

AWS RE:INVENT

re:Cap





ONLINE WEBINAR

AWS re:Invent Recap - ソリューション編

セッションは 10:00 から開始となります。
今しばらくお待ちください。

推奨ブラウザ

- Google Chrome (most recent 2 versions)
- Mozilla Firefox (most recent 2 versions)

お知らせとお願い

- 本日の資料はセミナー終了後にご案内するアンケートにダウンロードリンクがございます。
- 質問はセッション中、休み時間中、随時受け付けております。
GoToWebinar のチャット機能ではなく質問機能でご質問ください。

Time Table

- 10:00~10:20 Search/Streaming update
- 10:20~10:50 Data Lake update
- 10:50~11:00 休憩
- 11:00~11:25 DWH update
- 11:25~11:55 BI update

推奨ブラウザ

- Google Chrome (most recent 2 versions)
- Mozilla Firefox (most recent 2 versions)

お知らせとお願い

- 本日の資料はセミナー終了後にご案内するアンケートにダウンロードリンクがございます。
- 質問はセッション中、休み時間中、随時受け付けております。GoToWebinar のチャット機能ではなく質問機能でご質問ください。

注意事項

- アンケートに質問事項を記載頂いても回答をお送りできません。疑問点や不明点がある場合には GoToWebinar の質問機能をご利用ください。
- 頂いた質問事項は可能な限りセッション中に回答をお送りします。セッション中に回答出来なかった質問事項に関しては、セッション登録時に記入頂いたメールアドレス宛に後日送信致します。



AWS RE:INVENT RECAP - ソリューション編

Search & Streaming Updates

佐藤 祥多 (Shota Sato)

