Krpm NL-SAS x 410) 1基で構成
 - ▶ 10本のDiskを用いたRAID6(8D+2P)で構成し、1つのRAID6から1つのVolume(LUN)を切り出し構成
- ファイルシステム数は1つ (octfs)
- ファイルシステムにはDNEを使用しない
- Lustre ClientはInfiniband EDRを持つLinux Client

OCTPUS LNET構成

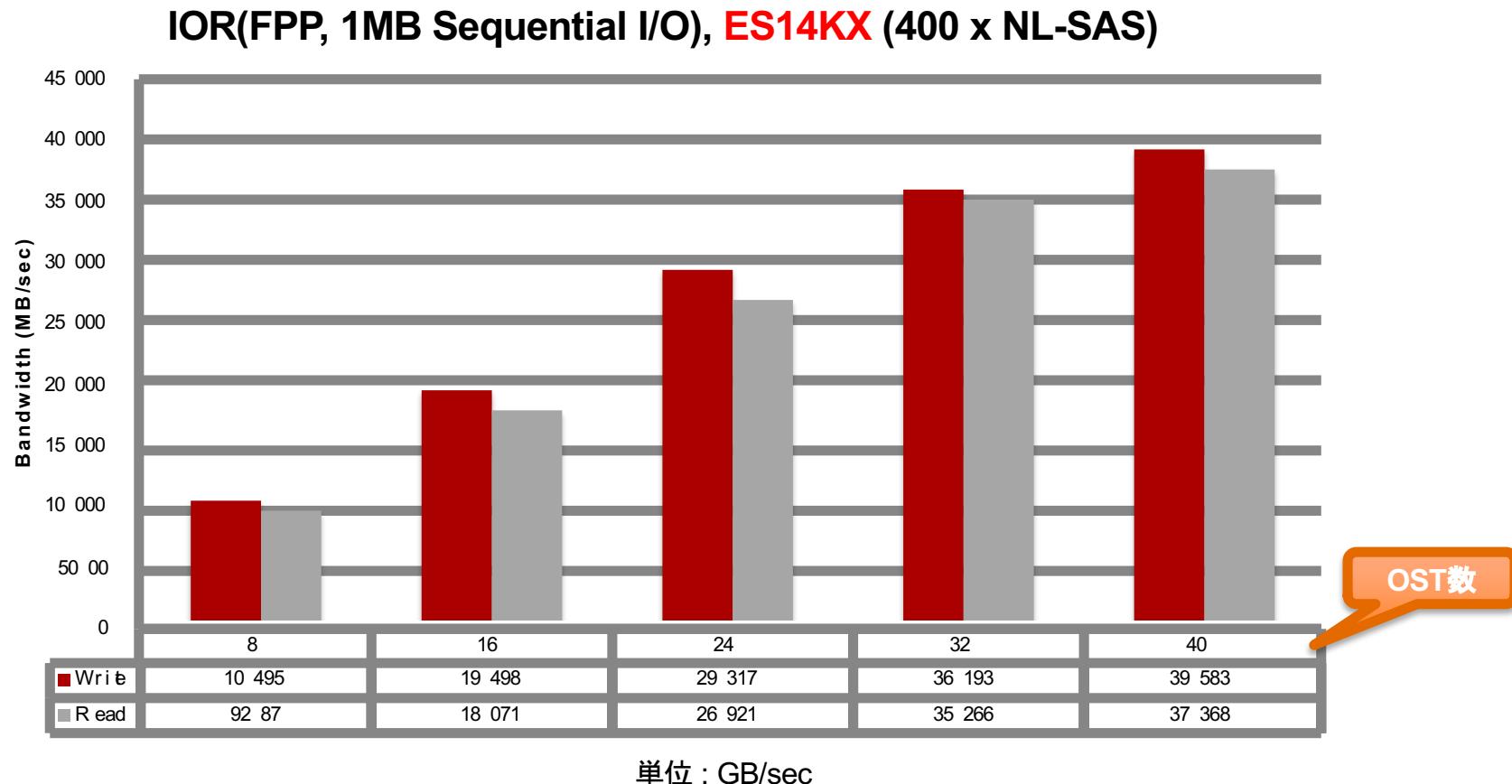


OCTOPUS LNET概要

- 1つのLNETを構成し、IB EDR Network(@o2ib)を使用
- Lustre FSを構成するMDS/OSSは2本のIB EDRを用いて LNET(@o2ib)を構成
- MDS/OSSの2本のIB EDR Cable/Portは通信時、Active/Activeで使用
- IB EDRを持つLinux Clientはo2ibを通して直接MDS/OSSへアクセス

