

TOSHIBA

Leading Innovation >>>

高速処理と高信頼性を両立し、
ペタバイト級の多種大量データを蓄積する、
ビッグデータ/ I o T時代のデータベースとは？？

株式会社 東芝
野々村 克彦

TOSHIBA

Leading Innovation >>>

高速処理と高信頼性を両立し、
ペタバイト級の多種大量データを蓄積する、
ビッグデータ/ I o T時代のデータベースとは？？

スケールアウト型DB GridDB

株式会社 東芝
野々村 克彦

プロフィール

- **名前：野々村 克彦**
- **経歴：**
 - 00年代 XMLデータベースTX1 開発メンバ
 - 10年代 スケールアウト型DB GridDB(旧GridStore) 開発メンバ
 - 15年～ GridDB オープンソース・チーム
- **昨年から始めたこと：**
 - 息子（7歳）とサッカー
 - GitHub

目次

1. はじめに

- ビッグデータ
- NoSQL
- IoTと既存NoSQLの課題

2. GridDB

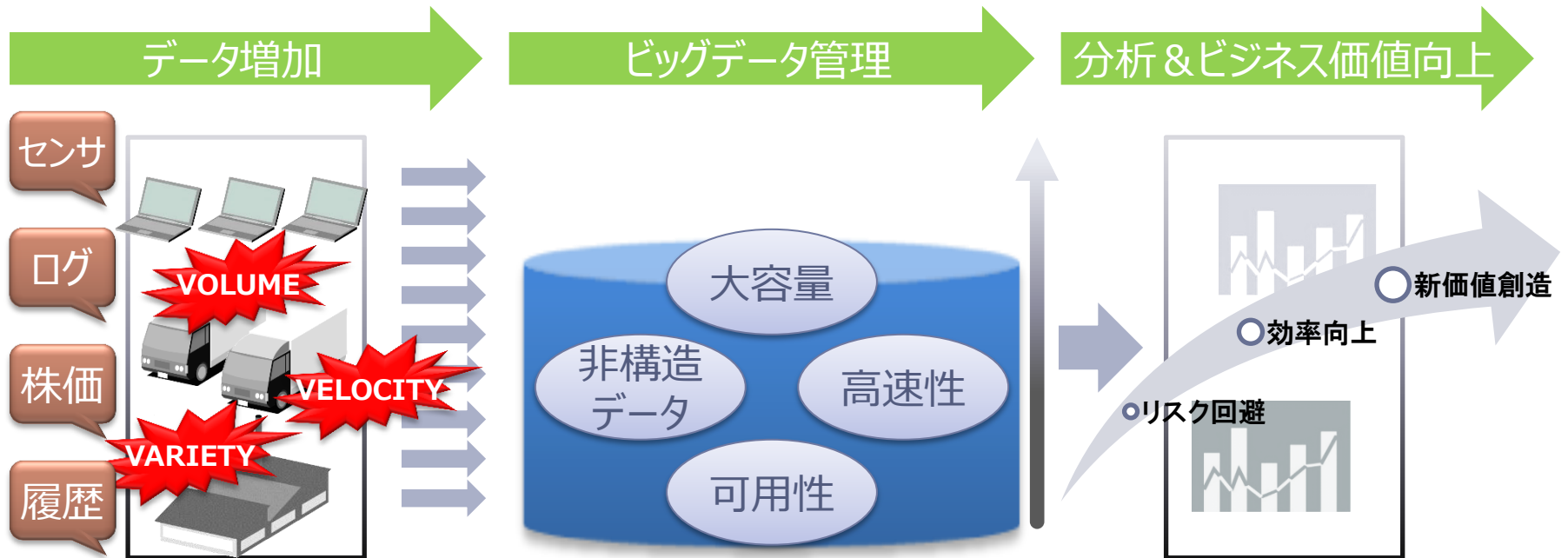
- オープンソース化
- 特長
- 適用事例
- 公開サイト

3. 利用方法

4. まとめ

ビッグデータ

- ビジネスの価値を向上させるビッグデータ活用が本格化
 - センサーデータ、履歴データなど多様なデータが日々増加
- ビッグデータ管理の要件に合わせたデータベースが必要



ビッグデータ管理は柔軟な拡張性が必須

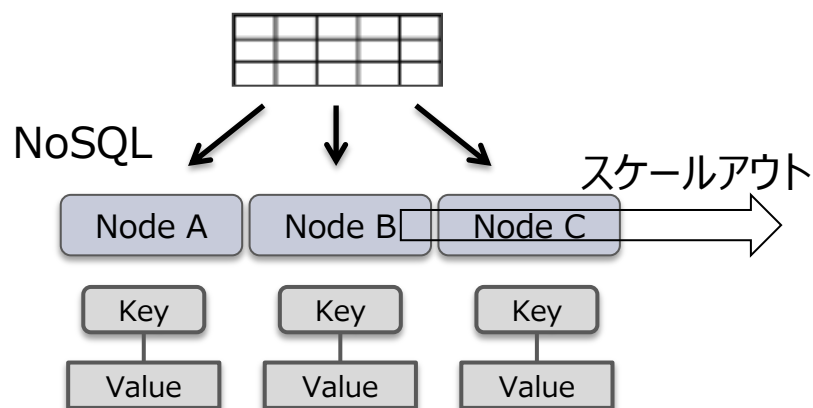
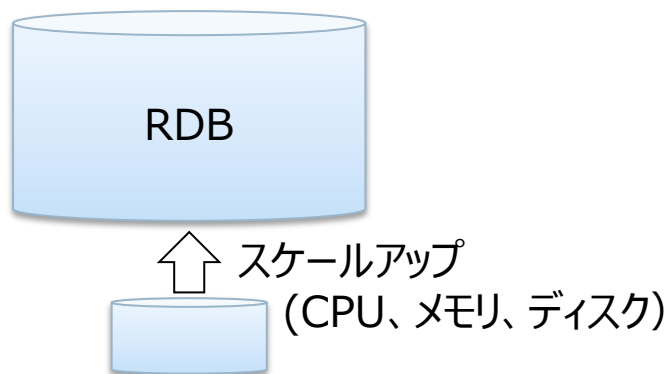
RDBとNoSQL

• RDB

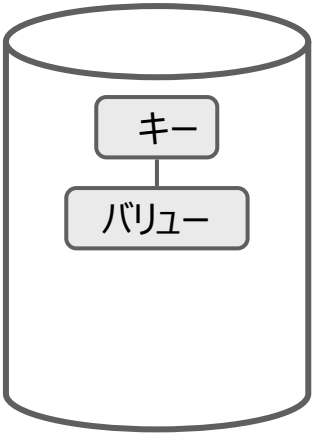
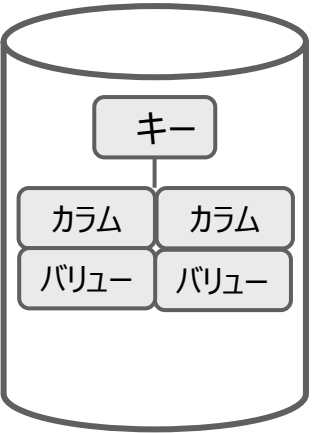
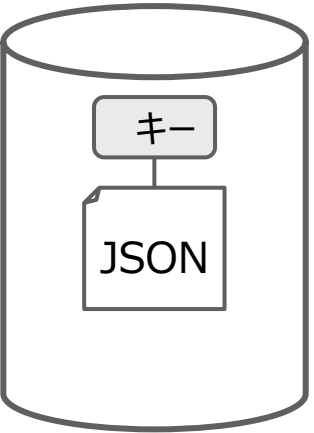
- スケールアップでは限界がある。ビッグデータを管理するのに適していない
- 一貫性を重視するため、スケールアウトは困難である

• NoSQL(Not Only SQL)

- キーによるPut/Getが基本I/F（キーバリュー型）
- スケールアウトによる性能向上で近年注目されている
- 一貫性を緩和する代わりにRDBでは対応できない規模の大容量データを管理可能である



