

# アーキテクチャを一新した IoT/ビッグデータ向けデータベースGridDB

**TOSHIBA**

東芝デジタルソリューションズ株式会社  
ソフトウェアシステム技術開発センター  
シニアフェロー 服部 雅一



# 自己紹介

- 東芝デジタルソリューションズ株式会社 シニアフェロー
- 入社以降、AIシステムやデータベースの研究開発に従事
- 2010年、(社)情報処理学会より「喜安記念業績賞」受賞
- 2012年、スケールアウト型DB “GridDB” の研究立ち上げ  
翌年、V1.0を上市。以降、チーフアーキテクトとして開発を主導
- 2020年、日本データベース学会よりGridDB事業化により  
弊社が「業績賞」受賞



# アジェンダ

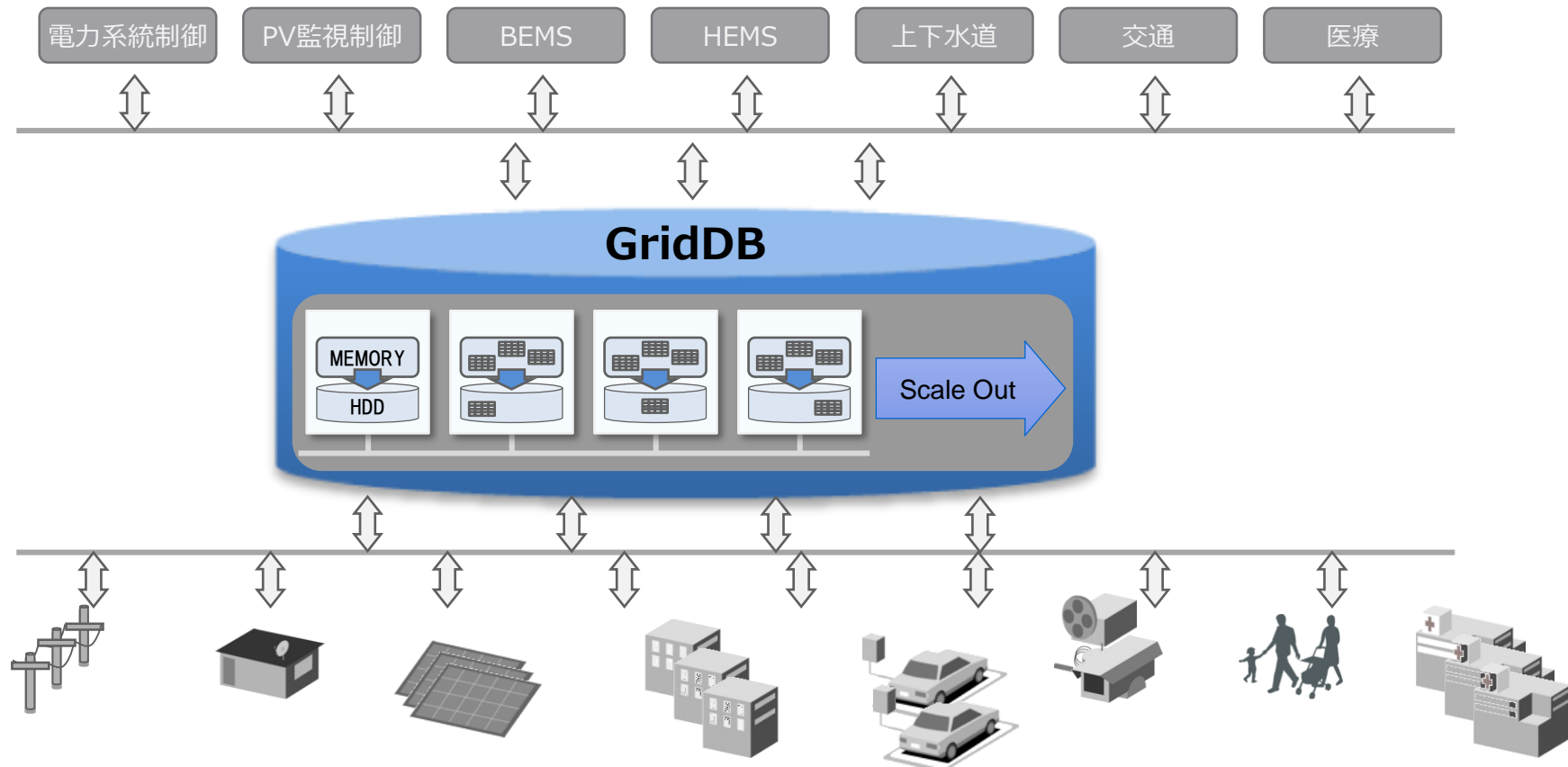
- スケールアウト型データベース GridDB
- SQLを補完するNoSQL
- NoSQLを超えたGridDB
- アーキテクチャを一新したV5



# ペタバイト級IoTデータを高速に処理するスケールアウト型データベース GridDB

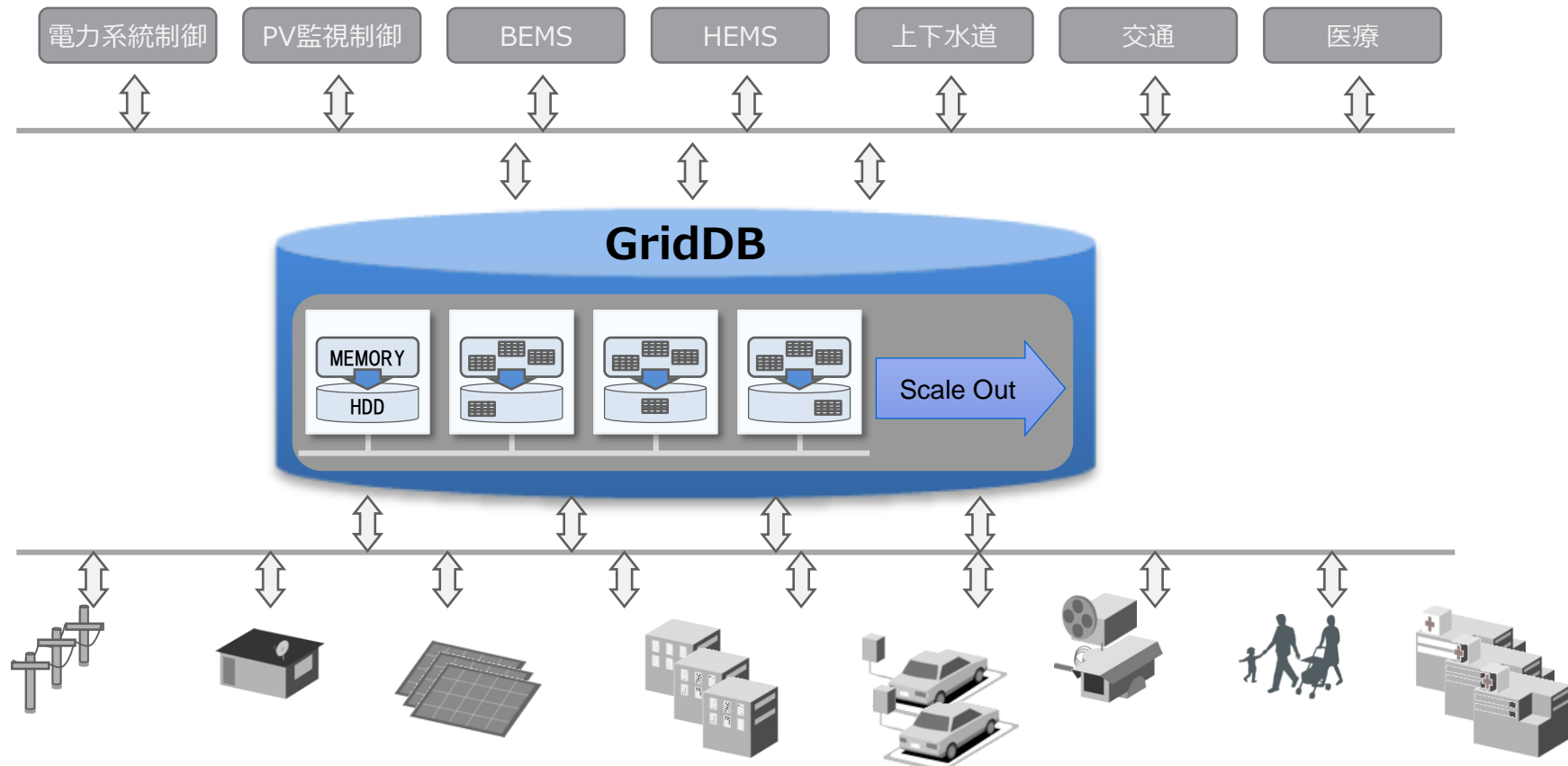
大規模IoT等、高い性能・可用性が求められるシステムに適用（現在バージョン4.6）

2013年、GridDB (NoSQLのみ)上市  
2015年、GridDB/NewSQL(SQLインターフェイス機能)  
2017年、SQLの並列分散処理化  
2019年、1ノードあたりペタバイト級のデータ管理に対応、など  
単体製品、DBaaSサービス、Meister DigitalTwinの基盤などを通してビジネス展開



# ペタバイト級IoTデータを高速に処理するスケールアウト型データベース GridDB

- 大規模IoT等、高い性能・可用性が求められるシステムに適用（現在バージョン4.6）
- 実適用データ規模 数TB～数PB





## GridDB Community Edition

高頻度・大量に発生する時系列データの蓄積とリアルタイムな活用をスムーズに実現する次世代の  
**オープンソースデータベース**



## GridDB Enterprise Edition

高頻度・大量に発生する時系列データの蓄積とリアルタイムな活用をスムーズに実現し、ビジネスを大きく成長させるために  
**最適化された次世代のデータベース**



【A23】

## GridDB Cloud

高頻度・大量に発生する時系列データの蓄積とリアルタイムな活用をスムーズに実現する  
**クラウドデータベースサービス**



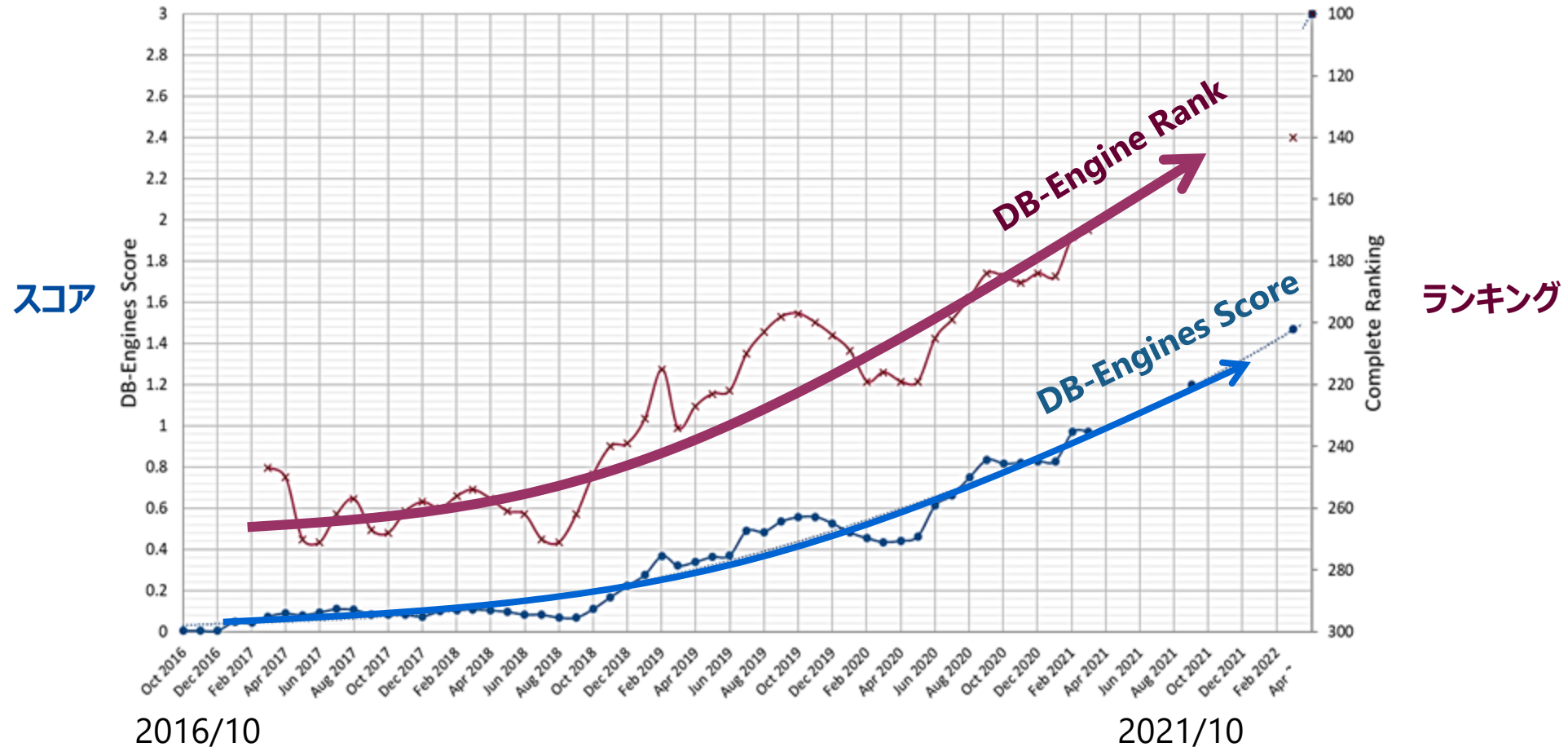
# 各エディションの違い

- インターフェイスはほぼ同じ
- クラスタ構成の有無の違い

項目	機能	Community Edition	Enterprise Edition	Cloud
	サポート		✓	✓
	プロフェッショナルサービス		✓	✓
データ管理	時系列コンテナ	✓	✓	✓
	コレクションコンテナ	✓	✓	✓
	索引	✓	✓	✓
	アフィニティ	✓	✓	✓
	テーブルパーティショニング	✓	✓	✓
クエリ言語	TQL	✓	✓	✓
	SQL	✓	✓	✓
NoSQLインタフェース	Java	✓	✓	✓
	C言語	✓	✓	✓
NewSQL(SQL) インタフェース	JDBC	✓	✓	✓
	ODBC		✓	✓
WebAPI		✓	✓	✓
時系列データ	時系列分析関数	✓	✓	✓
	期限付き解放機能		✓	✓
	長期アーカイブ		✓	✓
クラスタリング	機能クラスタ構成		✓	✓
	分散データ管理		✓	✓
	レプリケーション		✓	✓
運用管理	ローリングアップグレード		✓	
	オンラインバックアップ		✓	✓
	エクスポート / インポート		✓	✓
	運用管理GUI		✓	✓
セキュリティ	運用コマンド		✓	✓
	信暗号化 (TLS/SSL)		✓	✓
	認証機能 (LDAP)		✓	✓
オンプレミス環境	オンプレミス環境	✓	✓	
クラウドサービス	クラウドサービス			✓

