



AWS
Black Belt
Online Seminar

【AWS Black Belt Online Seminar】

データレイク入門：

AWS で様々な規模のデータレイクを分析する
効率的な方法

アマゾン ウェブ サービス ジャパン株式会社
ソリューションアーキテクト 川村誠

2018.06.19

自己紹介

□ 名前

川村 誠 (かわむら まこと)

□ 所属

アマゾン ウェブ サービス ジャパン 株式会社
技術統括本部 ストラテジックソリューション本部
ソリューション アーキテクト

□ 好きな AWS サービス

- ❖ Amazon EMR
- ❖ Amazon ECS
- ❖ Amazon SageMaker

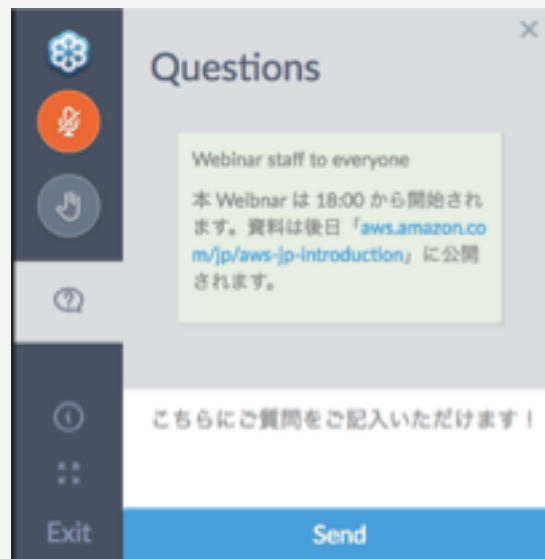


AWS Black Belt Online Seminarへようこそ！

質問を投げることができます！

- GoToWebinar の吹き出しマークから、質問を書き込んで下さい。
(書き込んだ質問は、主催者にしか見えません)
- 今後のロードマップに関するご質問は
お答えできませんのでご了承下さい。
- Twitterへツイートする際はハッシュタグ
#awsblackbelt をご利用下さい。

- ① 吹き出しをクリック
- ② 質問を入力
- ③ Sendをクリック



AWS Black Belt Online Seminarとは

AWSJのTechメンバがAWSに関する様々な事を紹介するオンラインセミナーです

【火曜 12:00～13:00】

主にAWSのソリューションや業界カットでの使いどころなどを紹介（例：IoT、金融業界向け etc.）

【水曜 18:00～19:00】

主にAWSサービスの紹介やアップデートの解説（例：EC2、RDS、Lambda etc.）

※開催曜日と時間帯は変更となる場合がございます。最新の情報は下記をご確認下さい。

オンラインセミナーのスケジュール&申し込みサイト <https://aws.amazon.com/jp/about-aws/events/webinars/>

内容についての注意点

- 本資料では2018年6月19日時点のサービス内容および価格についてご説明しています。最新の情報はAWS公式ウェブサイト (<http://aws.amazon.com>)にてご確認ください。
- 資料作成には十分注意しておりますが、資料内の価格とAWS公式ウェブサイト記載の価格に相違があった場合、AWS公式ウェブサイトの価格を優先とさせていただきます。
- 価格は税抜表記となっています。日本居住者のお客様が東京リージョンを使用する場合、別途消費税をご請求させていただきます。
- AWS does not offer binding price quotes. AWS pricing is publicly available and is subject to change in accordance with the AWS Customer Agreement available at <http://aws.amazon.com/agreement/>. Any pricing information included in this document is provided only as an estimate of usage charges for AWS services based on certain information that you have provided. Monthly charges will be based on your actual use of AWS services, and may vary from the estimates provided.

本セミナーの概要

□ 本セミナーで学習できること

- ❖ データレイクのデータを自動的／効率的に分析可能にする方法
- ❖ データレイクと **DWH** のデータを効率的に分析する方法

□ 対象者

- ❖ データ運用業務に関わるエンジニア、アナリスト、アーキテクトの方
- ❖ **DWH**・**DB Administrator** の方
- ❖ 次の **AWS** のサービスの概要レベルの知識が前提になります。
Amazon S3 / AWS Glue / Amazon Redshift など

Agenda

- データチャレンジ & データレイク
- データレイクのデータを自動的 / 効率的に分析可能にする方法
- データレイクとデータウェアハウスに入っている様々な規模のデータを効率的に分析する方法

データチャレンジ & データレイク

データチャレンジ

データの種類とデータ量が急増している



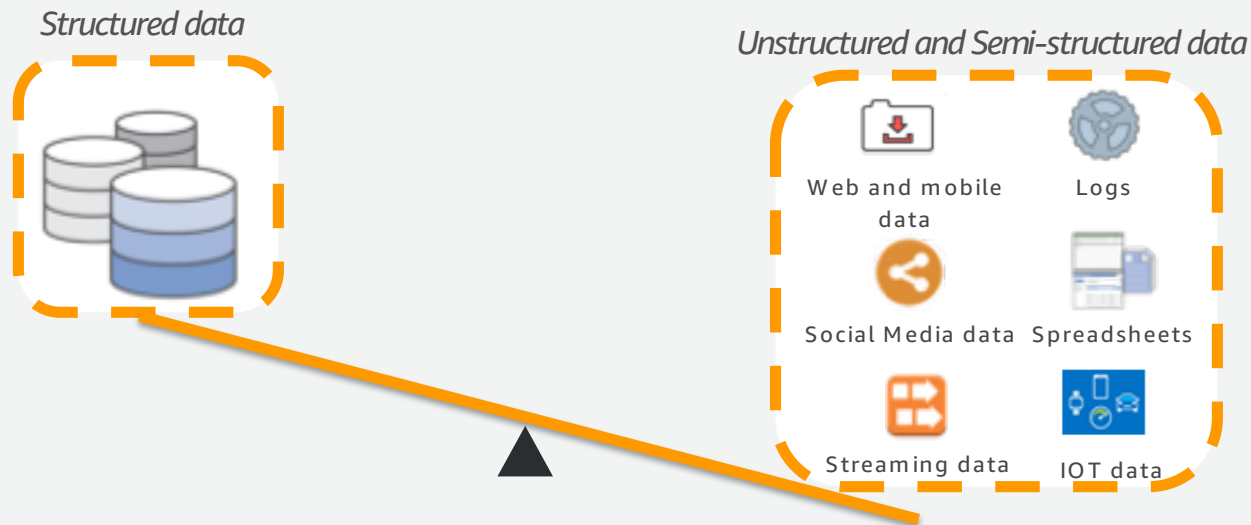
データを集め、理解し、データから価値を見出す

複数のデータ消費者とアプリケーション



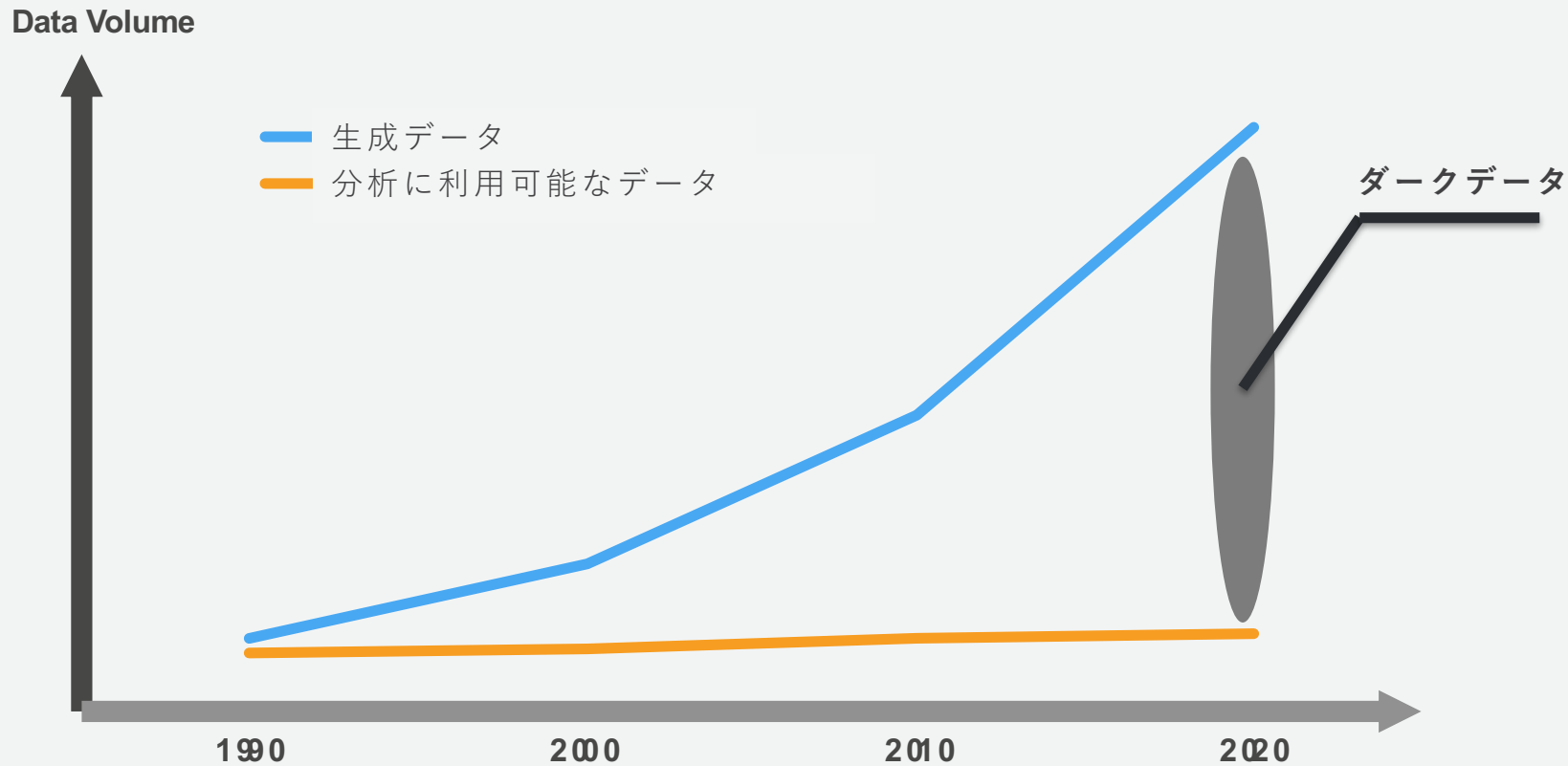
新しい洞察をデータからすばやく抽出し、ビジネスを加速する

何が新しいチャレンジなのか？

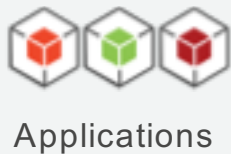
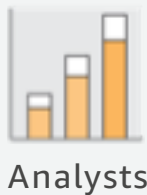