NoSQLのデータモデル

	キーバリュ型	列指向型	ドキュメント型
データモデル			
NoSQLの例	Riak	Cassandra	MongoDB

- **IoT(Internet of Things、モノのインターネット)**

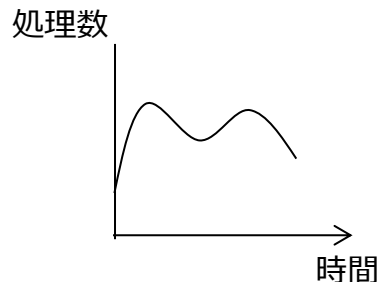
- 一意に識別可能な「もの」がインターネット/クラウドに接続され、情報交換することにより相互に制御する仕組み ※「IoT」『フリー百科事典 ウィキペディア日本語版』

- **特性**

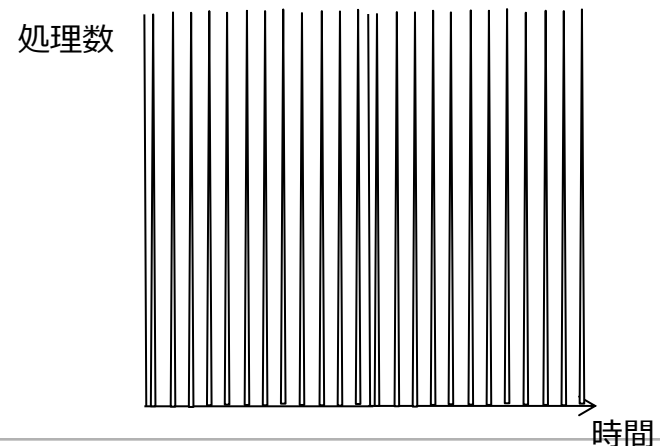
- 分、秒周期、さらにはそれ以下の周期で発生する膨大なセンサーデータを扱う必要がある。長時間に渡るデータを保持する必要がある。
- 各センサ内のデータ欠損や参照データの矛盾など、データ一貫性やデータ整合性を保つ必要がある。

人の活動で生成されるデータ

- SNS、ゲームなど
- テキスト、イメージデータ
- インメモリ前提



IoTデータ (センサー、ログ、履歴、株価等)

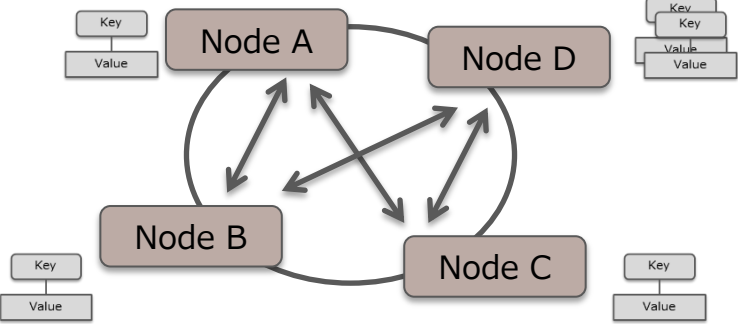
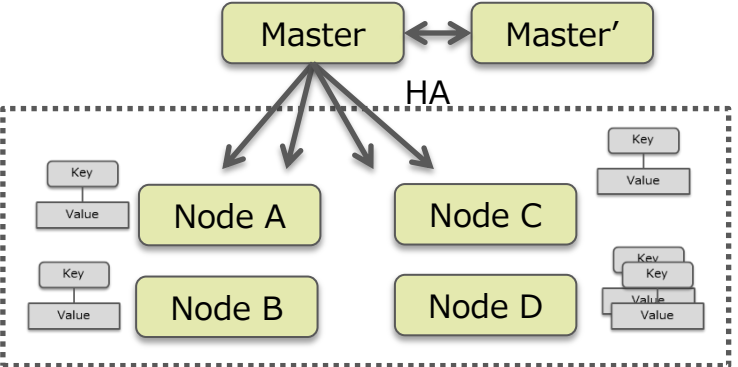


IoTにおける既存NoSQLの課題

(A) IoTのデータ管理が困難

- センサ単位の一貫性を保てない。時間範囲指定の検索ができない、メモリが足りない場合に性能が大幅に劣化、など

(B) 既存クラスタ管理方式に起因するトレードオフ問題

ピアツーピア(Peer to Peer)	マスタスレーブ(Master Slave)
	
<p>○ノード追加でのデータ再配置が容易 ×一貫性維持のためのノード間通信のオーバーヘッドが大⇒一貫性と処理速度がトレードオフ</p>	<p>○一貫性の維持は容易 ×マスタノードが単一障害点(SPOF) ×ノード追加でのデータ再配置が難しい</p>

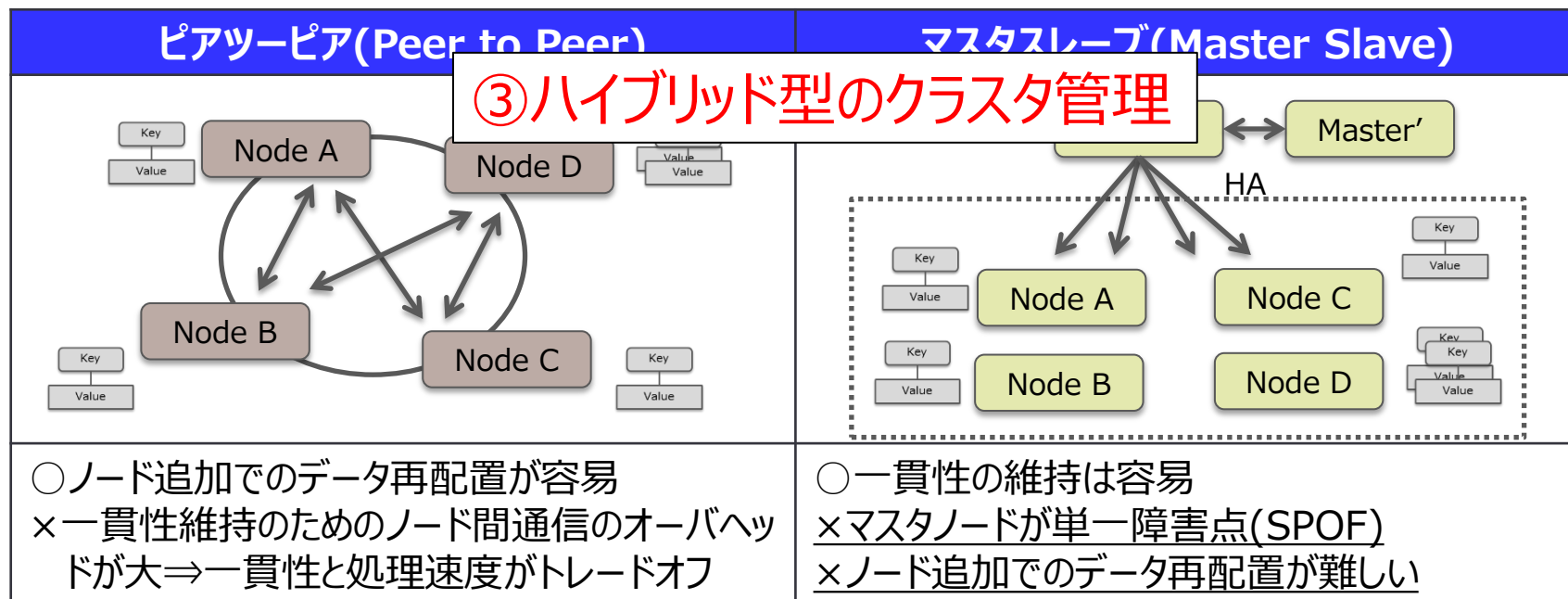
IoTにおける既存NoSQLの課題

(A) IoTのデータ管理が困難

- センサ単位の一貫性を保てない、時間範囲指定の検索ができない、メモリが足りない

①キーコンテナ型のデータモデル

(B) 既存クラスタ管理方式に起因するトレードオフ問題



目次

1. はじめに

- ビッグデータ
- NoSQL
- IoTと既存NoSQLの課題

2. GridDB

- オープンソース化
- 特長
- 適用事例
- 公開サイト

3. 利用方法

4. まとめ

GridDBのオープンソース化

- **GridDBとは**

- IoTデータ管理向けのスケールアウト型DB

- **GitHub上にNoSQL機能をソース公開(2016/2/25)**

- (英語)https://github.com/griddb/griddb_nosql/
- (日本語)https://github.com/griddb/griddb_nosql/README_ja.md