# 日本初OSS GridDB Community Edition

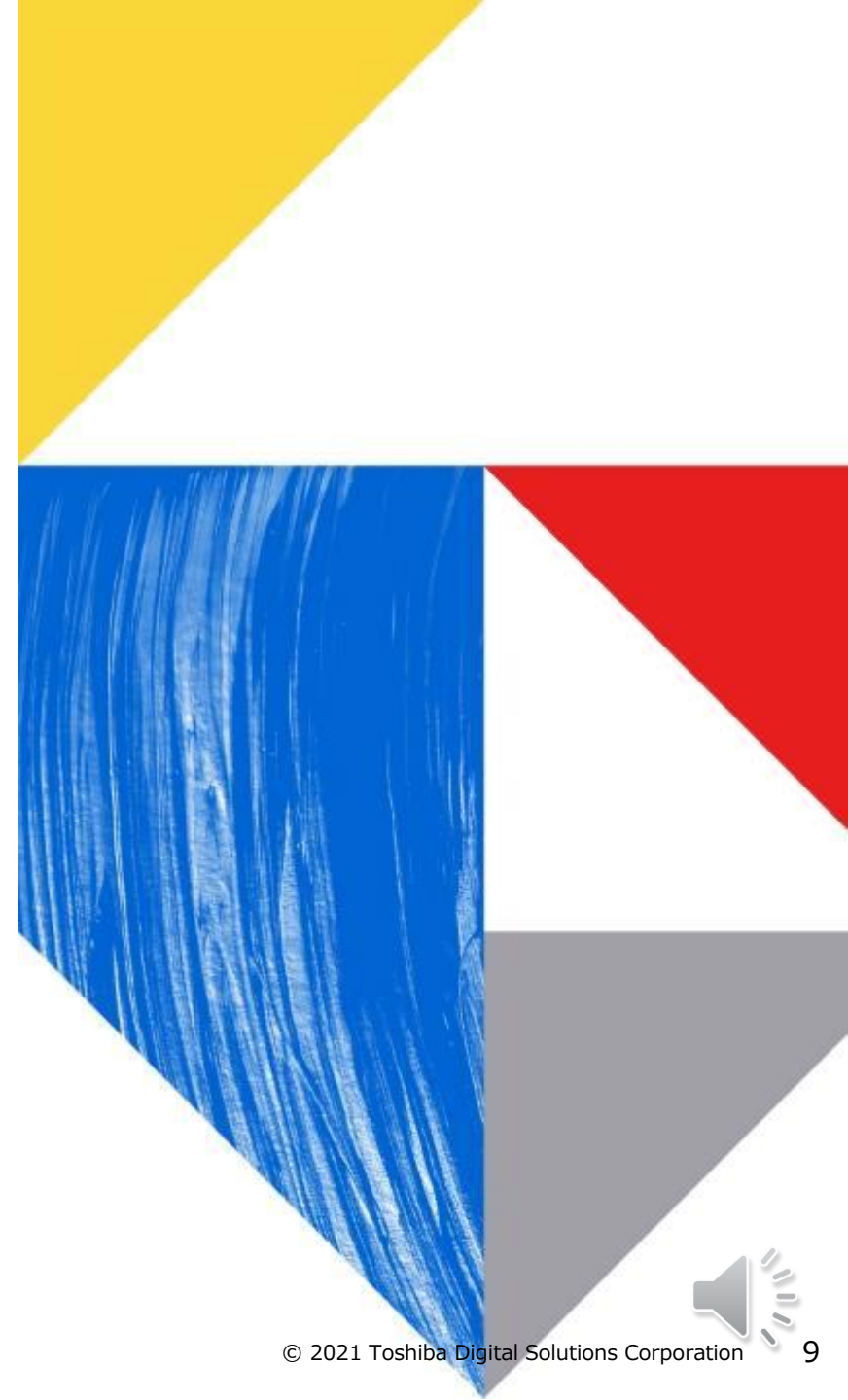
2016年公開。近年、DB-Engineランキングで急上昇





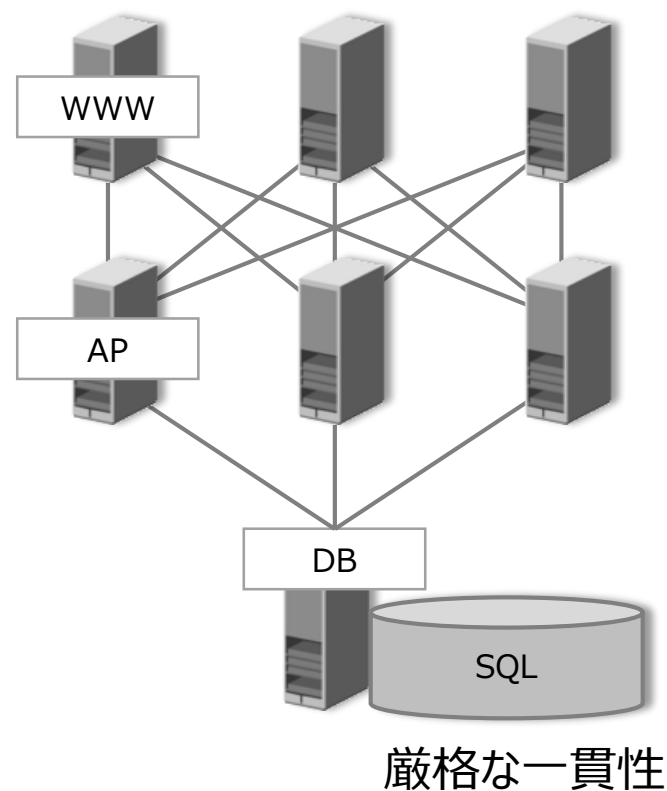
# 01

## SQLを補完するNoSQL



# NoSQL (Not Only SQL) の存在理由

スケールアップは物理的な限界やコスト高となる

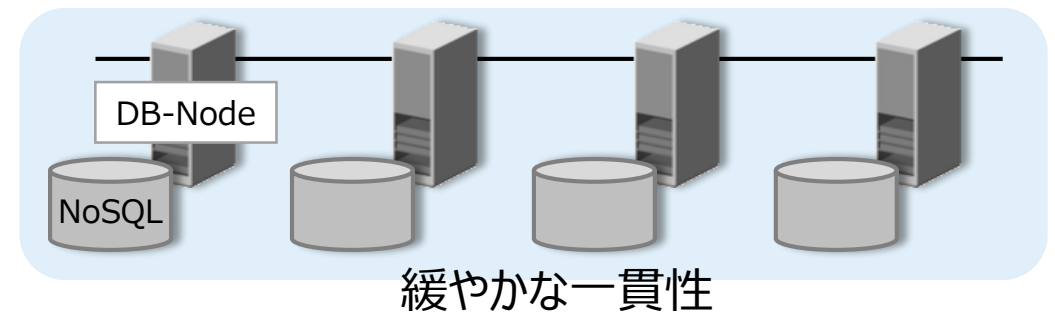
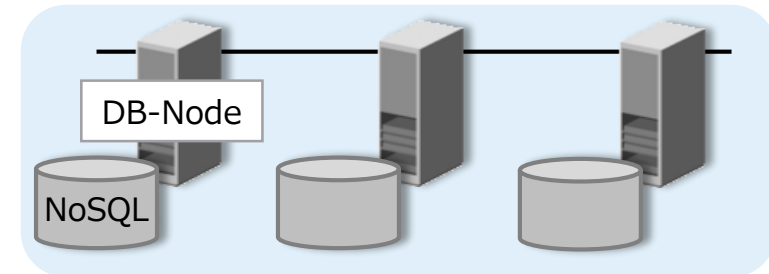
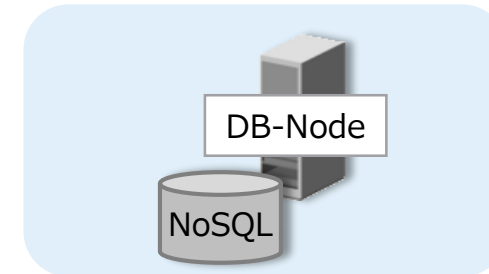
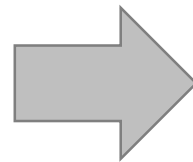
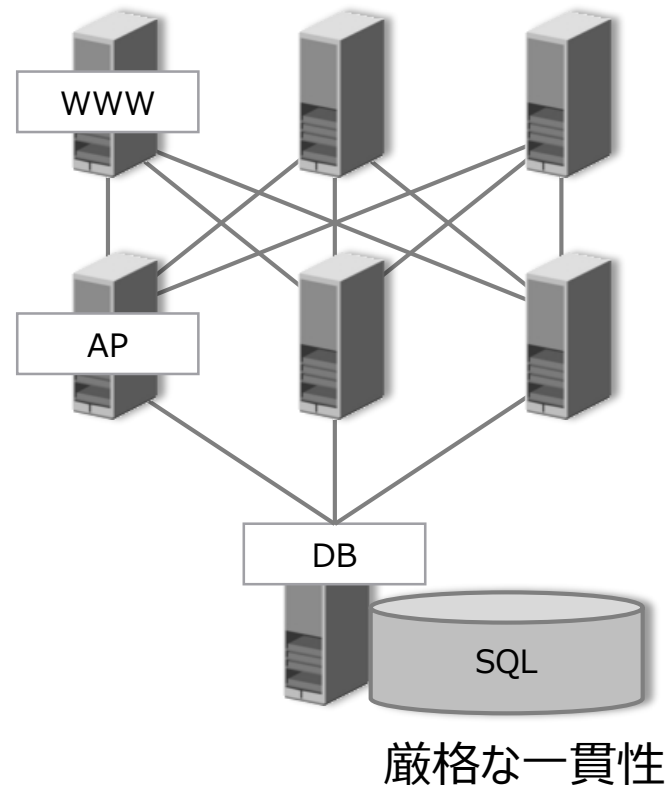


Webノード、APノードと異なり、  
DBのスケールアウトは非常に難しい。



# NoSQL (Not Only SQL) の存在理由

- ノード台数増で対応できるスケールアウトがお得
- その代わりに、データ一貫性の要件を緩和する必要



緩やかな一貫性

# NoSQLとSQLの関係

SQLとNoSQLは補完関係にある

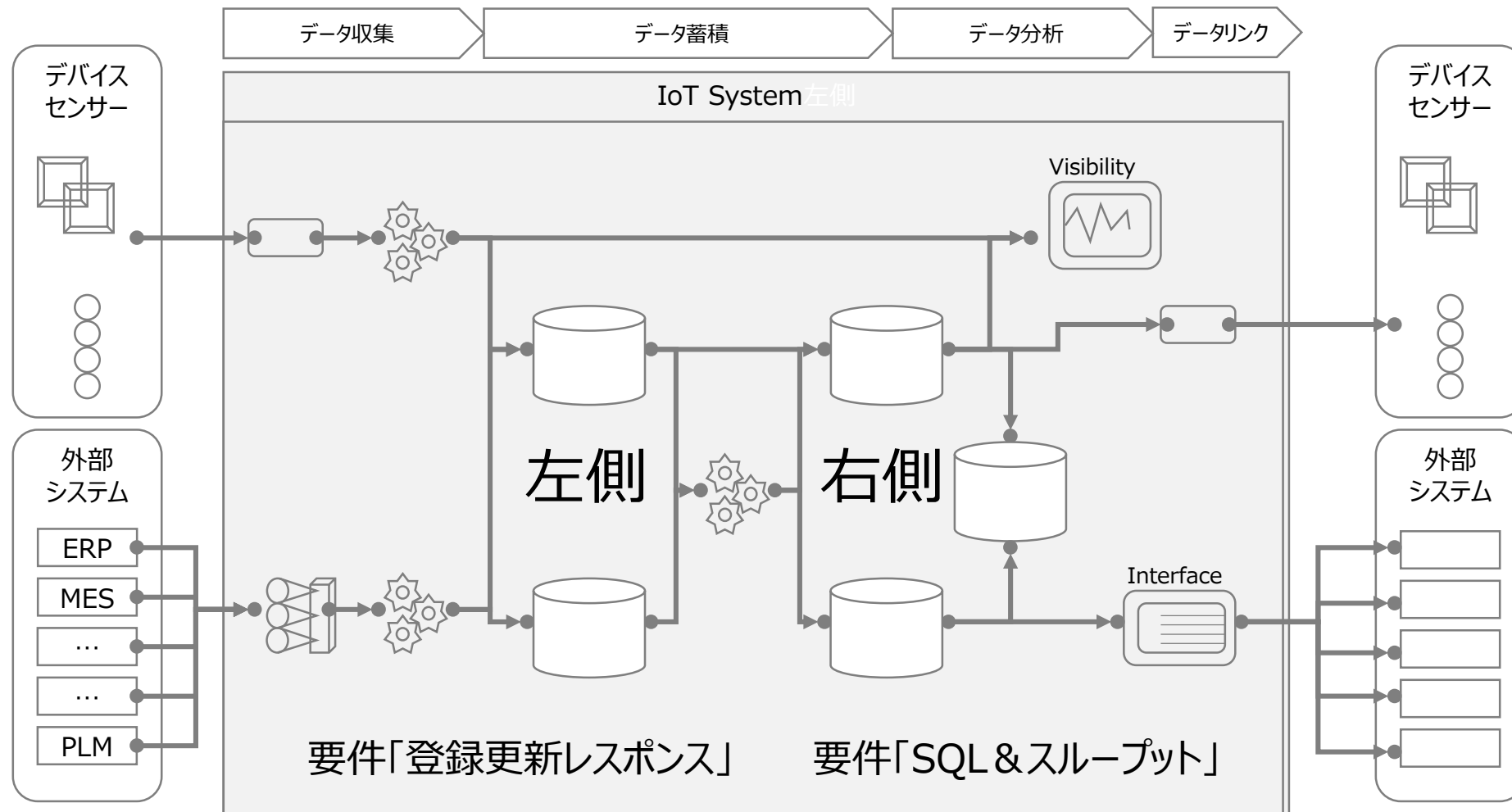
	某NoSQL	某NoSQL	SQL(RDB)
問い合わせ言語	独自I/F	独自I/F	SQL
スキーマ言語	無い	無い	SQL(DDL)
データ一貫性	参照一貫性 テーブル内一貫性	緩やかな一貫性 or 参照一貫性	厳格な一貫性
1台あたりの性能	普通	一貫性を上げると低速	実績あるRDBは ある程度
高可用性	ノード分散による 冗長性有り	ノード分散による 冗長性有り	別の仕掛けが必要
拡張性	スケールアウトする手段が不明	スケールアウト困難 特にオンラインスケールアウト	スケールアップ

# NoSQLをベースとしたスケールアウト型DB

SQLとNoSQLは補完関係にある

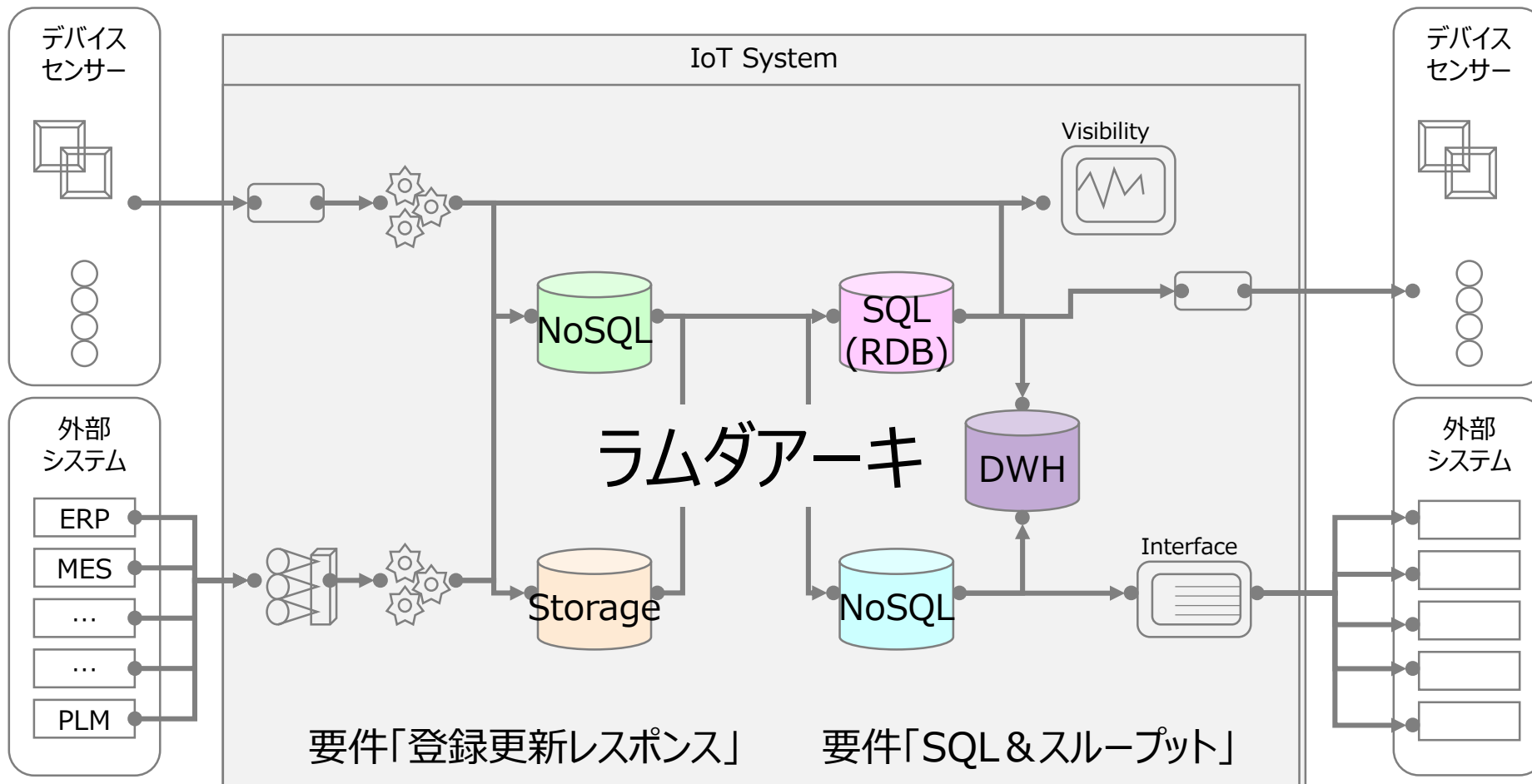
	某NoSQL	某NoSQL	SQL(RDB)
問い合わせ言語	独自I/F	独自I/F	SQL
スキーマ言語	無い	無い	SQL(DDL)
データ一貫性	参照一貫性 テーブル内一貫性	緩やかな一貫性 or 参照一貫性	厳格な一貫性
1台あたりの性能	普通	一貫性を上げると低速	実績あるRDBは ある程度
高可用性	ノード分散による 冗長性有り	ノード分散による 冗長性有り	別の仕掛けが必要
拡張性	スケールアウトする手段が不明	スケールアウト困難 特にオンラインスケールアウト	スケールアップ