ダークデータ

ダークデータチャレンジ



複数のデータ消費者と複数の要件



Agile

Real time

Flexible

Scale

データの複製が生成されてしまう！

伝統的なデータウェアハウス



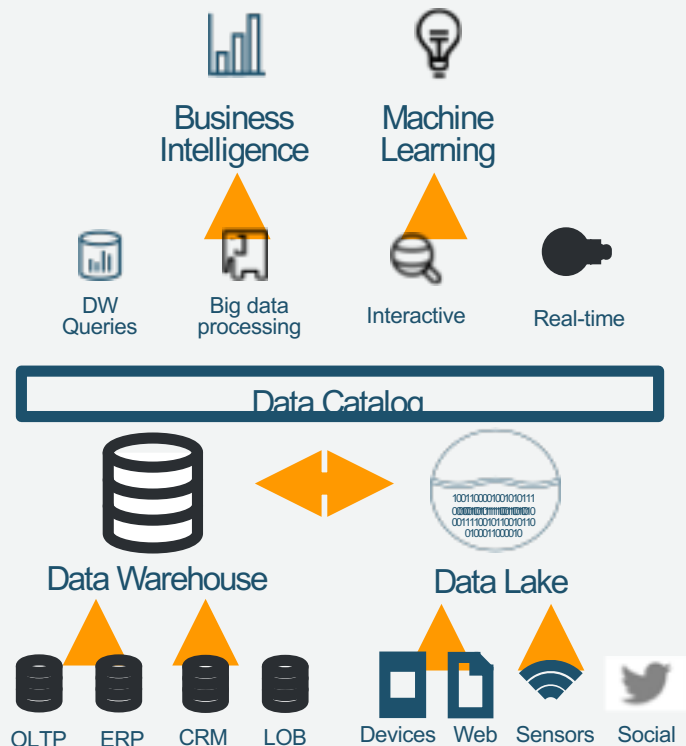
リレーショナルデータ

テラバイト～ペタバイト規模

データロード前に定義されるスキーマ

運用報告やアドホック分析

データウェアハウスを拡張するデータレイク



リレーショナルと非リレーショナル
データ

テラバイト～エクサバイト規模

分析中に定義するスキーマ
(Schema on Read)

インサイトを得るための多様な分析
エンジン

低コストストレージと分析用に設計

データレイク on AWS



データを取り込む様々な方法

エクサバイト規模での冗長性と可用性

セキュリティ、コンプライアンス、監査

同一データに対して移動無しでどんな分析も実行可能

バラバラに拡張可能なストレージと計算リソース

保存：0.025USD / GB-month ※

クエリ：0.005USD / GB scanned

データレイクとは何か？

どんな規模のデータも低コストで全て収集し、保存することが可能

データを配置し、価値を与え、セキュアに守ることが可能

組織内のデータへの民主化されたアクセスを提供する

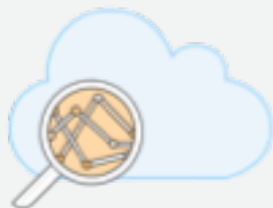
すばやく、簡単に新しいデータ分析形式を実行可能



データレイクの利点



どんな規模のデータも低コストで全て収集し、保存することが可能



信頼できる唯一の情報源(**single source of truth**)を持つことで、すばやく検索し、関連データを見つけることが可能



統一されたツール郡を使用し、データを簡単にクエリ可能

欠けている部品

- ① 格納場所に関係なくデータへの統合されたビュー
- ② 分析ツールとの統合
- ③ メタデータを自動的に構築し、進化するデータに合わせてメタデータを同期する方法

データレイクのデータを自動的／効率的に
分析可能にする方法

→ **AWS Glue**

AWS Glue とは何か？

発見

自動的にデータを発見して分類することで、様々なデータソースをまたいだ検索と参照がすぐに実行可能に

開発

様々なデータソース間でデータをクレンジングし、リッチ化し、確実に移動するコードが生成される。このコードを簡単にカスタマイズしたり、独自のコードを持ち込むことが可能

展開

サーバーレスで完全に管理されたスケールする環境でジョブが実行される。プロビジョニングまたは管理するための計算リソースは必要ない

AWS Glue 構成要素



Data Catalog

発見

Apache Hive メタデータ
ストア互換

AWSサービスと統合さ
れている

自動的にクローリング

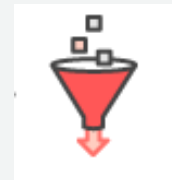


Job Authoring

開発

PySpark・Scalaに対応
したETLコードを自動
的に生成する

編集、デバッグ、共有
可能



Job Execution

展開

サーバレスで実行可能
柔軟なスケジューリング
モニタリングとアラート
実施可能

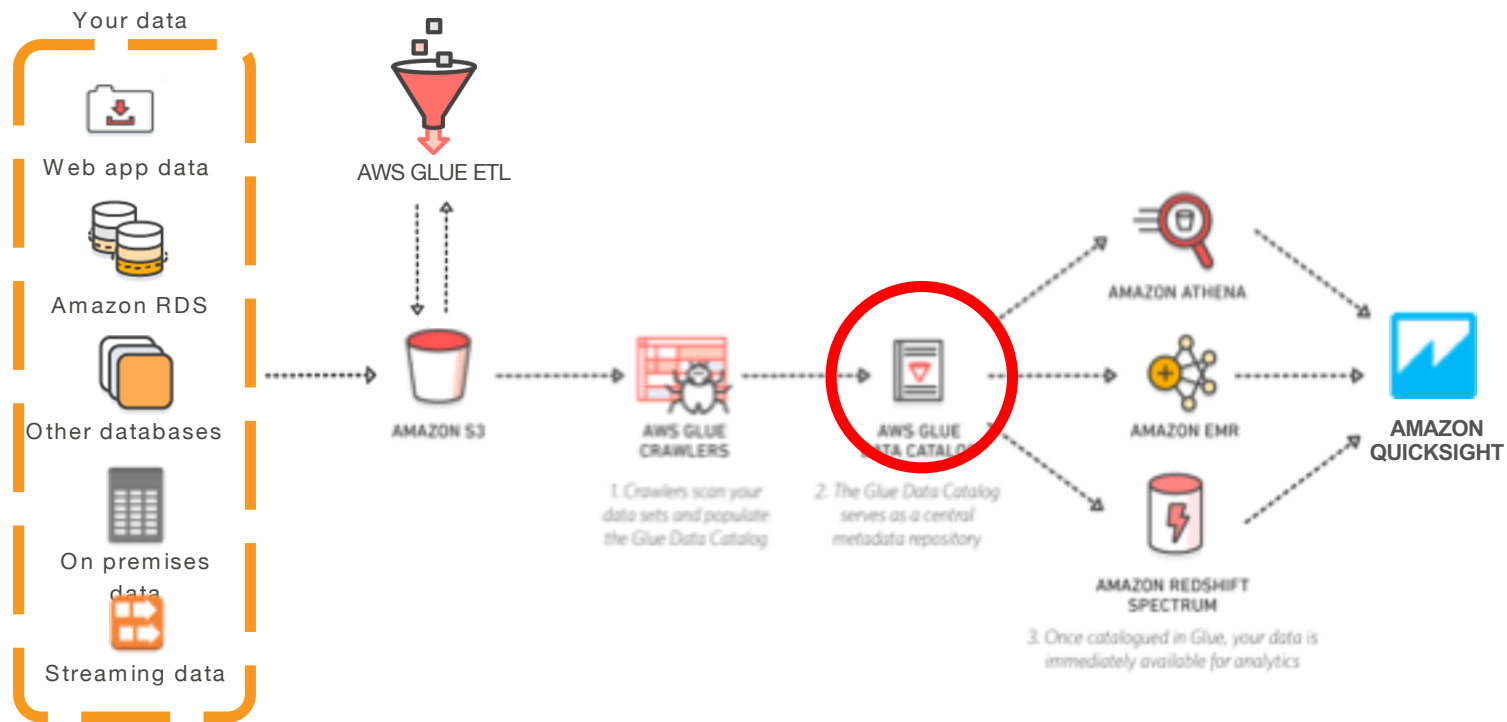
AWS Glue Data Catalog とは何か？

統合されたメタデータリポジトリ

リレーショナル・データベース、Amazon RDS、Amazon Redshift、Amazon S3 に至る(...対応製品は今後も増える予定！)まで

- ❑ データがどこに保存されているとしても**単一の View** を取得可能
- ❑ **検索可能な1つ**の中央リスト内にデータを自動的に**分類する**
- ❑ **スキーマを版管理**することでデータの変化を追跡する
- ❑ Amazon Athena もしくは Amazon Redshift Spectrum を利用してデータを**参照する**
- ❑ **Apache Hive metastore 互換**なので、Amazon EMR で実行するアプリケーションの外部メタストアとして利用可能