Analytics Specialist Solution Architect
Amazon Web Service Japan

自己紹介

佐藤 祥多

アマゾンウェブサービスジャパン
アナリティクススペシャリスト
ソリューションアーキテクト

ゲームバックエンドエンジニア
→ データコンサルタント
→ 現職

好きな AWS サービス
AWS Glue, AWS Lake Formation



アジェンダ

1. Amazon OpenSearch Service Updates
2. Streaming Updates (MSK/Kinesis)

Amazon OpenSearch Service Updates





Amazon OpenSearch Service

OpenSearch を簡単にデプロイ・管理、
スケール可能なフルマネージドサービス



フルマネージド: リソースのデプロイ、
管理に費やす時間を削減



セキュリティ: 認証、認可、暗号化、監査、
およびコンプライアンスのための高度な
セキュリティを維持

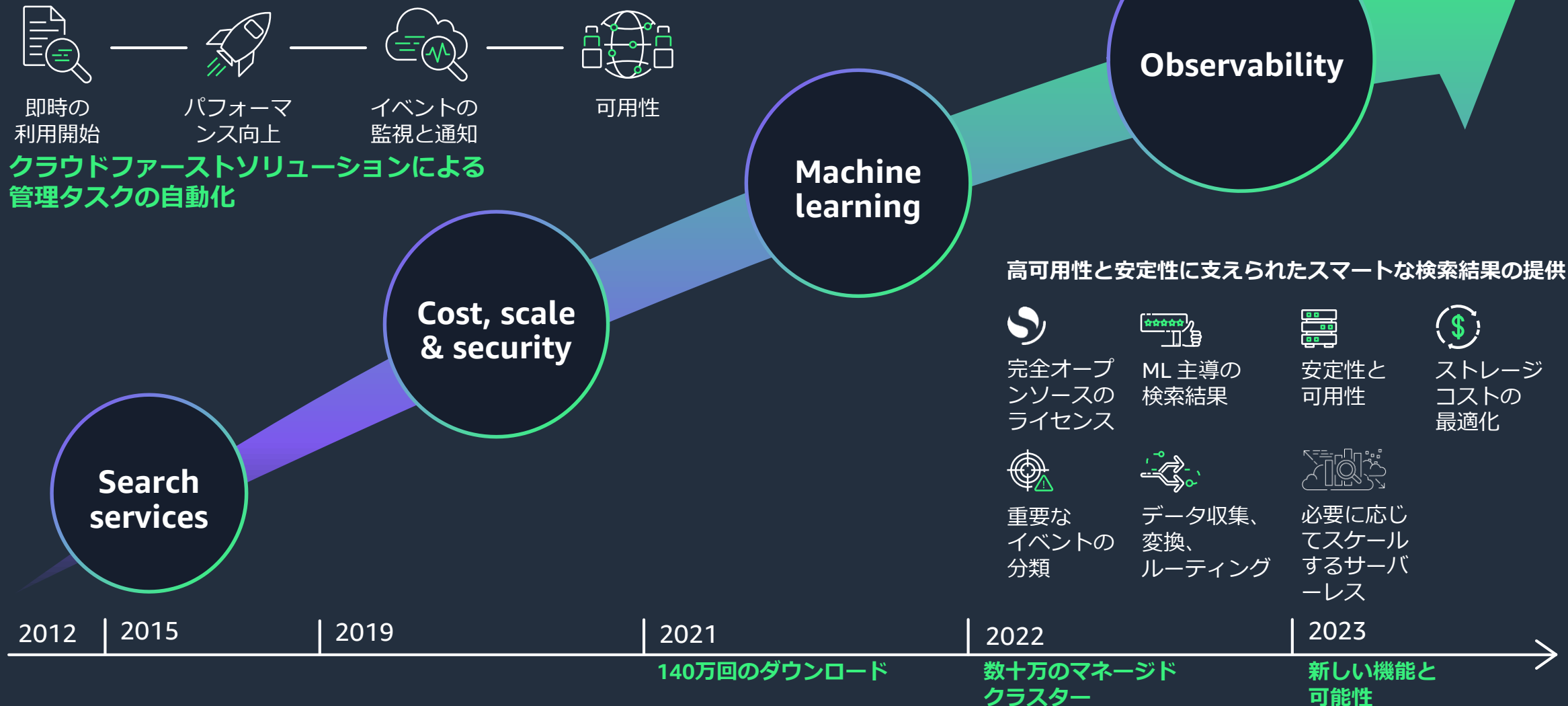


データ分析・オブザーバビリティ:
潜在的な脅威を体系的に検出し、機械学習、
アラート、可視化を活用して対処



コスト最適化: 各種リソースを最適化し、
戦略的な作業に注力

カスタマーの成長とユースケースの変遷



利用を阻む課題

- OpenSearch にデータを投入するのが大変
- 利用者が期待する検索精度を担保したい
- 複数の運用ログを横串で分析したい
- コストパフォーマンスよくクラスターを運用する必要がある

Getting data into OpenSearch



Data Engineer

"I want to process and ingest data into OpenSearch for use."


CHALLENGES

- Build custom tools and pipelines to fetch, prep, and ingest data

(from AWS Lambda, Logstash, FluentD, Amazon Glue, AWS Database Migration Service, DynamoDB Streams)

aws © 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Real-time analytics at scale



Developer / DevOps

"I want to ensure systems are monitored and data is available for forensics."

CHALLENGES

- Write queries to perform forensic analysis
- Builds visualizations, alerts, anomaly detection
- Manually correlates issues across sources (logs, metrics, traces)
- Builds tools to inspect logs for security events

aws © 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Easy and cost-effective platform



Platform / Infra Engineer


"I am responsible for running and managing clusters."

CHALLENGES

- Manually sizing clusters
- Monitors cluster health
- Manages uptime and SLAs

aws © 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Providing highly relevant search



Search Engineer

"I want to ensure my users have the best search experience."

CHALLENGES

- Manually fine tuning
- Boosting search results
- Updating synonyms
- Managing business requirements

aws © 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.

re:Invent 付近の主要なアップデート

re:Invent
前後

- Search Pipelines for OpenSearch の発表
- Sudachi Analyzer のサポート
- Neural Search のサポート
- Neural Sparse Retrieval のサポート
- OpenSearch Ingestion が Elasticsearch からの移行をサポート
- OpenSearch Ingestion が永続的バッファリングのサポート

re:Invent

- **Vector engine for OpenSearch Serverless**
- **Amazon DynamoDB との zero-ETL インテグレーション**
- **Amazon S3 との zero-ETL インテグレーション ※プレビュー**
- **OR1 インスタンスの発表**

Vector engine for OpenSearch Serverless

OpenSearch Serverless で Vector engine が一般利用開始に

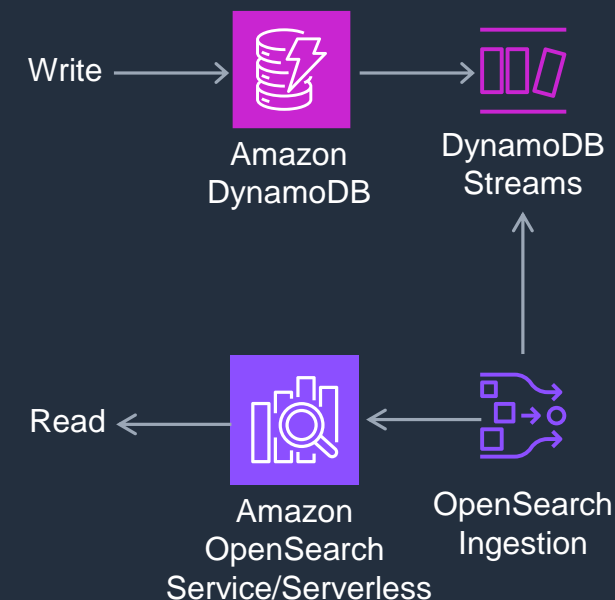
- 複雑なベクトルデータからなる数十億のデータをミリ秒単位の所要時間で保存・更新・検索が可能
- MLモデルで Embedding 生成をしたいときは外部で実施
- 検索は HNSW をサポート
- インフラ管理の手間なくベクトル検索機能を実現可能
- アプリケーションの検索ワークロードや生成 AI の RAG 基盤として利用