- **目的**

- ビッグデータ技術の普及促進
 - 多くの人に知ってもらいたい、使ってもらいたい。
 - いろんなニーズをつかみたい。
- 他のオープンソースソフトウェア、システムとの連携強化

GridDBの特長

① IoT向けデータモデル

- キーコンテナ型のデータモデル

② 高パフォーマンス(High Performance)

- メモリ指向アーキテクチャ

③ 高信頼性(High Reliability)

- (P2Pとマスタスレーブの) ハイブリッド型のクラスタ管理技術

④ 高スケーラビリティ(High Scalability)

- 自律データ再配置 (ADDA) 技術

① IoT向けのデータモデル

キーコンテナ型のデータモデル

- キーバリューをグループ化するコンテナ（テーブル）
- コンテナのスキーマ定義が可能。カラムにインデックスを設定可能
SQLライクなクエリ(TQL)が利用可能
- レコード単位でトランザクション操作（コンテナ単位でACID保証）

※ACID : Atomicity、Consistency、Isolation、Durability

IoTデータ



機器1のレコード

日時	センサA	センサB
2015/01/01 0:00	7.788683	0.648364

キー

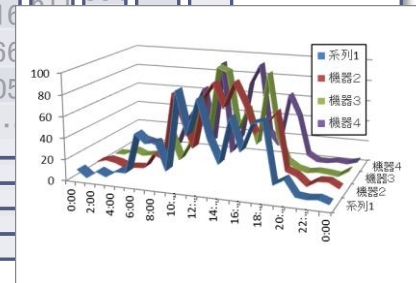
機器 1

対象毎にIoTデータを格納

データ格納

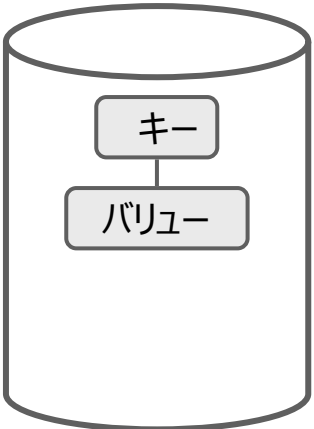
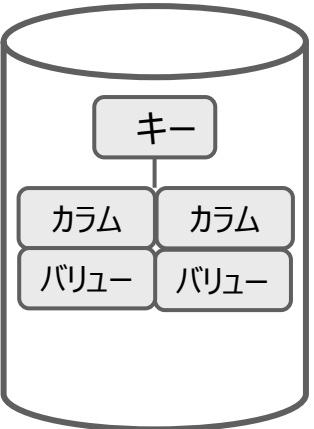
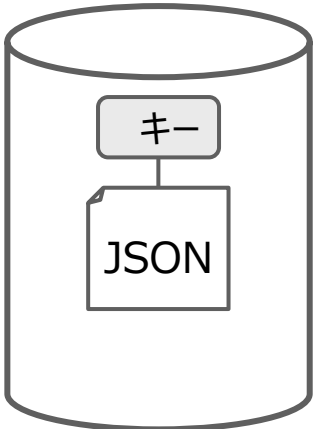
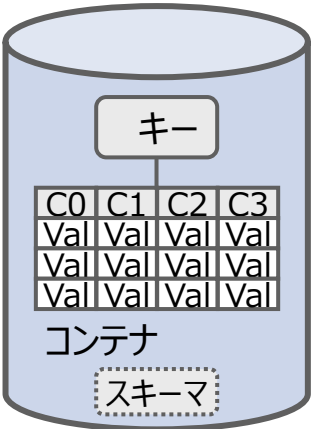
日時	センサA	センサB
2015/01/01 0:00	7.788683	0.648364
2015/01/01 1:00	0.68874	0.353611
2015/01/01 2:00	7.677135	5.881216
2015/01/01 3:00	3.1816	2.511166
2015/01/01 4:00	9.739242	0.655805
...

テーブル表現で管理



単純なキーバリュー型とは異なり、使い慣れたRDBに近いモデリングが可能

データモデルの比較

	キーバリュー型	列指向型	ドキュメント型	キーコンテナ型
データモデル				
NoSQLの例	Riak	Cassandra	MongoDB	GridDB

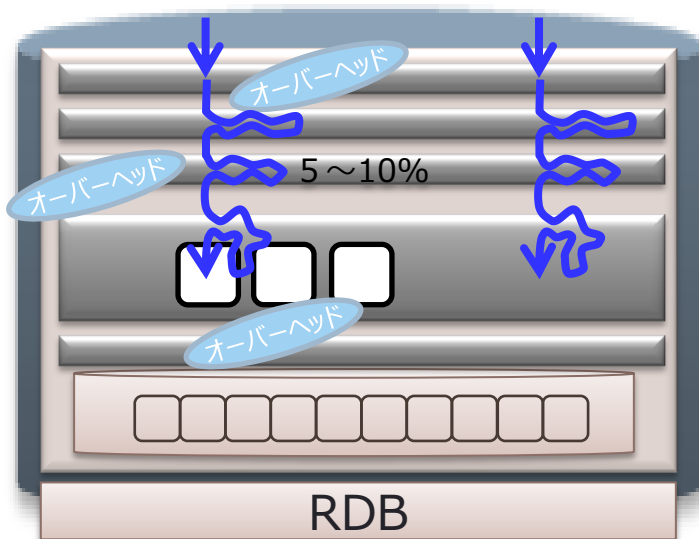
• コンテナの種類

- コレクションコンテナ：レコード管理用
- 時系列コンテナ：時刻で並べられたレコード集合。時系列データ管理用
 - 期限解放機能、サンプリング機能など

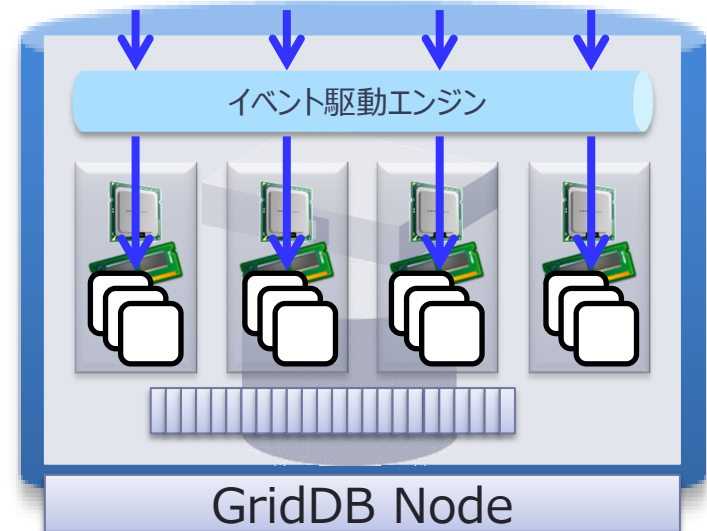
② 高パフォーマンス

メモリ指向アーキテクチャ

- イベント駆動エンジンであるため、少ないリソースで多くの要求を無駄なく処理
- メモリ、ディスクアクセスの排他処理や同期待ちを無くした、オーバーヘッドの無いデータ処理
- GB超級のメモリ搭載を前提とし、読み書きサイズを最適化しI/O効率を改善