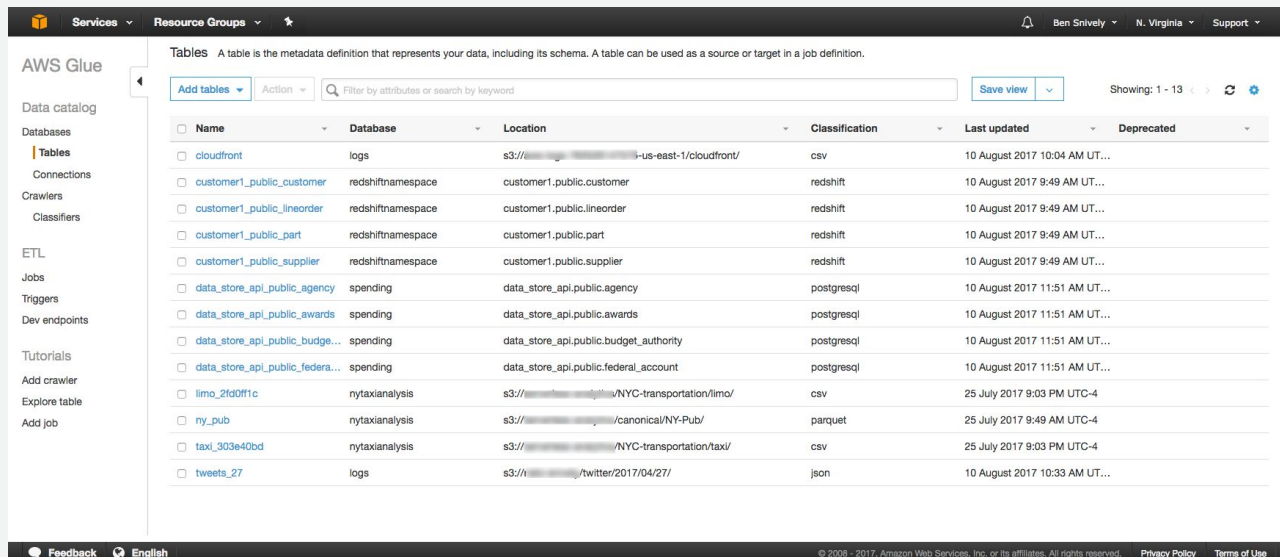
データレイク on Amazon S3 with AWS Glue



Glue Data Catalog を構築するより簡単な方法

1. データが保存されている場所を指定する
2. アップデートを確認する頻度を指定する

これだけで、**Data Catalog** を検索と参照に利用する準備が完了です！



The screenshot shows the AWS Glue Data Catalog console. The left sidebar contains navigation options: Databases, Tables (selected), Connections, Crawlers, Classifiers, ETL, Jobs, Triggers, Dev endpoints, and Tutorials. The main area displays a table of registered tables with columns for Name, Database, Location, Classification, Last updated, and Deprecated. The table lists various tables such as 'cloudfront', 'customer1_public_customer', and 'data_store_api_public_agency'.

Name	Database	Location	Classification	Last updated	Deprecated
cloudfront	logs	s3://[redacted]-us-east-1/cloudfront/	csv	10 August 2017 10:04 AM UT...	
customer1_public_customer	redshiftnamespace	customer1_public.customer	redshift	10 August 2017 9:49 AM UT...	
customer1_public_lineorder	redshiftnamespace	customer1_public.lineorder	redshift	10 August 2017 9:49 AM UT...	
customer1_public_part	redshiftnamespace	customer1_public.part	redshift	10 August 2017 9:49 AM UT...	
customer1_public_supplier	redshiftnamespace	customer1_public.supplier	redshift	10 August 2017 9:49 AM UT...	
data_store_api_public_agency	spending	data_store_api.public.agency	postgresql	10 August 2017 11:51 AM UT...	
data_store_api_public_wards	spending	data_store_api.public.ards	postgresql	10 August 2017 11:51 AM UT...	
data_store_api_public_budg...	spending	data_store_api.public.budget_authority	postgresql	10 August 2017 11:51 AM UT...	
data_store_api_public_federa...	spending	data_store_api.public.federal_account	postgresql	10 August 2017 11:51 AM UT...	
limo_2fd0ff1c	nytaxianalysis	s3://[redacted]/NYC-transportation/limo/	csv	25 July 2017 9:03 PM UTC-4	
ny_pub	nytaxianalysis	s3://[redacted]/canonical/NY-Pub/	parquet	25 July 2017 9:49 AM UTC-4	
taxi_303e40bd	nytaxianalysis	s3://[redacted]/NYC-transportation/taxi/	csv	25 July 2017 9:03 PM UTC-4	
tweets_27	logs	s3://[redacted]/twitter/2017/04/27/	json	10 August 2017 10:33 AM UT...	

Crawlerとは何か？

Crawler は自動的に Data Catalog を構築し、最新状態を維持する

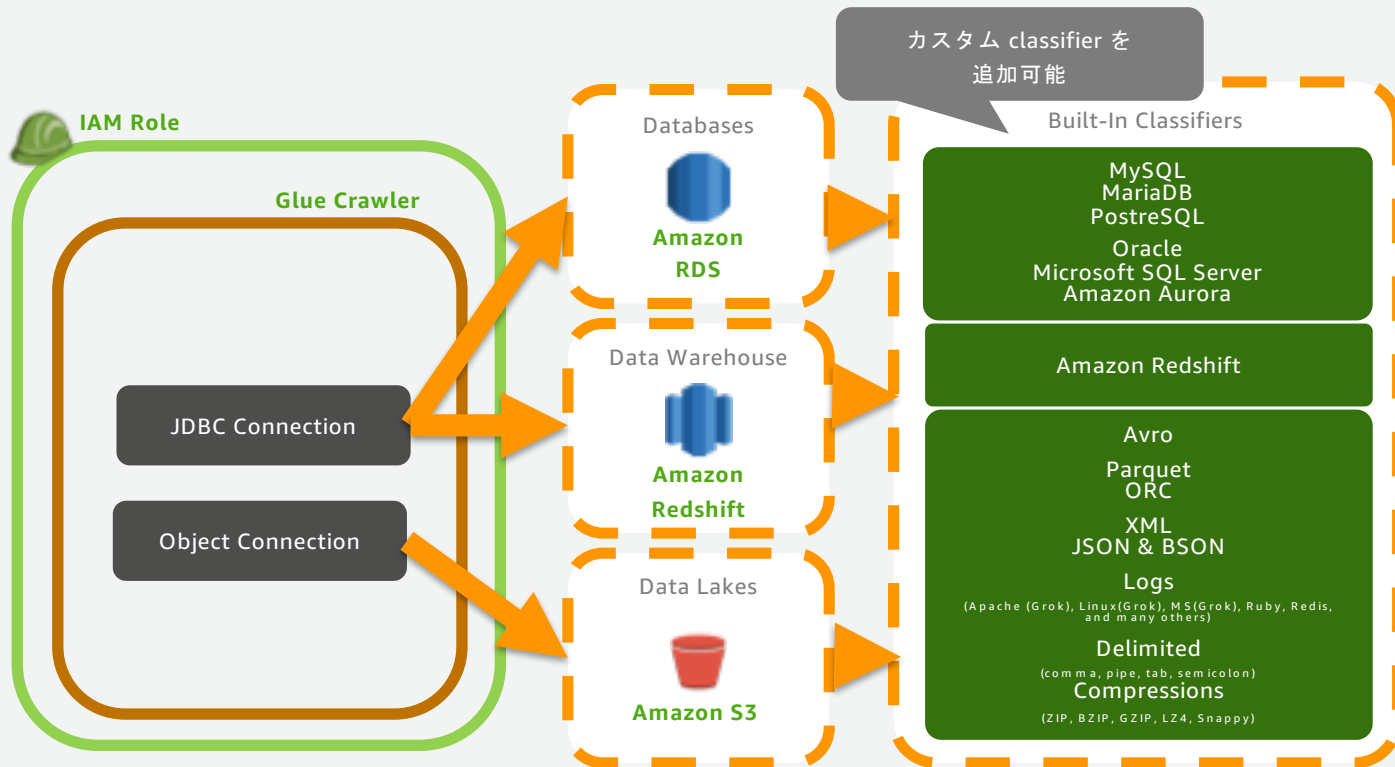
- 様々なデータストアに保存されているデータをスキャンし、メタデータとデータの統計量を抽出し、Data Catalog にテーブル定義を追加
 - ❖ 組み込み/カスタム classifier を利用して、データを分類
 - ❖ Grok 式を利用して独自の classifier を定義することができる
- 新しいデータを発見し、スキーマ定義を抽出
 - ❖ スキーマの変化を検出し、テーブルのバージョンを更新する
 - ❖ Amazon S3 にあるデータから Hive 形式のパーティションを検出する
- 要求に応じて、もしくは、スケジュールに基づいて実行
 - ❖ サーバーレスなので、Crawler が実行されたときだけ課金される