# IoTシステムの内部構成

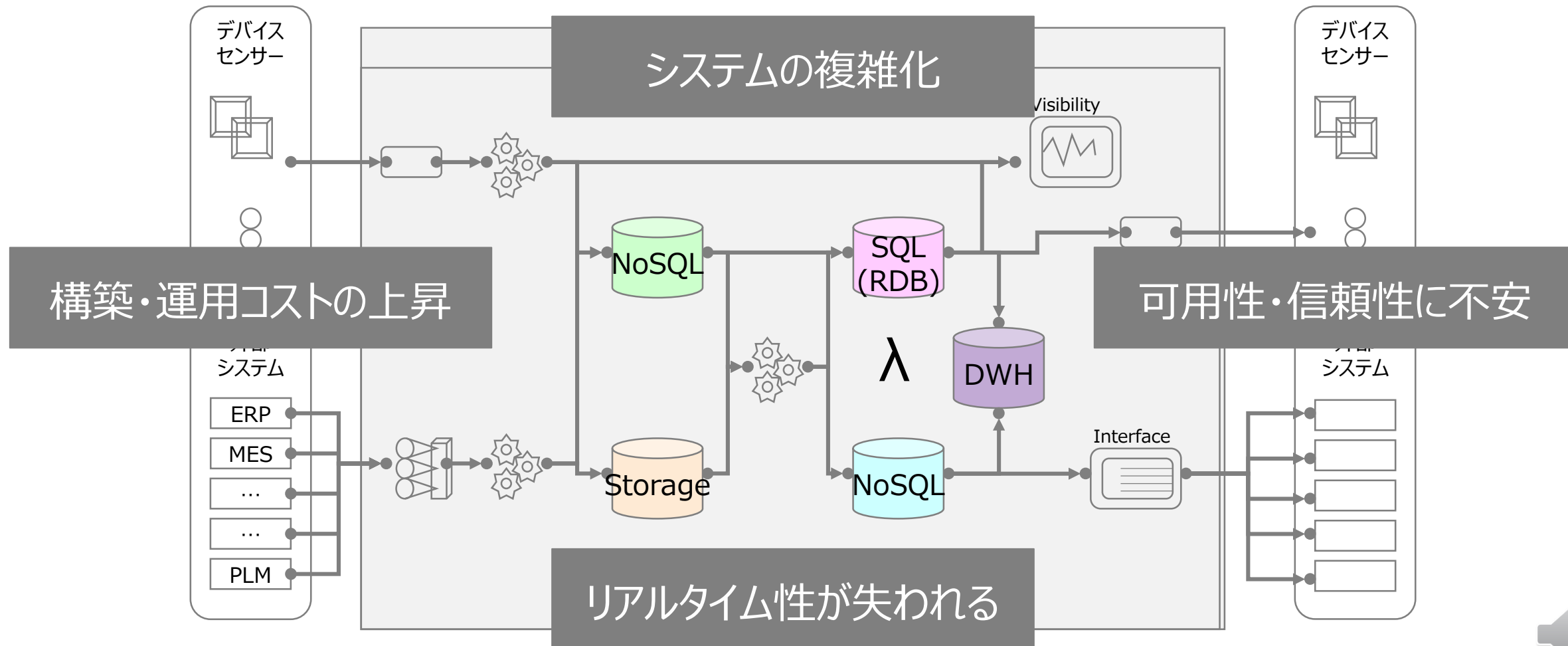


# IoTシステムにおけるDB

## 異なる特性を持つ複数のDBの使い分け

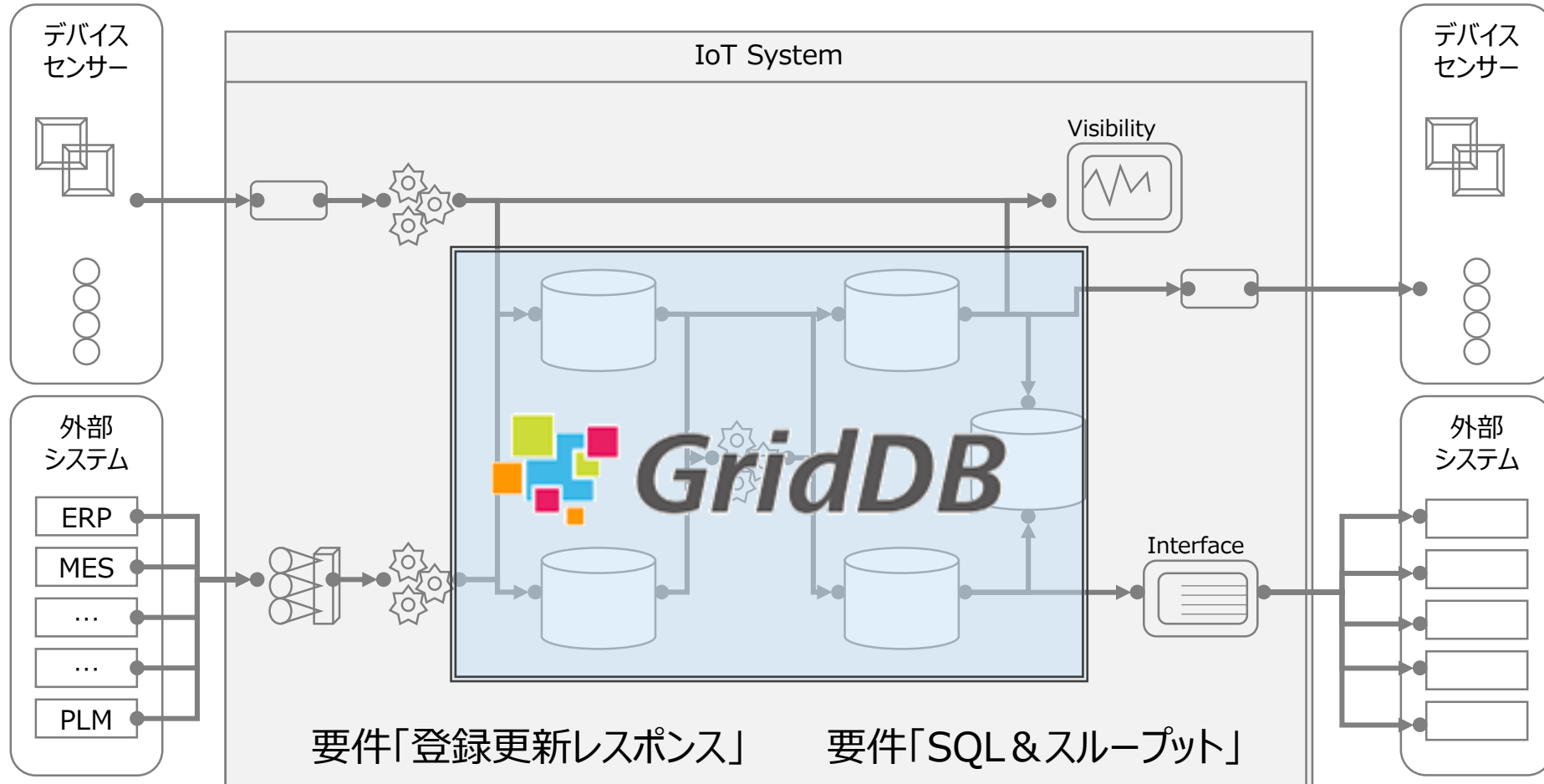


DBの使い分けがなぜ問題なのか？





GridDBならば、複数のDBを使わなくても良い。

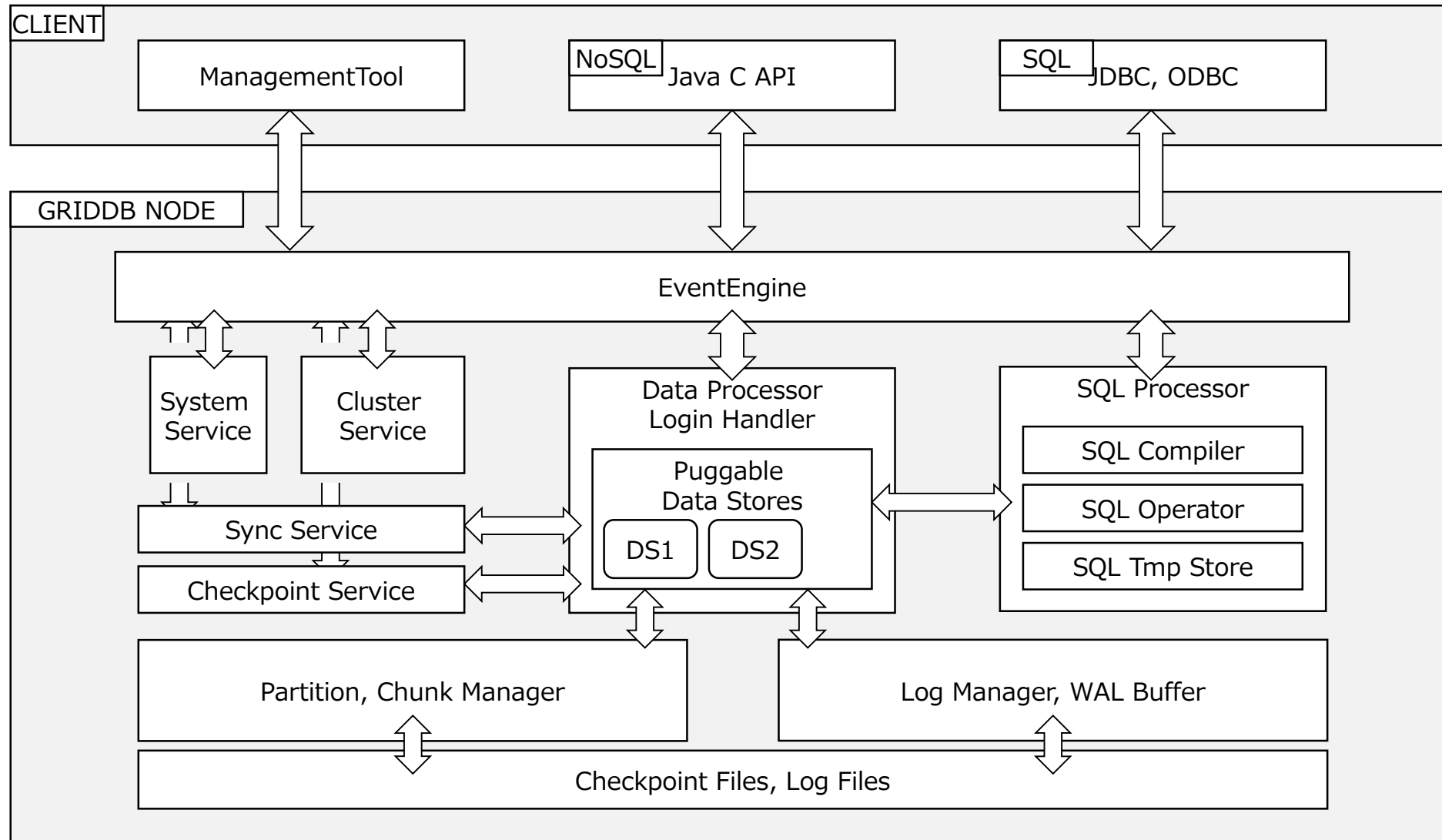


# 02

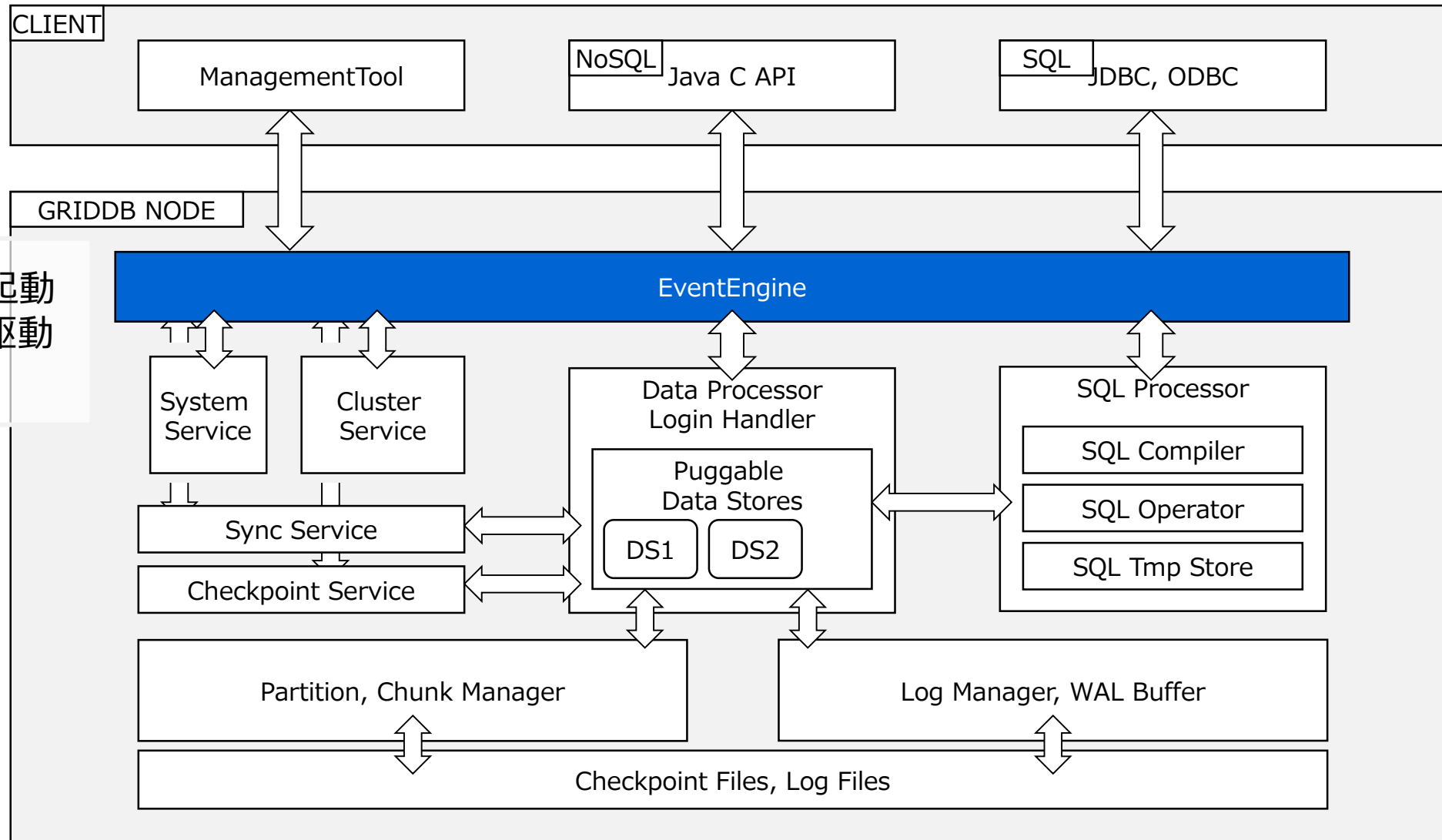
## NoSQLを超えたGridDB



# V5アーキテクチャ



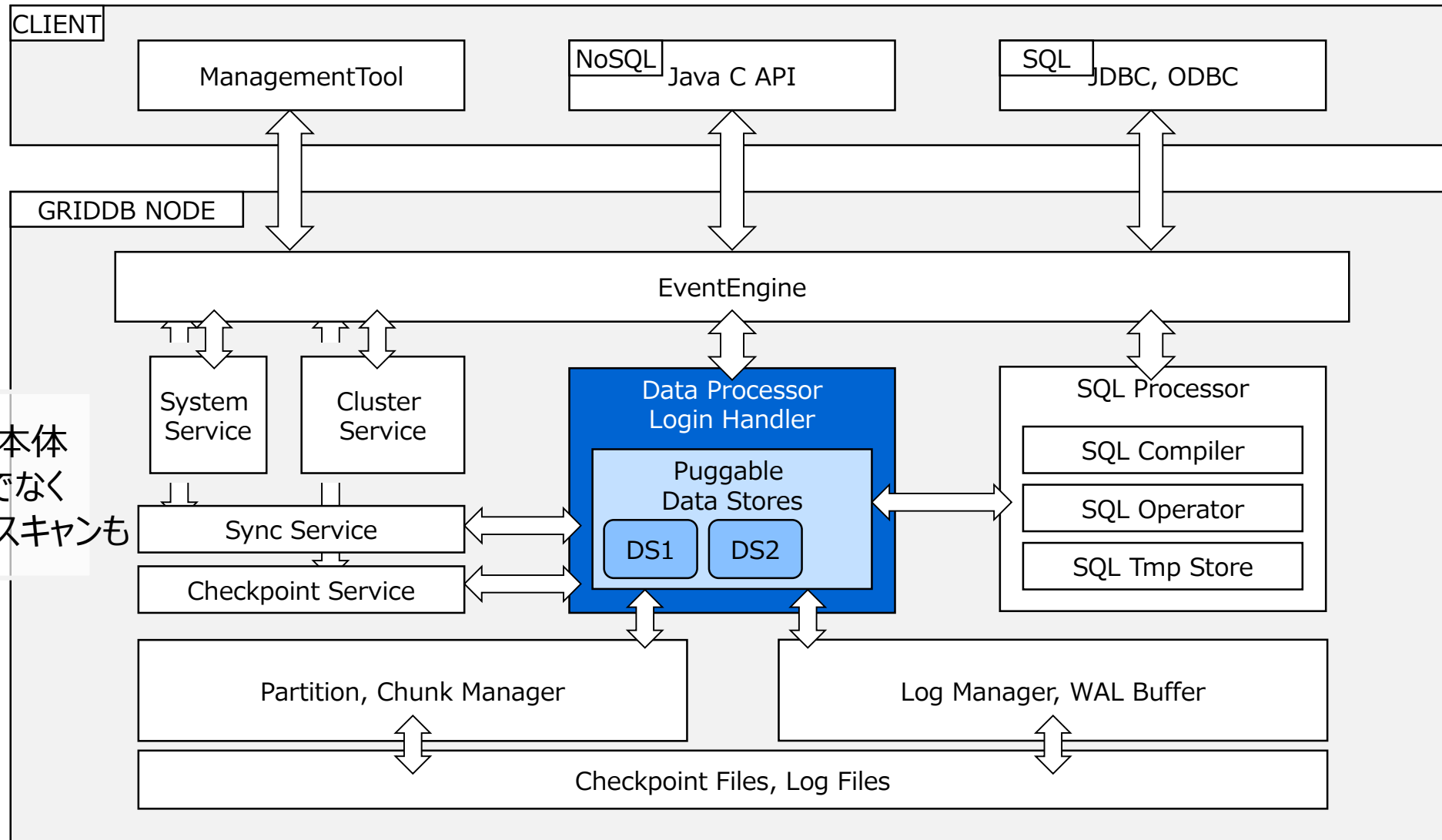
# V5アーキテクチャ



各種処理の起動  
全てイベント駆動  
一種のOS

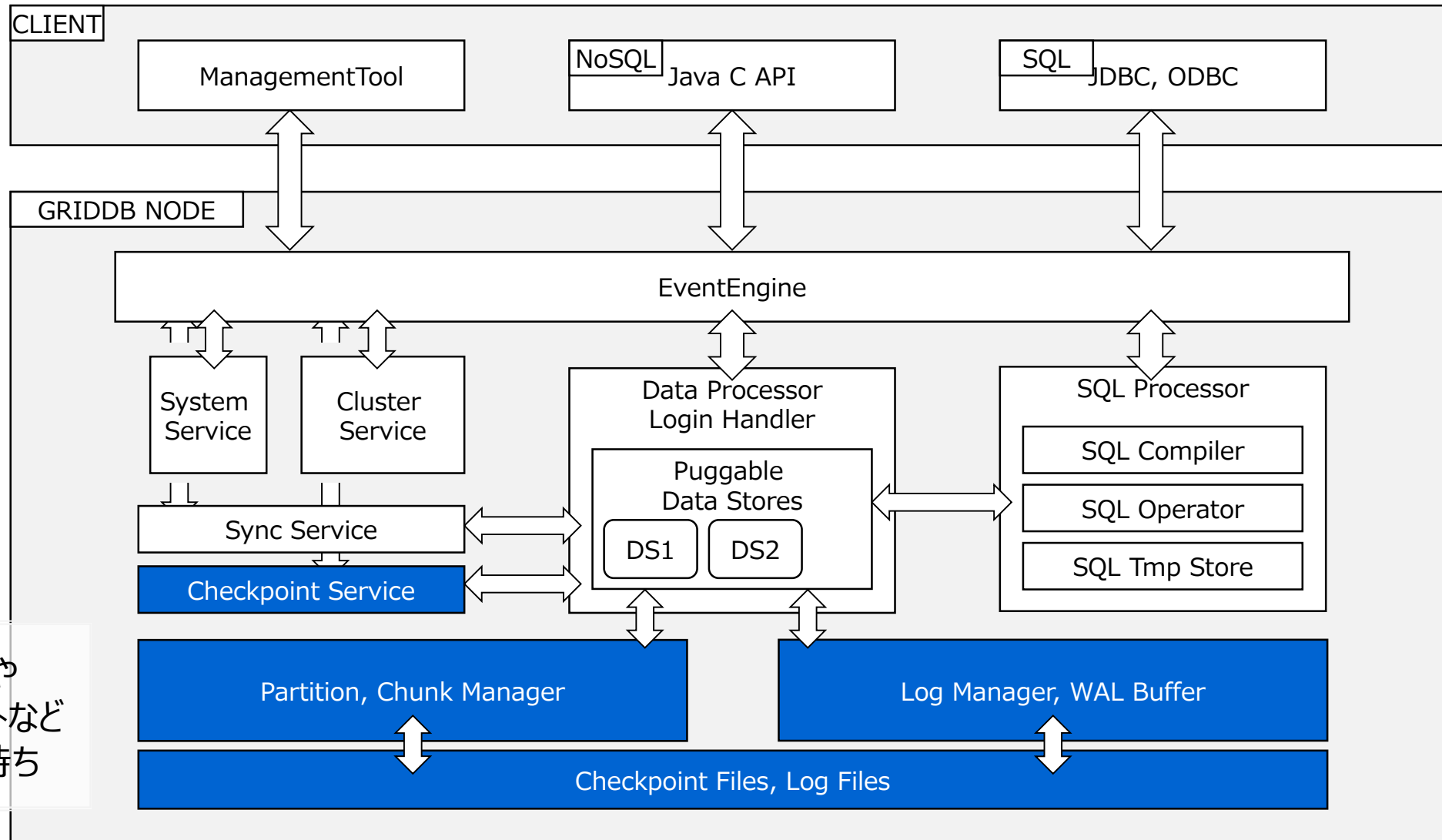