SFA14KXE(ES14KX)の性能

- OCTOPUSストレージと同構成のピーク実効スループット
 - ▶ ベンチマーク参考値



Agenda

- DDN会社/製品紹介
- Lustreについて
- Lustreの特徴
- Lustreのアーキテクチャ
- Lustreの各コンポーネント
- OCTOPUSストレージ概要
- OCTOPUSファイルシステムの利用
- Q&A

OCTOPUSファイルシステム

- マウントポイント : /octfs
- home領域 : /octfs/home/(利用者番号)
 - ▶ 1アカウントに付き10GBの領域を割当(拡張不可)
 - ▶ Lustre Project Quotaを利用
- work領域 : /octfs/work/(グループ名)/(利用者番号)
 - ▶ 1グループ(1申請)に付き1TBの領域を割当(拡張可能)
 - ▶ 使用容量はhome領域との合算
 - ▶ Lustre Group Quotaを利用

OCTOPUSファイルシステム

```
-bash-4.2$ df -H
```

ファイルシステム

/dev/mapper/rhel-root

devtmpfs

tmpfs

tmpfs

tmpfs

/dev/sda1

/dev/sda2

10.10.20.11@o2ib:10.10.20.12@o2ib:/octfs 3.2P 53T 3.1P 2% /octfs

	サイズ	使用	残り	使用%	マウント位置
/dev/mapper/rhel-root	237G	17G	220G	8%	/
devtmpfs	101G	0	101G	0%	/dev
tmpfs	101G	29M	101G	1%	/dev/shm
tmpfs	101G	3.7G	97G	4%	/run
tmpfs	101G	0	101G	0%	/sys/fs/cgroup
/dev/sda1	1.1G	206M	747M	22%	/boot
/dev/sda2	210M	0	210M	0%	/boot/efi

Lustreコマンド

- **lctl**
 - ▶ 管理者向けコマンド
 - ▶ Lustre FS Configuration Utility
 - ▶ Lustre Parameterの設定・変更などを行うことが可能
- **lfs**
 - ▶ 一部一般ユーザ権限で利用可能なコマンド
 - ▶ Lustre Utility
 - ▶ Stripe管理
 - ▶ Quota管理
 - ▶ **find** : 特定のOST上のファイル、stripeされたファイルなどを指定可能

Lustreの状態確認

```
lfs > mdts
MDTS:
0: octfs-MDT0000_UUID ACTIVE
lfs > osts
OBDS::
0: octfs-OST0000_UUID ACTIVE
1: octfs-OST0001_UUID ACTIVE
2: octfs-OST0002_UUID ACTIVE
3: octfs-OST0003_UUID ACTIVE
4: octfs-OST0004_UUID ACTIVE
5: octfs-OST0005_UUID ACTIVE
6: octfs-OST0006_UUID ACTIVE
7: octfs-OST0007_UUID ACTIVE
8: octfs-OST0008_UUID ACTIVE
9: octfs-OST0009_UUID ACTIVE
10: octfs-OST000a_UUID ACTIVE
11: octfs-OST000b_UUID ACTIVE
12: octfs-OST000c_UUID ACTIVE
13: octfs-OST000d_UUID ACTIVE
14: octfs-OST000e_UUID ACTIVE
15: octfs-OST000f_UUID ACTIVE
16: octfs-OST0010_UUID ACTIVE
17: octfs-OST0011_UUID ACTIVE
18: octfs-OST0012_UUID ACTIVE
19: octfs-OST0013_UUID ACTIVE
20: octfs-OST0014_UUID ACTIVE
21: octfs-OST0015_UUID ACTIVE
22: octfs-OST0016_UUID ACTIVE
23: octfs-OST0017_UUID ACTIVE
24: octfs-OST0018_UUID ACTIVE
25: octfs-OST0019_UUID ACTIVE
26: octfs-OST001a_UUID ACTIVE
27: octfs-OST001b_UUID ACTIVE
28: octfs-OST001c_UUID ACTIVE
29: octfs-OST001d_UUID ACTIVE
30: octfs-OST001e_UUID ACTIVE
31: octfs-OST001f_UUID ACTIVE
32: octfs-OST0020_UUID ACTIVE
33: octfs-OST0021_UUID ACTIVE
34: octfs-OST0022_UUID ACTIVE
35: octfs-OST0023_UUID ACTIVE
36: octfs-OST0024_UUID ACTIVE
37: octfs-OST0025_UUID ACTIVE
38: octfs-OST0026_UUID ACTIVE
39: octfs-OST0027_UUID ACTIVE
```