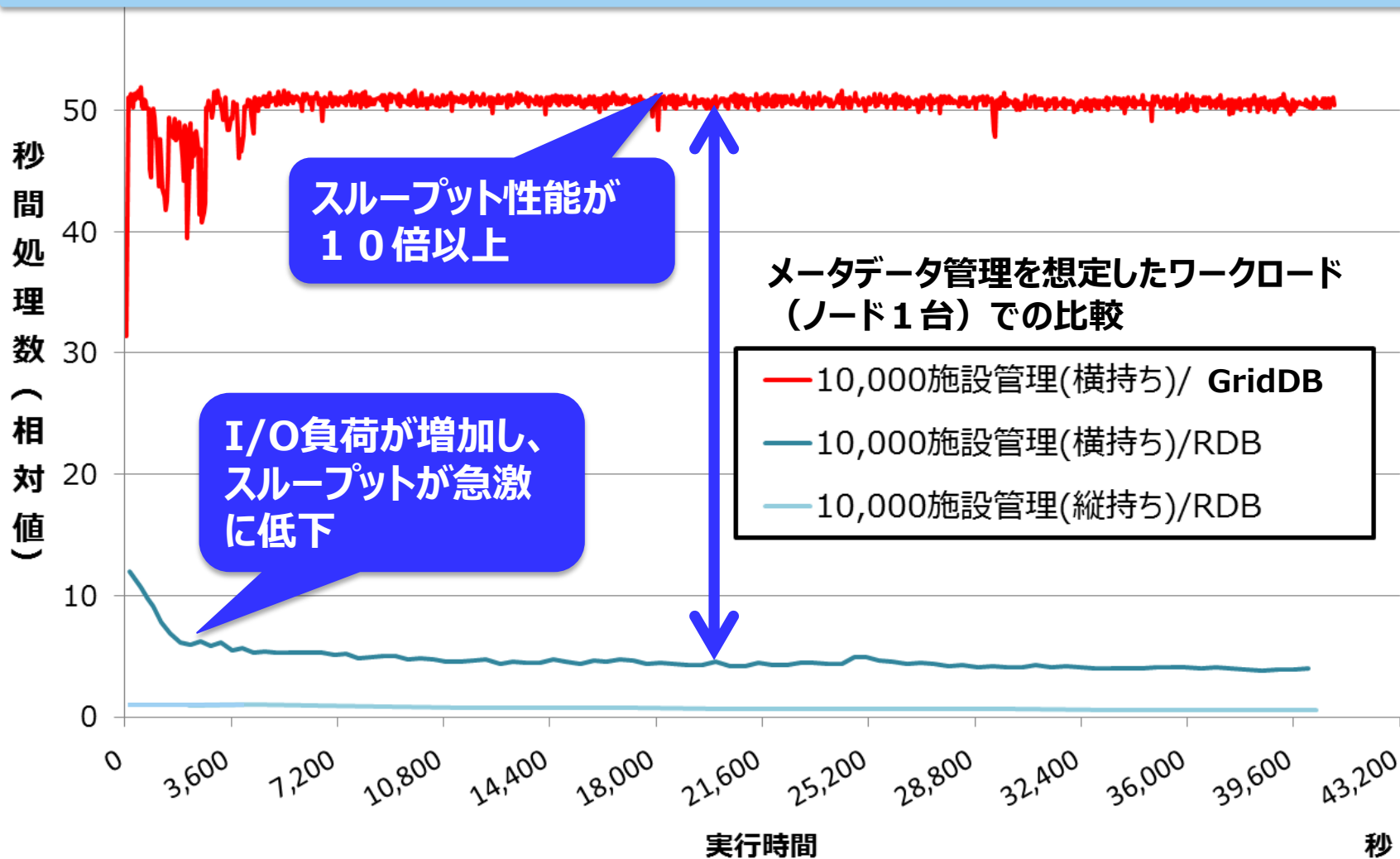
要求処理
トランザクション管理
クエリ処理
バッファ処理
I/O処理



H/Wのスペックを最大限に生かすインメモリ指向DB

RDBとの性能比較

登録件数が10億件超の規模になっても、安定した性能を維持



※RDBの10,000施設処理ケースでの最大値を1とした場合の相対値
横持ち: 1時刻で機器毎のメータをまとめて(50メータで)1レコード作成
縦持ち: 1時刻で1メータごとに1レコード作成

※1施設当たり50メータ

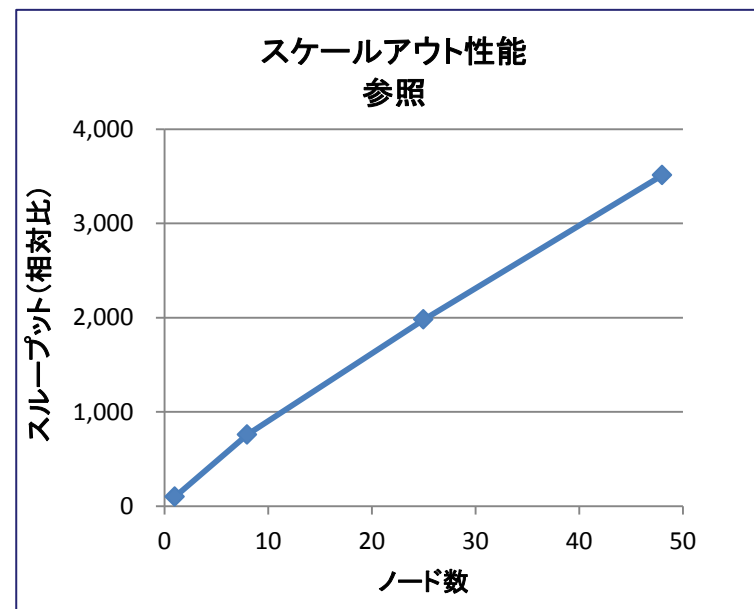
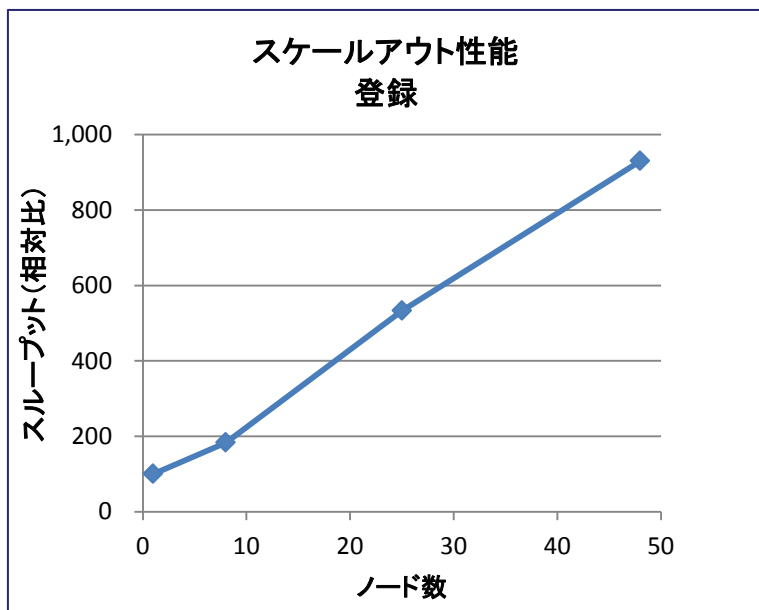
スケールアウト性能