# V5アーキテクチャ



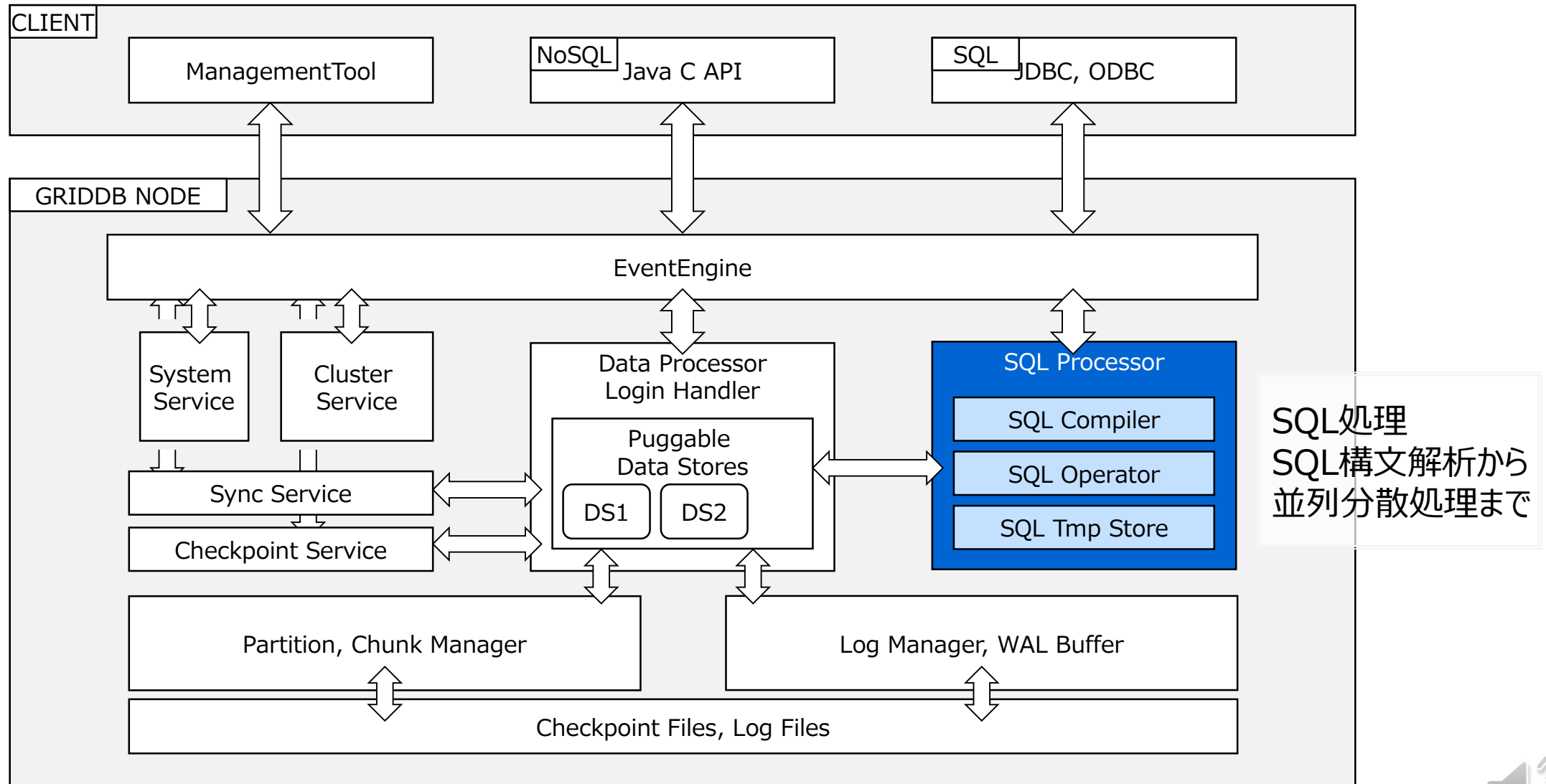
データ処理の本体  
NoSQLだけでなく  
SQLのデータスキャンも

# V5アーキテクチャ

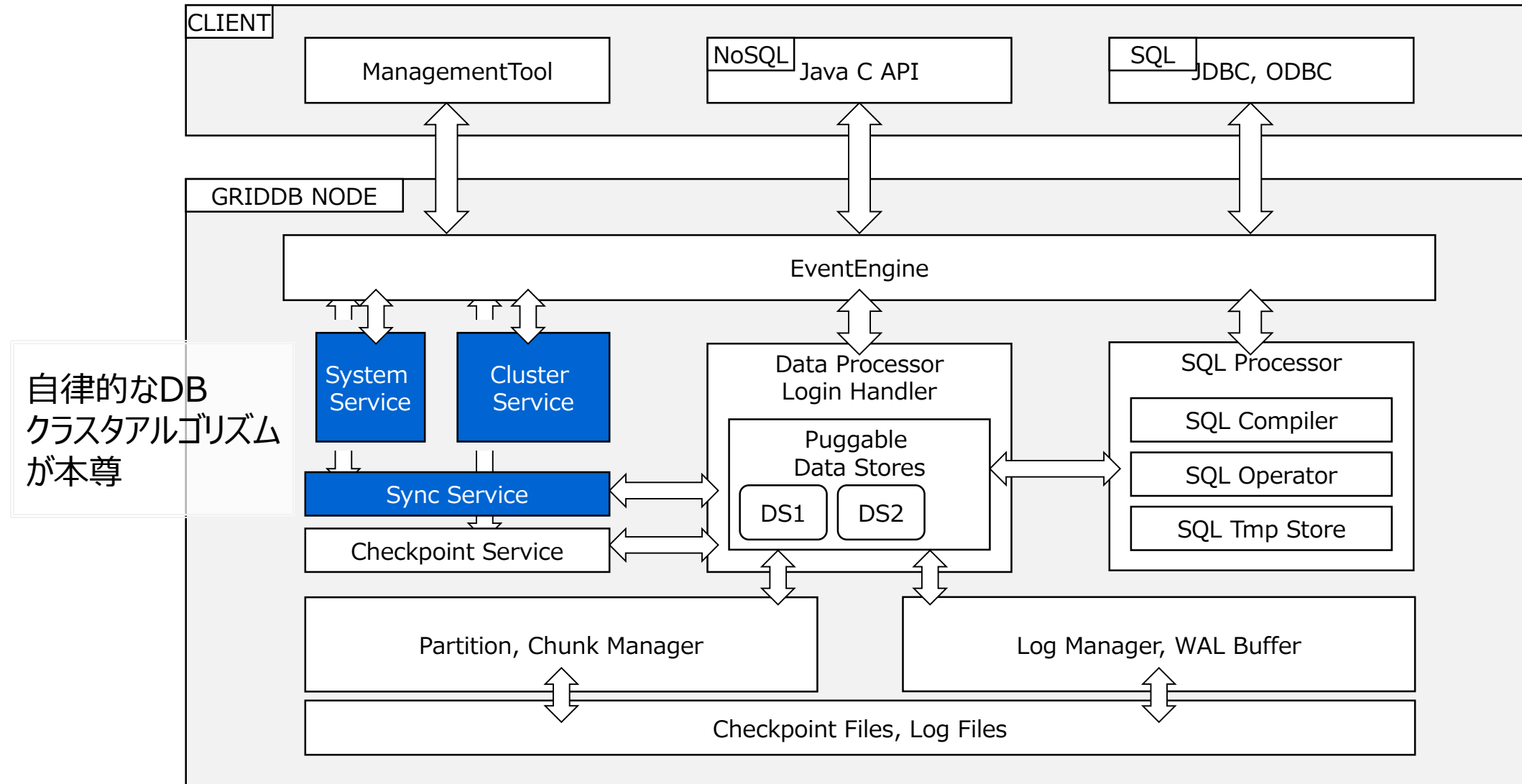


WALリカバリや  
スナップショットなど  
縁の下の力持ち

# V5アーキテクチャ

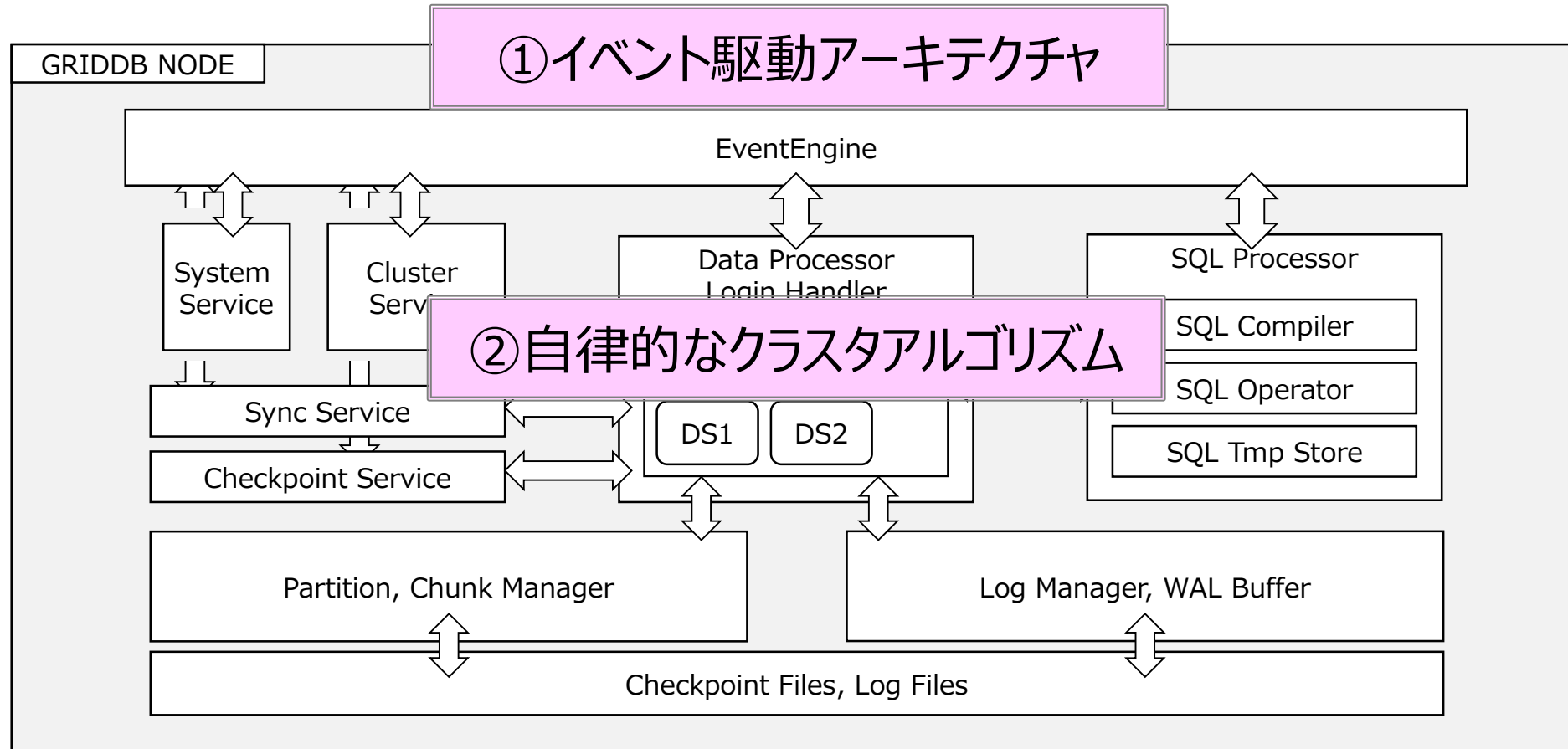


# V5アーキテクチャ



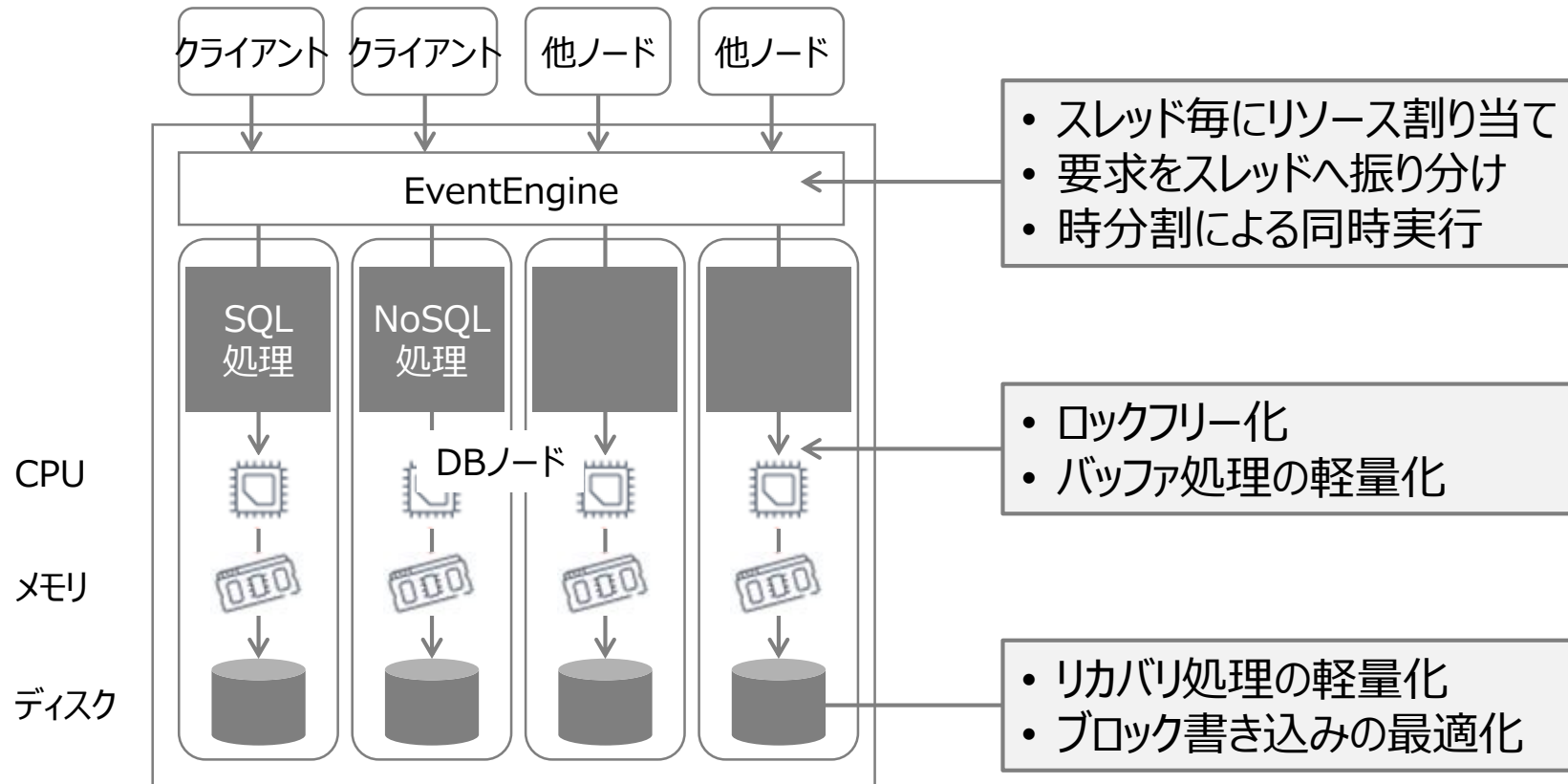


これまで培ってきた2つのコア技術



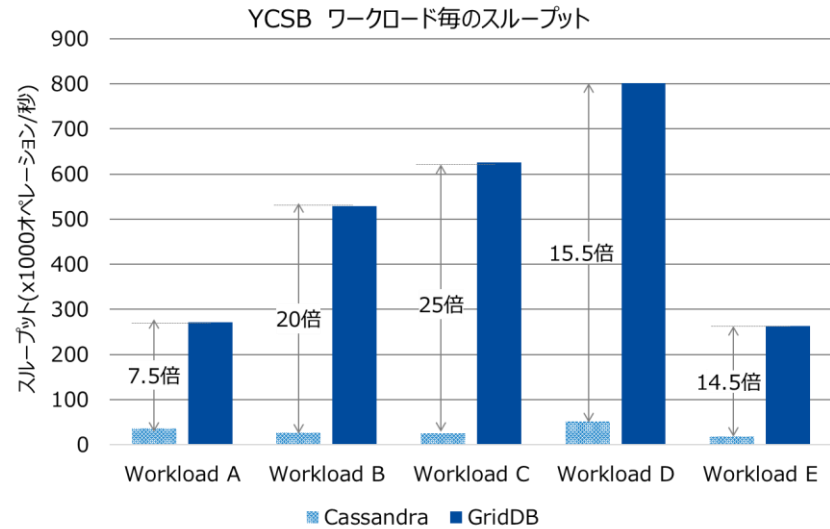
# ① イベント駆動アーキテクチャ

- CPUのマルチコア, メニーコア化を前提
- 非同期的なデータ処理を絶え間なく実行するイベント駆動方式

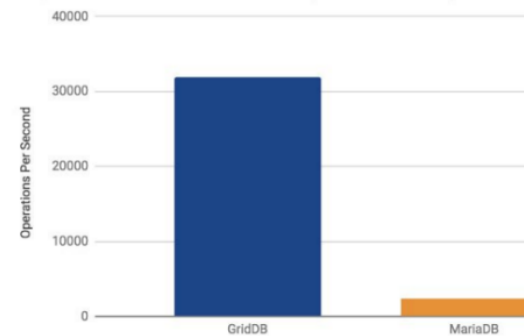