MDTの状態

OSTの状態

Lustreの使用状況確認 – 容量(block)

```
lfs > df -h
UUID      bytes  Used  Available Use% Mounted on
octfs-MDT0000_UUID    4.0T   19.1G   4.0T  0% /octfs[MDT:0]
octfs-OST0000_UUID    71.3T   1.3T   69.3T  2% /octfs[OST:0]
octfs-OST0001_UUID    71.3T   2.0T   68.7T  3% /octfs[OST:1]
octfs-OST0002_UUID    71.3T   1.0T   69.6T  1% /octfs[OST:2]
octfs-OST0003_UUID    71.3T   1.3T   69.4T  2% /octfs[OST:3]
octfs-OST0004_UUID    71.3T  1023.9G   69.6T  1% /octfs[OST:4]
octfs-OST0005_UUID    71.3T   1.1T   69.5T  2% /octfs[OST:5]
octfs-OST0006_UUID    71.3T  934.4G   69.7T  1% /octfs[OST:6]
octfs-OST0007_UUID    71.3T  968.7G   69.7T  1% /octfs[OST:7]
octfs-OST0008_UUID    71.3T   1.6T   69.0T  2% /octfs[OST:8]
octfs-OST0009_UUID    71.3T   1.6T   69.1T  2% /octfs[OST:9]
octfs-OST000a_UUID    71.3T   1.8T   68.8T  3% /octfs[OST:10]
octfs-OST000b_UUID    71.3T   1.2T   69.4T  2% /octfs[OST:11]
octfs-OST000c_UUID    71.3T   1.4T   69.2T  2% /octfs[OST:12]
octfs-OST000d_UUID    71.3T  1011.6G   69.6T  1% /octfs[OST:13]
octfs-OST000e_UUID    71.3T  863.2G   69.8T  1% /octfs[OST:14]
octfs-OST000f_UUID    71.3T  1019.1G   69.6T  1% /octfs[OST:15]
octfs-OST0010_UUID    71.3T  884.8G   69.8T  1% /octfs[OST:16]
octfs-OST0011_UUID    71.3T  985.2G   69.7T  1% /octfs[OST:17]
octfs-OST0012_UUID    71.3T   1.4T   69.2T  2% /octfs[OST:18]
octfs-OST0013_UUID    71.3T   1.3T   69.3T  2% /octfs[OST:19]
octfs-OST0014_UUID    71.3T   2.2T   68.5T  3% /octfs[OST:20]
octfs-OST0015_UUID    71.3T   1.1T   69.5T  2% /octfs[OST:21]
octfs-OST0016_UUID    71.3T   1.3T   69.3T  2% /octfs[OST:22]
octfs-OST0017_UUID    71.3T  1020.3G   69.6T  1% /octfs[OST:23]
octfs-OST0018_UUID    71.3T   1.4T   69.2T  2% /octfs[OST:24]
octfs-OST0019_UUID    71.3T  949.3G   69.7T  1% /octfs[OST:25]
octfs-OST001a_UUID    71.3T  954.2G   69.7T  1% /octfs[OST:26]
octfs-OST001b_UUID    71.3T   1.1T   69.5T  2% /octfs[OST:27]
octfs-OST001c_UUID    71.3T  992.3G   69.6T  1% /octfs[OST:28]
octfs-OST001d_UUID    71.3T   1.4T   69.2T  2% /octfs[OST:29]
octfs-OST001e_UUID    71.3T  969.4G   69.7T  1% /octfs[OST:30]
octfs-OST001f_UUID    71.3T  926.3G   69.7T  1% /octfs[OST:31]
octfs-OST0020_UUID    71.3T   1.1T   69.5T  2% /octfs[OST:32]
octfs-OST0021_UUID    71.3T   1.5T   69.1T  2% /octfs[OST:33]
octfs-OST0022_UUID    71.3T   1.1T   69.5T  2% /octfs[OST:34]
octfs-OST0023_UUID    71.3T   1.1T   69.5T  2% /octfs[OST:35]
octfs-OST0024_UUID    71.3T  979.9G   69.7T  1% /octfs[OST:36]
octfs-OST0025_UUID    71.3T   1.0T   69.6T  1% /octfs[OST:37]
octfs-OST0026_UUID    71.3T   1.4T   69.2T  2% /octfs[OST:38]
octfs-OST0027_UUID    71.3T   1.2T   69.4T  2% /octfs[OST:39]

filesystem summary: 2.8P 48.0T 2.7P 2% /octfs
```