Amazon DynamoDB zero-ETL integration

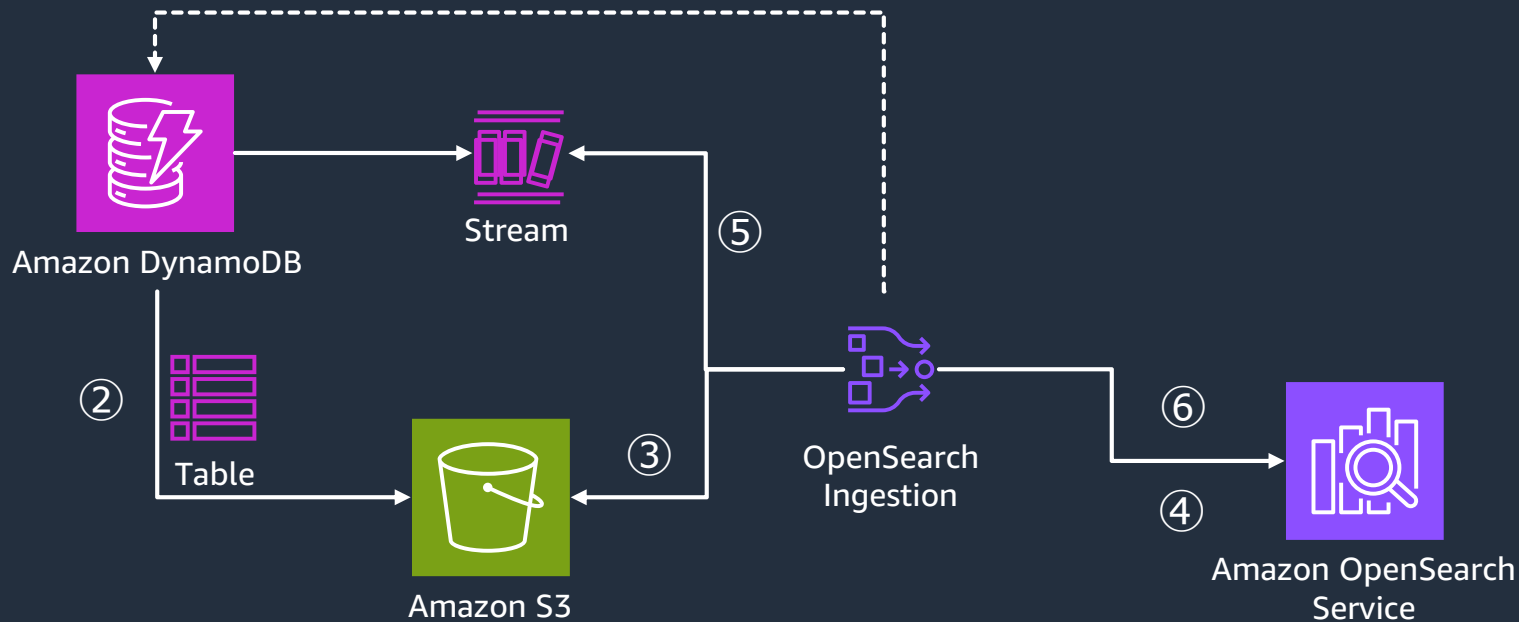
DynamoDB から OpenSearch へ簡単にデータをレプリケーションできる zero-ETL 機能

- ETL コードやインフラの管理が不要になり、OpenSearch の柔軟なクエリが利用しやすくなる
- ETL に OpenSearch Ingestion を利用するため、取り込み時のサンプリング、フィルタリング、データ変換をノーコードで実現
- 料金は DynamoDB Streams, PITR リカバリ, OpenSearch Ingestion の OCU 料金のみ
- Neural Search Plugin (OpenSearch Service のみ) と連携することで、Dynamo DB のデータをベクトル検索可能