どの様にデータが分類されるのか？

Crawler はデータをスキャンするときにデータに対して classifier の集合を適用し、結果、テーブルとしてメタデータを Data Catalog に追加する

- ❑ **classifier** はデータのフォーマットを認識し、スキーマを生成
- ❑ classifier が返却する分類に合っているかどうかを示す確からしさの値 (0.0 ~ 1.0) を元に、Crawler は分類できたかどうかを判断
- ❑ 順序付けられた classifier の集合を Crawler にセット、Crawler はマッチするまで、提供された順に classifier を実行する
- ❑ Glue はクロールの際に組込済みの classifier を利用できるだけでなく、**独自のカスタム classifier を定義することも可能！**

Crawler が分類できるものは何か？



カスタム classifier を定義する方法

- ❑ **Grok** パターンと一致したスキーマに割り当てる分類ラベルを設定することで、独自のカスタム **classifier** を定義することが可能
- ❑ **Grok** パターンは、一度に1行ずつデータを照合するために使用される正規表現の名前付き集合
- ❑ **Example:**
`%{TIMESTAMP_ISO8601:timestamp}`
`¥[%{MESSAGEPREFIX:message_prefix}¥]`
`%{CRAWLERLOGLEVEL:loglevel} :`
`%{GREEDYDATA:message}`

Classifier name

Classification

Describes the format of the data classified or a custom label.

Grok pattern

Built-in and custom named patterns used to parse your data into a structured schema. For more information, see the [list of built-in patterns](#).

Custom patterns

1	CRAWLERLOGLEVEL (BENCHMARK ERROR WARN INFO TRACE)
2	MESSAGEPREFIX .*-.*-.*-.*-.*
3	

Optional custom building blocks for the grok pattern.

カスタム classifier

1. カスタム classifier を定義する

Classifier name

Classifier type

Grok XML

Classification

Describes the format of the data classified or a custom label.

Grok pattern



2. Crawler にカスタム classifier を追加する

Classifiers infer the schema of your data. The first classifier in the list of custom classifiers to recognize your data is used. Subsequent classifiers are skipped. Built-in classifiers are used if you do not supply a classifier that matches.

Custom classifiers

Showing: 1 - 2 < >

Classifier	Classificatic	
Id Crawler...	crawlerlogs	Add
MyCusto...	MyLogFo...	Add

Selected classifiers

Showing: 1 - 1 < >

Classifier	Classificatic	
Id Crawler logs	crawlerlogs	x

Tables A table is the metadata definition that represents your data, including its schema. A table can be used as a source or target in a job definition.

<input type="checkbox"/> Name	Database	Location	Classification	Last updated	Deprecated
<input type="checkbox"/> exportedlogs	mys3	s3	/exportedlogs/	27 October 2017 12:22 PM U...	

Crawlers: 自動的にスキーマを推測する

enumerate
S3 objects

file 1

file 2

...

file N

identify file type
and parse files

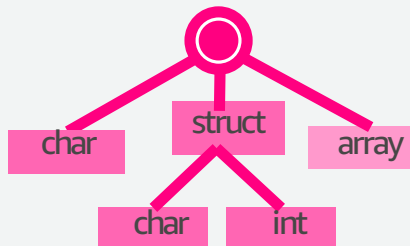
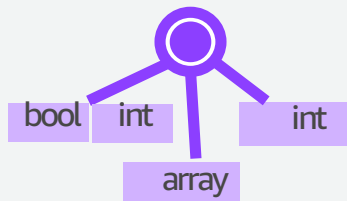
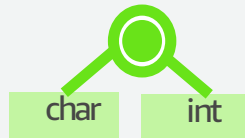
custom classifiers

Grok based parser

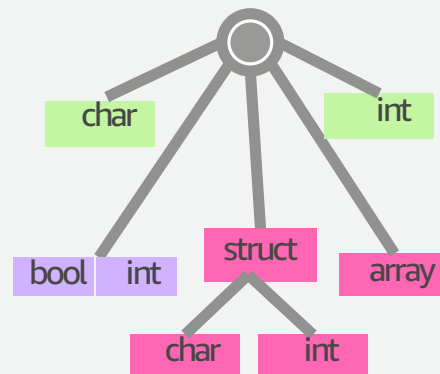
built-in classifiers

JSON parser
CSV parser
Parquet parser
...

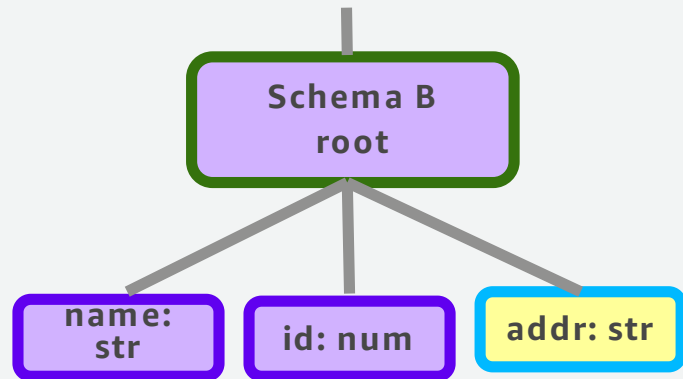
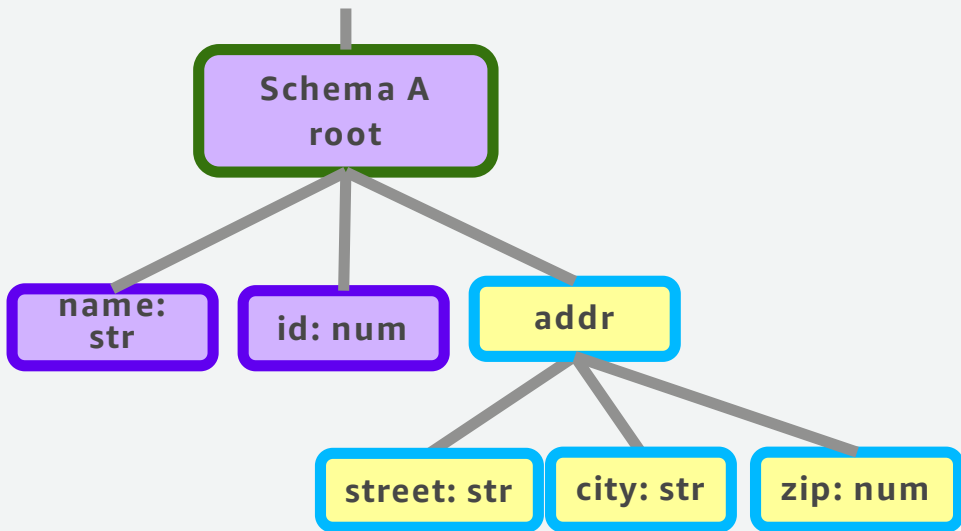
semi-structured
per-file schema



semi-structured
unified schema



スキーマの類似度の検出



Schema similarity heuristic

- 名前が一致したら +1 point
- データ型が一致したら +1 point
- $sim > 0.7$ だったらマッチ

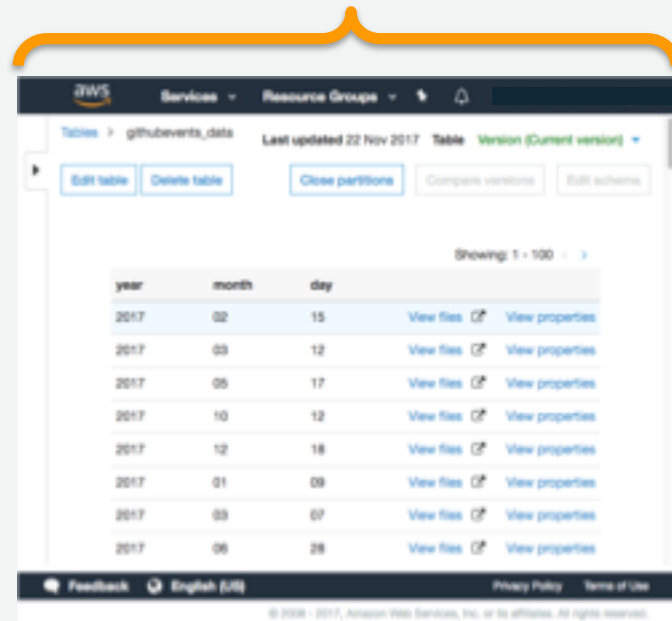
$$sim = \frac{\text{intersection}}{\min(A,B)} = \frac{7}{8} = .875$$

→ マッチ！

自動パーティション検出

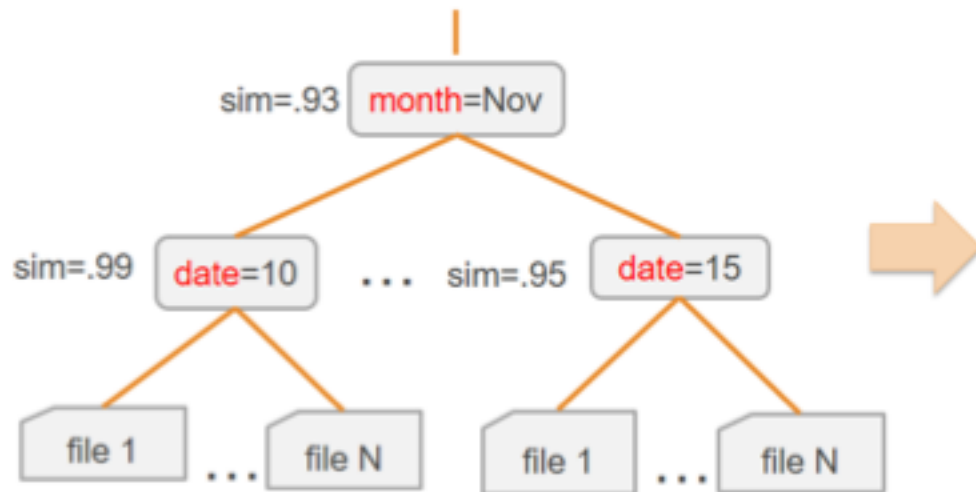


利用可能なパーティション



パーティションを検出する仕組み

S3 バケット階層構造



テーブル定義

Column	Type
month	str
date	str
col 1	int
	float
⋮	⋮

半構造化ログ／スキーマの進化を処理するために、各レベルのファイル間のスキーマ類似性を見積る

自動スキーマ版管理

データが進化すると自動的にテーブルのバージョンが更新される

The screenshot displays two versions of a table schema in the AWS Glue console. An orange bracket at the top indicates the transition from Version 1 to Version 2.