Lustreの使用状況確認 – inode

lfs > df -i -h

UUID	Inodes	IUsed	IFree	IUse%	Mounted on
octfs-MDT0000_UUID	2.7G	47.2M	2.7G	2%	/octfs[MDT:0]
octfs-OST0000_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:0]
octfs-OST0001_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:1]
octfs-OST0002_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:2]
octfs-OST0003_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:3]
octfs-OST0004_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:4]
octfs-OST0005_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:5]
octfs-OST0006_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:6]
octfs-OST0007_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:7]
octfs-OST0008_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:8]
octfs-OST0009_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:9]
octfs-OST000a_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:10]
octfs-OST000b_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:11]
octfs-OST000c_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:12]
octfs-OST000d_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:13]
octfs-OST000e_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:14]
octfs-OST000f_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:15]
octfs-OST0010_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:16]
octfs-OST0011_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:17]
octfs-OST0012_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:18]
octfs-OST0013_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:19]
octfs-OST0014_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:20]
octfs-OST0015_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:21]
octfs-OST0016_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:22]
octfs-OST0017_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:23]
octfs-OST0018_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:24]
octfs-OST0019_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:25]
octfs-OST001a_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:26]
octfs-OST001b_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:27]
octfs-OST001c_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:28]
octfs-OST001d_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:29]
octfs-OST001e_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:30]
octfs-OST001f_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:31]
octfs-OST0020_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:32]
octfs-OST0021_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:33]
octfs-OST0022_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:34]
octfs-OST0023_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:35]
octfs-OST0024_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:36]
octfs-OST0025_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:37]
octfs-OST0026_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:38]
octfs-OST0027_UUID	1.1G	1.2M	1.1G	0%	/octfs[OST:39]

filesystem summary: 2.7G 47.2M 2.7G 2% /octfs

Lustre Quota

- Lustre 2.4.1以降、block/inodeに対してそれぞれUser/Group Quotaを設定可能
- Lustre 2.7.21以降、Project Quotaをサポート
 - ▶ 同一Project IDを持つディレクトリ/ファイルに対するQuota
 - ▶ 複数のディレクトリを纏めてQuota設定可能
- OCTOPUSでは全てのQuota機能が有効な状態
 - ▶ Quotaの設定、変更、解除は管理者のみ可能

Lustre Quota設定(1)