Amazon DynamoDB zero-ETL integration による連携フロー

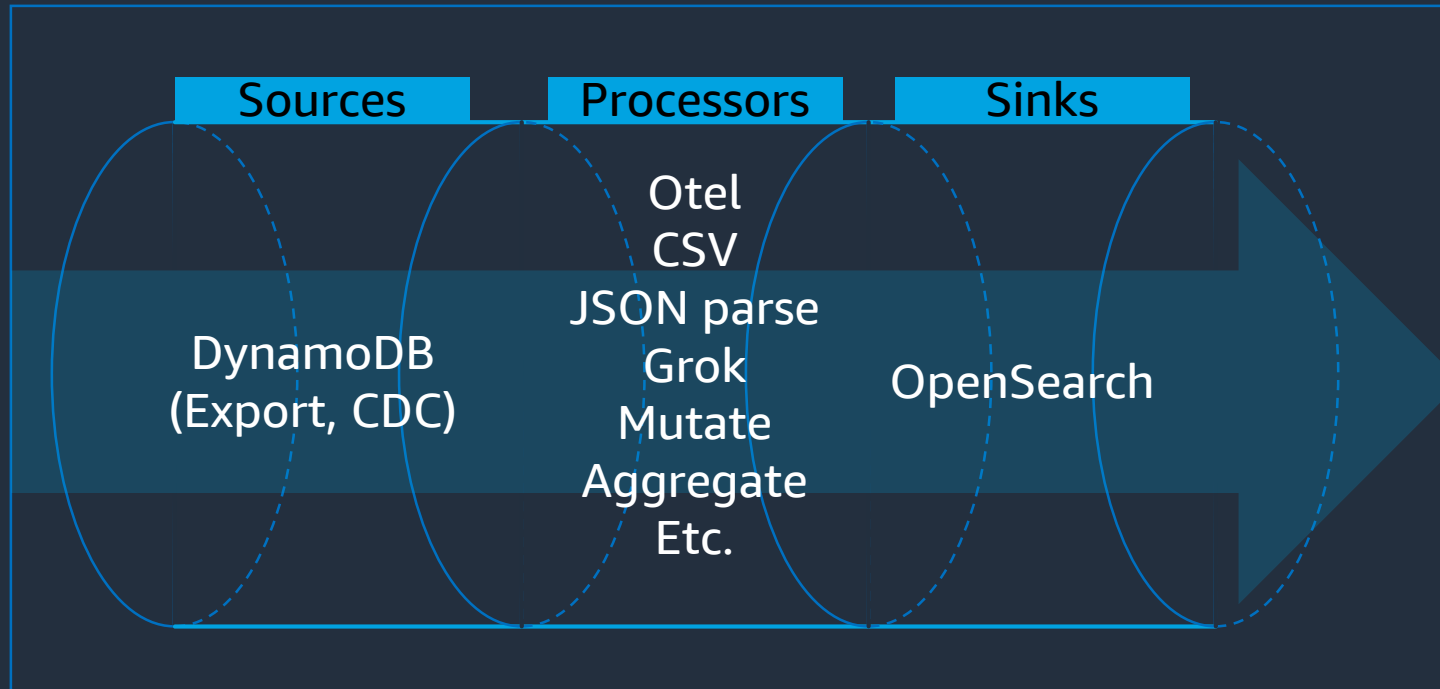
① ExportTableToPointInTime



1. テーブルエクスポートをリクエスト
2. S3 にエクスポートデータを出力
3. エクスポートされたデータを取得
4. OpenSearch ドメインへ書き込み
5. DynamoDB Streams からレコードを取得し、OpenSearch ドキュメントに変換
6. OpenSearch ドメインへ書き込み
7. 以降、パイプラインが中断されるまで 5 と 6 を繰り返し実行

blueprint からパイプライン設定を作成し、パイプラインを起動するだけ

zero-ETL Integration におけるパイプラインの処理フロー



- OpenSearch Ingestion のパイプラインは Sources, Processors, Sinks で構成されている
- Sink に複数の OpenSearch クラスターを登録することでデータを複数のクラスターに同期可能
- Processors の設定次第でデータ変換やフィールド追加などのデータ処理が可能

