

【SS研オブジェクトストレージ座談会】
大量データ利活用に向けたデータ管理
～テープ活用によるコスト最適化～

2021年3月23日

(株)富士通研究所 プラットフォーム革新PJ
田村 雅寿

タムラ

マサヒサ

田村 雅寿



- 1999年 富士通研究所入社
- 専門：分散ストレージシステム
 - 分散ファイルシステム
 - 分散ブロックストレージ
 - **分散オブジェクトストレージ**
 - ブロックチェーン, パーソナルデータストア (PDS)

- オブジェクトストレージとは
- オンプレのオブジェクトストレージ
 - オブジェクトストレージの特徴 (例: Ceph)
- オブジェクトストレージでのデータ活用
- テープ活用によるコスト最適化
 - テープの特徴
 - テープ階層化ストレージ

オブジェクトストレージとは

■ Webスケール時代の新しいストレージ

- Webアプリケーションからダイレクトにアクセス

■ オブジェクト単位のデータ保存

- ID（識別子）で、データとメタデータ（属性情報）を一括管理

■ シンプルインターフェース（REST API）

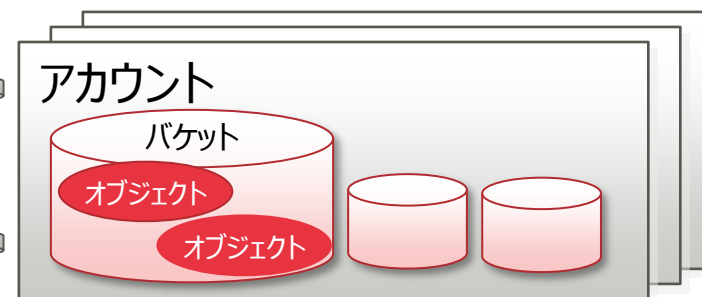
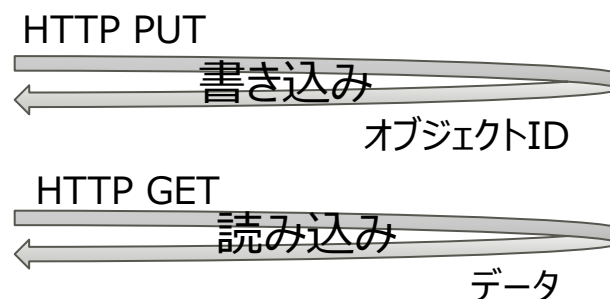
- HTTP(s) 通信でクラウドストレージサービス互換（S3 API）のアクセス
 - GET（読み込み）、PUT（書き込み）、DELETE（削除）、他

■ 広大でフラットなストレージ空間

- 階層構造の排除により上限の無い柔軟な拡張性



ユーザー
(Webアプリケーション)



■ クラウドストレージ

- Amazon S3 (Amazon)
- Azure Blob Storage (Microsoft)
- Google Cloud Storage (Google)

■ オンプレストレージ

- StorageGRID Webscae (NetApp), Elastic Cloud Storage (Dell EMC), ...
- Scality RING (Scality), HyperStore (Cloudian), ...
- OSS: **Ceph**, OpenStack Swift, MinIO, ...

オンプレでもオブジェクトストレージの製品やOSS活用が本格化

特徴 (例: Ceph)