Quota	対象	単位	説明
User block soft limit	User	Kbytes	対象ファイルシステムへ、あるユーザが書き込めるサイズの制限値 この値を超えていている場合でも、grace timeの間は書き込み可能
User block soft limit grace time	User	時間	User block soft limitを超えていられる時間
User inode soft limit	User	Inode数	対象ファイルシステムへ、あるユーザが書き込めるファイル数の制限値 この値を超えていている場合でも、grace timeの間は書き込み可能
User inode soft limit grace time	User	時間	User inode soft limitを超えていられる時間
User block hard limit	User	Kbytes	対象ファイルシステムへ、あるユーザが書き込めるサイズの制限値 この値を超えて書き込みを行うことは不可能
User inode hard limit	User	Inode数	対象ファイルシステムへ、あるユーザが書き込めるファイル数の制限値 この値を超えて書き込みを行うことは不可能
Group block soft limit	Group	Kbytes	対象ファイルシステムへ、あるグループが書き込めるサイズの制限値 この値を超えていている場合でも、grace timeの間は書き込み可能
Group block soft limit grace time	Group	時間	Group block soft limitを超えていられる時間
Group inode soft limit	Group	Inode数	対象ファイルシステムへ、あるグループが書き込めるファイル数の制限値 この値を超えていている場合でも、grace timeの間は書き込み可能
Group inode soft limit grace time	Group	時間	Group inode soft limitを超えていられる時間
Group block hard limit	Group	Kbytes	対象ファイルシステムへ、あるグループが書き込めるサイズの制限値 この値を超えて書き込みを行うことは不可能
Group inode hard limit	Group	Inode数	対象ファイルシステムへ、あるグループが書き込めるファイル数の制限値 この値を超えて書き込みを行うことは不可能

Lustre Quota設定(2)

Quota	対象	単位	説明
Project block soft limit	Project	Kbytes	対象ファイルシステムへ、あるプロジェクトが書き込めるサイズの制限値 この値を超えている場合でも、grace timeの間は書き込み可能
Project block soft limit grace time	Project	時間	Project block soft limitを超えていられる時間
Project inode soft limit	Project	Inode数	対象ファイルシステムへ、あるプロジェクトが書き込めるファイル数の制限値 この値を超えている場合でも、grace timeの間は書き込み可能
Project inode soft limit grace time	Project	時間	Project inode soft limitを超えていられる時間
Project block hard limit	Project	Kbytes	対象ファイルシステムへ、あるプロジェクトが書き込めるサイズの制限値 この値を超えて書き込みを行うことは不可能
Project inode hard limit	Project	Inode数	対象ファイルシステムへ、あるプロジェクトが書き込めるファイル数の制限値 この値を超えて過去込みを行うことは不可能

Quotaの確認方法

- OCTOPUSではusage_viewコマンドによってwork領域利用状況の確認ができますが、ここではlfsコマンドによる確認方法を示します
- lfs quota [-u <user_name | uid> | -g <group_name | gid> | -p project_id] mount_point
- lfs quotaその他のオプション
 - ▶ -t [-u | -g | -p] mount_point
 - User, Group, Projectそれぞれに対するblock およびinodeのgrace timeを表示
 - ▶ -v
 - MDT, OSTごとの使用状況を確認可能
 - ▶ -h
 - human readable format

Quotaの確認 (User Quota)

- OCTOPUSではUser Quotaは有効ですが、使用されていません

```
lfs > quota -u v60333 /octfs
Disk quotas for usr v60333 (uid 1938):
  Filesystem  kbytes  quota  limit  grace  files  quota  limit  grace
    /octfs      10240024      0      0       -       6      0      0       -
```

見出し	意味
Filesystem	ファイルシステム名(mount_point)
kbytes(-hを使用した場合はused)	現在の使用量、超過している場合数字の後ろに * が付与される
quota	block soft quota
limit	block hard quota
grace	block grace time
files	現在のファイル数、超過している場合数字の後ろに * が付与される
quota	inode soft quota
limit	inode hard quota
grace	inode grace time

Quotaの確認 (Group Quota)

- OCTOPUSではGroup Quotaによって全体使用量の制限を行っています

```
lfs > quota -h -g G14547 /octfs
Disk quotas for grp G14547 (gid 14547):
  Filesystem  used  quota  limit  grace  files  quota  limit  grace
    /octfs     20k    0k    1T      -      5      0      0      -
```

