



データレイクの活用に必要な運用のポイント

Noritaka Sekiyama

Big Data Architect, AWS Glue & Lake Formation

2020/6/11

AWSオンラインセミナーへようこそ

お気軽にご質問ください！

- 書き込んだ質問は主催者にしか見えません
- 最後のQ&A時間で、いただいたご質問からピックアップしてご回答をさせていただきます

- ① 吹き出しをクリック
- ② 質問を入力
- ③ Sendをクリック



「AWSではじめるデータレイク」 2020年7月発売予定

上原誠 / 志村誠
下佐粉昭 / 関山宜孝

AWSではじめる データレイク

Data Lake starting with AWS

クラウドによる
統合型
データリポジトリ
構築入門

クラウドを軸に 統合データ基盤を作る

企業に関わるすべてのデータを統合する
レポジトリを作るためのノウハウを結集、
構築から運用に関するAWSの
各サービスに対応させて詳細に解説

ログ ネット
メトリック カスタマー
Data Lake
AI サマリ 分析 予測

データメディア

自己紹介

関山 宜孝

Big Data Architect

AWS Glue & Lake Formation

- 5年間 AWS サポートにて技術支援を担当
- 2019年からプロダクト開発チームにジョイン
- GlueとLake Formationに関するユーザーに近い部分の開発を担当