データ変換の例 – 埋め込み

- OpenSearch Service の Neural Search と連携し、DynamoDB レコードの特定フィールドからベクトルを生成・格納することが可能
- OpenSearch Index にデフォルトの Ingest Pipeline を設定することでノーコードでの埋め込みを実装
- 事前に OpenSearch 側にコネクタやモデルを設定しておく必要はある
- DynamoDB のデータを活用したベクトル検索が可能に
- ベクトル検索は OpenSearch のエンドポイントに発行する必要あり

```
processor:  
- add_entries:  
  entries:  
  - key: "embedding_chunk"  
    format: "${item_name} ${product_description} ${question_text} ${answer_text}"  
    overwrite_if_key_exists: true  
  template_type: "index-template"  
  template_content: |  
{  
  "template": {  
    "settings": {  
      "index.knn": true,  
      "default_pipeline": "nlp_pipeline",  
      "number_of_shards": 2,  
      "number_of_replicas": 0  
    }  
    "mappings": {  
      "properties": {  
        "pk": { "type": "keyword", "ignore_above": 256 },  
        "sk": { "type": "keyword", "ignore_above": 256 },  
        "doc_type": { "type": "keyword", "ignore_above": 256 },  
        "embedding": {  
          "type": "knn_vector",  
          "dimension": 1536,  
          "method": {
```


パイプラインのサイジング

- シンプルなパイプラインについては、1000 WCU ごとに 1 ingestion OCU をベースとしてサイジング
 - 1 ingestion OCU ごとに 2 vCPU + 8 GiB memory を割り当て。
 - 1 OCU につき 1時間ごとに \$0.326 のコストが発生 (東京リージョン)
- OpenSearch Ingestion はオートスケールに対応しているため、必要に応じてスケールさせることが可能。OCU 上限も設定可

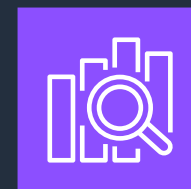
“For simple pipelines, a single OpenSearch Compute Unit (OCU) can process about 1 MB per second of writes. This is the equivalent of about 1000 write request units (WCU). Depending on your pipeline's complexity and other factors, you might achieve more or less than this.”

<https://docs.aws.amazon.com/amazondynamodb/latest/developerguide/OpenSearchIngestionForDynamoDB.html>

zero-ETL integration with Amazon S3

S3 にあるデータに対して OpenSearch Service で直接クエリして可視化できる Zero-ETL 機能

- S3 にある膨大なログデータ (VPC Flow ログ、ELB ログなど) に対してインシデントがあった時に OpenSearch 上で分析と可視化が簡単に実施可能に
- 料金は直接クエリするための計算リソースである OpenSearch Compute Unit に対して発生 (分単位)



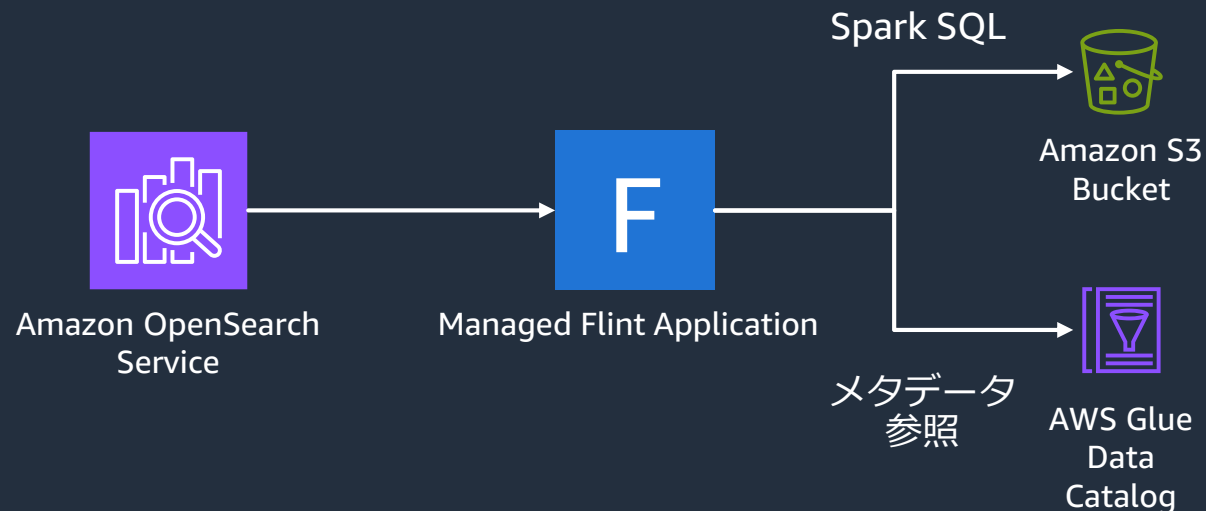
Amazon OpenSearch Service



Amazon S3

zero-ETL integration with Amazon S3 の仕組み

- Glue Data Catalog をメタデータストアとして使用
- クエリは OpenSearch クラスターから OpenSearch Flint アプリケーション実行基盤に。Flint アプリケーションが Spark SQL を実行することで S3 バケット内のデータ分析を実施