- Soft Limitが設定されていない為、Hard Limitに達した時点で書き込めなくなります
- inode数(ファイル数)に対する制限はありません

Quotaの確認 (Project Quota)

- OCTOPUSではProject Quotaによってhome領域の使用量制限を行っています
 - ▶ Project IDはフロントエンドノードでは確認できません

```
lfs > quota -h -p 1938 /octfs
Disk quotas for prj 1938 (pid 1938):
  Filesystem  used  quota  limit  grace  files  quota  limit  grace
    /octfs     12k    0k   10G      -       3      0      0      -
```

- Soft Limitが設定されていない為、Hard Limitに達した時点で書き込めなくなります
- inode数(ファイル数)に対する制限はありません

Quotaの確認 (MDT/OSTごと)

- -vオプションによってMDT/OSTごとの使用状況の確認ができます
 - ▶ 以下はGroup Quotaの状態

```
lfs > quota -v -h -g G14547 /octfs
Disk quotas for grp G14547 (gid 14547):
  Filesystem    used   quota   limit   grace   files   quota   limit   grace
    /octfs     9.766G    0k      1T      -       6      0      0      -
octfs-MDT0000_UUID
          12k      -    0k      -       6      -      0      -
octfs-OST0000_UUID
          0k      -  13.76G      -      -      -      -      -
octfs-OST0001_UUID
          0k      -   16G      -      -      -      -      -
<省略>
octfs-OST000b_UUID
          9.766G      -   16G      -      -      -      -      -
<省略>
octfs-OST0027_UUID
          0k      -  16.25G      -      -      -      -      -
Total allocated inode limit: 0, total allocated block limit: 592.4G
```

File Stripingの利用

- Lustreでは、File Stripingをサポートしています
 - ▶ 1つのファイルを複数OSTに書き込む機能
- 各ファイルを格納するStripe Count(使用するOSTの数)およびStripe Sizeをディレクトリやファイルごとに設定可能
 - ▶ OCTOPUSのデフォルトはStripe設定無し
 - ▶ 各ファイルは1つのOSTへ格納
 - ▶ ディレクトリにStripeを設定した場合、以降そのディレクトリ以下に作られるファイルは全てStripeされる

Stripeの設定方法

```
lfs setstripe [--stripe-size | -S <stripe_size>] [--stripe-count | -c <stripe_count>] [--stripe-index | -I <start_ost_idx>] [--pool | -p <pool_name>] [--block | -b] <directory | filename> [---ost-list|-o <ost_indices>]
```

Option項目	意味
stripe_size	各OSTをStripeする際のStripe Sizeを指定 k, m, gでそれぞれKbytes, Mbytes, GBytes指定が可能 デフォルトは0でファイルシステムデフォルト(OCTOPUSでは1MiB)を意味する
start_ost_idx	Stripeを開始するOSTを指定 デフォルトは-1でランダム
stripe_count	使用するOSTの数を指定 -1で全OST デフォルトは0でファイルシステムデフォルト(OCTOPUSでは1)を意味する
pool_name	使用するプール名を指定 ※ OCTOPUSでは OST Poolを使用していません
block	data migration時にファイルへのアクセスをblockする
ost_indices	使用するOSTのインデックス番号を指定

Stripeの確認方法

```
lfs getstripe [--ost | -O <uuid>] [--quiet | -q] [--verbose | -v] [--recursive | -r] <directory | filename>
```

Option項目	意味
--ost -O <uuid>	uuidで指定したOSTに属するファイルを表示
--quiet -q	出力項目の制限
--verbose -v	Verbose Mode
--recursive -r	Recursive Mode

Stripe使用例(1)

- 作成したディレクトリに対し、2x OST, Stripe Size=1MBytesでStripeを設定

```
-bash-4.2$ mkdir /octlfs/work/G14547/v60333/s2
-bash-4.2$ lfs setstripe -c 2 -S 1m /octlfs/work/G14547/v60333/s2
-bash-4.2$ lfs getstripe /octlfs/work/G14547/v60333/s2
/octlfs/work/G14547/v60333/s2
stripe_count: 2 stripe_size: 1048576 stripe_offset: -1
```