Version 1 (Left): Last updated 21 Aug 2017. Serde parameters: paths, entities, id, retweeted, text, user. Table properties: mycustom abc, CrawlerSchemaSerializer/Version 1.0, recordCount 1001, averageRecordSize 456, CrawlerSchemaDeserializer/Version 1.0, compressionType none, typeOfData file.

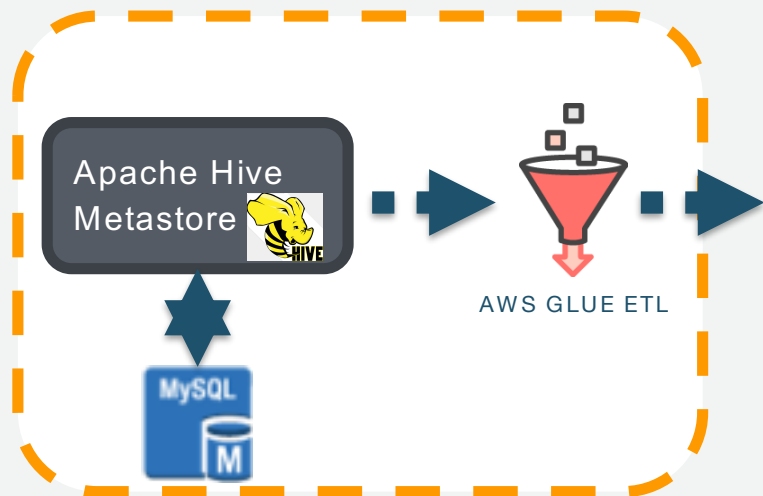
Change	Column name	Data type	Key
	id	bigint	
	retweeted	boolean	
	text	string	
	user	struct	

Version 2 (Right): Last updated 25 Nov 2017. Serde parameters: paths, entities, id, retweeted, text, user. Table properties: mycustom abc, CrawlerSchemaSerializer/Version 1.0, recordCount 1001, averageRecordSize 456, CrawlerSchemaDeserializer/Version 1.0, compressionType none, typeOfData file.

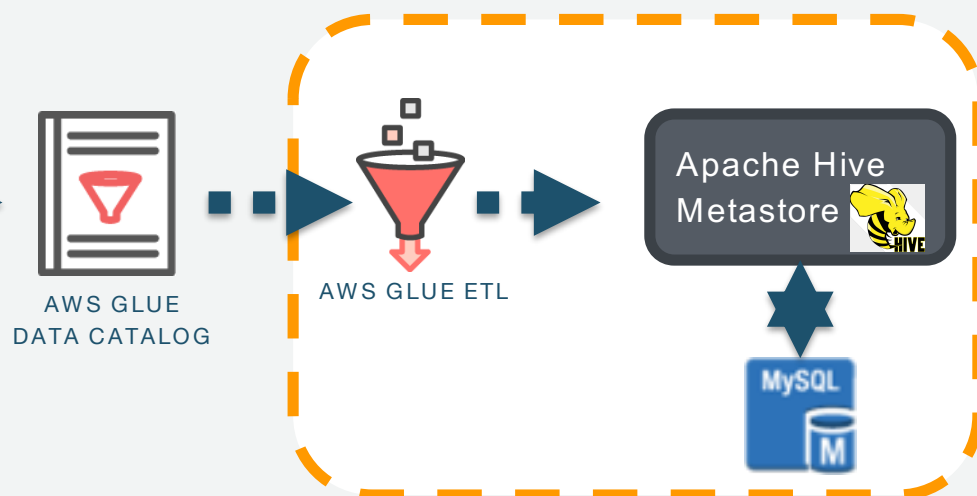
Change	Column name	Data type	Key
	id	bigint	
	retweeted	boolean	
	text	string	
	user	struct	
Added	url	string	

メタデータの Import/Export

Import from an external metastore



Export to an external metastore



import/export ETL スクリプト

https://github.com/aws-labs/aws-glue-samples/tree/master/utilities/Hive_metastore_migration

データレイクのデータが **AWS Glue** によって、
自動的／効率的に分析可能に

…その次は？

データをすばやく発見する

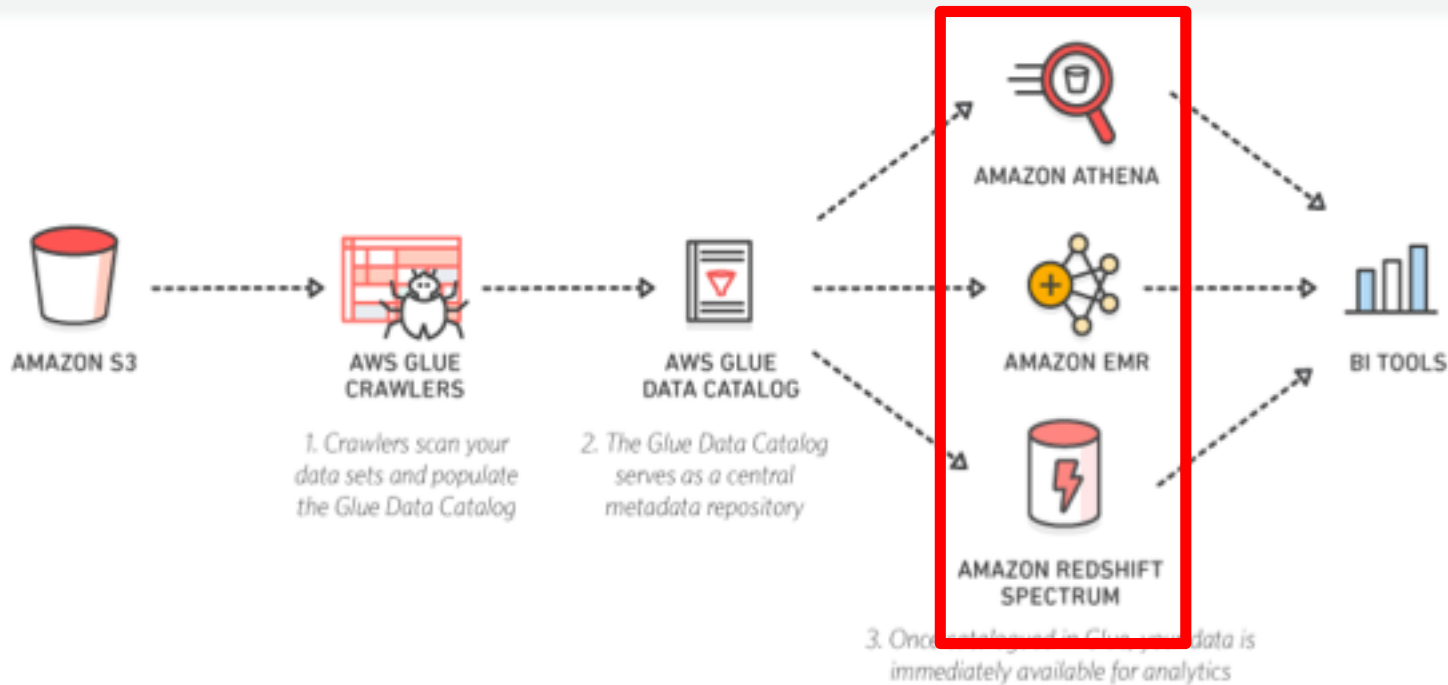
テーブル属性でのフィルタリング、もしくは、
キーワード検索を実行可能

検索結果を保存し、後から参照可能

The screenshot shows the Amazon Athena console interface. At the top, there's a header with the text "Tables A table is the metadata description that represents your data, including its schema. A table can be used as a source or target in a job definition." Below this, there are several buttons: "Add tables", "Action", "Search" (with a search box containing "log"), "Filter or search for tables...", "Save view", and "Showing: 1 - 3". A table list is displayed with columns for Name, Location, and other details. The "Action" menu is open, showing options like "Edit table details", "View details", "View data", and "Delete table". An orange arrow points from the "View data" option to a query editor and results view. The query editor shows a SQL query: `select * from cloudtrail.parquettrail where eventtime > '2017-10-23T12:00:00Z' AND eventtime < '2017-10-23T13:00:00Z' order by eventtime asc;`. The results view shows a table with columns: eventversion, eventId, eventtime, sharedeventid, requestparameters, durationseconds. The table contains 10 rows of data.