# ベンチマーク

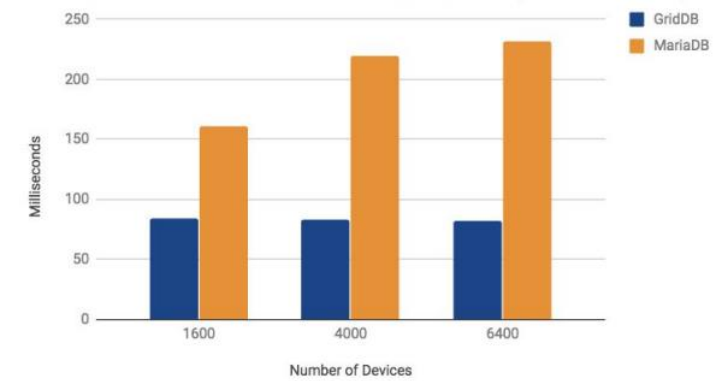
## ノードあたりの性能、クラスタの性能。各種ベンチで他DBを圧倒



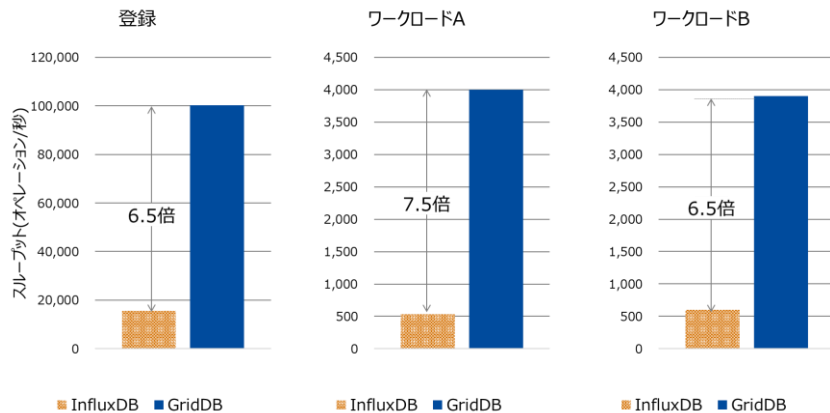
Ingest Operations Per Second (More is Better)



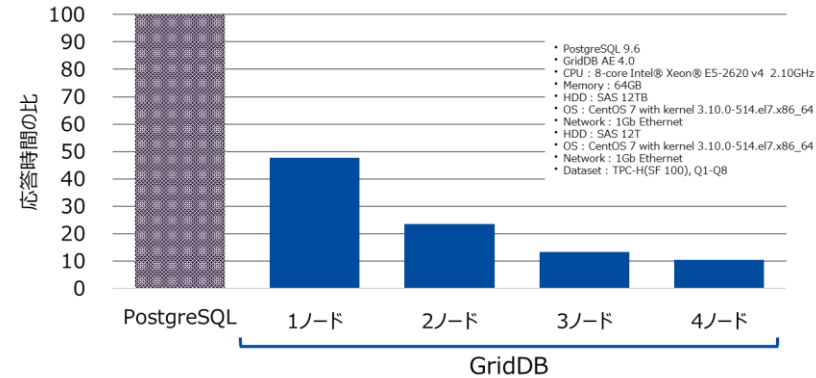
Average Milliseconds Per Billing Aggregation (Fewer is better)



[https://griddb.net/en/docs/Benchmarking\\_Application\\_GridDB\\_MariaDB.pdf](https://griddb.net/en/docs/Benchmarking_Application_GridDB_MariaDB.pdf)