NoritakaS-AWS

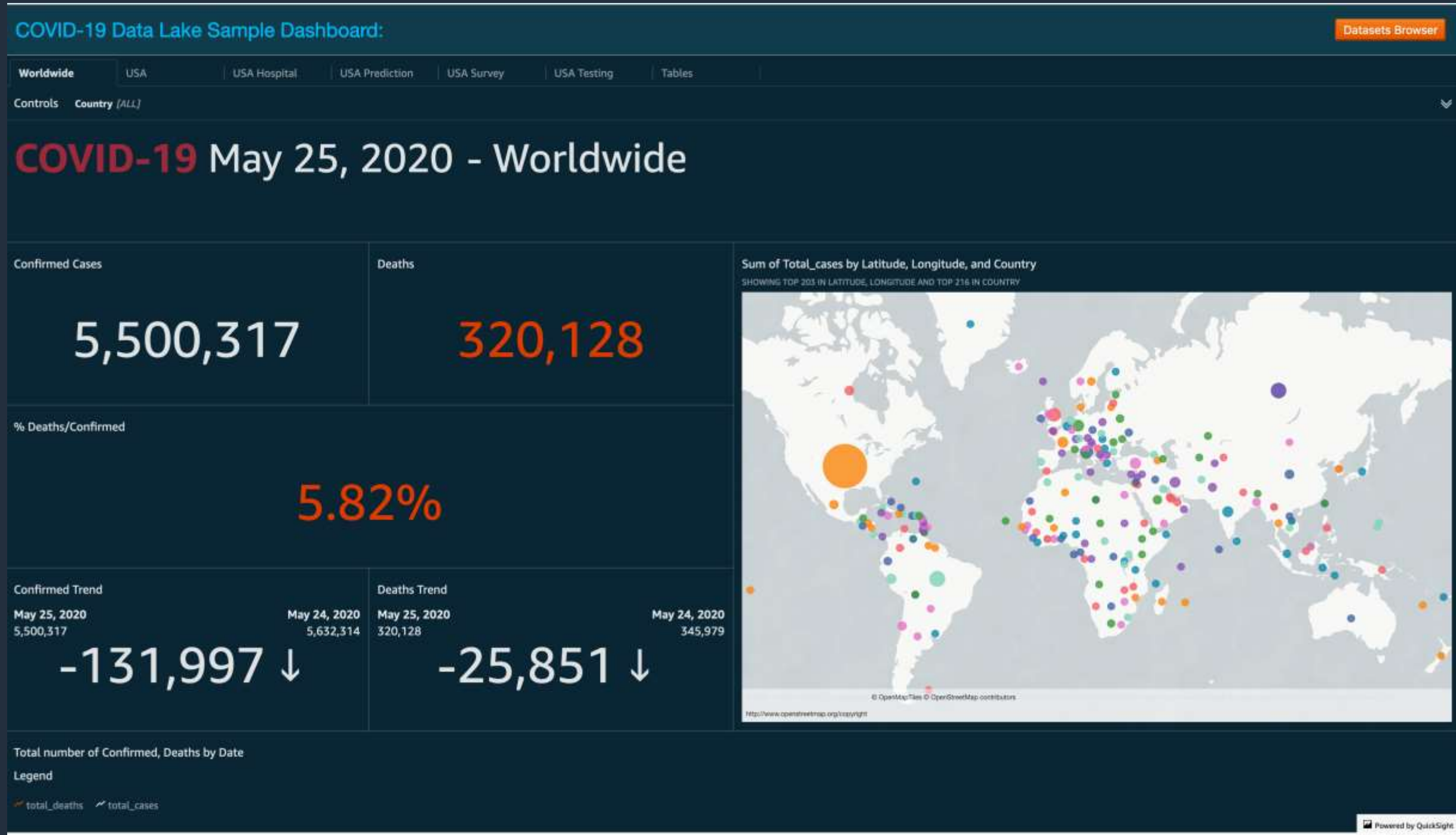


@moomindani

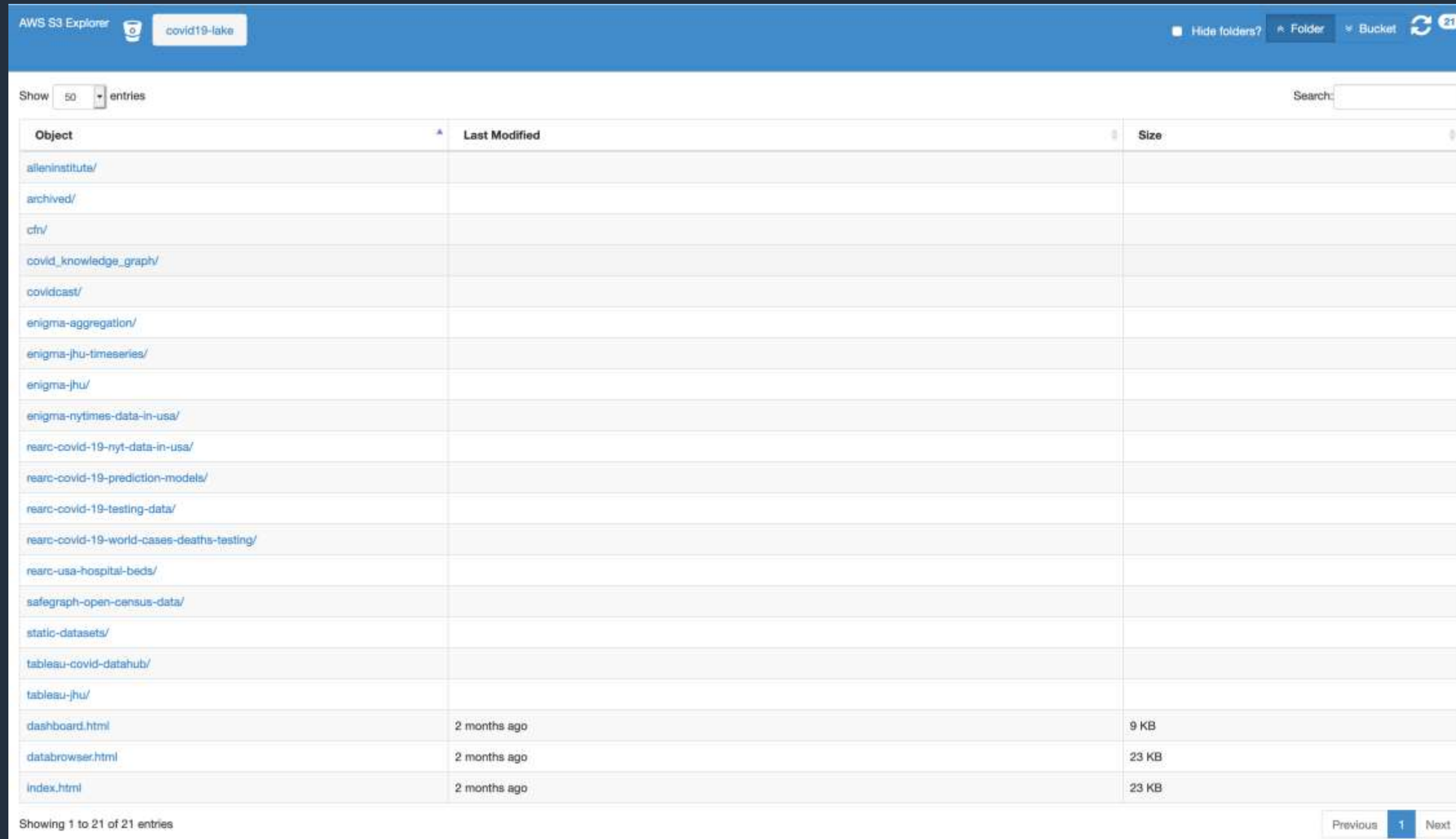


moomindani

サンプル: COVID-19 データレイク



サンプル: COVID-19 データレイク



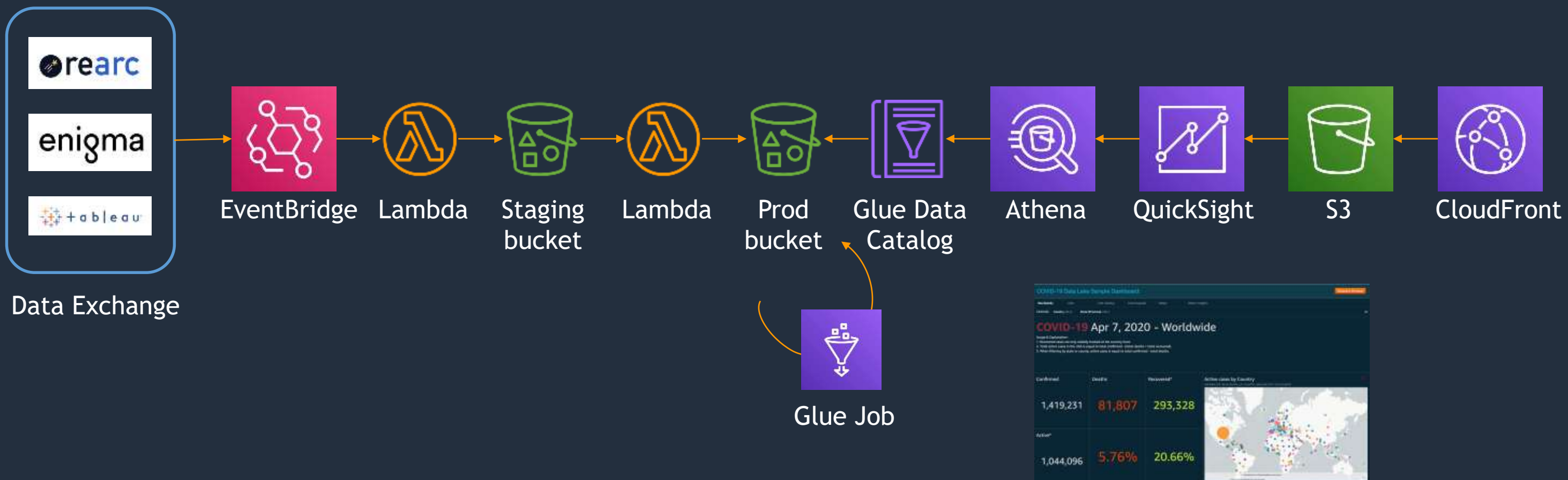
The screenshot shows the AWS S3 Explorer interface for a bucket named 'covid19-lake'. The interface includes a search bar, a 'Show 50 entries' dropdown, and a table of objects. The table has columns for 'Object', 'Last Modified', and 'Size'. The objects listed are folders and files, with the last three being HTML files: 'dashboard.html', 'databrowser.html', and 'index.html', all modified 2 months ago.