Amazon Athena でデータ参照可能
(Amazon Athena へのショートカット)

異なるエンジンで同じデータを分析する



Amazon Athena

- 標準 SQL を使用して Amazon S3 でデータを分析するインタラクティブなクエリサービス
- 設定または管理するインフラはなく、ロードするデータもない

Query Instantly



セットアップにコストは不要、Amazon S3 を指定するだけですぐにデータにクエリを実行可能

Pay per query



クエリ課金／データ圧縮、パーティショニング、列指向フォーマットを利用することで、1 クエリあたりの料金を 30～90% 節約とパフォーマンスの向上が可能

Open



ANSI SQL 準拠のインタフェース、JDBC/ODBC ドライバ、標準データフォーマット、圧縮、そして、複雑な join 処理に対応

Easy



サーバレスで、サーバの設定や管理は不要、Amazon QuickSight (BI) と統合されている

Amazon EMR

- ❑ 20 のオープンソースプロジェクトによるスケーラブルな分析と機械学習が可能
- ❑ Apache Spark、Apache Hive、Presto 用に AWS Glue Data Catalog と統合
- ❑ エンタープライズグレードのセキュリティ

Latest versions



リリース後30日
以内に最新のオー
プンソースフレー
ムワークで更新

Low cost



秒課金、EC2スポット
インスタンス、リザー
ブドインスタンス、
オートスケーリングを
利用した柔軟な課金体
系でコストを50-80%
削減可能

Use S3 storage



Amazon S3 に構築
したデータレイク
のデータを EMRFS
コネクタによるハ
イパフォーマンス
で直接、セキュア
に処理可能

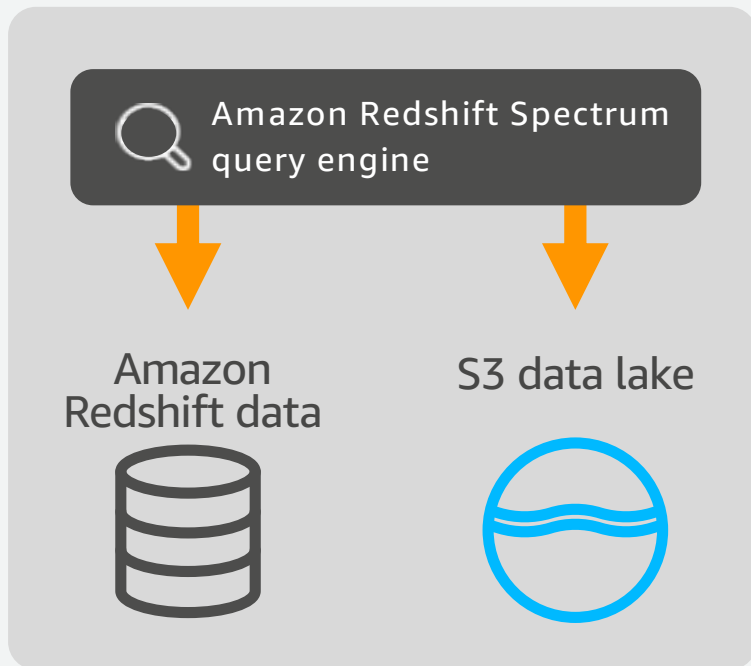
Easy



数分で完全マネージド
な Apache Hadoop &
Apache Spark を起動
可能。クラスタセット
アップ、ノードプロビ
ジョニング、クラスタ
チューニング不要

Amazon Redshift Spectrum

- S3 データレイクにデータウェアハウスを拡張する



S3 に対してエクサバイトクラスの Amazon Redshift SQL クエリを実行可能

Redshift と S3 をまたいだデータ結合

計算リソースとストレージを別々にスケール可能

安定したクエリのパフォーマンスと無制限の同時実行性

Parquet, ORC, Grok, Avro, CSV などのフォーマットに対応

スキャンしたデータ量に対するクエリ課金

データレイクとデータウェアハウス
に入っている様々な規模のデータを
効率的に分析する方法

→ **Amazon Redshift Spectrum**

Amazon Redshift – Data Warehousing

- 1/10 のコストで、高速で、強力で、シンプルで、完全に管理されたデータウェアハウス
- 大規模並列、ギガバイトからペタバイトまで拡大

Fast at any scale



I/O 効率を向上させる列指向ストレージテクノロジーの使用、および、複数ノード間のクエリ並列化により、高速クエリパフォーマンスを実現

Open file formats



最新 SSD 上で最適化されたデータフォーマット、Amazon S3 にある全てのオープンデータフォーマットを分析可能

Secure



すべてを監査; データをエンドツーエンドで暗号化; 豊富な認定とコンプライアンス

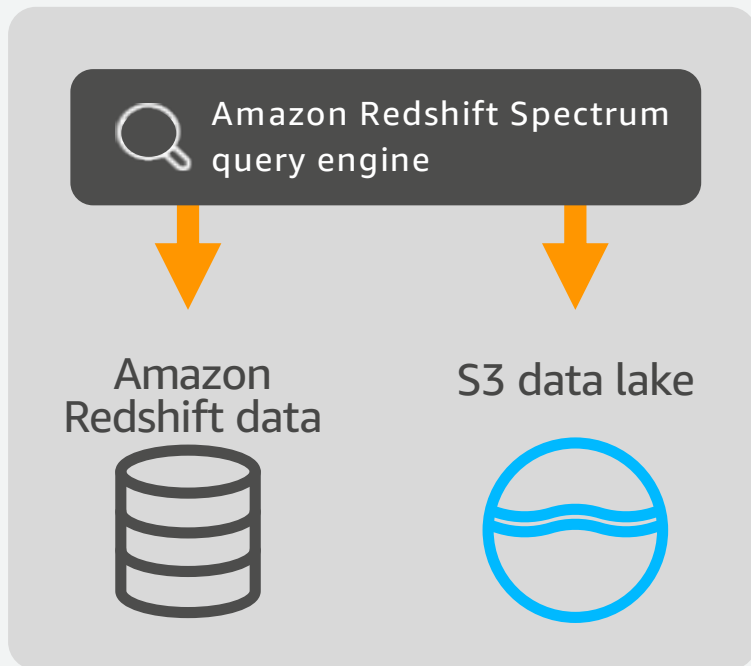
Inexpensive



年間 1 テラバイトあたりわずか **1,000 USD**。ウェアハウスソリューションのコストを従来の **1/10** に抑えることが可能

Amazon Redshift Spectrum (再掲)

- S3 データレイクにデータウェアハウスを拡張する



S3 に対してエクサバイトクラスの Amazon Redshift SQL クエリを実行可能

Redshift と S3 をまたいだデータ結合

計算リソースとストレージを別々にスケール可能

安定したクエリのパフォーマンスと無制限の同時実行性

Parquet, ORC, Grok, Avro, CSV などのフォーマットに対応

スキャンしたデータ量に対するクエリ課金

Amazon Redshift Spectrum アーキテクチャ

超並列、共有なしの列指向アーキテクチャ

❑ Leader Node

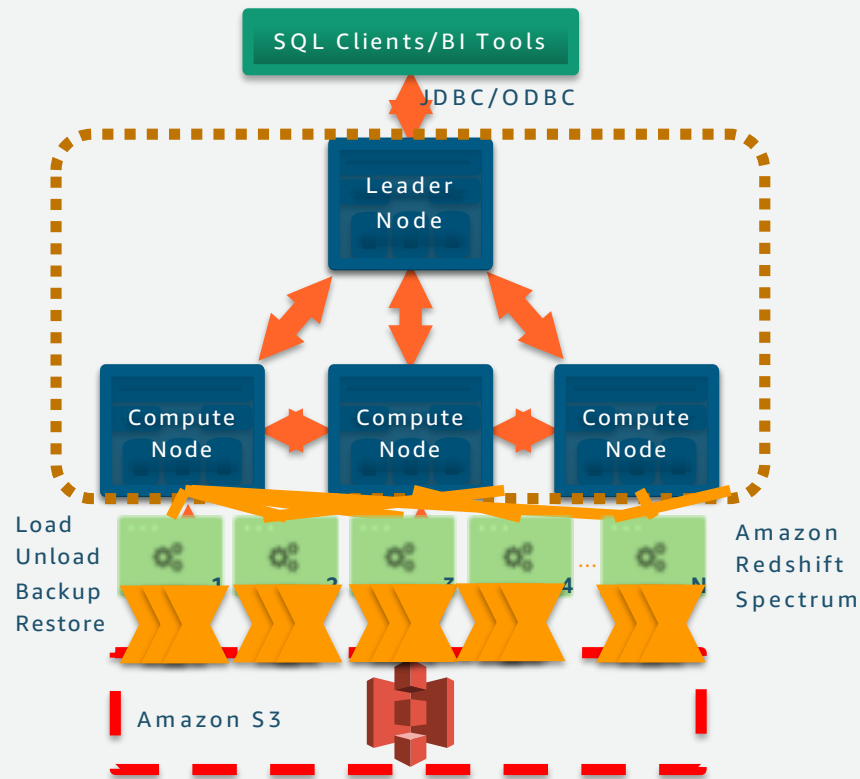
- ❖ SQL エンドポイント
- ❖ メタデータを保存
- ❖ 並列クエリ処理をコーディネート