- 作成したディレクトリに対し、全OST, Stripe Size=4MBytesでStripeを設定

```
-bash-4.2$ mkdir /octlfs/work/G14547/v60333/st_all
-bash-4.2$ lfs setstripe -c -1 -S 4m /octlfs/work/G14547/v60333/st_all
-bash-4.2$ lfs getstripe /octlfs/work/G14547/v60333/st_all
/octlfs/work/G14547/v60333/st_all
stripe_count: -1 stripe_size: 4194304 stripe_offset: -1
```

Stripe使用例(2)

- 作成したs2ディレクトリ以下にファイルを作成

```
-bash-4.2$ pwd  
/octfs/work/G14547/v60333/s2  
-bash-4.2$ dd if=/dev/zero of=testfile1 bs=1M count=1000  
1000+0 レコード入力  
1000+0 レコード出力  
1048576000 バイト (1.0 GB) コピーされました、0.693553 秒、1.5 GB/秒  
-bash-4.2$ lfs getstripe testfile1  
testfile1  
lmm_stripe_count: 2  
lmm_stripe_size: 1048576  
lmm_pattern: 1  
lmm_layout_gen: 0  
lmm_stripe_offset: 27  
obdidx objid objid group  
27 13564237 0xcef94d 0  
37 13564336 0xcef9b0 0
```

- testfile1はobdidxで示される2つのOSTにStripeされており、obdidxはlfs ostsで確認可能です

File Stripe使用上の注意

- 必ずしも、大きいStripe Count == 高性能ではありません
 - ファイルシステム全体のスループットはStripe無しが最も高い
- Stripe Sizeはデフォルトの1MiBが最適です
 - ストライプする場合
- 充分大きなファイルはStripeした方が良い
 - 1ファイルが1OSTに書ける最大サイズは16TiB
 - 各OSTの使用量がアンバランスになると性能が劣化する可能性あり
- シングルスレッドアプリケーションではStripeの性能に対する効果はほぼ無い
- OSTは手動で指定せずランダムが最適