Object	Last Modified	Size
alleninstitute/		
archived/		
cfv/		
covid_knowledge_graph/		
covidcast/		
enigma-aggregation/		
enigma-jhu-timeseries/		
enigma-jhu/		
enigma-nytimes-data-in-usa/		
rearc-covid-19-nyt-data-in-usa/		
rearc-covid-19-prediction-models/		
rearc-covid-19-testing-data/		
rearc-covid-19-world-cases-deaths-testing/		
rearc-usa-hospital-beds/		
safegraph-open-census-data/		
static-datasets/		
tableau-covid-datahub/		
tableau-jhu/		
dashboard.html	2 months ago	9 KB
databrowser.html	2 months ago	23 KB
index.html	2 months ago	23 KB

Showing 1 to 21 of 21 entries

Previous 1 Next

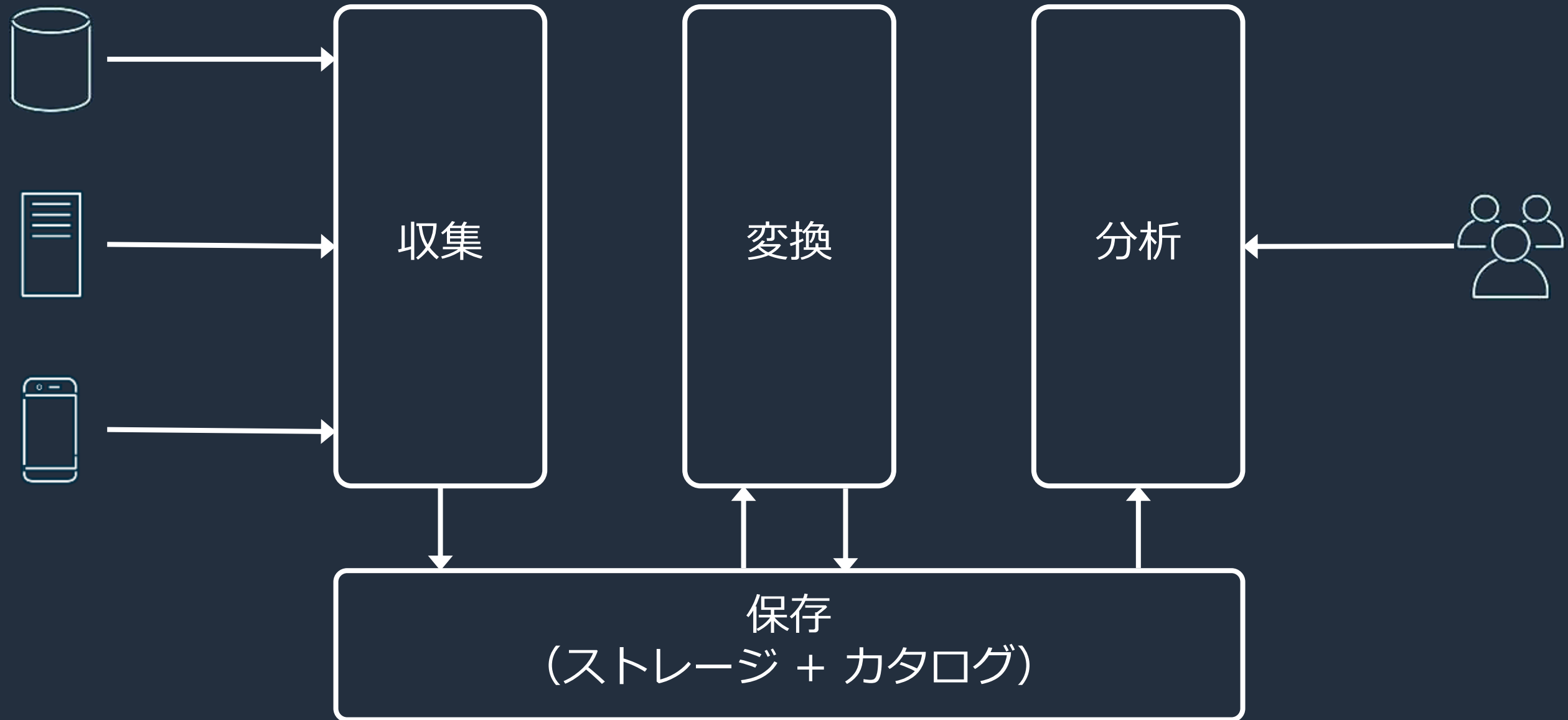
サンプル: COVID-19 データレイクのアーキテクチャ



アジェンダ

- データレイクでよくあるトラブル
- データレイクの運用
- 利用状況の計測とデータレイクの継続的改善
- Q&A

おさらい：データレイクの構築



データレイクでよくあるトラブル