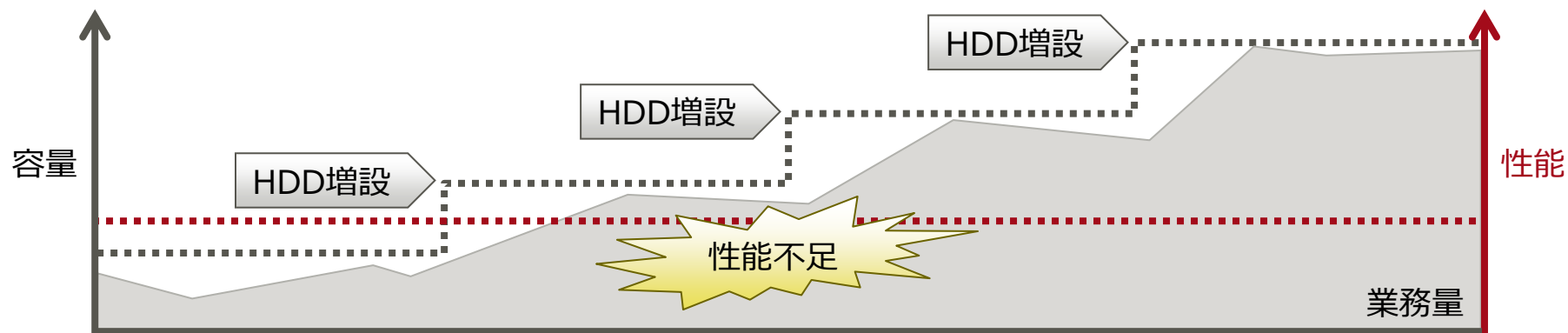
柔軟な拡張性(スケールアウト、自動調整)、新陳代謝、高可用性

規模に応じた容量/性能の拡張

柔軟な拡張

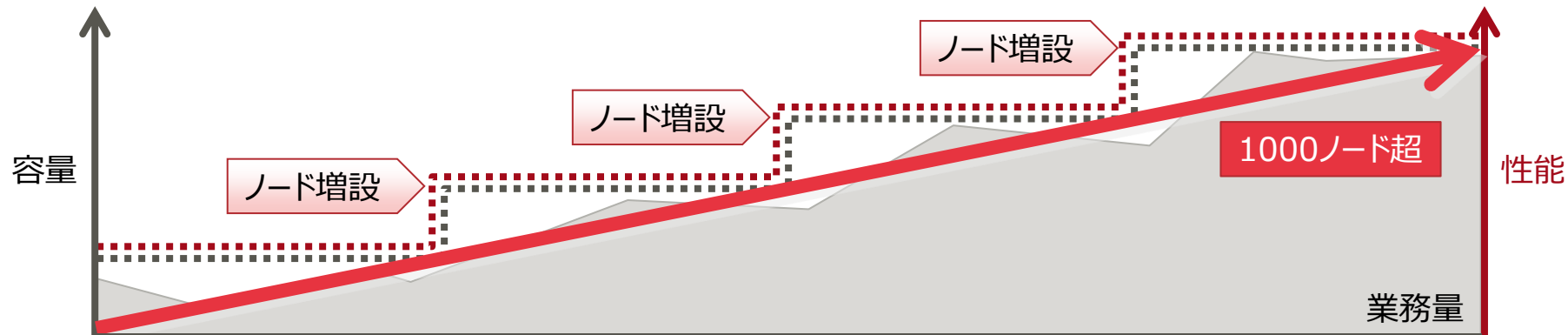
スケールアップ

HDD増設はできても、コントローラーの性能が頭打ちとなる



スケールアウト

HDDとコントローラーを組み合わせたノード増設により、容量/性能をリニアに拡張可能

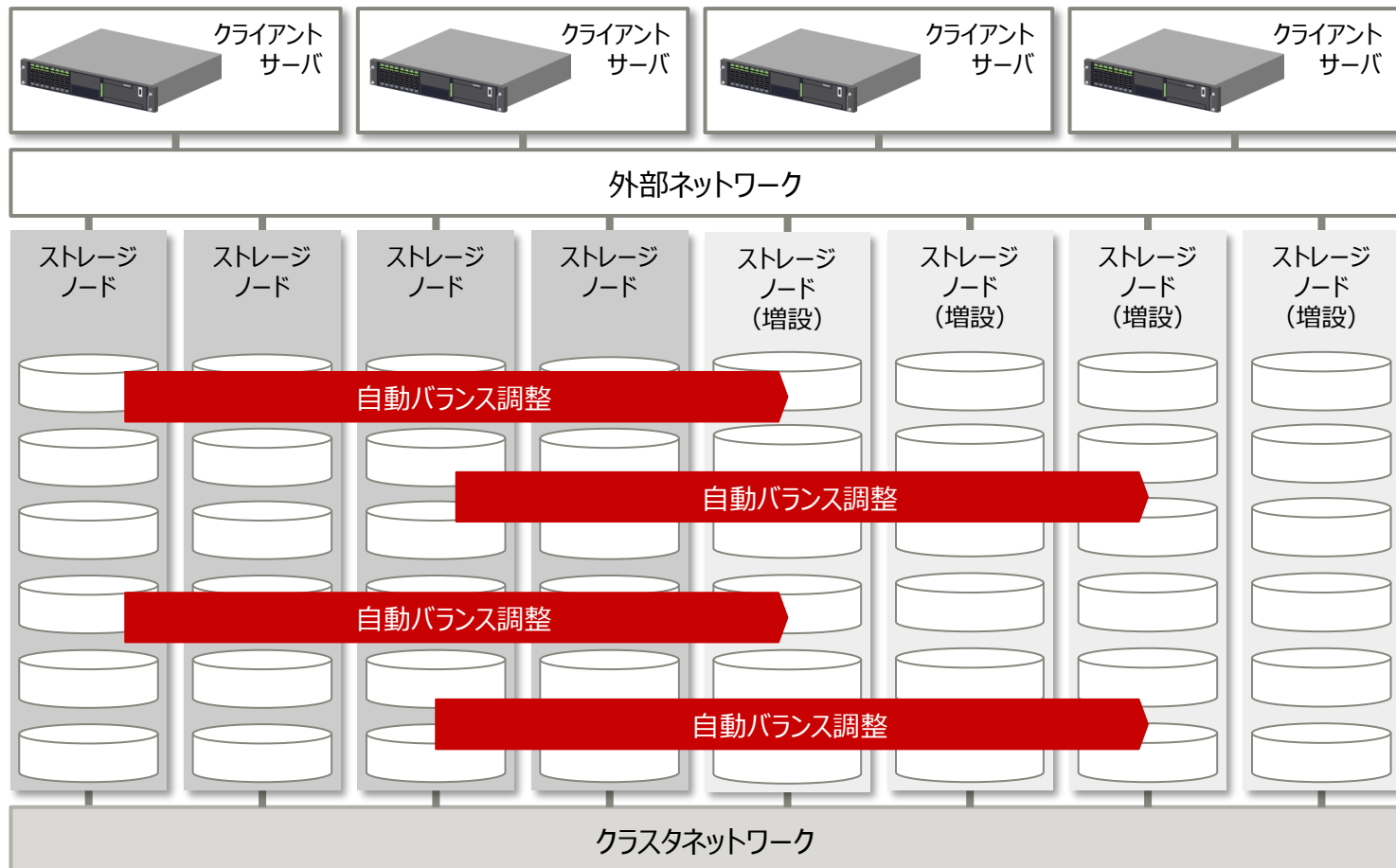


容易なシステム拡張

柔軟な拡張

自動調整

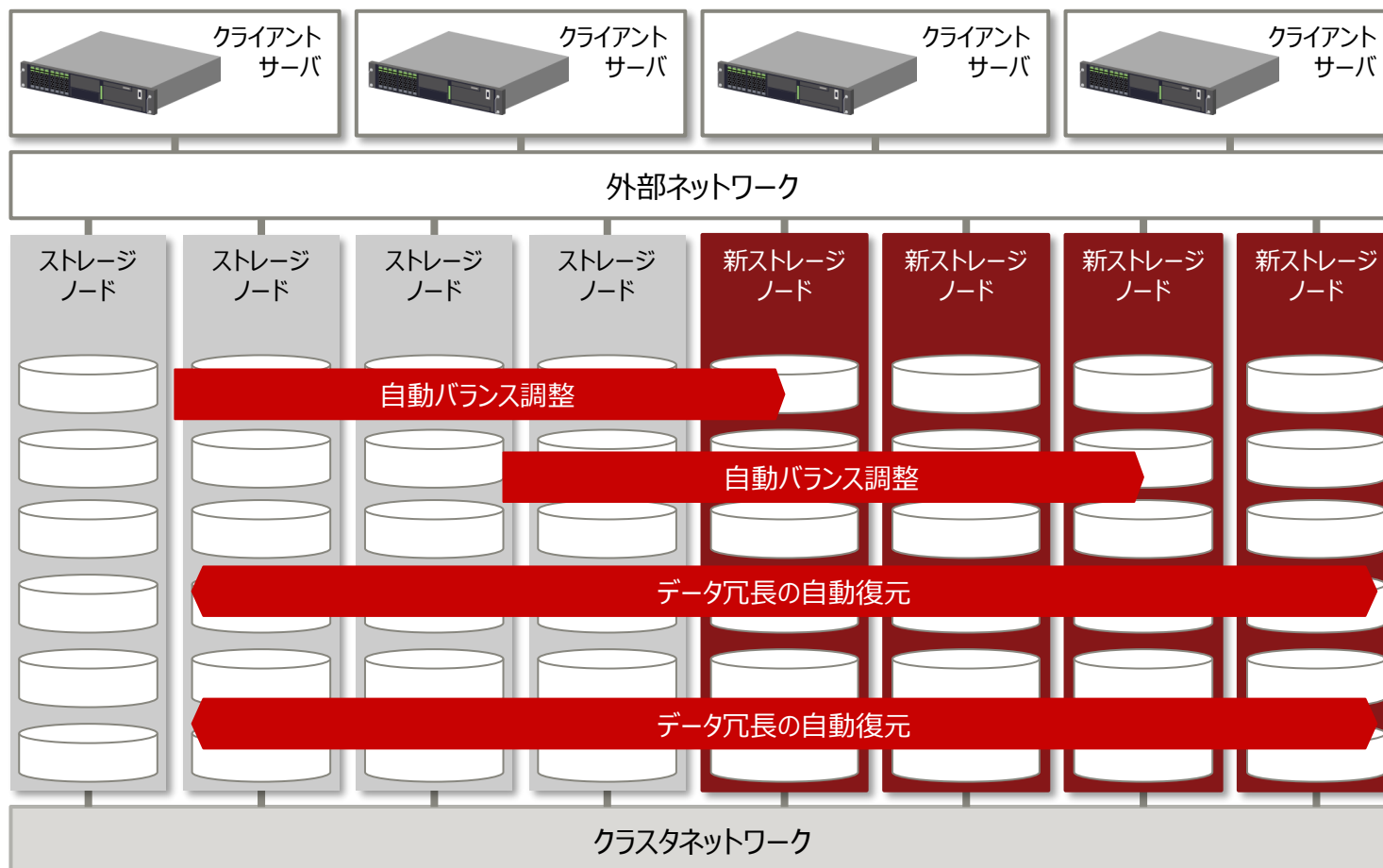
- オンラインでストレージノード単位の増設が可能 (4~1,000ノード超)
- システム全体でデータ配置のバランスを常に自動調整し、ホットスポット発生を防止