- Yahoo Cloud Serving Benchmark (YCSB) による測定

<http://labs.yahoo.com/news/yahoo-cloud-serving-benchmark>

- ノード数に対して
性能がリニアに向上

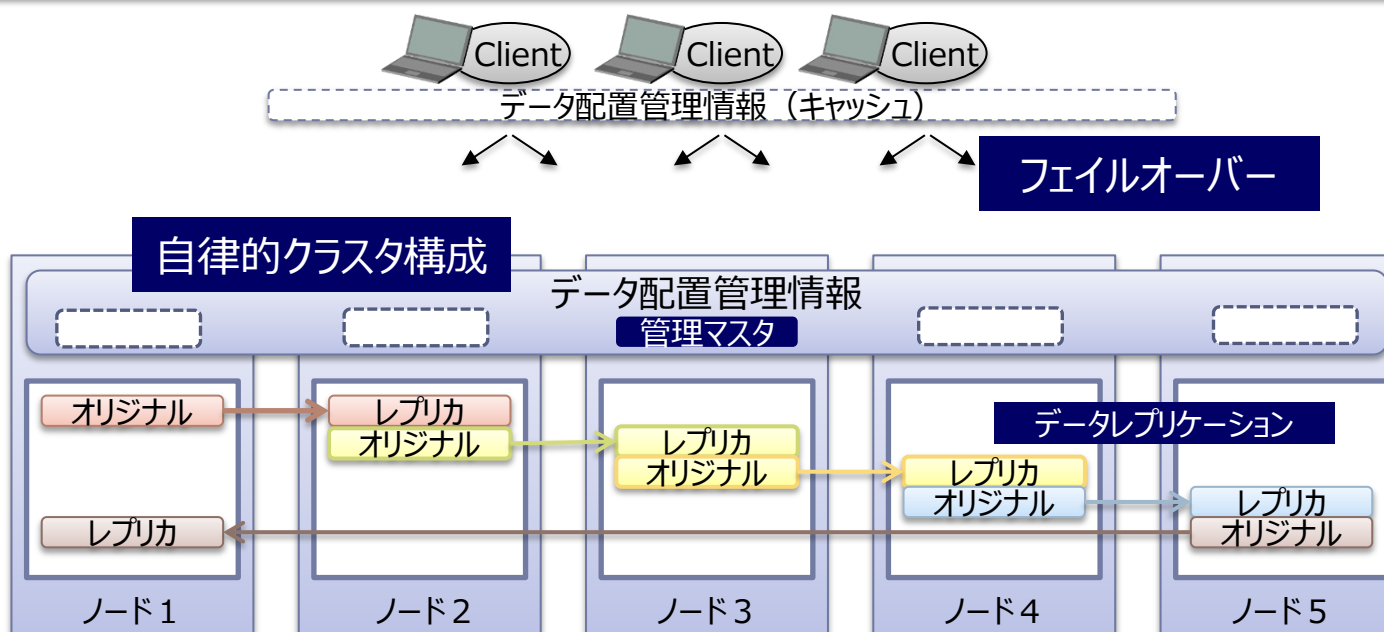
項目	値
台数	1~50台
ロウサイズ (レコードサイズ)	724B
件数	7.5億件
レプリカ数	3



③ 高信頼性

ハイブリッド型クラスタ技術

- ノード間で自律的、動的にマスタノードを決定。単一故障点（SPOF）を排除
- レプリケーションによるデータ多重化、フェールオーバーが可能
- 永続化（インメモリ／ディスク併用）

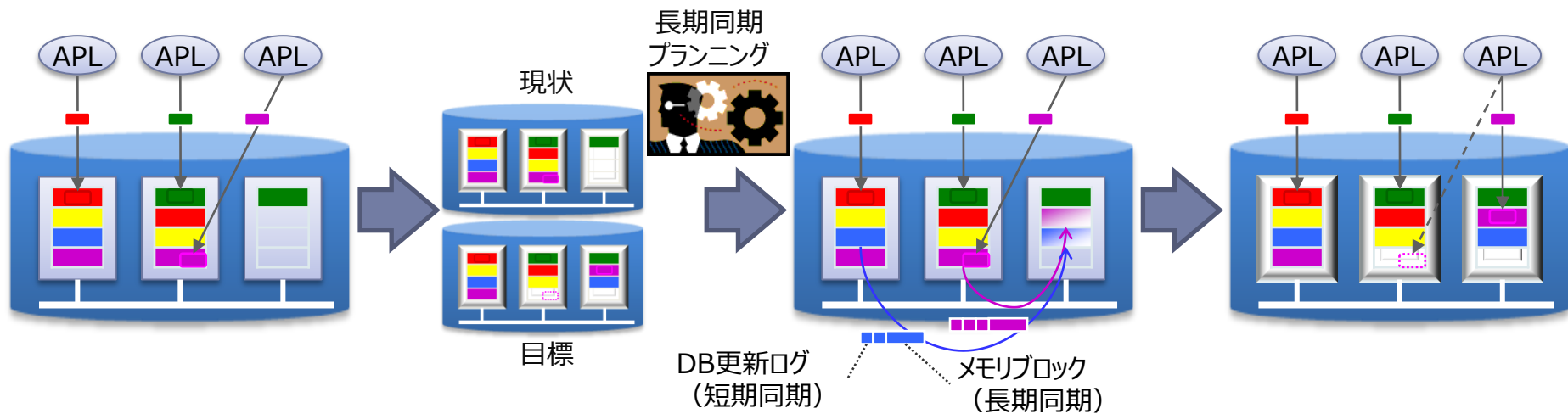


特別なスキルを必要とせずに、高可用なクラスタ構成が可能

④ 高スケーラビリティ

自律データ再配置技術 (A D D A : Autonomous Data Distribution Algorithm)

- インバランス状態を検知、長期同期プランニング
- 2種類のデータを使ってバックグラウンド高速同期、完了後切替

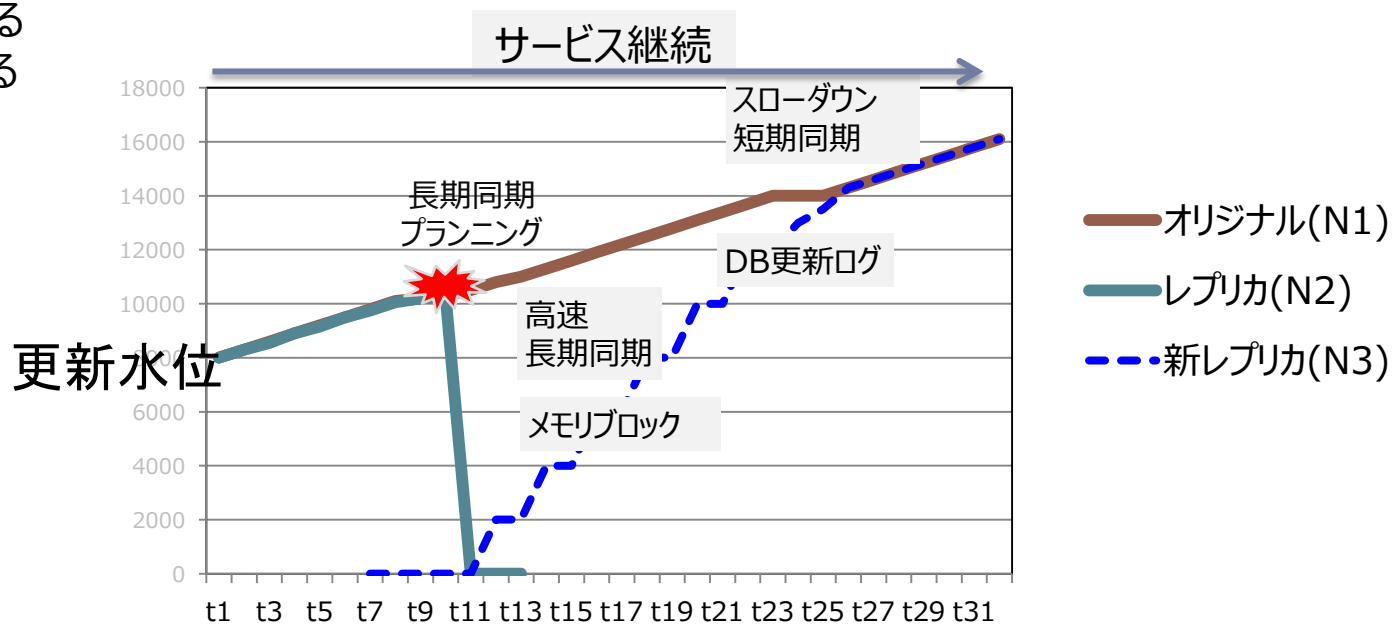
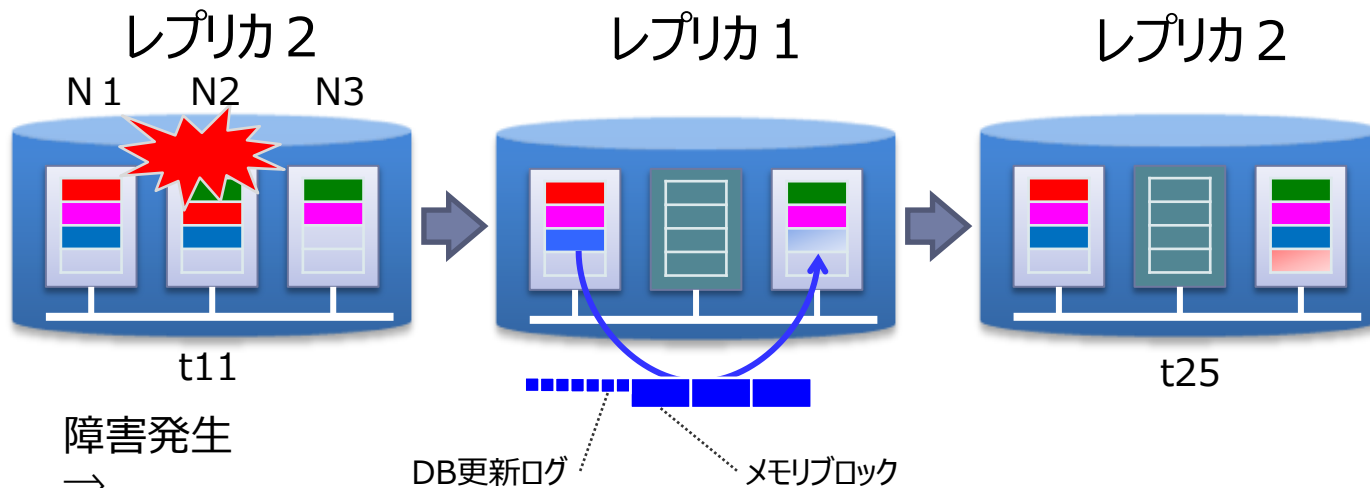