データレイクが完成したら・・・

待望のデータレイク完成！
みなさんたくさん使ってください

監視？運用？
まだあんまり考えてないけど
そのうちなんとかしますね



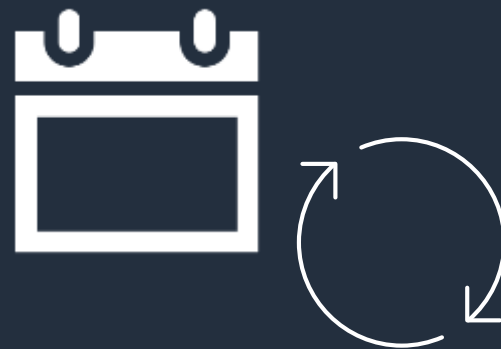
データレイク管理者
(初心者)

よくあるトラブル#1. データが古い

3日前からデータが更新されてなくて
新しいデータを分析できず困っています



データレイクユーザー



データレイク管理者
(初心者)

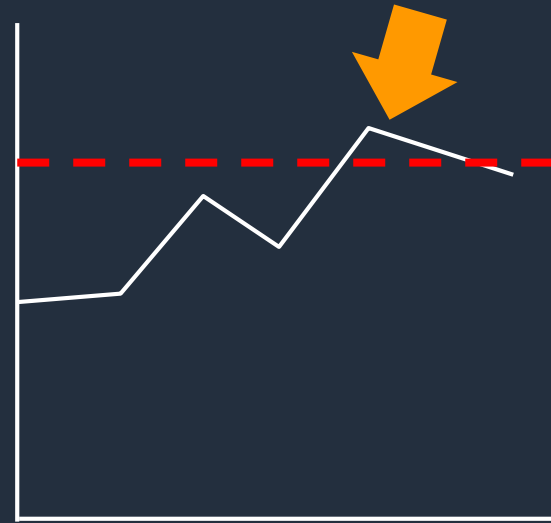
データが届いてない・・・？
ジョブが失敗した・・・？

よくあるトラブル#2. クエリが遅い

クエリしてみたら以前に比べて
50%長く待たされるようになりました



データレイクユーザー



データレイク管理者
● (初心者)

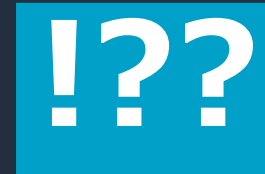
データが予想より増えてた・・・？
でも50%も遅くなる・・・？

よくあるトラブル#3. 想定してない結果が返ってきた

クエリしてみたら本来は日時が格納されるはずのカラムが NULL になっていました



データレイクユーザー



データレイク管理者
• (初心者)

ファイルの中身がおかしい・・・？
データの変換に失敗した・・・？

データレイクと運用

- データレイクは作って終わりではない
 - データレイクは活用されて初めて意味がある
- 適切に運用されていないデータレイクは、まともには使えない
 - 更新されていないデータは価値を失っていく
 - クエリの遅いデータレイクは不便で、急ぎの判断には使えない
 - 信頼性がチェックされていないデータは判断の根拠に使えない