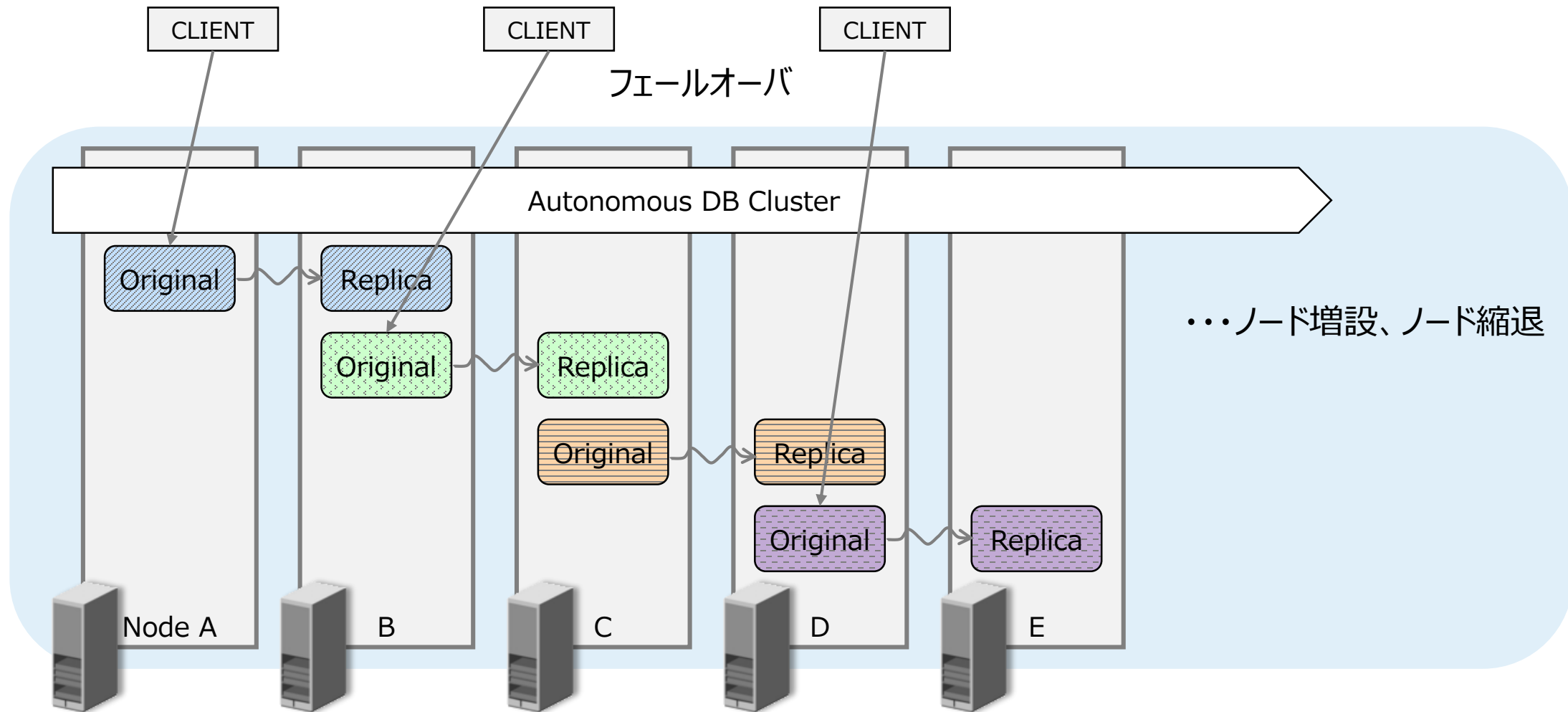


TPC-H(SF100) 応答時間



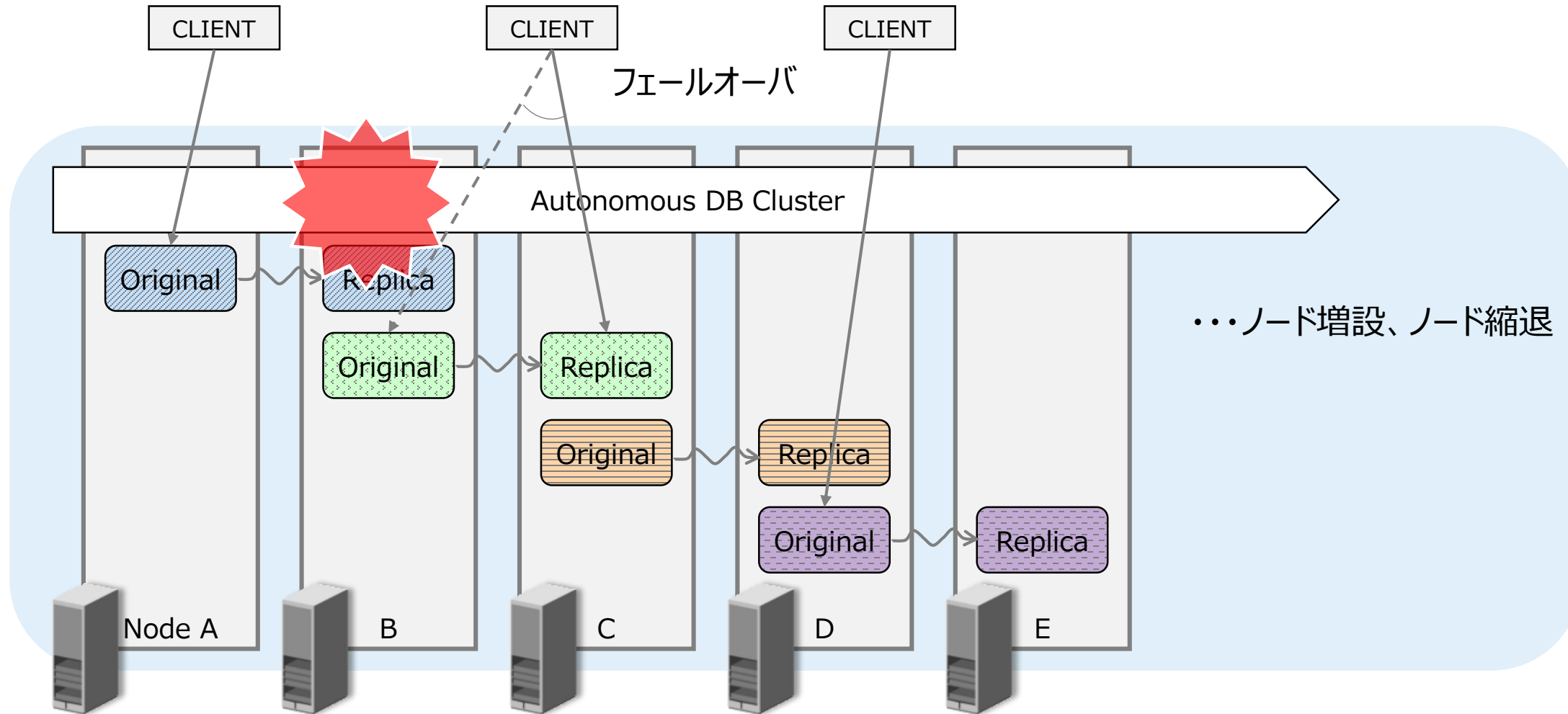
## ② 自律的なクラスタアルゴリズム

シャーディングとレプリケーションを制御する自律的なクラスタ

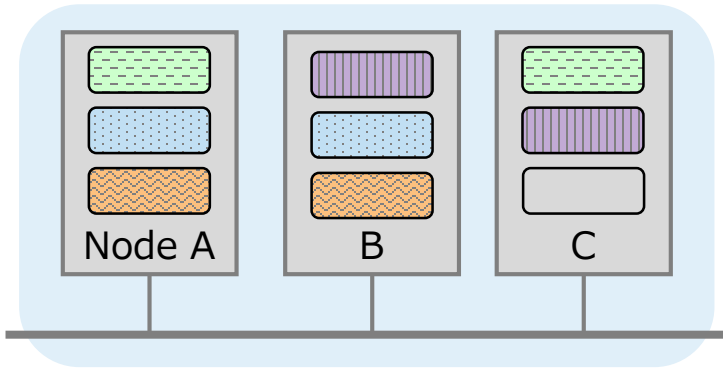


## ② 自律的なクラスタアルゴリズム

シャーディングとレプリケーションを制御する自律的なクラスタ



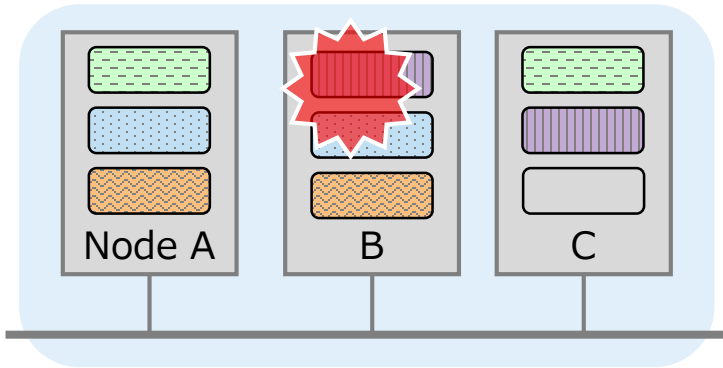
## ② 自律的なクラスタアルゴリズム



3ノード構成  
データ複製度：2



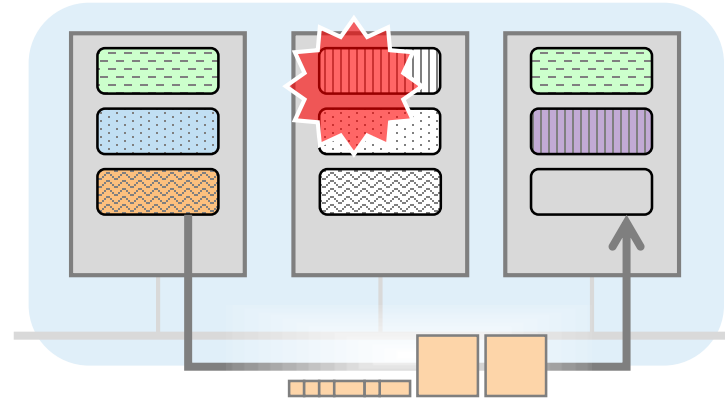
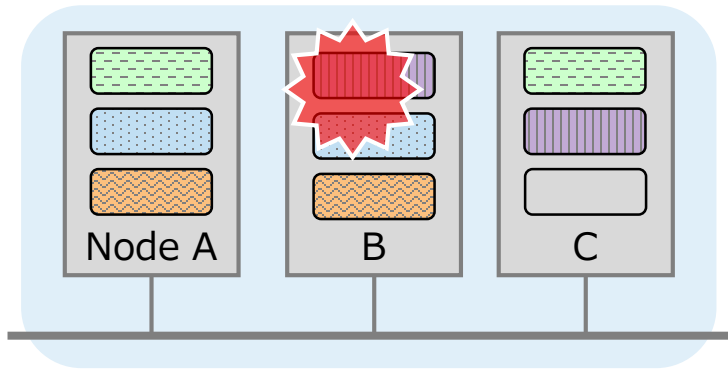
## ② 自律的なクラスタアルゴリズム



ノード障害発生！  
データ複製が減った！可用性が落ちる！



## ② 自律的なクラスタアルゴリズム

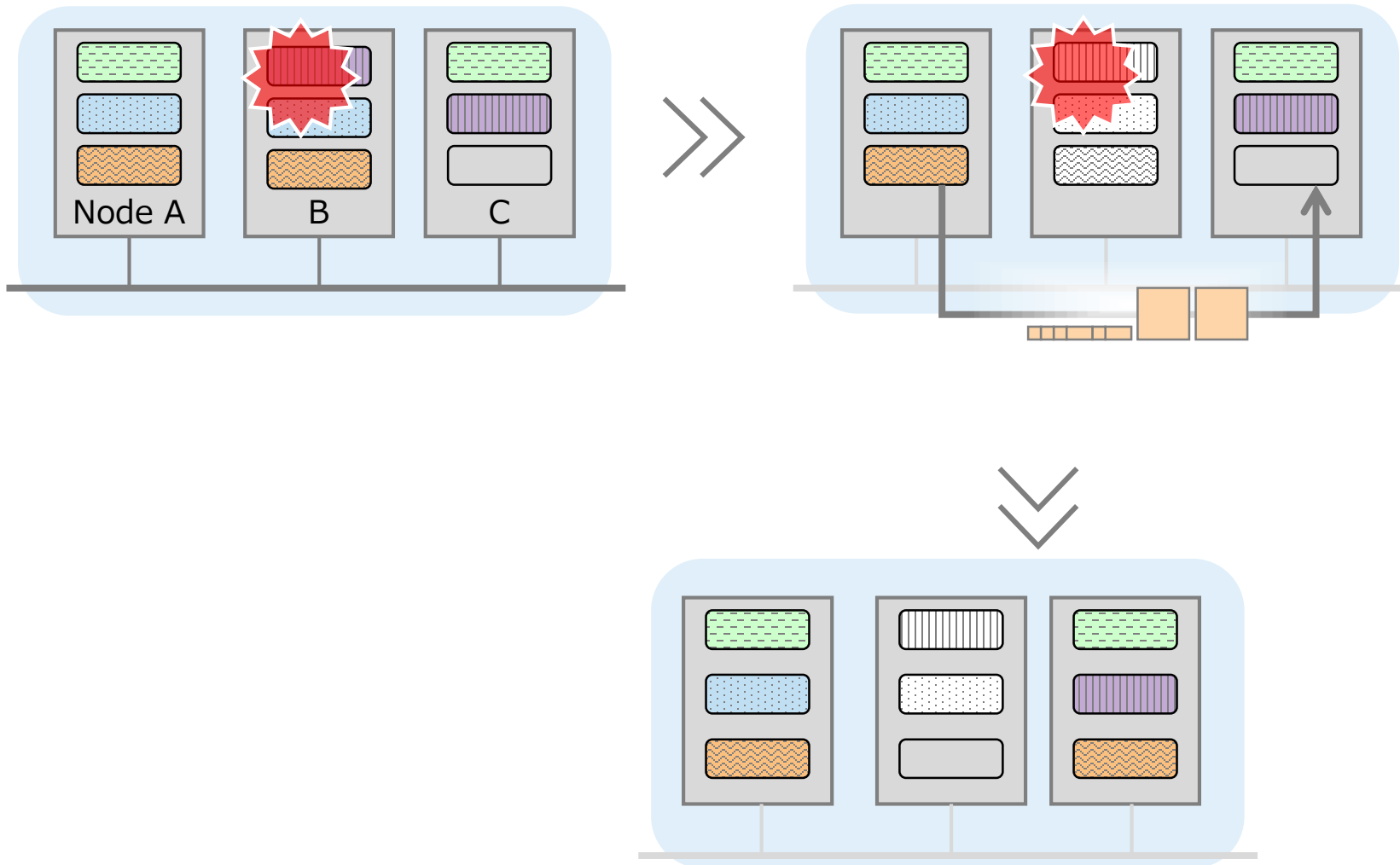


データの複製を作り始める。  
WALログ転送だけでは時間がかかる。  
メモリイメージとWALログの合わせ技で高速化する。





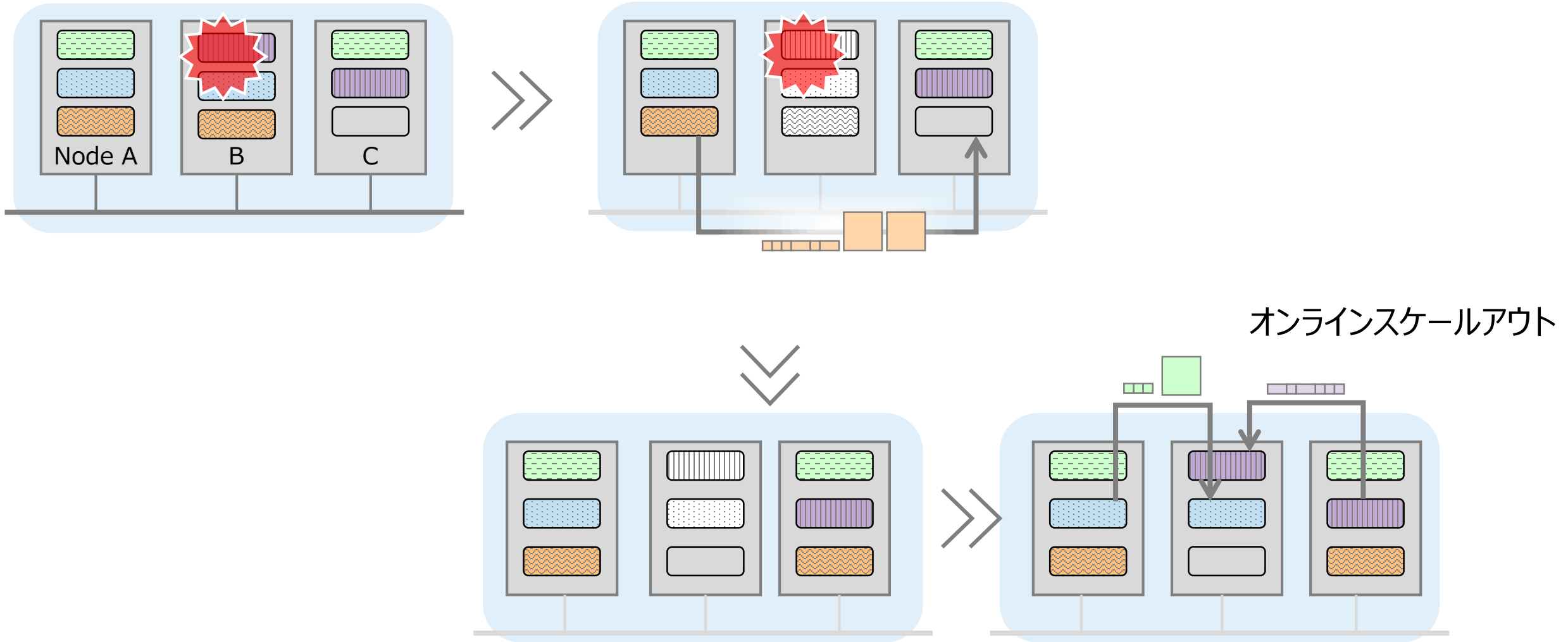
## ② 自律的なクラスタアルゴリズム



データ複製の成功  
クライアントのアクセスを切り替えさせる。




## ② 自律的なクラスタアルゴリズム



この動作の繰り返し。  
可用性を維持できた！


# NoSQLをベースとしたスケールアウト型DB

テーブル間一貫性を除けば、全てをカバーする

	 GridDB	NoSQL	SQL
問い合わせ言語	SQL + 独自I/F	独自I/F	SQL
スキーマ言語	SQL(DDL)	無い	SQL(DDL)
データ一貫性	参照一貫性 テーブル内一貫性	BASE理論 緩やかな一貫性	厳格な一貫性
1台あたりの性能	イベント駆動による 最大限の性能	SW次第	実績あるRDBは ある程度
高可用性	通常、高速の 2レプリケーション	ノード分散による 冗長性有り	別の仕掛けが必要
拡張性	自律的な スケールアウト	スケールアウト SW次第	スケールアップ

# NoSQLをベースとしたスケールアウト型DB

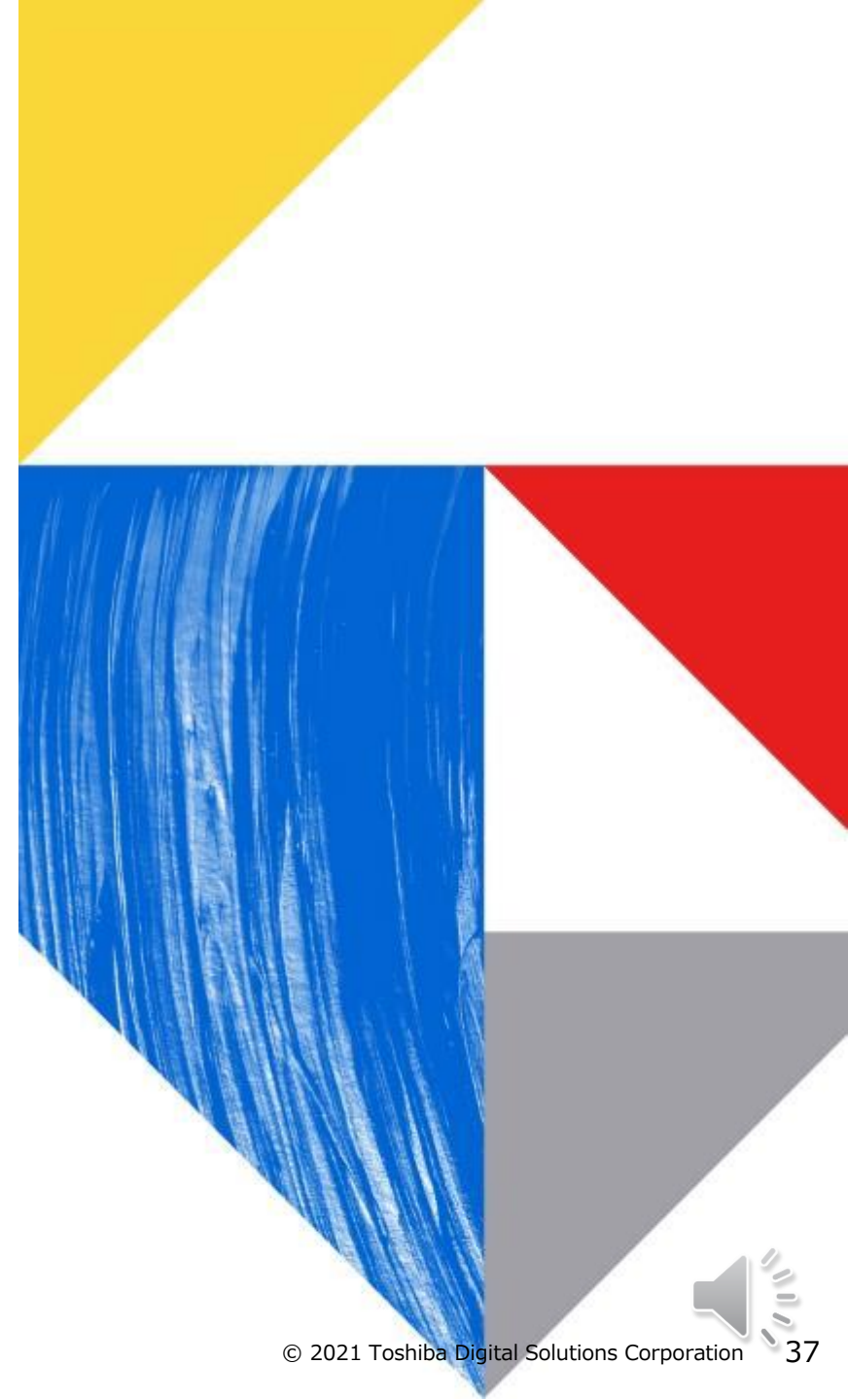
テーブル間一貫性を除けば、全てをカバーする

	 GridDB	NoSQL	SQL
問い合わせ言語	SQL + 独自I/F	独自I/F	SQL
スキーマ言語	SQL(DDL)	無い	SQL(DDL)
データ一貫性	参照一貫性 テーブル内一貫性	BASE理論 緩やかな一貫性	厳格な一貫性
1台あたりの性能	イベント駆動による 最大限の性能	SW次第	実績あるRDBは ある程度
高可用性	通常、高速の 2レプリケーション	ノード分散による 冗長性有り	別の仕掛けが必要
拡張性	自律的な スケールアウト	スケールアウト SW次第	スケールアップ