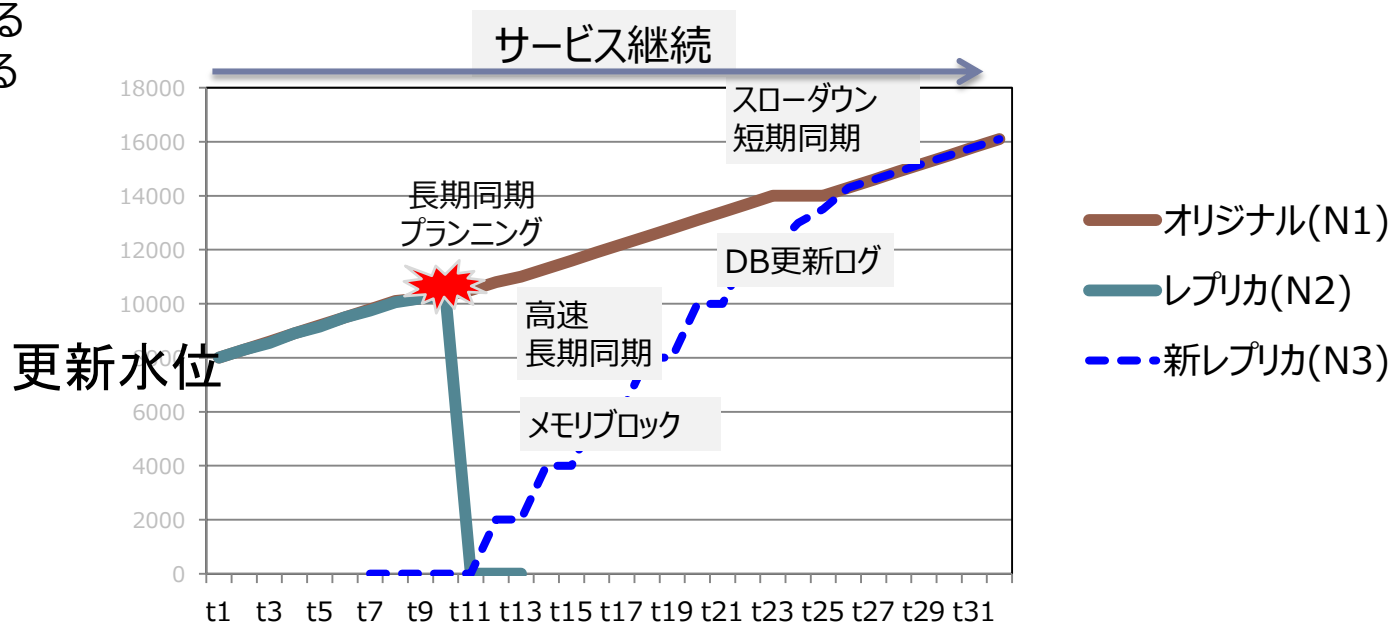
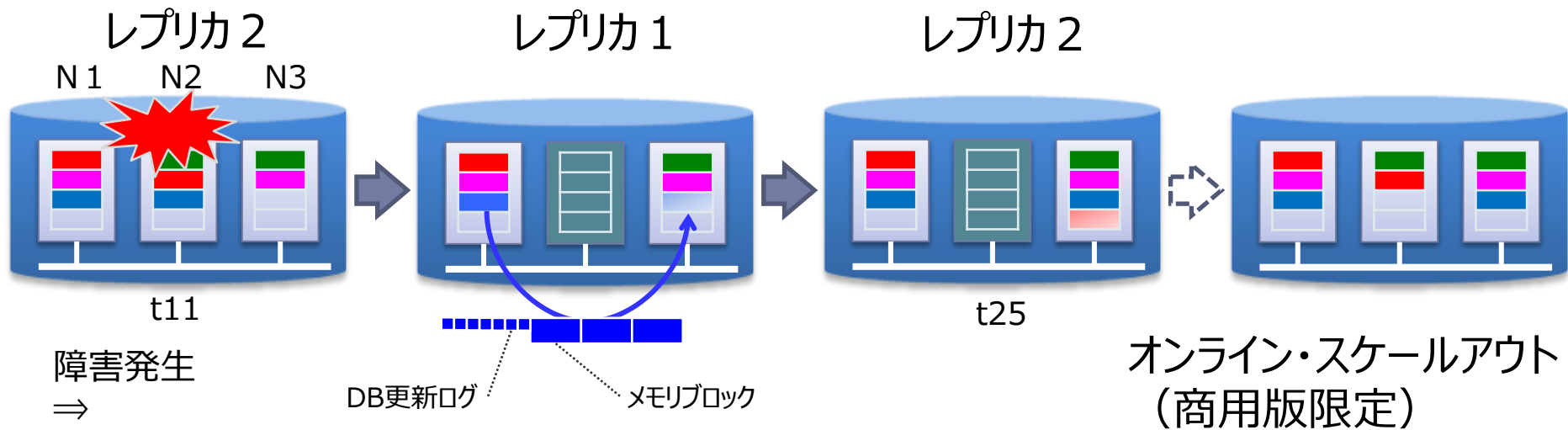
① 負荷インバランス検知