システム更改の負担軽減

新陳代謝

- ソフトウェア・デファインド技術が世代の異なるストレージノード混在を実現
- 新ノードの増設・旧ノードの減設により、ハードウェアを段階的に新陳代謝
- ゼロダウンタイムでシステム移行が可能であり、データ長期保存、サービス長期運用に最適



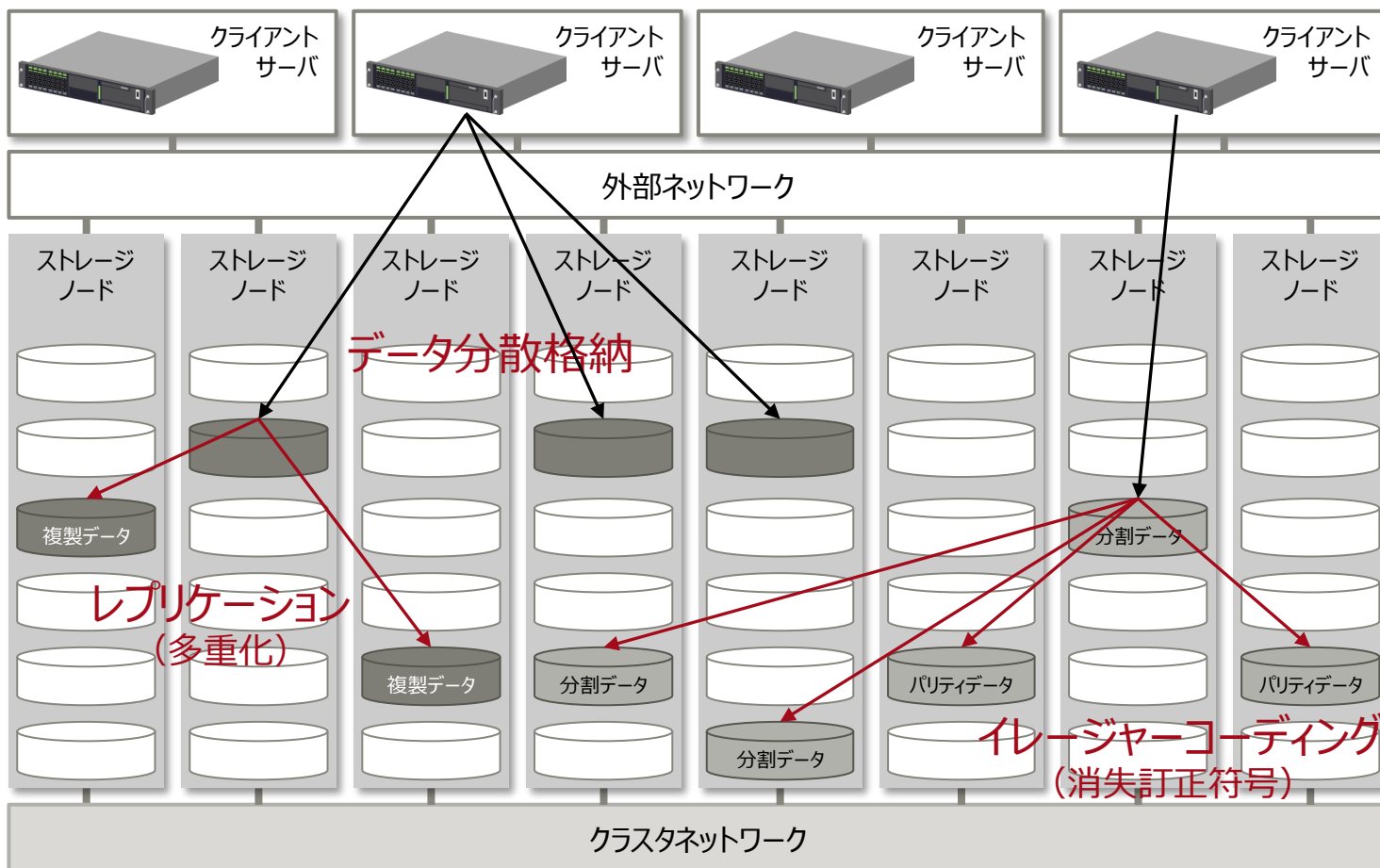
選択可能なデータ保護方式

高可用性

■ データの用途に応じた2種類のデータ保護

- レプリケーション：異なる複数のノードにデータを複製。高い可用性と高速なリカバリーが特長
- イレジャーコーディング：データ分割しパリティ情報を付与。高い容量効率が特長 (オブジェクトアクセス限定)