Query Acceleration

クエリ実行を高速化するための 3 つのオプションを提供

Skipping Index

- スキャン対象外とするオブジェクトの統計情報を OpenSearch 側で持ち、クエリを高速化
- Partition、MinMax(特定カラムの最小値 & 最大値)、ValueSet(特定の値を持つか)

Covering Index

- 特定カラムのデータを OpenSearch インデックスにインポートすることで、OpenSearch 上で高速・柔軟にデータ検索を行う
- 自動インデックスリフレッシュ機能をサポート

Materialized View

- クエリ実行結果を OpenSearch インデックスにインポートすることで、可視化を高速化
- Window 関数を一定間隔ごとに繰り返し呼び出すことで差分更新を行うことが可能

OR1 インスタンスの発表

OpenSearchに最適化されたインスタンスファミリーOR1を発表

- マネージドサービス内部でストレージシステムのアーキテクチャを見直して実現
- ベンチマーク結果では30%のコスト効率向上
- インデクシング処理負荷が重い運用分析ワークロード向き
- バージョンは 2.11 以降をサポート

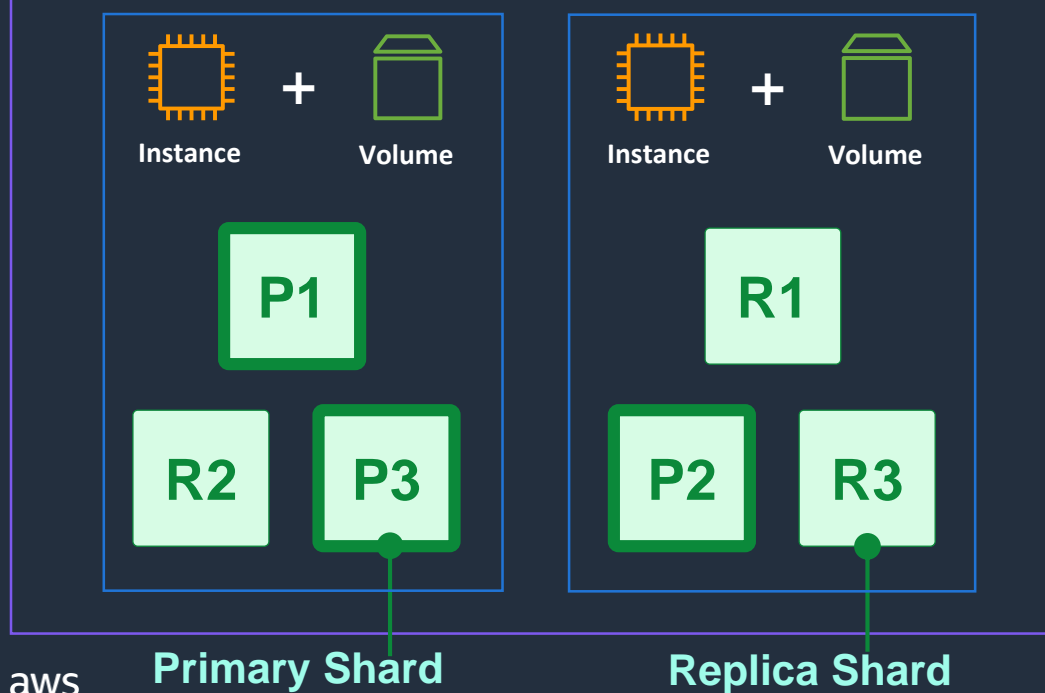
| Instance type | Instances |
|---------------|---|
| OR1 | or1.medium.search or1.large.search or1.xlarge.search or1.2xlarge.search or1.4xlarge.search or1.8xlarge.search or1.16xlarge.search |