② 長期同期プランニング

③ 長期同期実行

④ アクセス切替

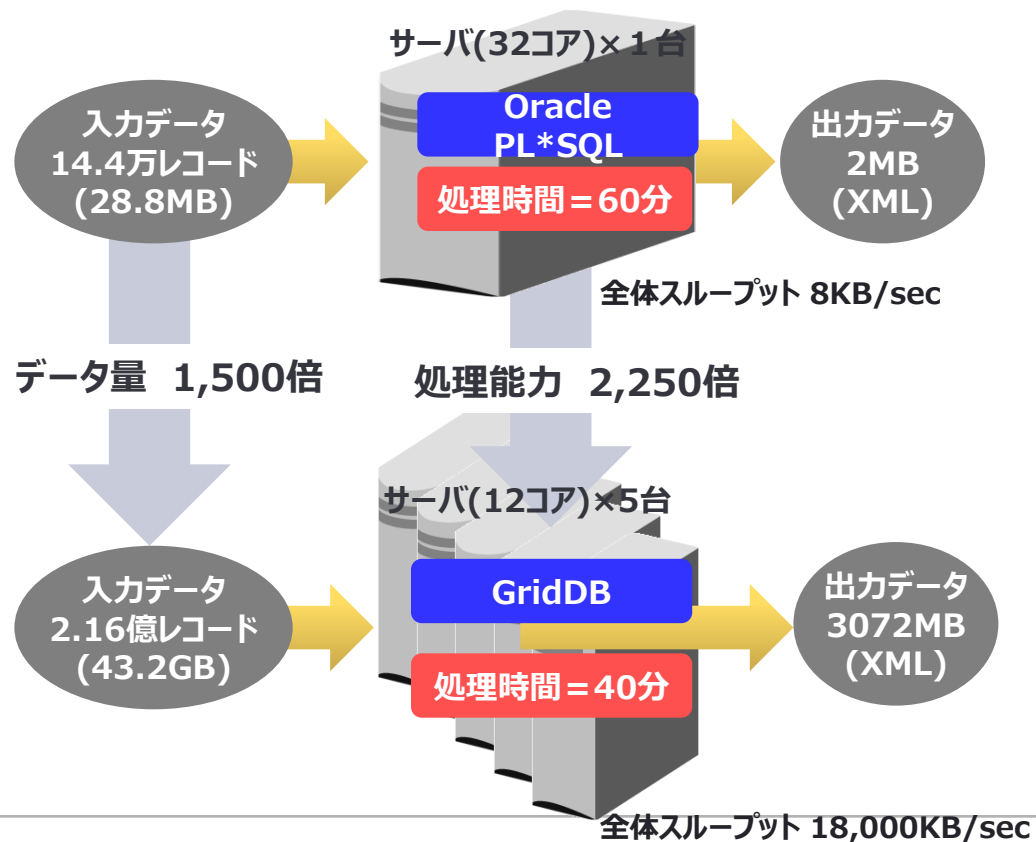




適用事例

GridDB（商用版）を適用し、1,500倍のデータ量を従来比の2,250倍の処理能力で対応

電力小売り事業者に対し、電力送配電網を提供し、契約ユーザの利用量に応じた料金を請求するシステム
電力の自由化に伴い、多数の電力小売り事業者が参入し、契約数の増加（3,000契約→450万契約）によるデータ量の爆発的増加へビッグデータ技術を適用し対応



公開サイト

(2016/2時点)

● 主な公開物件

– ソース :

- サーバ(C++)
- Javaクライアント
- 運用コマンド

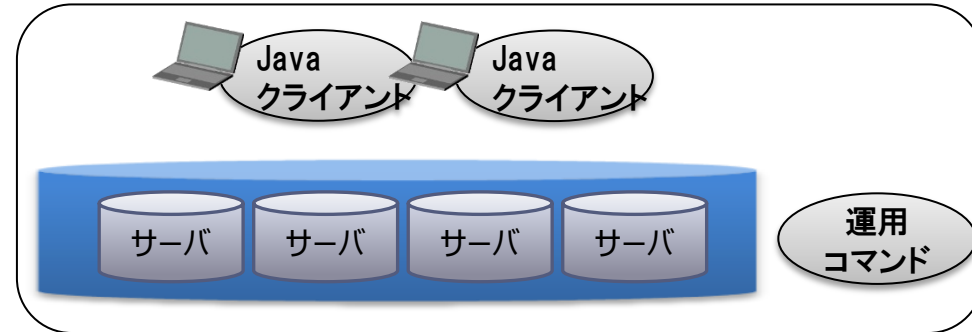
– ドキュメント (日/英) :

- スタートアップガイド
- 技術文書(クラスタ管理)
- APIリファレンス

● コミュニティ方法

– 質問、要望等は

GitHubのIssueを利用



README.md	first release	22 hours ago
README_ja.md	first release	22 hours ago
bootstrap.sh	first release	22 hours ago
configure.ac	first release	22 hours ago

Overview

GridDB has a KVS (Key-Value Store) that can be easily scaled.

- High Reliability
It is equipped with a replication mechanism to prevent node failure, and a failover mechanism to ensure high availability.
- High Performance
Even if the number of nodes increases, the performance is maintained around 10% of the original performance by lightening the data access.
- Advanced data management
In a traditional RDB, data management is done by functions greatly dependent on the RDB table name and Java API (Application Programming Interface).

Quick start

概要

GridDBは、時系列で蓄積されるセンサーデータに適したKVS(Key-Value Store)型のデータモデルを持ち、センサの数に応じて容易にスケールアウトすることができるデータベースです。

- 高信頼性
ノード間でキー・バリューデータの複製を互いのノードで持ち合う仕組みを備えており、万が一ノード障害が発生しても、他ノードの複製を使うことにより、数秒で自動フェイルオーバーが可能です。
- 高速性(インメモリ)
RDBではメモリを大容量化しても、バッファ管理等に大きなオーバーヘッドが発生するため本質的なデータ処理にCPUリソースの10%前後しか割当てられずCPUパワーを十分発揮できないことが知られています。GridDB/NoSQLは、大容量化されたメモリを前提に、バッファ処理の軽量化、リカバリ処理の軽量化、データ処理時のロックフリー化を行うことで、これまでのオーバーヘッドを最小化します。
- 高度なデータモデルと操作モデル
従来の分散KVSでは、Put/Get/Removeという操作によりデータを操作します。GridDB/NoSQLは、これを大幅に拡張し、構造化データの定義機能、SQLライクなクエリ機能、トランザクション機能、JavaのAPI(Application Programming Interface)をサポートしており、RDBユーザがスムーズに導入できるようになっています。キー・バリューをキー・コンテナと呼ぶレコードの集合体でデータを表現します。これはRDBのテーブル名とテーブルの関係に類似しています。また、センサーデータ管理向けの応用機能も備わっています。

クイックスタート

GridDBの機能概要

特徴	内容
データ構造	キーコンテナ型： コンテナ └レコード（ロウ） └カラム 型：文字列、数字、真偽値、小数、時刻、バイナリ、配列
クエリ・インデックス	クエリ：SELECT、集計、サンプリングなど インデックス：文字列、数字、真偽値、小数、時刻型のカラム
API	Java Native API、SQLライクな言語(TQL) HTTP(JSON)
水平分散	ハイブリッド型、自動負荷分散あり
レプリケーション	オーナ（オリジナル）、バックアップ（レプリカ）
その他	永続化（インメモリ／ディスク併用） コンテナ単位のトランザクション 時系列データ管理、期限解放 トリガ、など

目次

1. はじめに

- ビッグデータ
- NoSQL
- IoTと既存NoSQLの課題

2. GridDB

- オープンソース化
- 特長
- 適用事例
- 公開サイト

3. 利用方法

4. まとめ

起動、サンプル実行（1台構成）

● サーバの起動

```
$ export GS_HOME=$PWD
```

```
$ export GS_LOG=$PWD/log
```

```
$ bin/gpasswd admin
```

```
#input your_password
```

```
$ vi conf/gs_cluster.json
```

```
# "clusterName":"your_clustername" #<-- input your_clustername
```

```
$ bin/gstartnode
```

```
$ bin/gjoincluster -c your_clustername -u admin/your_password
```

● サンプルプログラムの実行

```
$ export CLASSPATH=${CLASSPATH}:$GS_HOME/bin/gridstore.jar
```

```
$ mkdir gsSample
```

```
$ cp $GS_HOME/docs/sample/program/Sample1.java gsSample/.
```

```
$ javac gsSample/Sample1.java
```

```
$ java gsSample/Sample1 239.0.0.1 31999 your_clustername
```

```
admin your_password
```

```
--> Person: name=name02 status=false count=2 lob=[65, 66, 67, 68,  
69, 70, 71, 72, 73, 74]
```

環境変数の設定

パスワードの設定

クラスタ名の設定

起動、クラスタ構成

ビルド

実行

ブロードキャスト用の
IPアドレス、ポートNo

ディレクトリ構成

```
bin/ //運用コマンド、モジュールディレクトリ
    gsserver // サーバ実行ファイル
    gridstore.jar // Javaクライアント
    gs_startnode // 運用コマンド
    ...
conf/ // 定義ファイルディレクトリ
    gs_node.json // ノード定義ファイル
    gs_cluster.json // クラスタ定義ファイル
    password // ユーザ定義ファイル
data/ // データベースファイルディレクトリ
    gs_cp_xxx // チェックポイントファイル
    gs_log_xxx.log // トランザクションログ(Redoログ)
log/ // イベントログ出力ディレクトリ
    gridstore_xxx.log // イベントログファイル
                                (トレース出力、エラー出力)
```