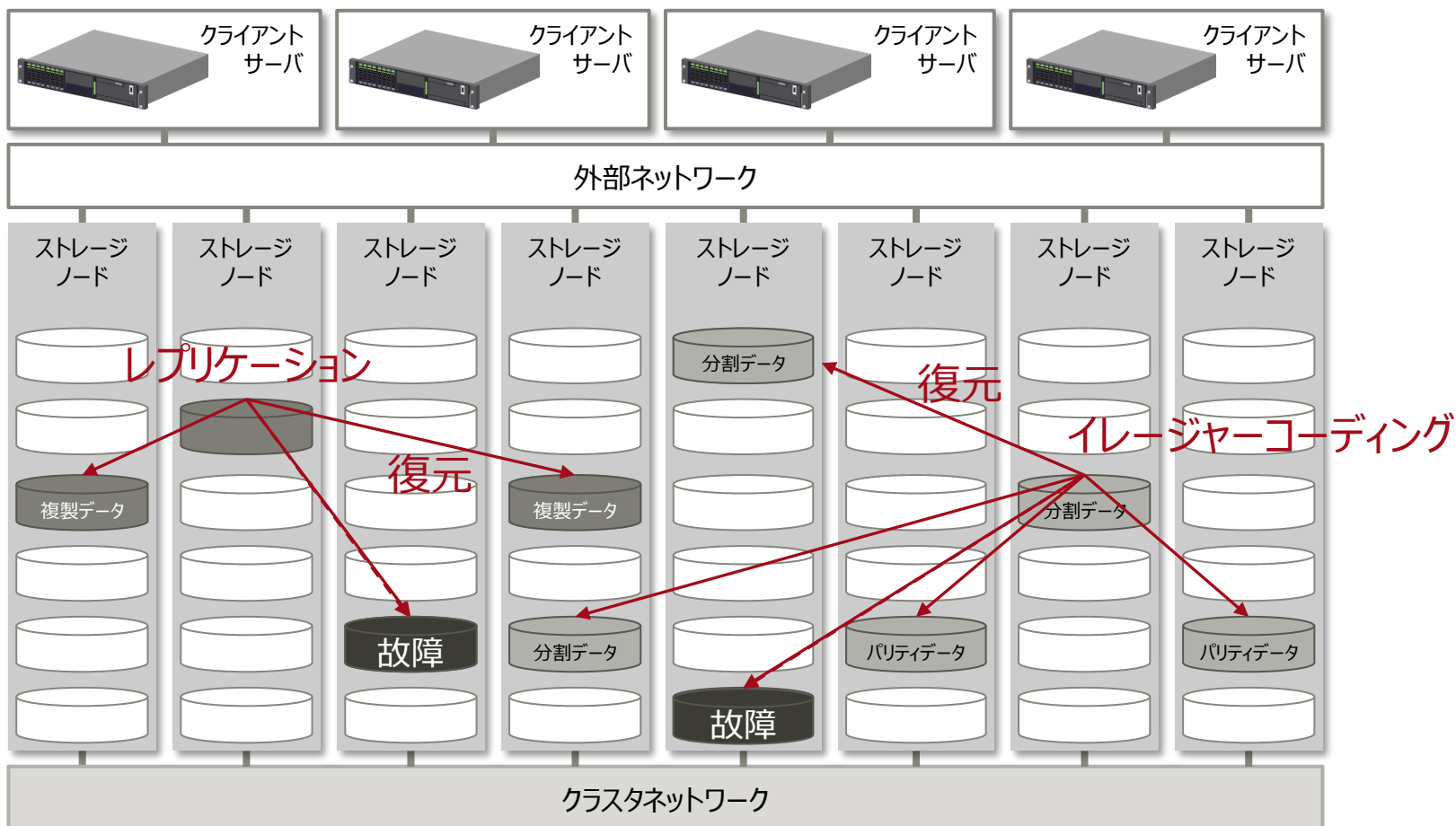
■ 複数ノード障害でも安全にデータを保護できる高い堅牢性



冗長性の自動復元

高可用性

- ノードやディスク故障時にデータ冗長性を自動復元 (セルフヒーリング)
- 故障したディスクやノードを除いたコンポーネントに、データを最適に分散再配置
 - 故障時の特定ノードへの負荷偏りとディスク故障によるデータロストのリスクを回避



蓄積と分析を一体化したデータ基盤

■ 従来型のデータ分析

- 様々なデータソースから発生するデータをオブジェクトストレージに蓄積しつつ、データをHadoopクラスタにコピーして分析

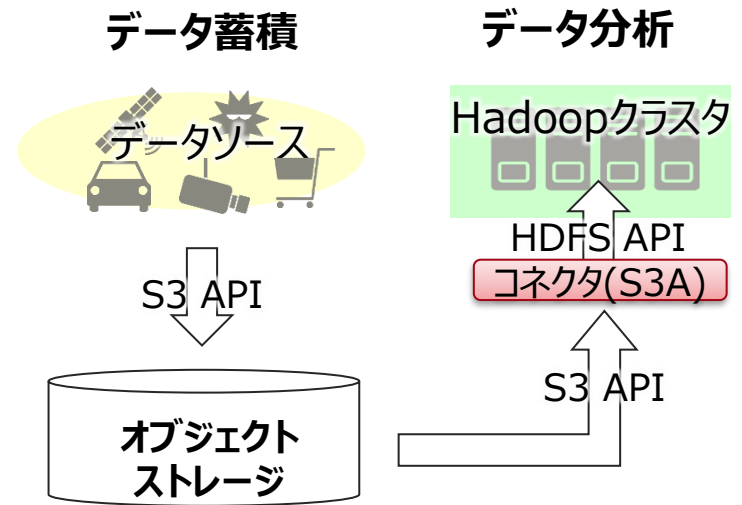
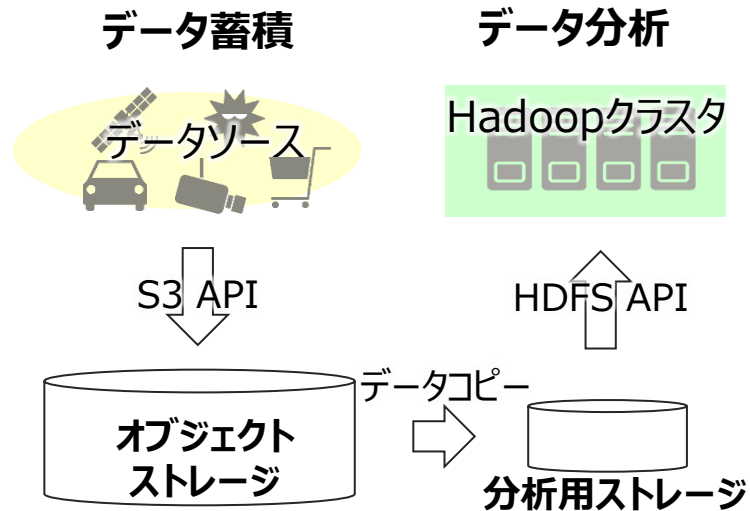
データ量が増えるにつれ、データコピー時間や分析用ストレージのコストが課題に