# 03

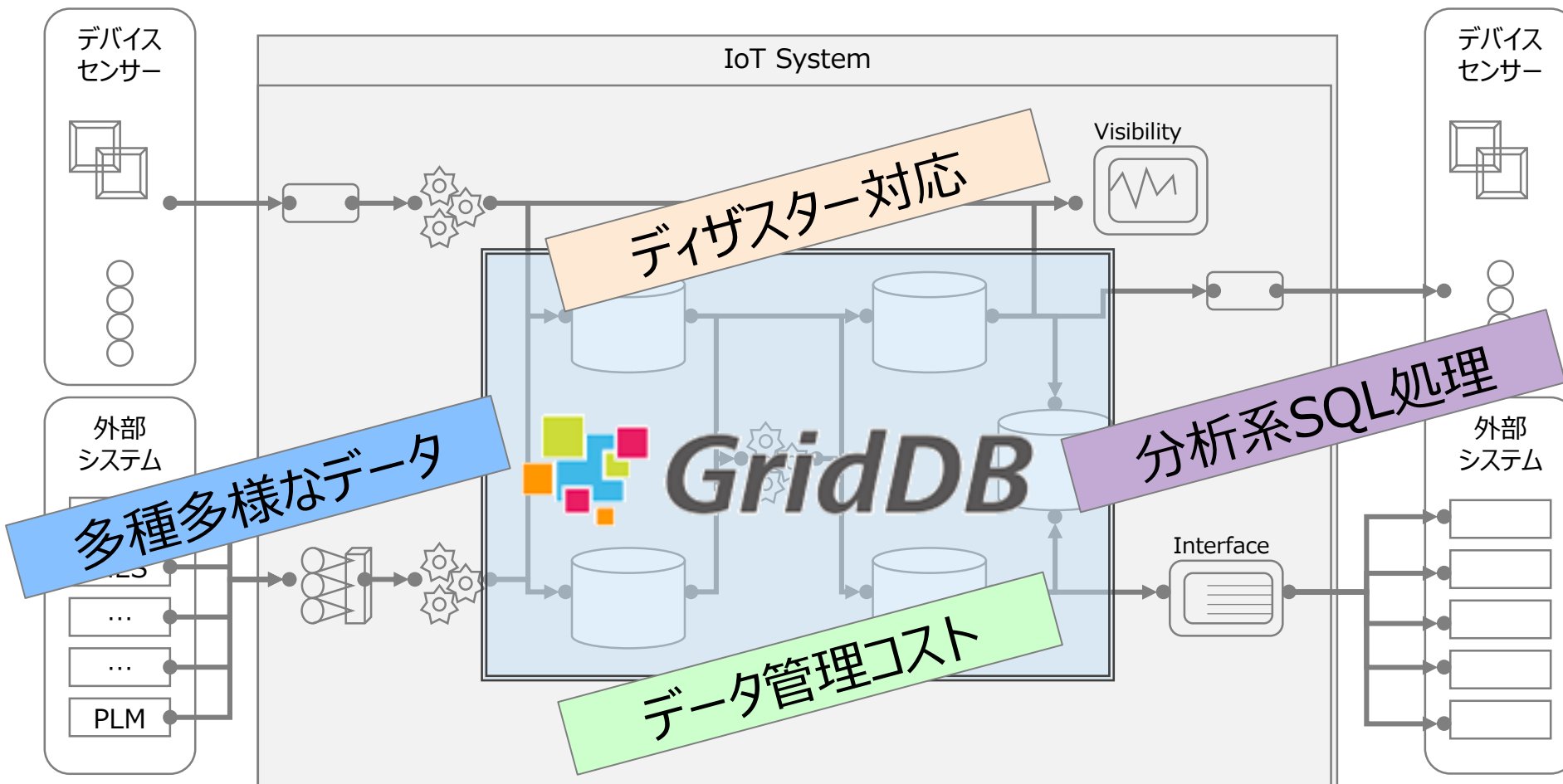
新アーキテクチャV5

－ V5シリーズの狙い －



# 適用領域拡大に伴う新たな要求

以下の要望に対応するため、アーキテクチャー新

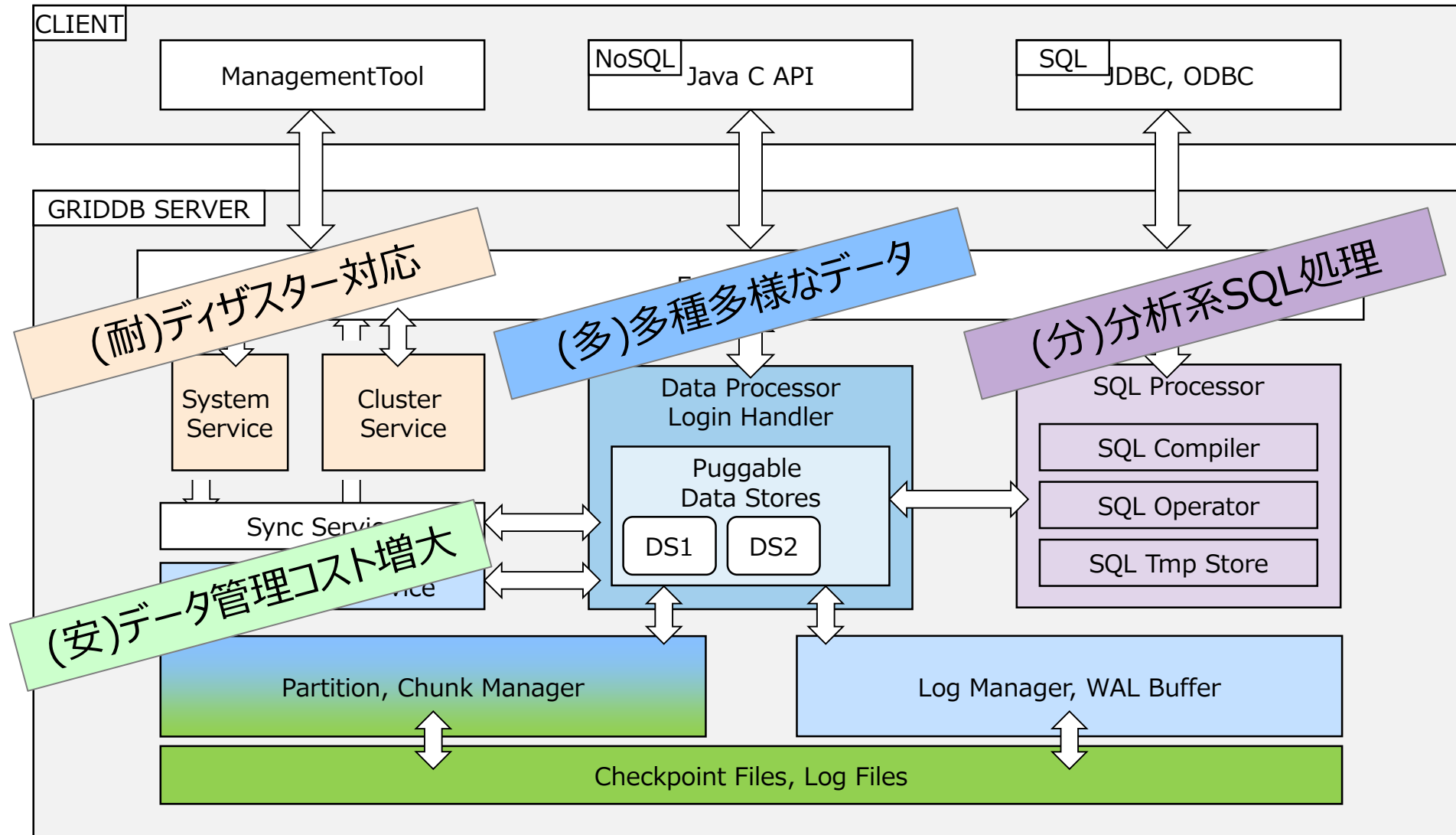


## GridDBへの要求の再掲 “多安分耐”

- (多) 多種多様なデータの最適管理
- (安) データ管理コストの削減
- (分) 分析系SQLの大幅な性能改善
- (耐) ディザスター機能の提供

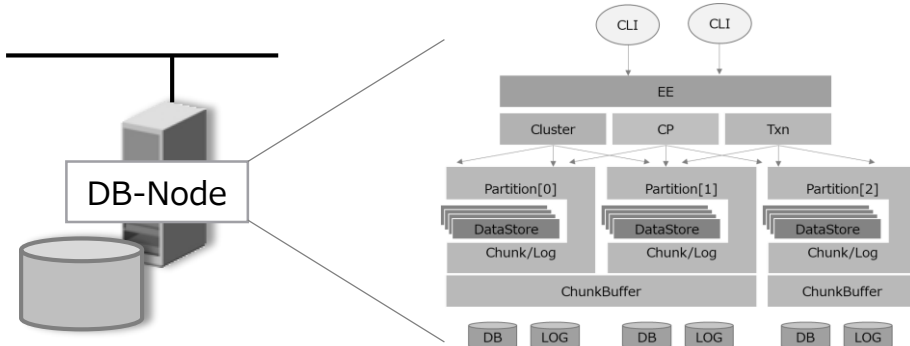
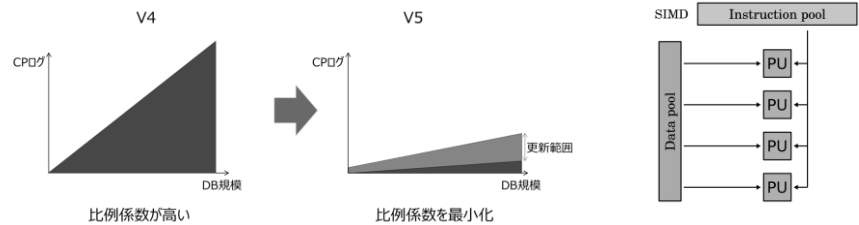
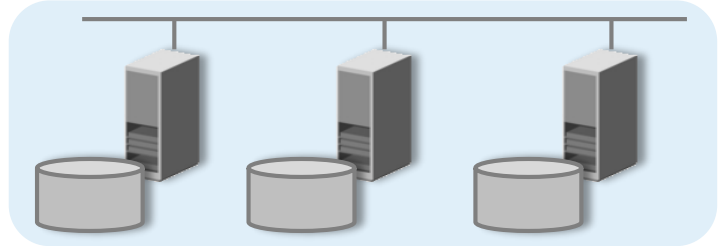


# アーキテクチャとの対応関係





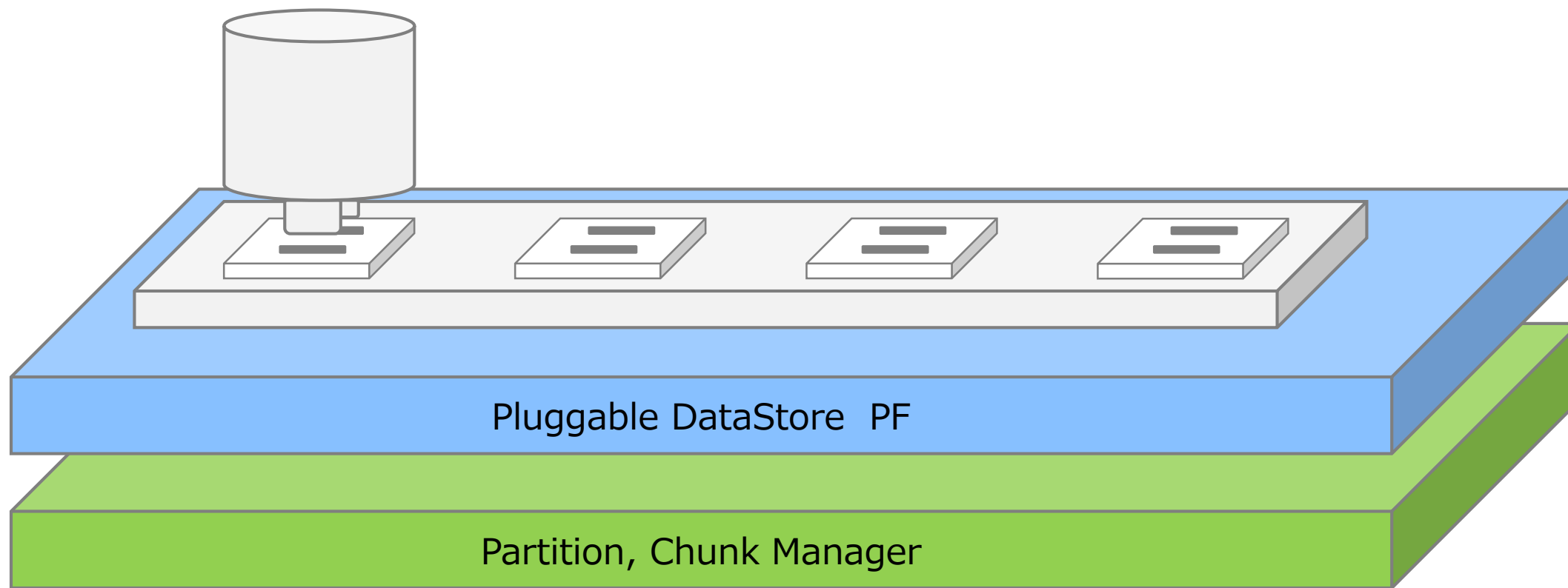
# (多) 多種多様なデータの最適管理

機能	達成レベル	イメージ図
<p>プラグブルデータストア</p> <p>5.0</p>	<ul style="list-style-type: none"> <li>リレーショナルモデル準拠</li> <li>例えば、カラムストア</li> </ul> <hr/> <ul style="list-style-type: none"> <li>任意データモデル</li> <li>例えば、オブジェクトストア</li> </ul>	
<p>ペタバイト対応強化</p> <p>5.0</p>	<ul style="list-style-type: none"> <li>ログサイズ削減(1/N倍)</li> <li>メモリサイズ削減</li> </ul>	
<p>各種高速化</p> <p>5.0</p>	<ul style="list-style-type: none"> <li>チェックポイント高速化</li> <li>DB削除、テーブル削除、スキャン(N倍)の高速化</li> <li>一部処理のSIMD化やアルゴリズム変更によるI/O高速化</li> </ul>	
<p>クラスタスナップショット</p>	<ul style="list-style-type: none"> <li>Copy-On-Writeに基づく</li> <li>クラスタ全体で瞬時バックアップ</li> </ul>	