❑ Compute Node

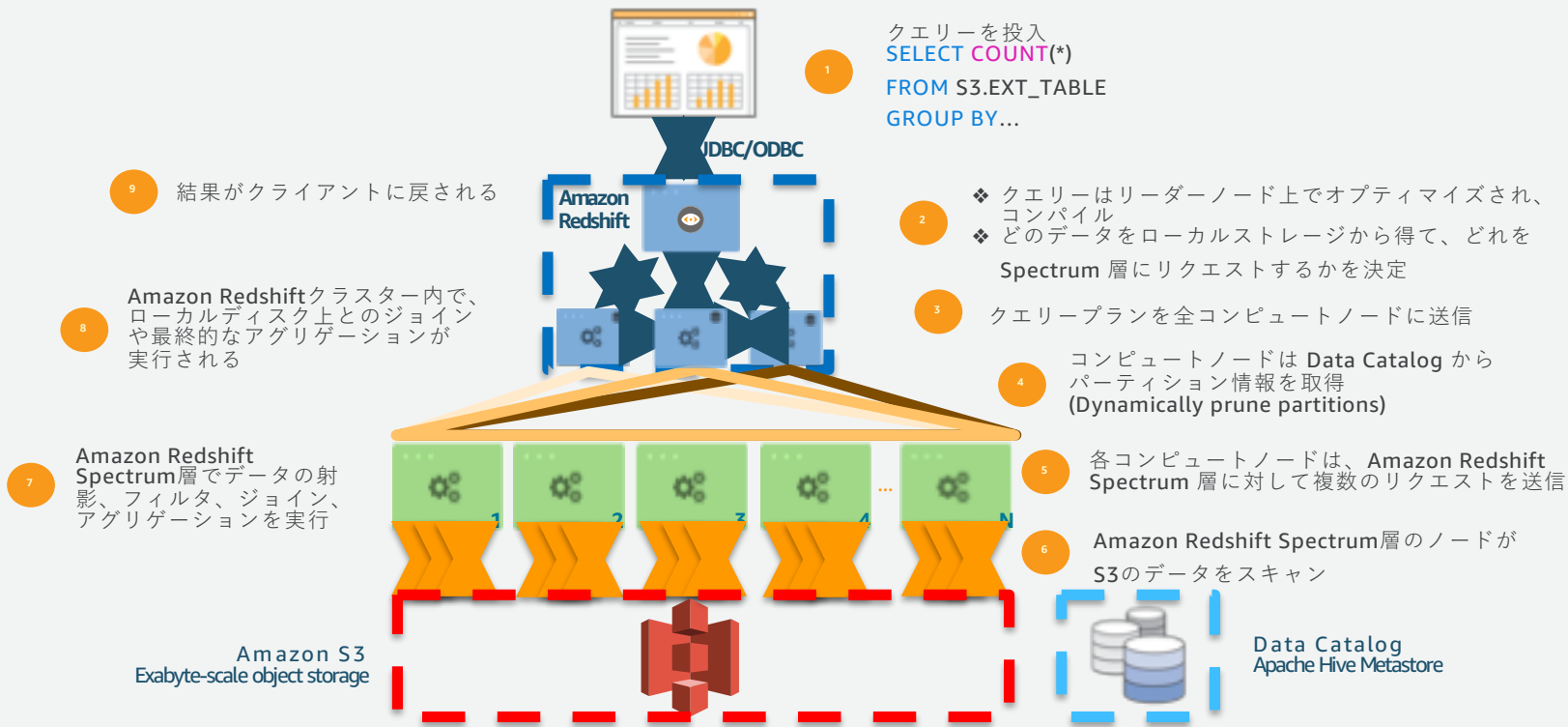
- ❖ ローカル列指向ストレージ
- ❖ 並列にクエリを実行
- ❖ データの load / unload / backup / restore

❑ Amazon Redshift Spectrum Node

- ❖ Amazon S3 に対して直接クエリを実行
- ❖ Redshift Spectrum は数千インスタンスにまで自動的に拡張し、エクサバイトのデータに対してさえもクエリは高速に動作



Amazon Redshift Spectrum クエリ実行の流れ



データレイクアーティファクトの定義 (Schema on Read)

Data Catalog を利用して Amazon Redshift に外部スキーマを定義する

```
CREATE external schema archived_trips
from data catalog database 'sampledb'
iam_role 'arn:aws:iam::123456789012:role/MySpectrumRole'
region 'us-east-2';
```

外部スキーマを参照する

```
select * from svv_external_schemas
```

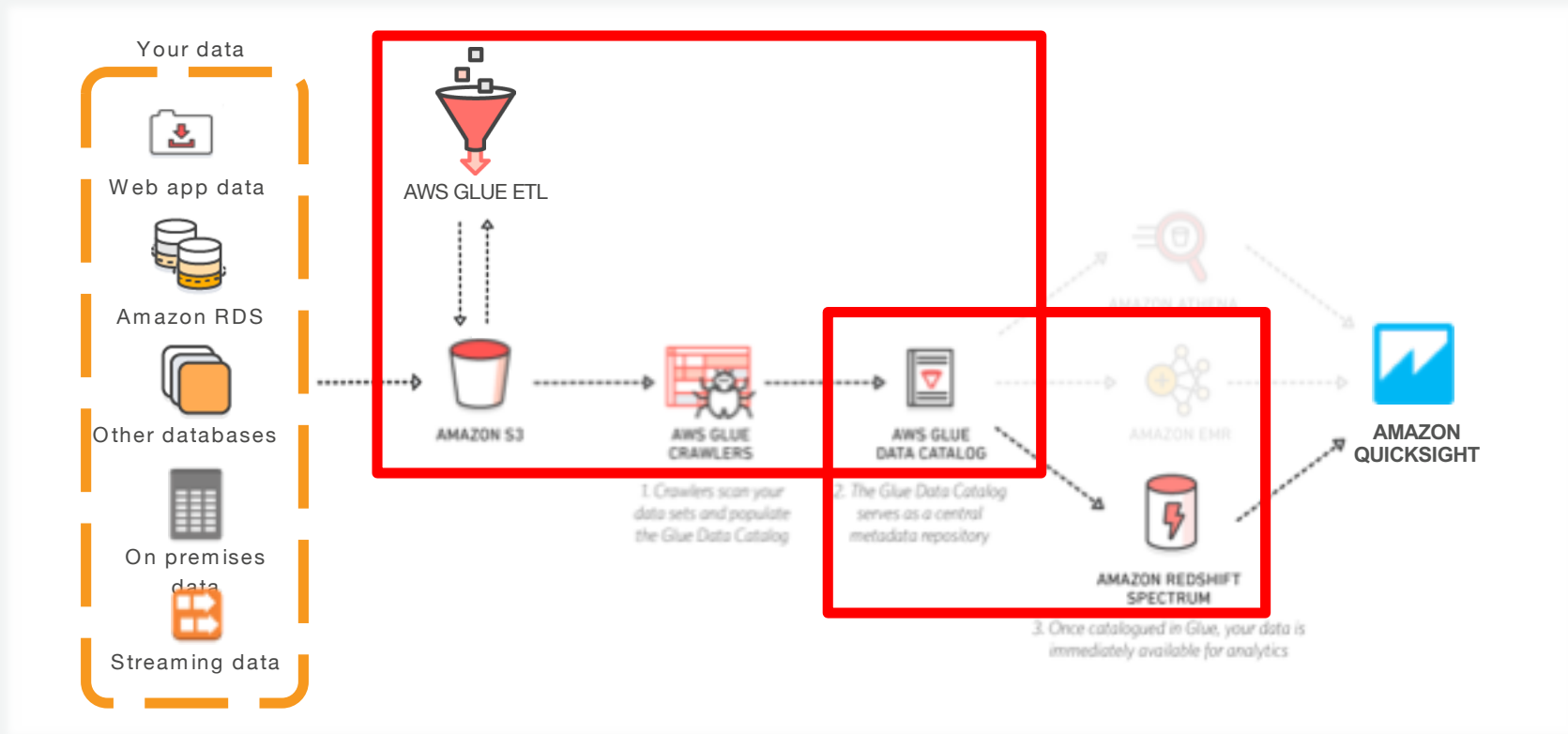
外部テーブルを参照する

```
select * from svv_external_tables
```

権限設定

- Amazon Redshift は AWS Glue にある Data Catalog と Amazon S3 にあるデータファイルにアクセスするための権限を必要とする
- 権限を与えるため、まず最初に AWS Identity and Access Management (IAM) ロールを生成する必要がある
- それから、クラスターにロールをアタッチし、Amazon Redshift の外部スキーマ生成文の中で、ARN(Amazon Resource Name) をロールに指定する

データレイク on Amazon S3 with AWS Glue (再掲)



データレイクとデータウェアハウスに入っている様々な規模のデータを **Redshift Spectrum** を利用して効率的に分析可能

・・・ベストプラクティスは？

ベストプラクティス - 1 / 5

Amazon Redshift Spectrum を使用して、スキャン集約的な同時作業負荷を改善する

- ❑ Redshift Spectrum は、利用している Redshift クラスタとは独立した専用のサーバー群にある
- ❑ フィルター処理や集約処理といった、多くのコンピュータインテンシブな処理を Redshift Spectrum 層で行うことで、クエリが使用する Redshift クラスタの処理キャパシティが大きく削減される