OR1 アーキテクチャー

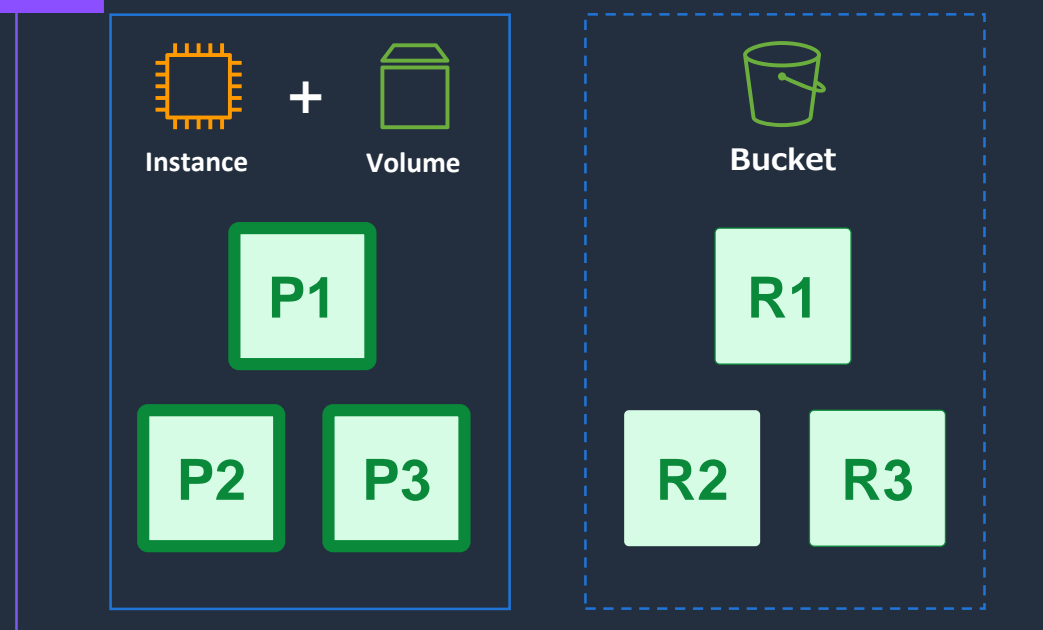
- 従来の OpenSearch ドメインはノード間でシャードを分散配置
- OR1 はシャードのレプリカを S3 バケットに配置。レプリカシャードはクエリされない



従来型のドメイン



OR1



OR1 インスタンスの特徴



パフォーマンス向上

- S3 からのスナップショット取得による負荷低減
- ノード障害からの復旧時、構成変更時にも S3 からデータを取得するためプライマリノードの負荷が低減
- レプリケーション仕組みの変更による書き込みスループットの改善



データ耐久性

- レプリカは Amazon S3 に配置することで、11 9s の耐久性と安価なストレージを活用できるように



自動復旧

- インデックスが利用不可能になった場合、S3 から自動で即時リストアを実行

OR1 インスタンスのスペックと料金体系 (東京リージョン)

- vCPU:memory 比が 1:8 となっており、r6g や r5 と同じ比率
- ノードごとにインスタンス料金と EBS ボリューム料金が発生するほか、レプリカを配置するマネージドストレージの料金が発生する
- 東京リージョンでは、マネージドストレージの料金は、1 GB ごとに 0.026 USD 実際に格納されたインデックスサイズに応じて課金される

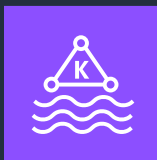
| type | vCPU | mem (GB) | EBS Volume Limit | Pricing (USD/hour) * on-demand | Pricing (USD/hour) * RI 1yr all |
|---------------------|------|----------|------------------|-----------------------------------|------------------------------------|
| or1.medium.search | 1 | 8 | 400 GiB | 0.126 | 0.081 |
| or1.large.search | 2 | 16 | 800 GiB | 0.253 | 0.164 |
| or1.xlarge.search | 4 | 32 | 1.5 TiB | 0.505 | 0.328 |
| or1.2xlarge.search | 8 | 64 | 3.0 TiB | 1.009 | 0.655 |
| or1.4xlarge.search | 16 | 128 | 6.0 TiB | 2.019 | 1.312 |
| or1.8xlarge.search | 32 | 256 | 12.0 TiB | 4.038 | 2.624 |
| or1.12xlarge.search | 48 | 384 | 18.0 TiB | 6.056 | 3.936 |
| or1.16xlarge.search | 64 | 512 | 24.0 TiB | 8.062 | 5.24 |

| type | vCPU | mem (GB) | EBS Volume Limit | Pricing (USD/hour) * on-demand | Pricing (USD/hour) * RI 1yr all |
|---------------------|------|----------|------------------|-----------------------------------|------------------------------------|
| r6g.large.search | 2 | 16 | 2.0 TiB | 0.202 | 0.139 |
| r6g.xlarge.search | 4 | 32 | 3.0 TiB | 0.404 | 0.279 |
| r6g.2xlarge.search | 8 | 64 | 6.0 TiB | 0.807 | 0.557 |
| r6g.4xlarge.search | 16 | 128 | 12.0 TiB | 1.615 | 1.114 |
| r6g.8xlarge.search | 32 | 256 | 16.0 TiB | 3.23 | 2.229 |
| r6g.12xlarge.search | 48 | 384 | 24.0 TiB | 4.845 | 3.343 |

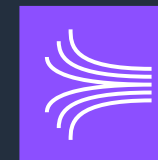
Streaming Updates (MSK/Kinesis)



AWS Streamings の主要なアップデート



Amazon Managed Streaming
for Apache Kafka
(Amazon MSK)