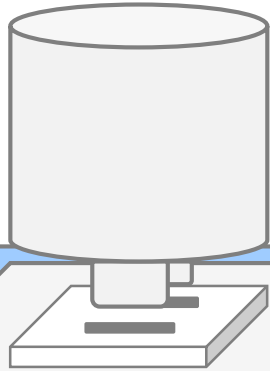
# プラグブルデータストア

登録更新に  
適した  
ロウ指向ストア

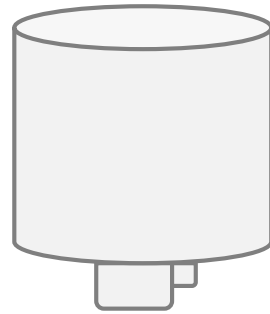


# プラグブルデータストア

登録更新に  
適した  
ロウ指向ストア



データ分析に  
適した  
カラム指向ストア

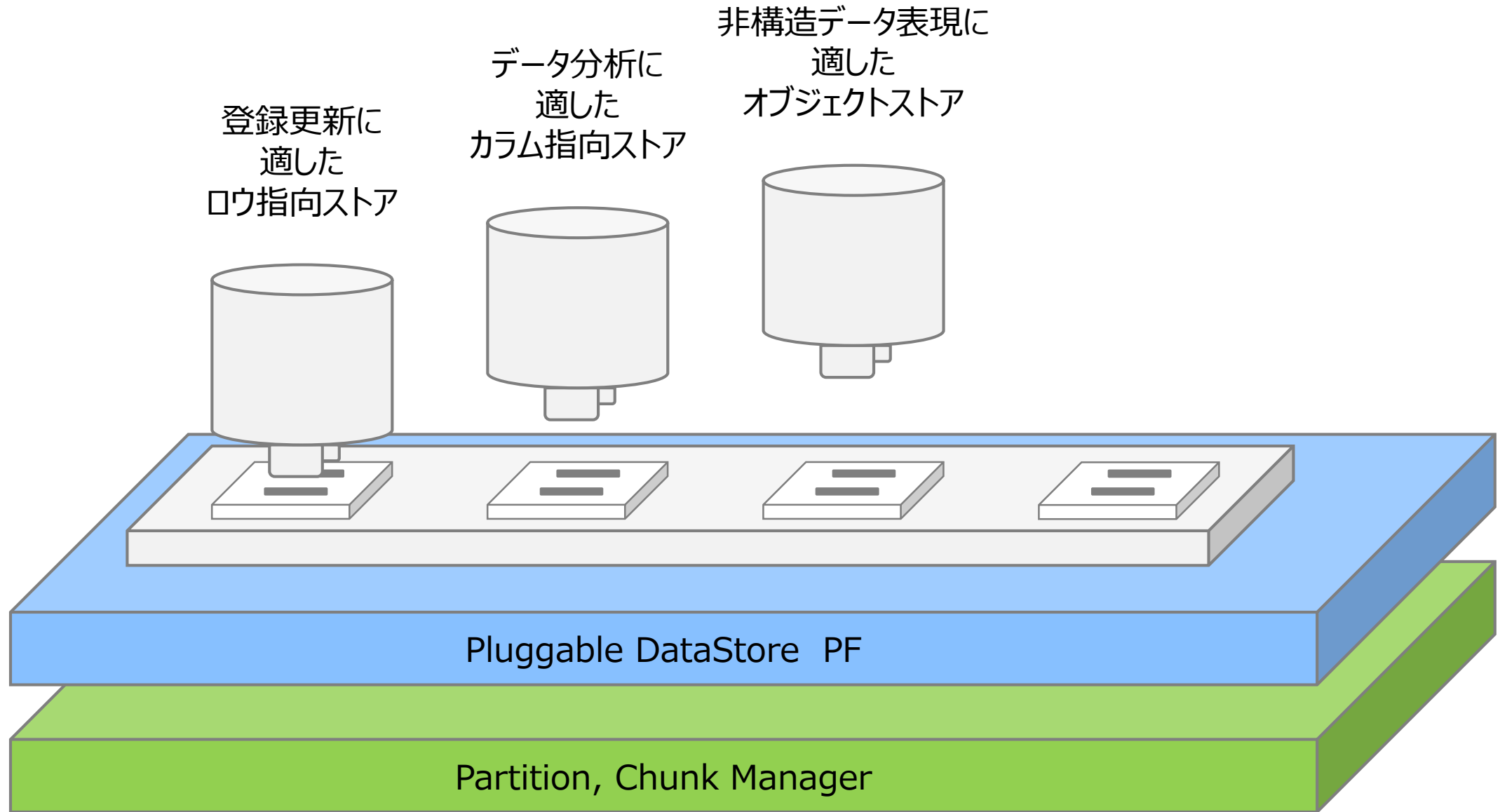


Pluggable DataStore PF

Partition, Chunk Manager

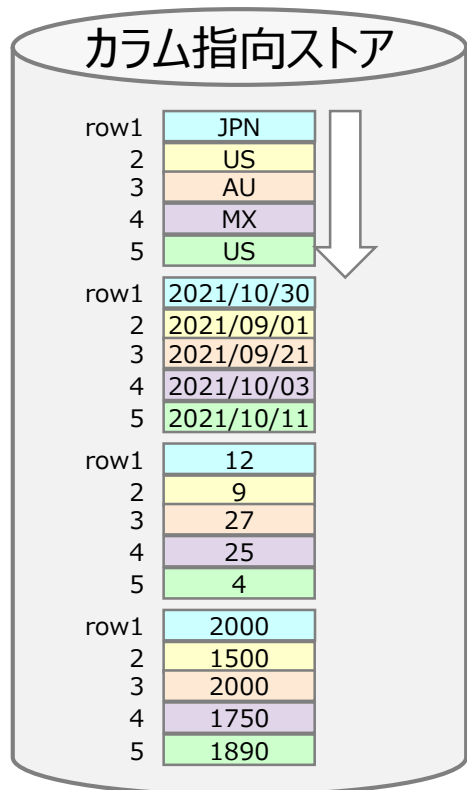
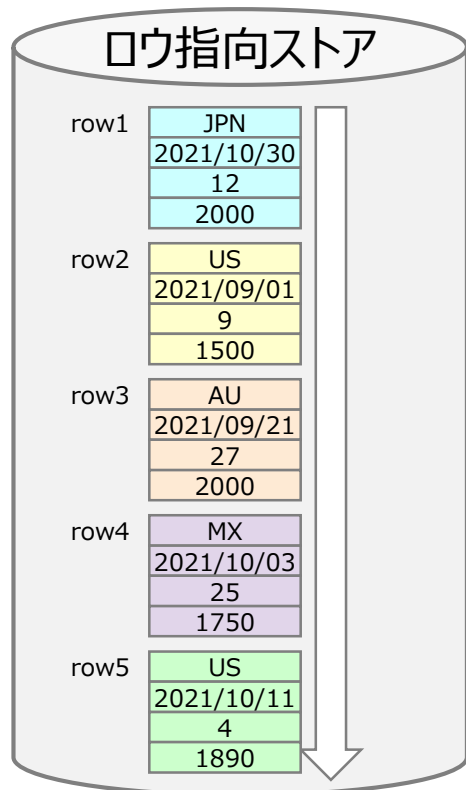


# プラグブルデータストア



# プラグブルデータストア データストア選択のメリット

	Country	Date	Quantity	Price
row1	JPN	2021/10/30	12	2000
2	US	2	入カテータブル	1500
3	AU	2		2000
4	MX	2021/10/03	25	1750
5	US	2021/10/11	4	1890

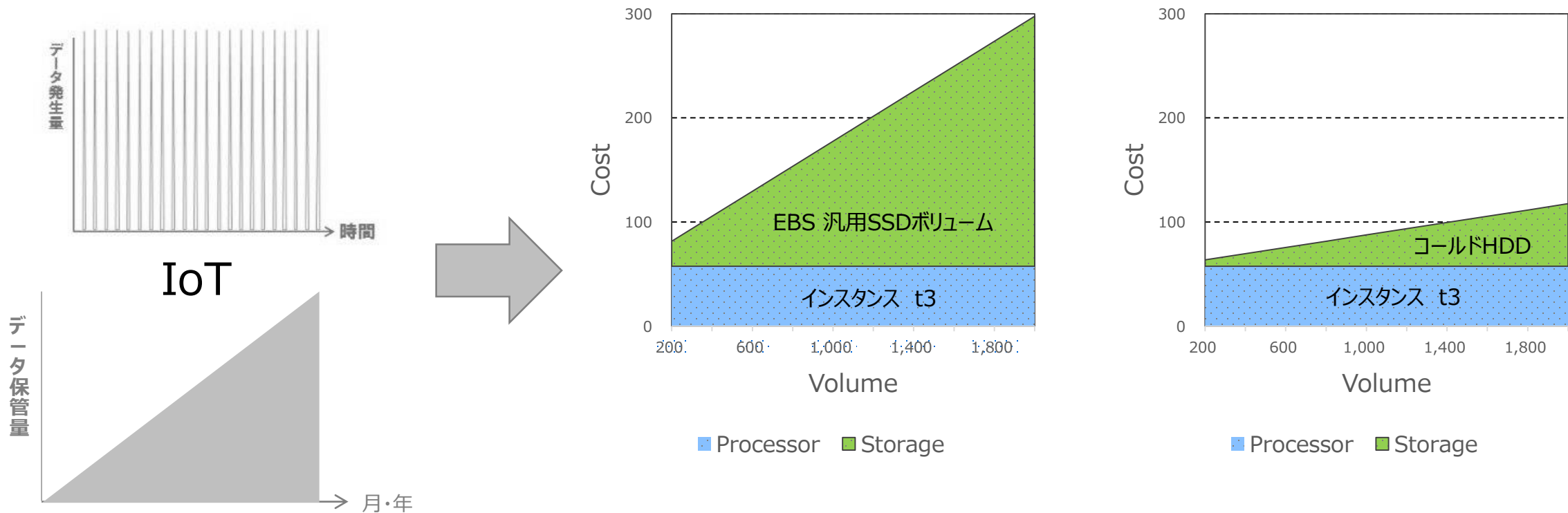


	ロウ指向ストア	カラム指向ストア
登録更新	○	×
レコード取得	○	△
カラム参照	△	○
集計	△	○
サイズ	△	○

ロウ指向とカラム指向の優劣

# (安) データ管理コストの削減

- × IoTにおいて、ストレージコストが全コストを支配
- × 一方、ストレージの性能とコストはトレードオフ

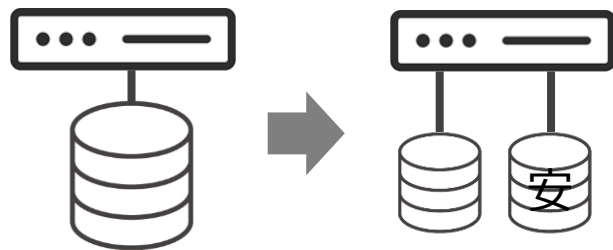


# (安) データ管理コストの削減

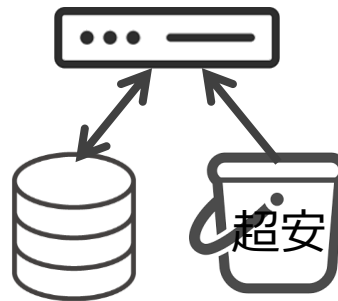
高速ストレージと格安ストレージの使い分けが必要

→ いろいろなタイプの使い分けを提供

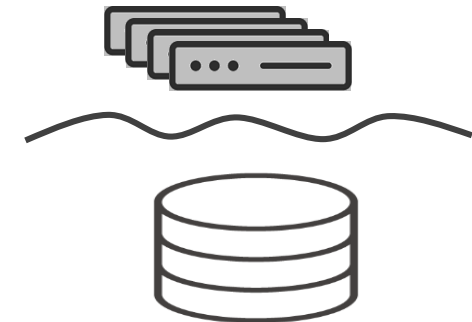
① GridDBデータファイル仮想化  
によるストレージミックス機能



② 外部テーブル参照機能



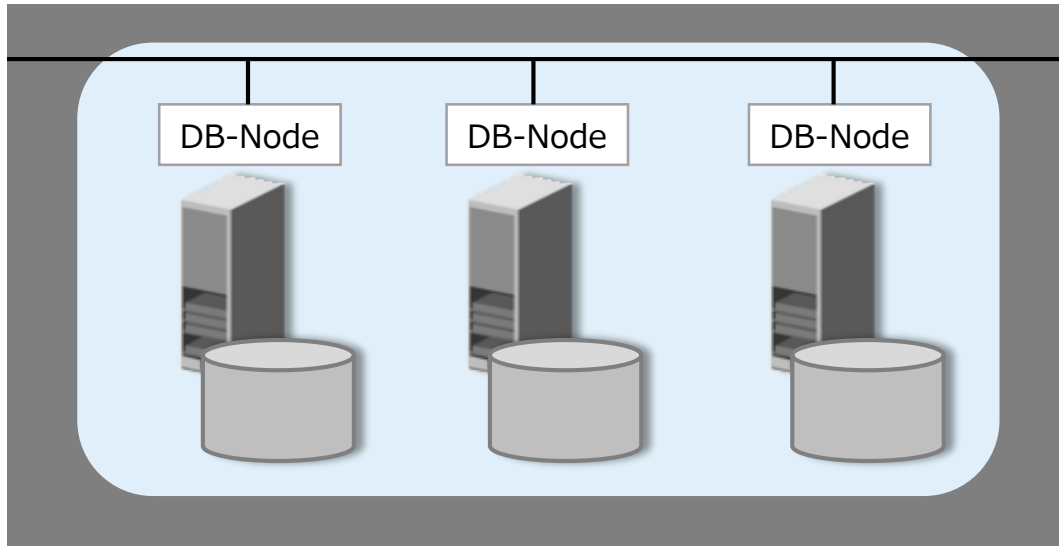
③ ストレージ分離機能  
(サーバレス一種)



### ③ストレージ分離機能

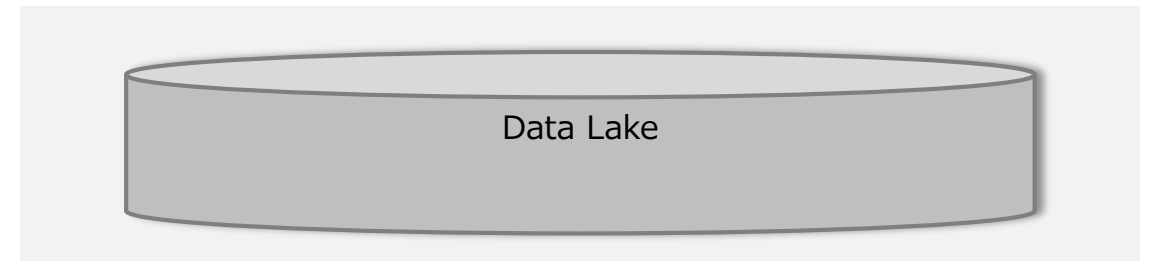
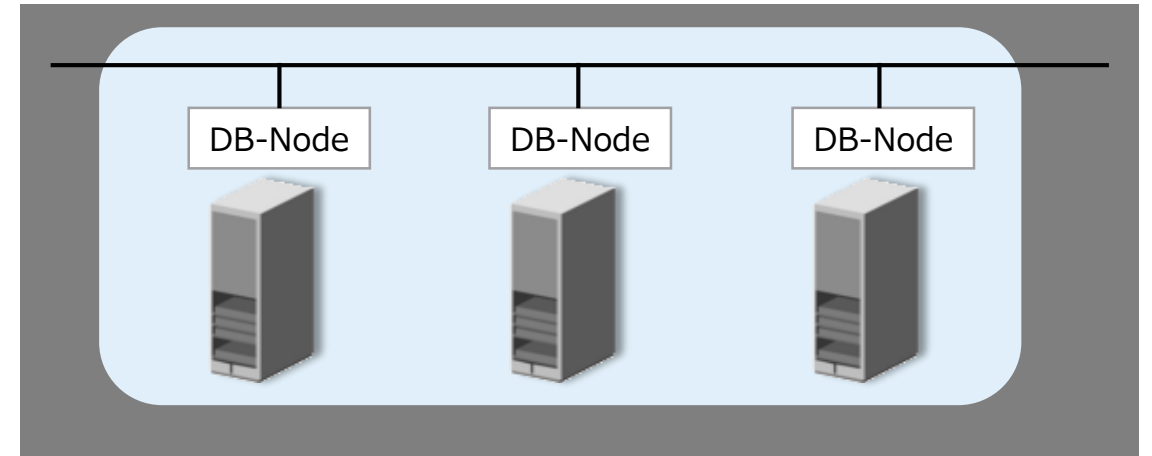
#### 非共有アーキテクチャ

- 高いパフォーマンス



#### ストレージ分離 ディスク共有アーキテクチャ

- 低コスト
- 負荷変動に対する弾力性



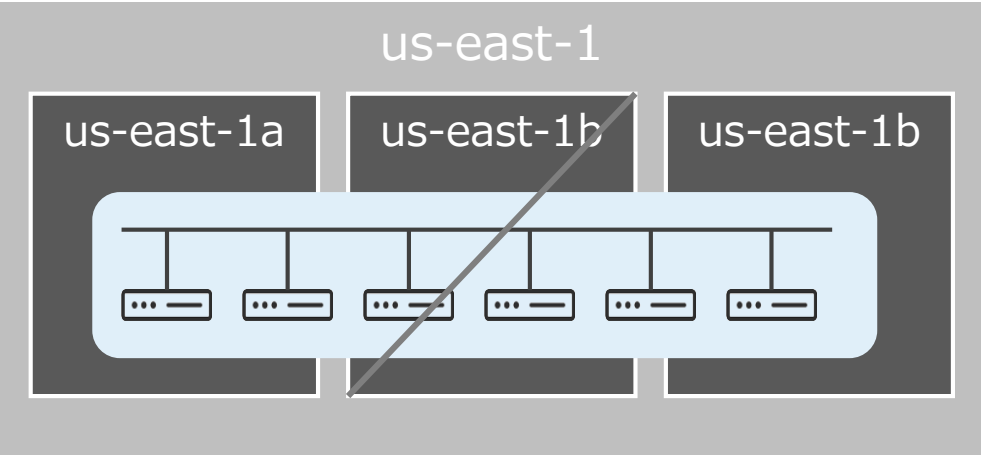
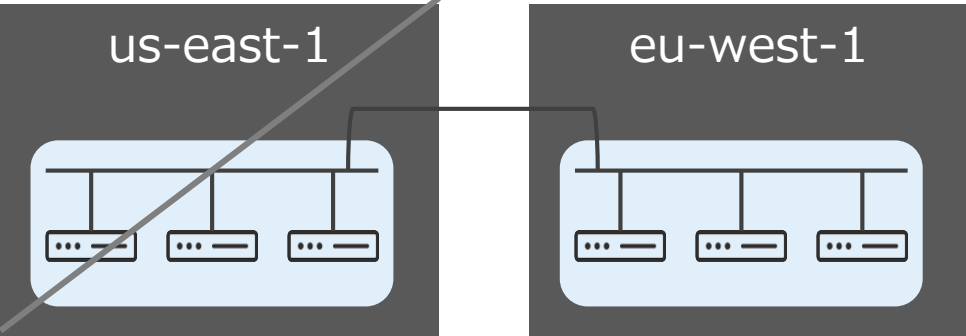


# (分) 分析系SQLの大幅な性能改善

機能	達成レベル	イメージ図
SQLコンパイラ コストベース最適化 (CBO)	<ul style="list-style-type: none"> <li>ジョイン最適化</li> </ul>	
	<ul style="list-style-type: none"> <li>オペレータ間最適化</li> </ul>	
	<ul style="list-style-type: none"> <li>統計情報精細化</li> </ul>	
並列分散SQL処理 高速化	<ul style="list-style-type: none"> <li>分散ハッシュジョイン (x10~)</li> <li>分散ハッシュグループ (x10~)</li> </ul>	
パイプライン制御強化	<ul style="list-style-type: none"> <li>ノード間データ流量制御</li> </ul>	



# (耐) ディザスター機能の提供

機能	達成レベル	イメージ図
耐ゾーン障害	<ul style="list-style-type: none"> <li>レプリカ増加</li> </ul>	
	<ul style="list-style-type: none"> <li>レプリカ不変</li> <li>ゾーンアウェアなノード配置</li> </ul>	
耐サイト障害	<ul style="list-style-type: none"> <li>アクティブクラスタ&amp;コールドクラスタ</li> <li>停止期間：～1Day</li> </ul>	
	<ul style="list-style-type: none"> <li>アクティブクラスタ&amp;アクティブクラスタ</li> <li>停止時間：～1Hour</li> </ul>	



## おわりに

- V5アーキテクチャー新により、開発効率アップ
- Community Editionも含めて、リリースサイクルを短縮
- 今後、3エディションでリリースされるV5シリーズにご期待ください。



# TOSHIBA

ご清聴ありがとうございました。