ベストプラクティス - 2 / 5

クエリはデータレイクを最適化する – Apache Parquet を使う

- ❑ Apache Parquet は、データ処理フレームワークやデータモデル、プログラミング言語に依らず利用可能な列指向フォーマット
- ❑ SVL_S3QUERY_SUMMARY テーブルを調べることで、パーティション分けされた Parquet ファイルを使う際の、S3 に関する様々な興味深いメトリクスを確認できる
- ❑ 特に s3_scanned_rows と s3query_returned_rows という 2 つのメトリクスに注目してみると、CSV ファイルを処理するときと比べて、Redshift Spectrum から Redshift クラスターに送られるデータ総量が驚異的に削減されていることがわかる

ベストプラクティス - 3 / 5

クエリはデータレイクを最適化する - Parquet ファイルでパーティションする

- 次のSQLは、パーティションプルーニングの有効性を分析する
- クエリが少数のパーティションにしか触れない場合は、すべてが期待通りに動作しているかどうかを確認できる:

```
SELECT query, segment, max(assigned_partitions) as total_partitions,  
max(qualified_partitions) as qualified_partitions FROM svl_s3partition  
WHERE query=<Query-ID> GROUP BY 1,2;
```

ベストプラクティス - 4 / 5

データレイクに投入するクエリを最適化する

□ Amazon Redshift Spectrum クエリの同時実行性能は、以下の 2 つのレベルで制御可能

- ❖ クエリレベル（クエリごと 1 スライスにつき最大 **10** の同時実行数）
 - いくつのクエリが同時に実行されているかによって、同時実行数が変わる
 - 割り当てられた同時実行数によって、**S3** をスキャンするスレッド数が制限される
- ❖ ノードレベル（ノード上で動作するすべての **S3** をスキャンするクエリに適用される。ノードタイプによって数が異なる）
 - より大きなノードタイプを選択するほど、上限数も高くなる

ベストプラクティス - 5 / 5

Predicate pushdown によるデータレイククエリのパフォーマンスの向上

- AmazonのRedshift Spectrumレイヤーにプッシュダウンできる特定のSQL操作があるので、可能であれば、これらの機能を利用する

例) :

- ❖ GROUP BY 句やいくつかの文字列関数
- ❖ 等価述語や LIKE のようなパターンマッチ条件
- ❖ COUNT/SUM/AVG/MIN/MAX/その他多くの共通集約関数
- ❖ Regex_replace 等の関数

ベストプラクティス - 5 / 5 (Cont.)

Predicate pushdown によるデータレイククエリのパフォーマンスの向上

- ❑ DISTINCT や ORDER BY のような特定の SQL 操作は、Amazon Redshift Spectrum にプッシュダウンできないため、Amazon Redshift で実行される。それらの使用を最小限に抑え、できるだけ使用を避ける

例) :

- ❖ DISTINCT を GROUP BY で置き換える

Amazon Redshift Spectrum 10 のベストプラクティス

<https://aws.amazon.com/jp/blogs/news/10-best-practices-for-amazon-redshift-spectrum/>

まとめ

まとめ

- データレイクにデータを集め始めてから直面する課題の解決に利用可能な**2つ**の方法
 1. データレイクのデータを自動的／効率的に分析可能にする方法
→ **AWS Glue**
 2. データレイクとデータウェアハウスに入っている様々な規模のデータを効率的に分析する方法
→ **Amazon Redshift Spectrum**

Data Lake on AWS



データレイクがビッグデータ
ストレージソリューションと
して最大限の柔軟性を提供！

<https://d1.awsstatic.com/whitepapers/Storage/data-lake-on-aws.pdf>

参考資料

- ❑ AWS Glue
 - ❖ <https://aws.amazon.com/jp/glue/>
- ❑ AWS Glue ドキュメント
 - ❖ <https://aws.amazon.com/jp/glue/details/>
- ❑ AWS Glue 開発者用リソース
 - ❖ <https://aws.amazon.com/jp/glue/developer-resources/>

- ❑ Amazon Redshift
 - ❖ <https://aws.amazon.com/jp/redshift/>
- ❑ Amazon Redshift 開発者用リソース
 - ❖ <https://aws.amazon.com/jp/redshift/developer-resources/>

Q&A



オンラインセミナー資料の配置場所

AWS クラウドサービス活用資料集

- <https://aws.amazon.com/ip/aws-ip-introduction/>



Amazon Web Services ブログ

- 最新の情報、セミナー中のQ&A等が掲載されています。
- <https://aws.amazon.com/ip/blogs/news/>

公式Twitter/Facebook AWSの最新情報をお届けします



@awscloud_jp



検索

もしくは

<http://on.fb.me/1vR8vWm>

最新技術情報、イベント情報、お役立ち情報、
お得なキャンペーン情報などを日々更新しています！

AWSの導入、お問い合わせのご相談

AWSクラウド導入に関するご質問、お見積、資料請求をご希望のお客様は以下のリンクよりお気軽にご相談下さい。

<https://aws.amazon.com/ip/contact-us/aws-sales/>

お問い合わせ

日本担当チームへのお問い合わせ >

関連リンク

フォーラム

日本担当チームへのお問い合わせ

AWS クラウド導入に関するご質問、お見積り、資料請求をご希望のお客様は、以下のフォームよりお気軽にご相談ください。平日営業時間内に日本オフィス担当者よりご連絡させていただきます。

※ご請求金額またはアカウントに関する質問はこちらからお問い合わせください。
※Amazon.com または Kindle のサポートに問い合わせはこちらからお問い合わせください。

アスタリスク (*) は必須情報となります。

姓*

名*

※ 「AWS お問い合わせ」 で検索して下さい。

AWS Well Architected 個別技術相談会お知らせ

- Well Architectedフレームワークに基づく数十個の質問項目を元に、お客様がAWS上で構築するシステムに潜むリスクやその回避方法をお伝えする個別相談会です。

<https://pages.awscloud.com/well-architected-consulting-ip.html>

- 参加無料
- 毎週火曜・木曜開催

【毎週火、木曜開催】AWS Well-Architected 個別技術相談会

AWS 上で構築するシステムのリスクの把握・回避方法をご希望のお客様

この度 AWS をご活用しているお客様を対象に「AWS Well-Architected 個別技術相談会」を開催致します。

Well-Architected 個別技術相談会では、リスクの把握・回避を目的として、セキュリティ・信頼性・パフォーマンス・コスト・運用の5つの観点で、お客様の AWS 活用状況や構成についてお話しします。AWS のベストプラクティスに基づき作成された Well-Architected フレームワークを元に、今までお客様がお気づきでなかったリスクやAWS活用の改善点を見つけることができます。例えば、自動化においては納品前段階、業務を定期的に行うのと同時に、本相談会はおお客様の AWS 上のシステムをよりよく活用頂くことを目的としております。

▶ 説明資料(PDF) [AWS Well-Architected Framework -クラウド設計 -運用ベストプラクティスの活用-]

Well-Architected 個別技術相談会にご参加頂くには、本ページにてお申し込み後、弊社担当者からお送りするアサインシートにご記入・担当者にご送付頂く必要があります。その内容を元に、当日の相談会では AWS のソリューションアーキテクトと共に技術的なディスカッションをさせて頂きます。また、遠方のお客様、アマゾン東京オフィスへのご来社が困難の場合は録画配信は、Web のプレイングマネージャーツールや、双方向共有のチャットでの開催となります。



下記のフォームよりお申込みください。

• 姓:

• 名:

AWS Black Belt Online Seminar 配信予定



2018年6月の配信予定

~~6/06 (水) 18:00-19:00 AWS 認定取得に向けて~~

~~6/12 (火) 12:00-13:00 AWS で実現するライブ動画配信とリアルタイムチャットのアーキテクチャパターン →~~

~~6/13 (水) 18:00-19:00 AWS Cloud9入門~~

6/19 (火) 12:00-13:00 データレイク入門：AWS で様々な規模のデータレイクを分析する効率的な方法

6/20 (水) 18:00-19:00 AWS AWS Support

6/26 (火) 12:00-13:00 AWS Summit 2018 振り返り

6/27 (水) 18:00-19:00 Amazon Alexa Skills

お申し込みサイト：<https://aws.amazon.com/jp/about-aws/events/webinars/>

「AWS セミナー」で検索

AWS Black Belt Online Seminar 配信予定



2018年7月の配信予定

7/03 (火) 12:00-13:00 Amazon Neptune

7/04 (水) 18:00-19:00 Amazon Elastic File System (Amazon EFS)

7/10 (火) 12:00-13:00 AWSで実現するウェブサイトホスティング

7/11 (水) 18:00-19:00 Trusted Advisor

7/17 (火) 12:00-13:00 大阪ローカルリージョンの活用とAWSで実現するDisaster Recovery

7/18 (水) 18:00-19:00 AWS Service Catalog

7/24 (火) 12:00-13:00 Amazon Elastic Container Service for Kubernetes (Amazon EKS) / AWS Fargate

7/25 (水) 18:00-19:00 AWS Systems Manager

7/31 (火) 12:00-13:00 S3ユースケース紹介及びサービスアップデート解説

お申し込みサイト: <https://aws.amazon.com/jp/about-aws/events/webinars/>
「AWS セミナー」で検索

簡単なアンケートにご協力下さい

画面に表示されるアンケートフォームに入力をお願いします。

皆様のご意見は、今後の改善活動に活用させていただきます。

コメント欄には1行で自由な内容を書き込み下さい。

例)

- 本オンラインセミナーへのご意見
- 今後オンラインセミナーで取り上げて欲しい題材
- 発表者への激励

等々

※Q&A同様に書き込んだ内容は主催者にしか見えません

ご参加ありがとうございました