■ 一体型のオブジェクトストレージ活用

- Hadoopクラスタから、コネクタを介して直接オブジェクトストレージにアクセスし、必要なデータのみを取り出して分析

オブジェクトストレージのスケールアウト性を活かし蓄積と分析を一体化したデータ基盤を実現

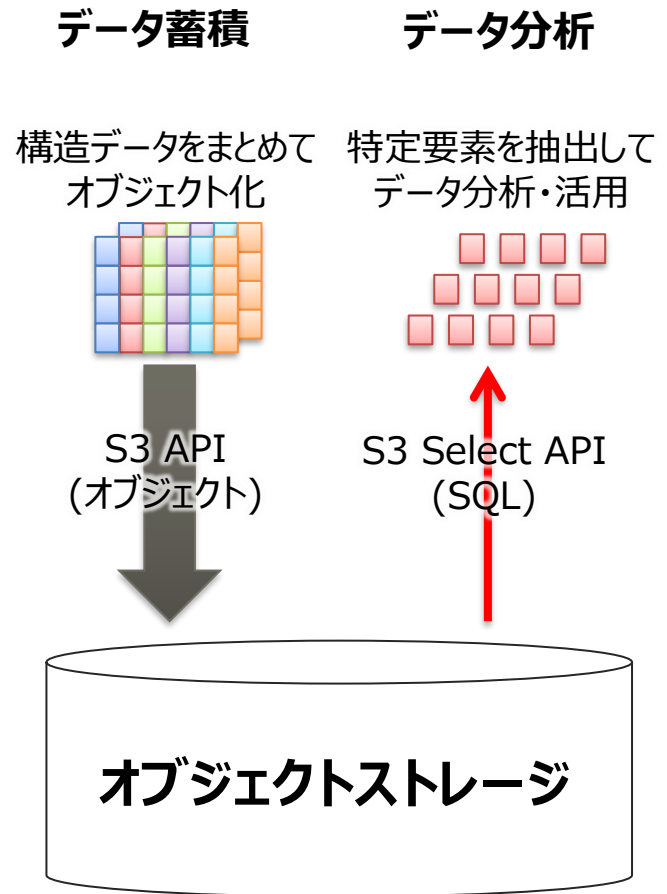


データレイク的なデータ活用

■ SQLによるオブジェクト内検索

- センサデータや構造データをまとめてオブジェクトとして蓄積しておき、活用時にはSQL等で、必要となる特定要素のみを抽出する方法
- クラウドストレージで一般的に利用できるが、近年オンプレストレージでも活用可能に
 - AWS: S3 Select, Amazon Athena, ...
 - GCP: BigQuery
 - OSS: Ceph, MinIO, ...

データ蓄積時のスキーマ設計が不要となり、データ基盤の導入/運用が簡単に



■ データアーカイブ市場の動向

- データ利活用の本格化に伴い、大量データの効率的管理がより重要に
- 現在、年間約4.5ZBのデータが新たに創出。その約3.5~4%程度が長期保存が必要なデータであり、その市場は約1兆円になると予測

■ 大容量・低コストなテープストレージが再注目

- 従来のバックアップ用途ではなく、アーカイブ用途としての利用に注目
 - バックアップ・・・データ保護を目的として、**原本データの一時的な複製を作成**すること
 - アーカイブ・・・保管/利用を目的として、**原本データを永久/長期に保存**すること
- アーカイブ用途の**大規模テープストレージはプラス成長を継続**すると予測 (IDC)

■ メガクラウド各社もアーカイブクラスのラインナップを強化 (システム構成は不明)

	ストレージサービス	アーカイブクラス
AWS	Amazon S3	Glacier, Deep Archive
Azure	Azure Blob Storage	Cold, Archive
Google	Cloud Storage	Nearline, Coldline, Archive