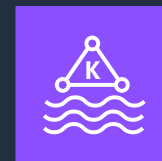
Amazon Kinesis

- Graviton 3 搭載の M7g インスタンスを **新規クラスター**でサポート
- Amazon MSK Serverless は全ての開発言語をサポート
- Amazon MSK は AWS IAM サポートを **新規クラスター**で全ての開発言語をサポート
- AWS Trusted Advisor において Amazon MSK のパーティション数が過多でないかのチェックを追加
- Amazon MSK が ストレージ容量のアラートをサポート
- Tiered Storage が利用可能な Apache Kafka version 3.6.0 サポートを発表
- Amazon Kinesis Data Streams が AWS Lambda によるクロスアカウントアクセスをサポート
- Kinesis Data Streams のコンソールが Amazon EventBridge Pipes との統合がリリース

Amazon MSK が M7g インスタンスをサポート

Amazon Managed Streaming for Apache Kafka (MSK) が M7g インスタンスをサポート

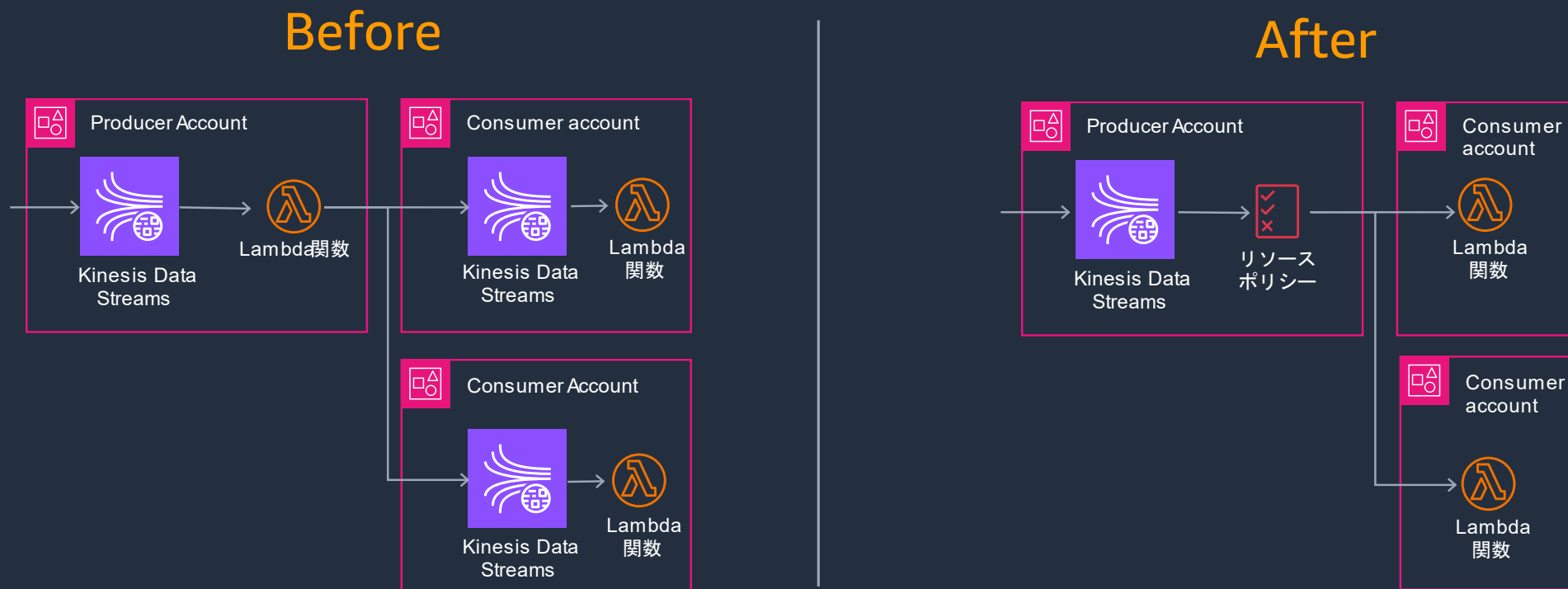
- **新規クラスタ**で利用できる
- M7g インスタンスは AWS Graviton 3 プロセッサを搭載しており、同等サイズの M5 と比較してコストが下がり、パフォーマンスも向上
- ストレージスループットが**最大 25% 向上**
- ネットワークスループットが**最大 88% 向上**



Amazon Managed Streaming
for Apache Kafka
(Amazon MSK)



Amazon Kinesis Data Streams が AWS Lambda によるクロスアカウントアクセスをサポート



- Kinesis Data Streams でリソースポリシーがサポートされるようになった
- データ重複保持コストを削減可能に



Thank you!