データレイクの活用には適切な運用が欠かせない

データレイクの運用

データレイク、どう運用したらいい？

データレイク、うまいこと運用しといてね



上司



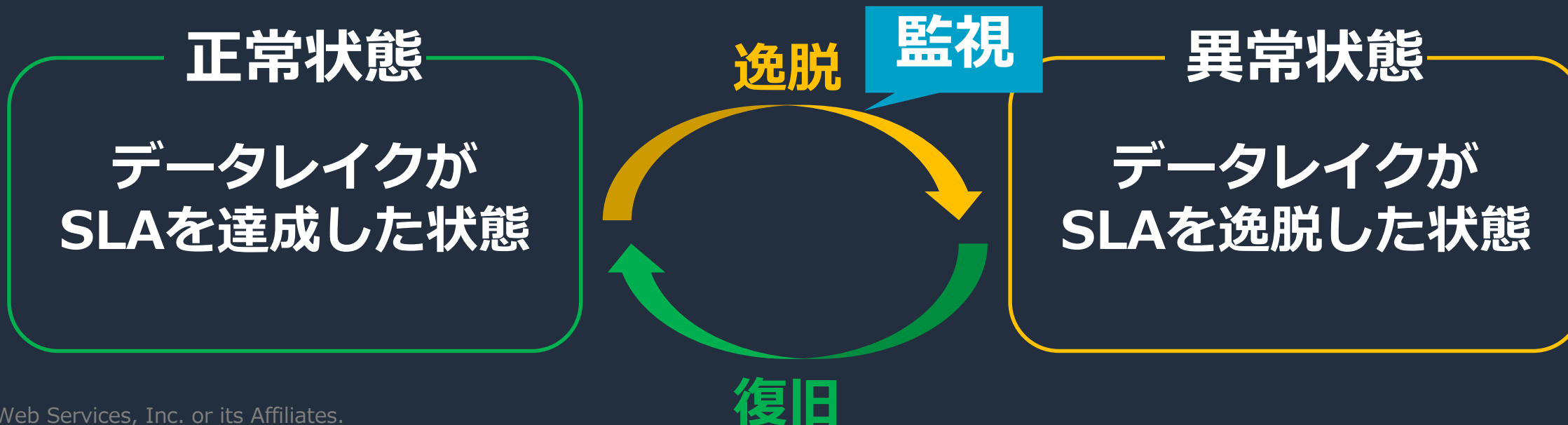
- データレイク管理者 (初心者)

データレイクといってもITシステムだし
CPU使用率とかディスク使用率とか
従来通りのサーバー監視を
設定しておけばいいよね・・・？

データレイクの運用のポイント

正常状態を定義し、その状態を維持する

- データレイクがどのような状態であれば正常かを定義する
 - 正常性指標 (SLA: Service Level Agreement)
- データレイクが正常である状態を維持し、逸脱した場合は正常に戻す



データレイクの正常性定義と SLA の例

管理対象	正常状態	SLA の例
データの鮮度	データ発生から利用可能になるまでの時間が規定内	アプリ上で出力されたログデータが XX時間以内に利用可能 となっている
	所定の日時までに対象データが利用可能となる	レポートニング用のデータの前処理が 午前xx時まで に完了している
データの信頼性	データ起因のエラーや結果不正の発生率が規定内	分析クエリの エラーレート が XX%以下 におさまっている
データレイクのコスト	データレイク全体のコストが規定内	ストレージ、ETL、クエリ等の 総コストが\$XX以下 におさまっている