テープの特徴(1) ～信頼性～

■ 一昔前とは全然違う、新しいテープの時代

■ テープストレージ技術の進化による品質向上

テープにまつわる不安のほとんどは過去の話

切れる、絡む？



テープメディアとドライブ
双方の技術革新により
物理ダメージ発生は
大幅減。

定期的な巻き直しが必要？



テープ素材の改良により
テープ貼りつきや磁気転写の
心配は全くありません。

カビが生える？



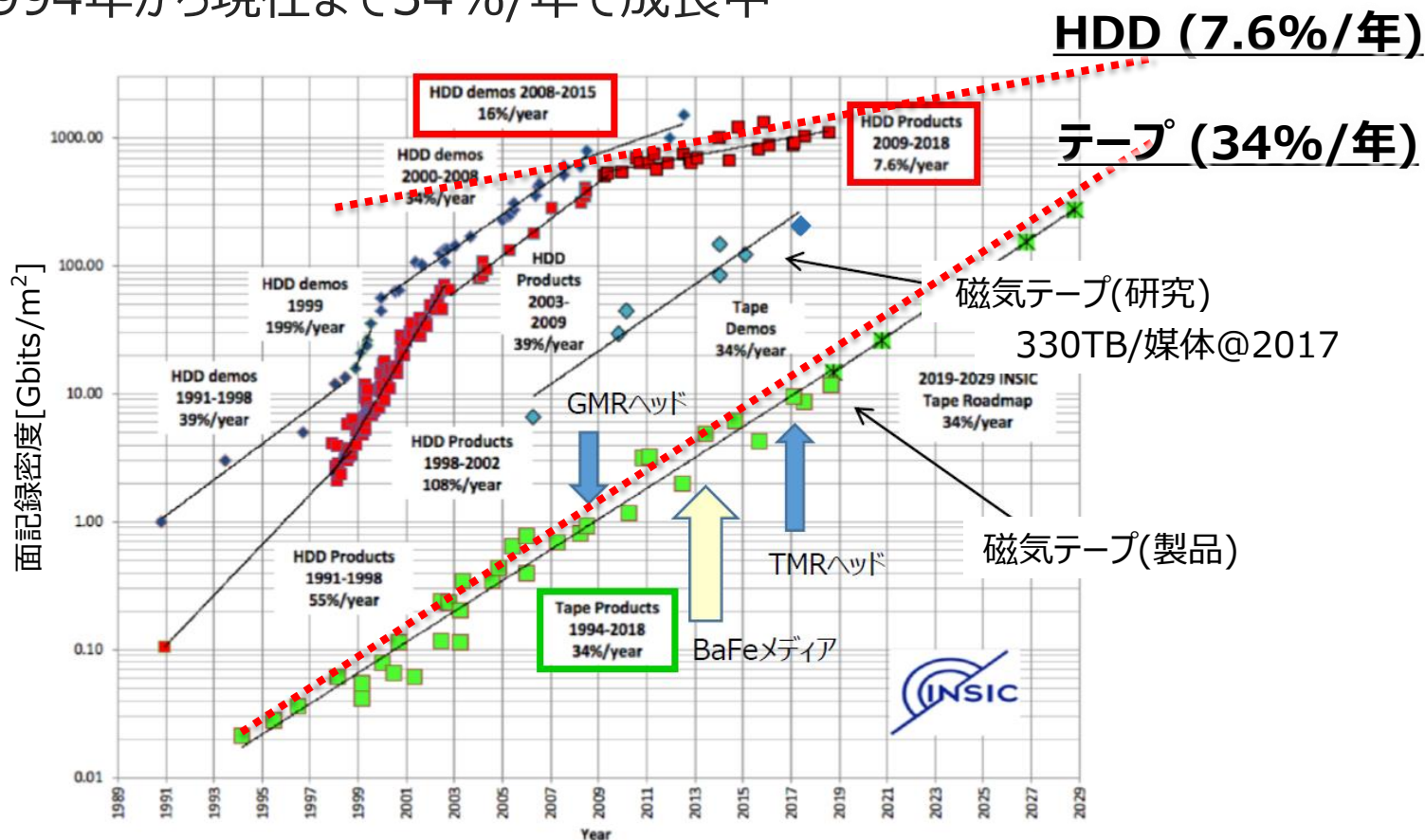
密閉構造のカートリッジ、
テープ素材の改良によって
カビの心配もありません。

出典：JEITAテープストレージ専門委員会 (CEATEC 2019 発表資料)

テープの特徴(2) ～大容量～

■ 記録密度の高い成長

- HDDは2009年頃から大容量化の速度が大幅に鈍化 (7.6%/年)
- テープは1994年から現在まで34%/年で成長中

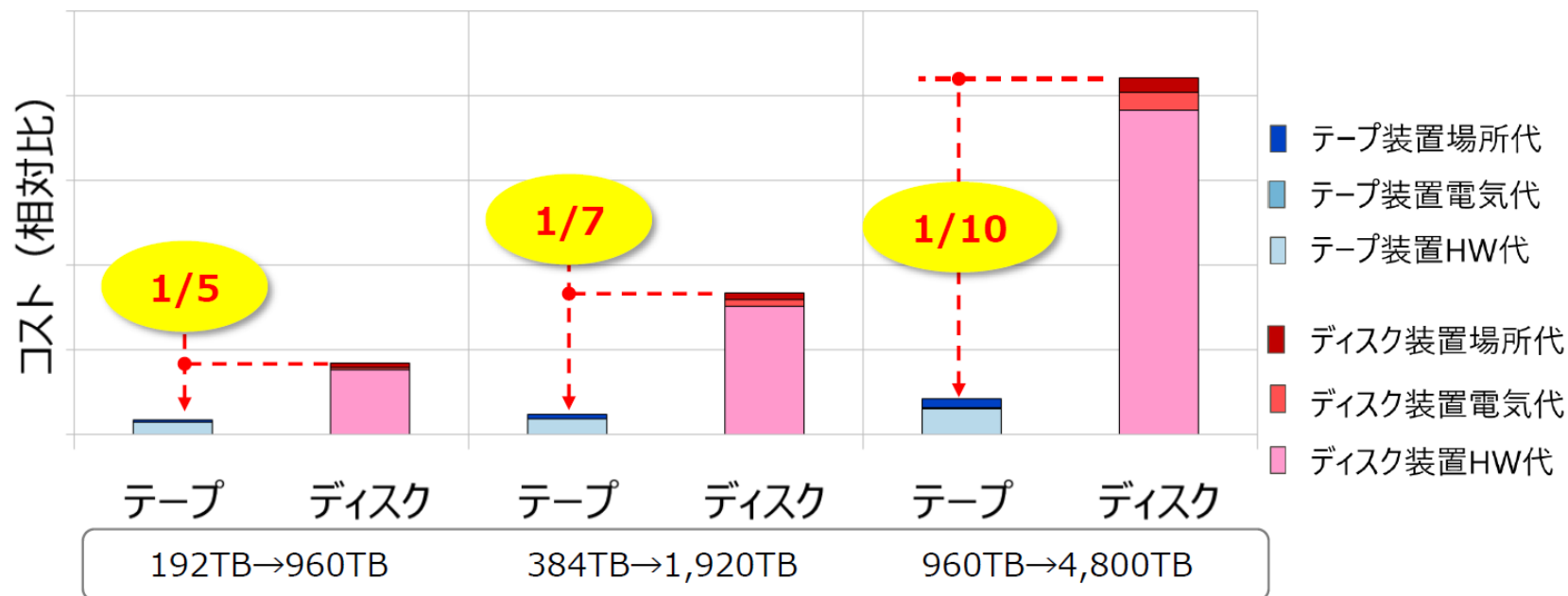


出典: INSIC 2019 Areal Density Trends. Hard Disk Drive, Tape Product and Tape Technology Roadmap

テープの特徴(3) ～低コスト～

■ ディスク装置と比較して数分の1のコスト

■ 容量が多いほど差が大きくなる



5年間の容量増加量

- ※ テープ装置 : 80巻テープライブラリ、LTO 8ドライブ搭載 (非圧縮12TB)
- ※ ディスク装置 : RAID 6構成、高密度実装タイプ、エコモード、Near Line 12TB HDD

出典: JEITAテープストレージ専門委員会 (CEATEC 2019 発表資料)

テープの特徴(4) ～利便性の向上～

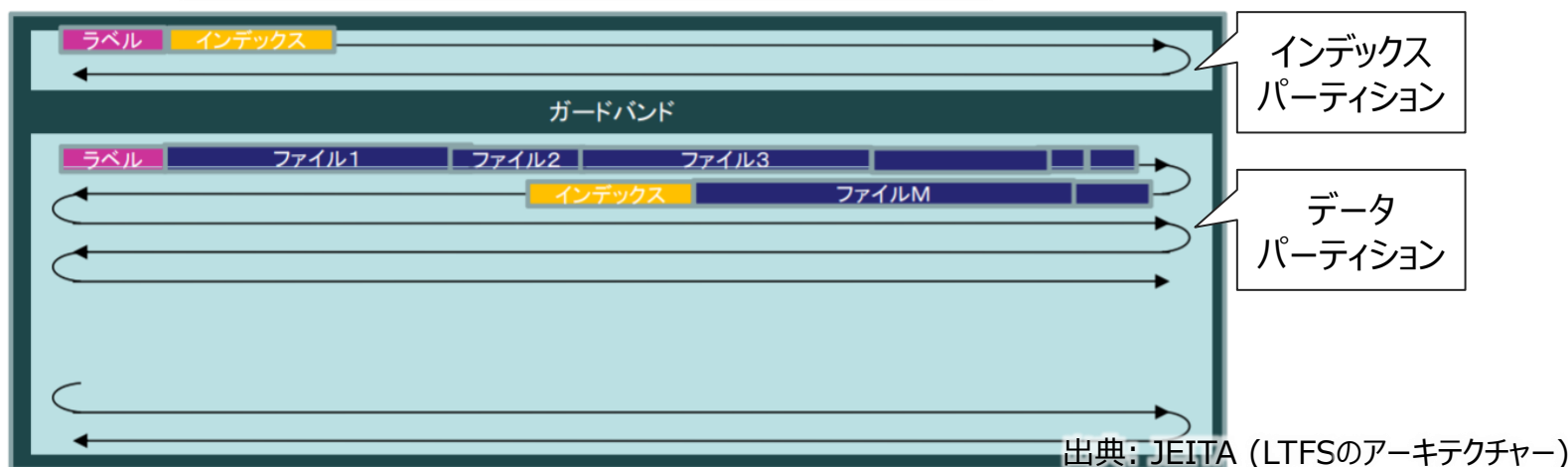
■ LTO (Linear Tape-Open)