設定手順

(A) 初期設定

1. 環境変数(GS_HOME, GS_LOG)
2. ユーザアカウントの作成(gs_adduser, gs_passwd, gs_deluser)
3. クラスタ名の設定(conf/gs_cluster.jsonのclusterName)

(B) パラメータ設定(conf/gs_node.json)

1. メモリ上限値(dataStore/storeMemoryLimit)
2. 並列度(dataStore/concurrency)

(C) クラスタ構成時の設定

1. ノード数(gs_joinclusterの-nオプション)
2. ブロードキャスト用のIPアドレス、ポートNo
(conf/gs_cluster.jsonのtransaction/notificationAddress,
transaction/notificationPort)

運用コマンドの実行手順

1. 起動(各ノードについて)

- gs_startnode

2. クラスタ構成(各ノードについて。1台構成でも必要)

- gs_joincluster [-s 接続サーバ:ポート] -n 構成ノード数 -c クラスタ名
-u ユーザ名/パスワード

3. アプリ実行

4. クラスタ停止(クラスタを構成している1ノードに対して)

- gs_stopcluster [-s 接続サーバ:ポート] -u ユーザ名/パスワード

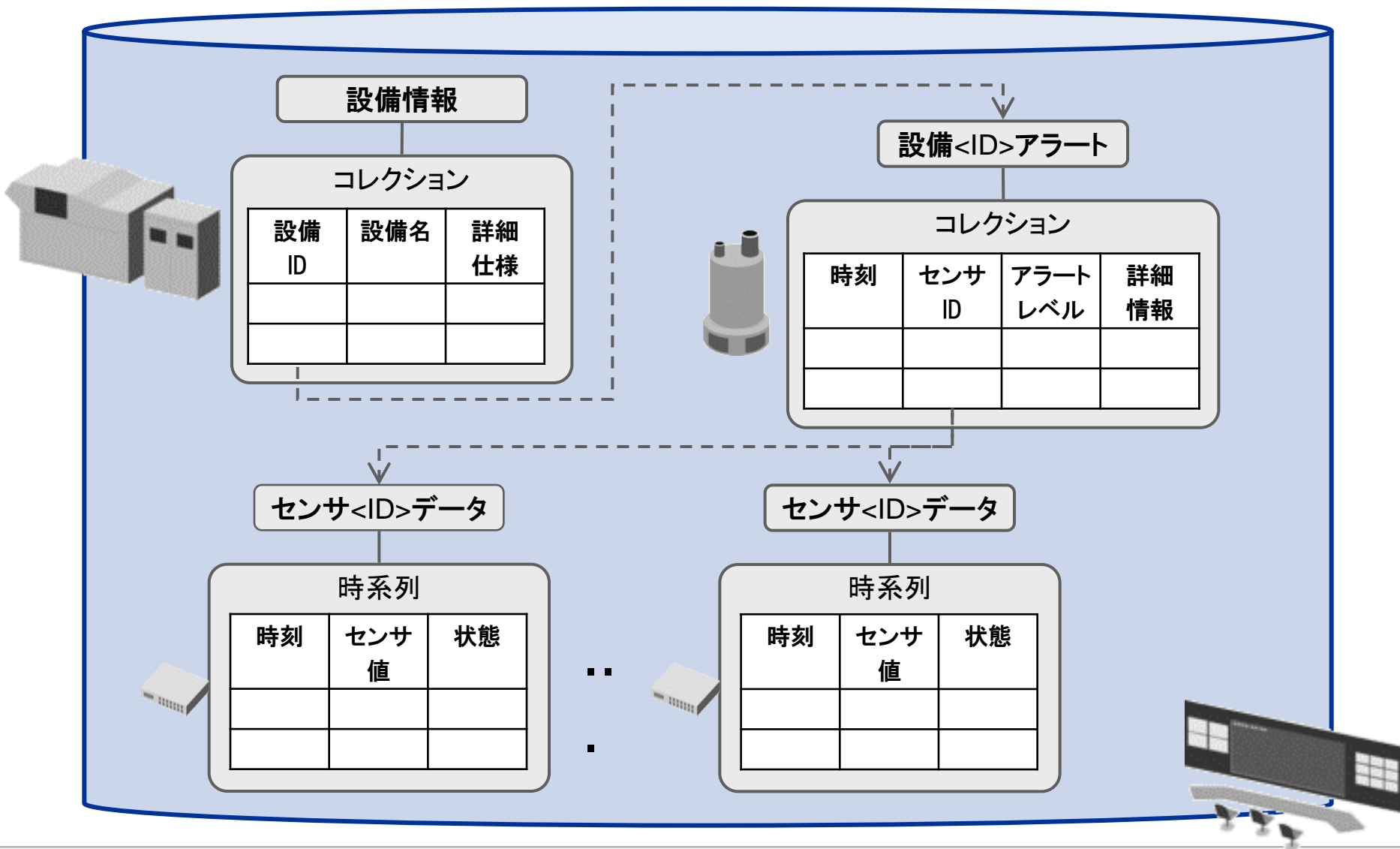
5. ノード停止(各ノードについて)

- gs_stopnode [-w [WAIT_TIME]][-s 接続サーバ:ポート] [-f]
-u ユーザ名/パスワード

※STAT取得(各ノードについて)

- gs_stat [-s 接続サーバ:ポート] -u ユーザ名/パスワード

プログラミングの例



レコードデータの登録

```
// アラート情報
class Alert {
    @RowKey long id;
    Date    timestamp;
    String  sensorId;
    int     level;
    String  detail;
}
```

```
final GridStore store = GridStoreFactory.getInstance().getGridStore(prop);
// コレクション登録
Collection<Long,Alert> alertCol
    = store.putCollection(alertColName, Alert.class);

// 索引設定
alertCol.createIndex("timestamp");
alertCol.createIndex("level");
alertCol.setAutoCommit(false);
...
Alert alert = new Alert();
while ((nextLine = reader.readNext()) != null) {
    ...
    alert.timestamp = new Date(datetime);
    alert.sensorId  = nextLine[2];
    alert.level     = Integer.valueOf(nextLine[3]);
    alert.detail    = nextLine[4];

    // 値登録
    alertCol.put(alert);
}
```


24時間以内に発生した重大アラートについて直前の時系列データを検索

```
Collection<String,Alert> alertCol
    = store.getCollection(alertColName, Alert.class);
...
// 24時間以内の重大アラート(level>3)を検索
String from = TimeStampUtil.format( TimeStampUtil.add(new Date(now), -24, TimeUnit.HOUR) );
Query<Alert> alertQuery
    = alertCol.query("select * where level > 3 and time > " + from );

RowSet<Alert> alertRs = alertQuery.fetch();
while( alertRs.hasNext() ) {
    // 重大アラート発生したセンサ
    Alert alert = alertRs.next();
    // 直前の時系列データを検索
    TimeSeries<Point> ts = store.getTimeSeries(alert.sensorId, Point.class);

    Date endDate = alert.timestamp;
    Date startDate = TimeStampUtil.add(alert.timestamp, -10, TimeUnit.MINUTE);
    RowSet<Point> rowSet = ts.query(startDate, endDate).fetch();

    while (rowSet.hasNext()) {
        Point ret = rowSet.next();
    }
}
```

目次

1. はじめに

- ビッグデータ
- NoSQL
- IoTと既存NoSQLの課題

2. GridDB

- オープンソース化
- 特長
- 適用事例
- 公開サイト

3. 利用方法

4. まとめ

まとめ

- **GridDBは高速処理と高信頼性を両立し、ペタバイト級の多種大量データを蓄積する、ビッグデータ/ I o T時代のデータベースです。**
 - High Performance
 - High Scalability
 - High Reliability

GridDBはオープンソースになったので、是非ともご覧ください、使ってみてください。

● 本資料に掲載の製品名、サービス名には、各社の登録商標または商標が含まれています。

TOSHIBA

Leading Innovation >>>