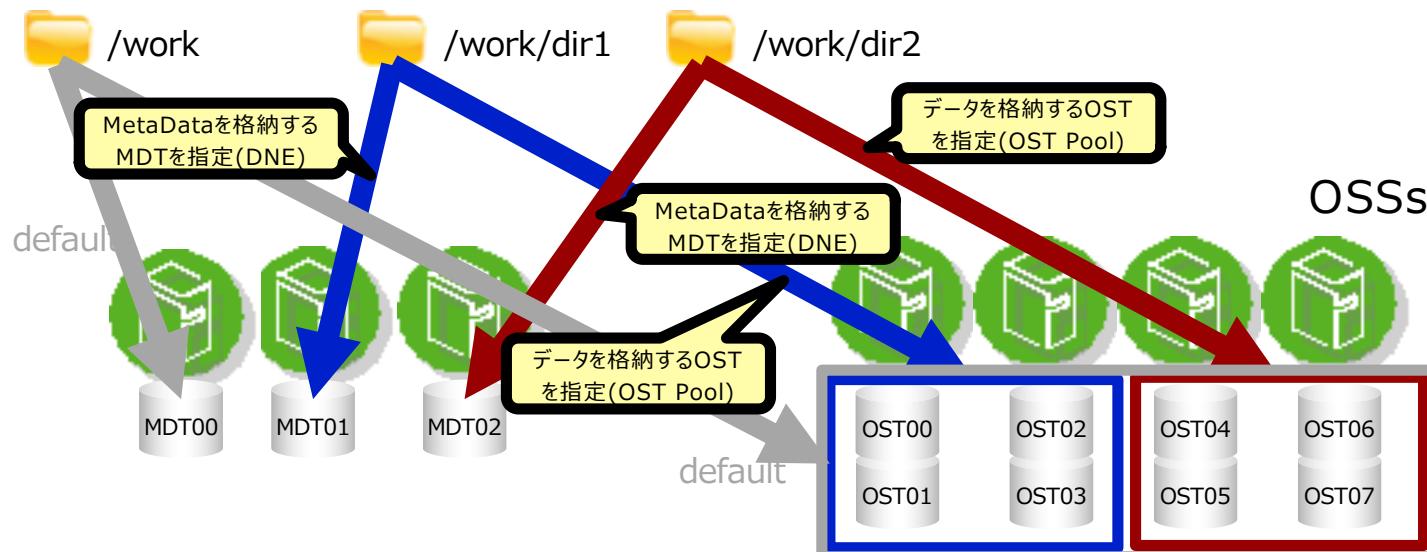
Lustre Tips

- 苦手なこと
 - ▶ 大量の小さいファイル
 - ▶ 小さいIO Size
 - ▶ ランダムアクセス
- メタデータ性能
 - ▶ 複数プロセスによる同一ディレクトリへの操作は遅くなる
 - ▶ プロセス毎に異なるディレクトリを操作する方が良い
- スループット性能
 - ▶ プロセス毎に異なるファイルを操作
 - ▶ できる限り各プロセスにIOを行わせる
 - ▶ 1 Clientのスループット性能は最高でも10GB/s
 - 単一プロセスでは2-3GB/sでれば良い方、Xeon Phiならもっと出ない
 - 1クライアント/複数IOプロセス + 複数ノード = スループット向上

DNEとOST Poolのご紹介(1)

- OCTOPUSでは利用されていませんが、Lustreが提供しているDNEとOST Poolを紹介します
 - 次回調達のご参考までに、、、
- DNE(Distributed Namespace Environment)
 - Directoryごとに使用するMDS/MDTを指定可能
 - メタデータオペレーションのスケーリング
 - 他ジョブによる性能への影響を回避
- OST Pool
 - Directoryごとに使用するOSS/OSTを指定可能
 - Disk性能が異なる際に有用
 - 他ジョブによる性能への影響を回避

One Namespaceで
複数FSのような利用が可能



DNEとOST Poolのご紹介(2)

- DNEのメリット・デメリット

メリット	デメリット
複数台のMDS/MDTを使用することによる全体でのメタデータオペレーション性能の向上	2つのRAIDを1つにまとめて使用する場合に比べて、DNEとして構成するMDTを意図的に使用する必要がある
複数台のMDS/MDTを使用することでジョブ間での影響を抑える	DNE間を跨る場合、Directory間におけるfile/directoryのmv時にはdataのcopy処理も実施される MDT用RAIDの数に依存するが、後でやり直すことがほぼ不可能。但し、動的な追加は可能

⇒ 構築前に構成FIXが必要

- Stripeのメリット・デメリット

メリット	デメリット
1 FileへのIO性能向上	FS全体でのスループットではStripe無しが高い
1 Fileあたりの最大サイズ向上	大きなファイルを作成した場合、特定OSTの使用量が大幅に増加し、そのOSTが使われにくくなる

⇒ 使いながら調整可能であり、アプリやデータごとに変更することがBEST

DNEとOST Poolのご紹介(3)

- 例えばOCTOPUSであれば、
 - ▶ home領域とwork領域を別MDTにする
 - ▶ work領域を複数MDTで構成し、専用MDSを持つ領域を作る
 - ▶ ただし、MDTの追加が必要です
- OST Pool
 - ▶ Sub Directory Mountと組み合わせた専有OST領域
 - セキュリティを考慮したコンテナからの利用など
 - ▶ ただし、スループットは限定される
 - ▶ OCTOPUSでの利用は、現時点では必要ないと考えます

Agenda

- DDN会社/製品紹介
- Lustreについて
- Lustreの特徴
- Lustreのアーキテクチャ
- Lustreの各コンポーネント
- OCTOPUSストレージ概要
- OCTOPUSファイルシステムの利用
- Q&A

Thank You!

Keep in touch with us



Team-JPSales@ddn.com



@ddn_limitless



company/datadirect-networks



Tokyu Bancho Bldg. 8F
6-2 Yonbancho Chiyoda-ku,
Tokyo 102-0081



+81-3-3261-9101
+81-3-3261-9140