- Seagate社、HP社、IBM社が共同で開発した**磁気テープの統一規格**
- 2000年に第1世代(LTO-1)を発売。最新世代はLTO-8。

■ LTFS (Linear Tape File System)

- LTO-5以降でサポートされた**テープ用ファイルシステム**
- インデックスパーティションとデータパーティションに2分割
 - インデックスパーティション・・・**最新のインデックス情報**。テープ先頭部分から上書き
 - データパーティション・・・**ユーザデータを追記保管**。インデックス情報も**随時書き込まれる**

分割されたテープとファイル保存のイメージ



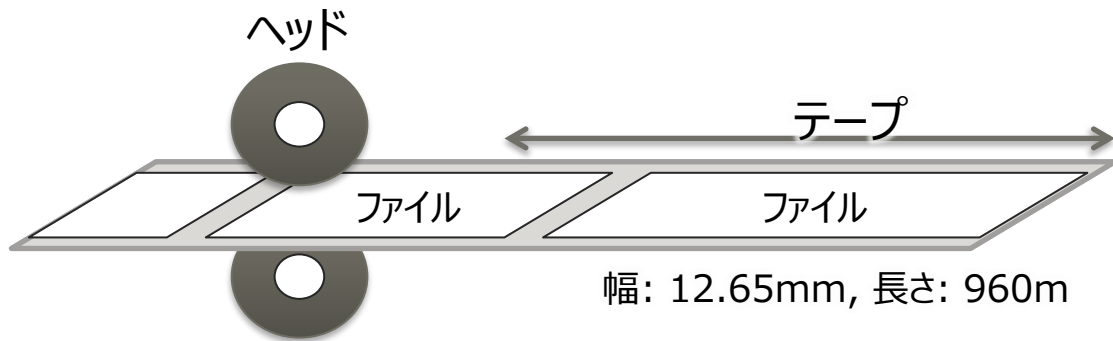
テープの特徴(5) ～アクセス性能～

■ シーケンシャルアクセスは高速

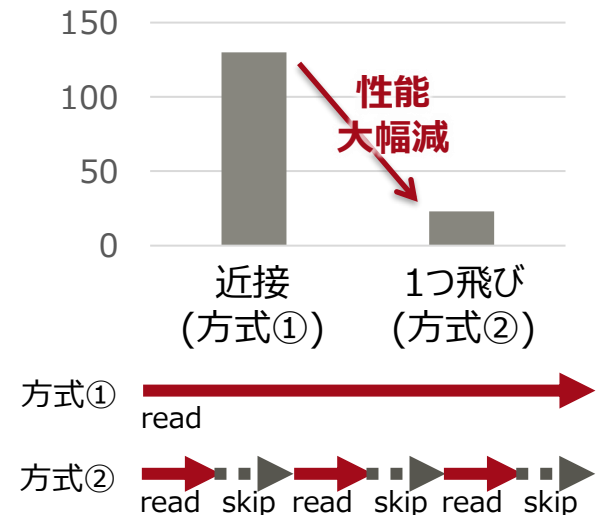
規格	LTO-6	LTO-7	LTO-8
発売日	2012	2015	2017
容量(非圧縮)	2.5TB	6.0TB	12TB
速度(非圧縮)	160MB/s	300MB/s	360MB/s

■ ランダムアクセスは非常に低速

■ ヘッドの移動と位置合せに時間がかかる



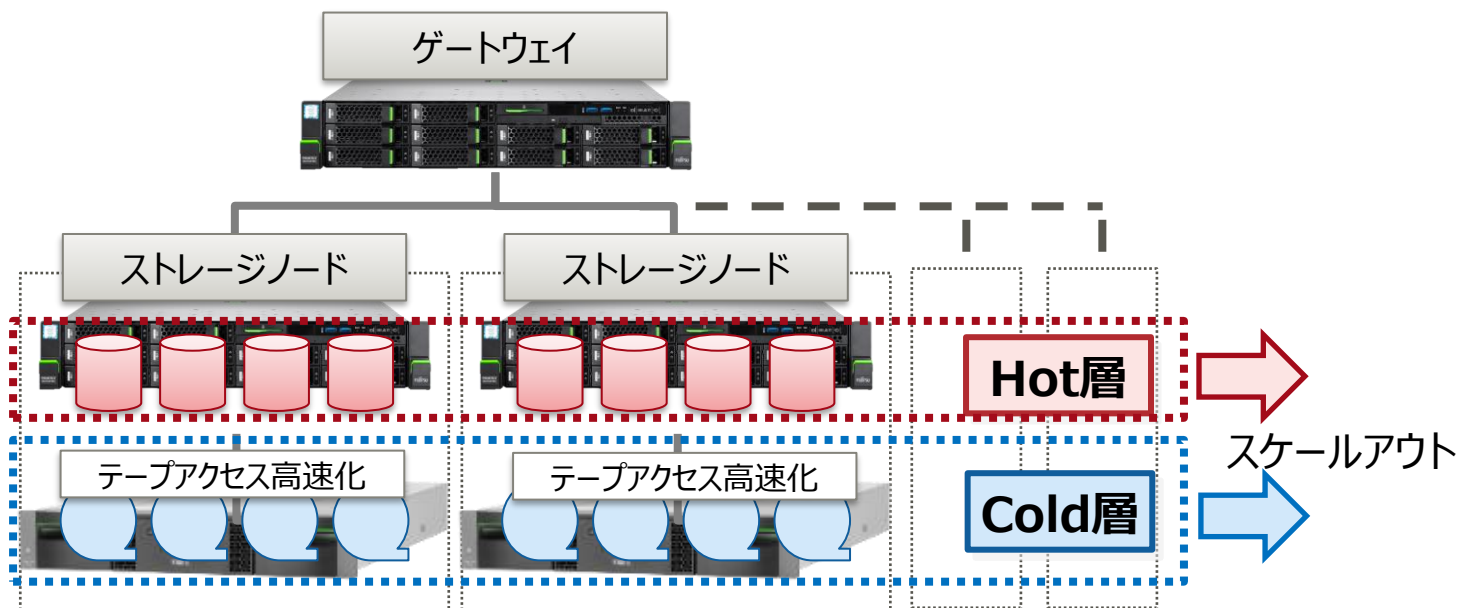
読出性能 [MB/s]
(LTFSを用いた4MBファイルの読出)



富士通研究所では、オブジェクトストレージCephを基盤としたテープ階層化ストレージの研究開発を進めています

■ 特徴

- 蓄積場所に関わらず、ユーザからはシームレスにアクセス
- 独自のテープアクセス高速化技術を組み込み、高い性能コスト比を実現
- ストレージ階層毎に独立にスケールアウト
- オブジェクト単位で階層移動を制御



Wrap-aware列指向データ配置技術

テープのデータレイアウトを利用し、時系列センサデータのアクセス高速化

■ テープ上のデータレイアウト

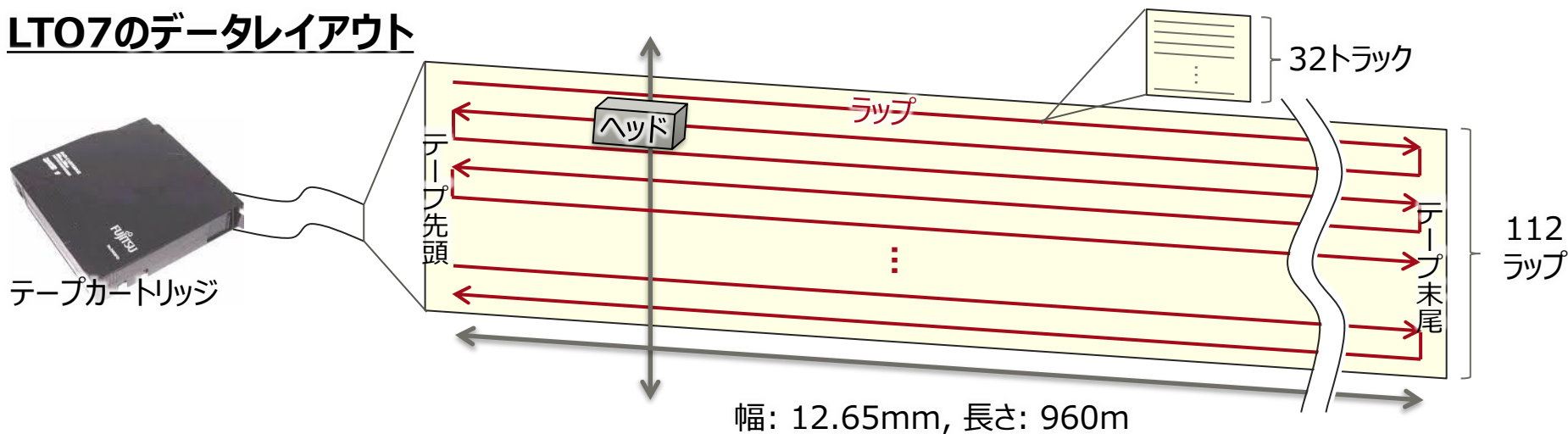
■ ラップ(Wrap)と呼ばれる多数の帯から構成

- LTO7の場合、112本のラップで構成 (テープ容量6TBに対し、各ラップ容量は約53GB)

■ データアクセスの特徴

- データ書き込みは、各ラップを折り返しながら追記的に書き込んでいく
- ラップを跨いだアクセスはヘッドシークのみで、テープ送りすることなく高速実行

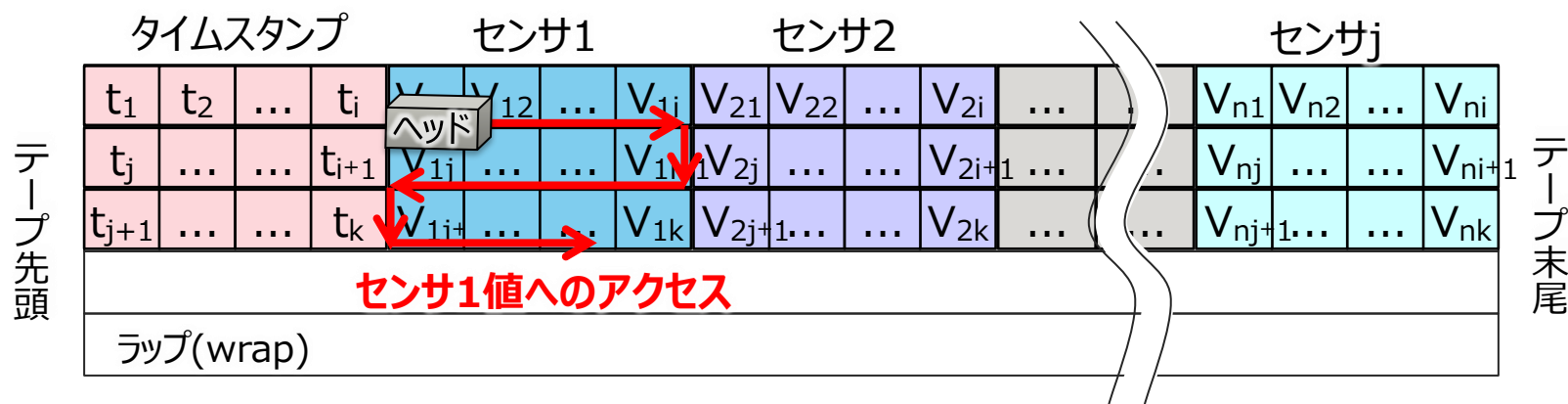
LTO7のデータレイアウト



Wrap-aware列指向データ配置技術

■ 技術概要

- **1ラップ分**(約53GB)のデータをバッファして、**列指向形式**に変換してテーブルに書く
- ラップ方向に合わせて**並び順を反転** → 同じセンサ値は**2次元的に近くに配置**




■ 評価

- 特定のセンサデータを抽出する際に、従来の時間順にセンサデータを並べた方式と比較して、約3倍の高速化を実現
 - データ量 = 350GB, カラム数 = 100カラム, Zipf分布に合わせて平均6カラムにアクセス

- オブジェクトストレージは、クラウドストレージだけでなく、オンプレでも製品やOSS活用が本格化
- 拡張性に優れるだけでなく、データ活用に向けた取り組みも進んでいる
- 大量データを蓄積・活用するにはコスト最適化が重要
 - クラウドストレージ ⇒ アーカイブクラスのストレージサービスの活用
 - オンプレストレージ ⇒ テープ等のコールドストレージの活用

データは「21世紀のオイル」：
「データが価値を生む」、「上手く活用しないと価値が生まれない」

増加し続けるデータを蓄積・活用する基盤として、
オブジェクトストレージは、高いポテンシャルを持っている



FUJITSU

shaping tomorrow with you