COVID-19 データレイクの場合

データレイクの各データが **直近 2日以内** のデータに更新された状態



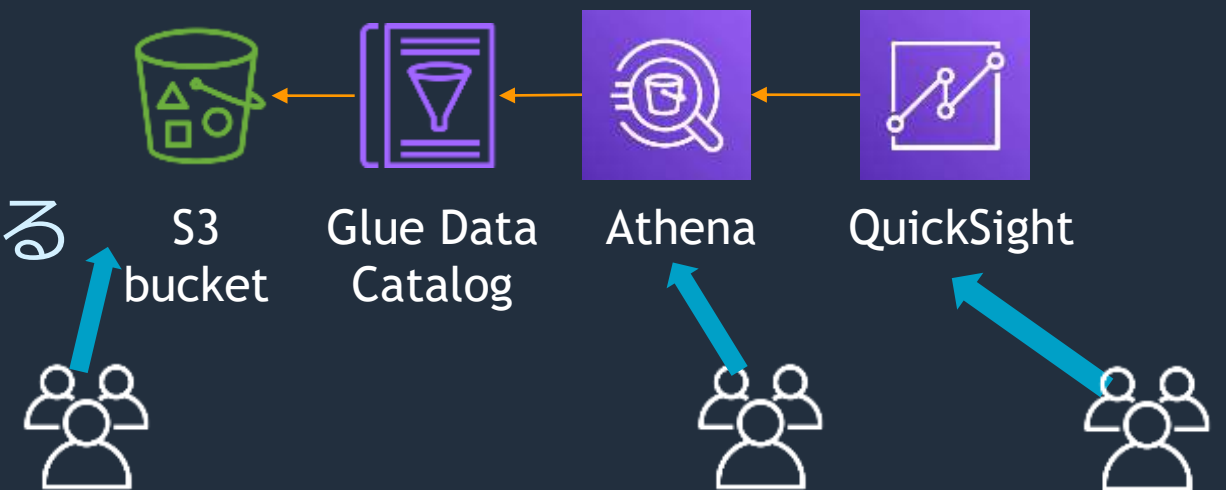
データレイクの正常性 SLA の監視

- ユーザーからの End-to-End のアクセスパターンに対するデータレイクの動作の正常性をチェックする
 - **実ユーザーのアクセス**の記録をチェックして判定
 - 監視用にユーザーアクセスを**定期的にシミュレート**して判定

COVID-19 データレイクの場合

異なるアクセスパターンがある

- S3バケットに直接アクセスする
- Athena (EMR/Redshift) 経由でクエリする
- QuickSight ダッシュボードを閲覧する



データレイクの正常性 SLA の監視

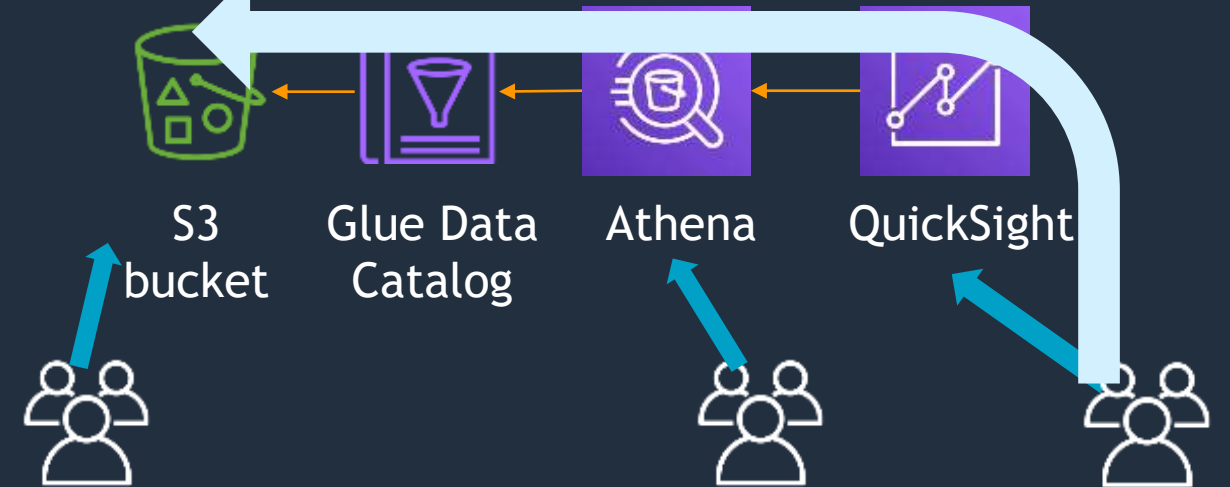
- ユーザーがデータレイクにアクセスする
- 実ユースケースをシミュレートして監視

```
New query 1 +
1 SELECT s.date
2 FROM "covid-19"."covid_testing_us_daily" u
3 left outer join "covid-19"."covid_testing_states_daily" s on u.date=s.date
4 left outer join "covid-19"."us_state_abbreviations" abb on s.state=abb.abbreviation
5 GROUP BY s.date
6 ORDER BY s.date DESC
7 LIMIT 1
Run query Save as Create (Run time: 1.86 seconds, Data scanned: 1.62 MB) Format query Clear
```

COVID-19 データレイクの場合

QuickSight へのデータロード時に Athena に対して発行されるクエリと同等のクエリをシミュレート実行してデータの更新日時と現在日時を比較

サービスレベル監視



Tips : 複数の Athena クエリの定期実行

The screenshot shows the AWS Glue console interface. On the left is a navigation menu with categories like 'データカタログ', 'データベース', 'ETL', and 'ワークフロー'. The main area displays a workflow named 'ワークフロー (2)' with a job named 'monitor-covid19-worldwide'. The job's code is shown in a text editor, which is a Python script using boto3 to execute an Athena query and check its status.

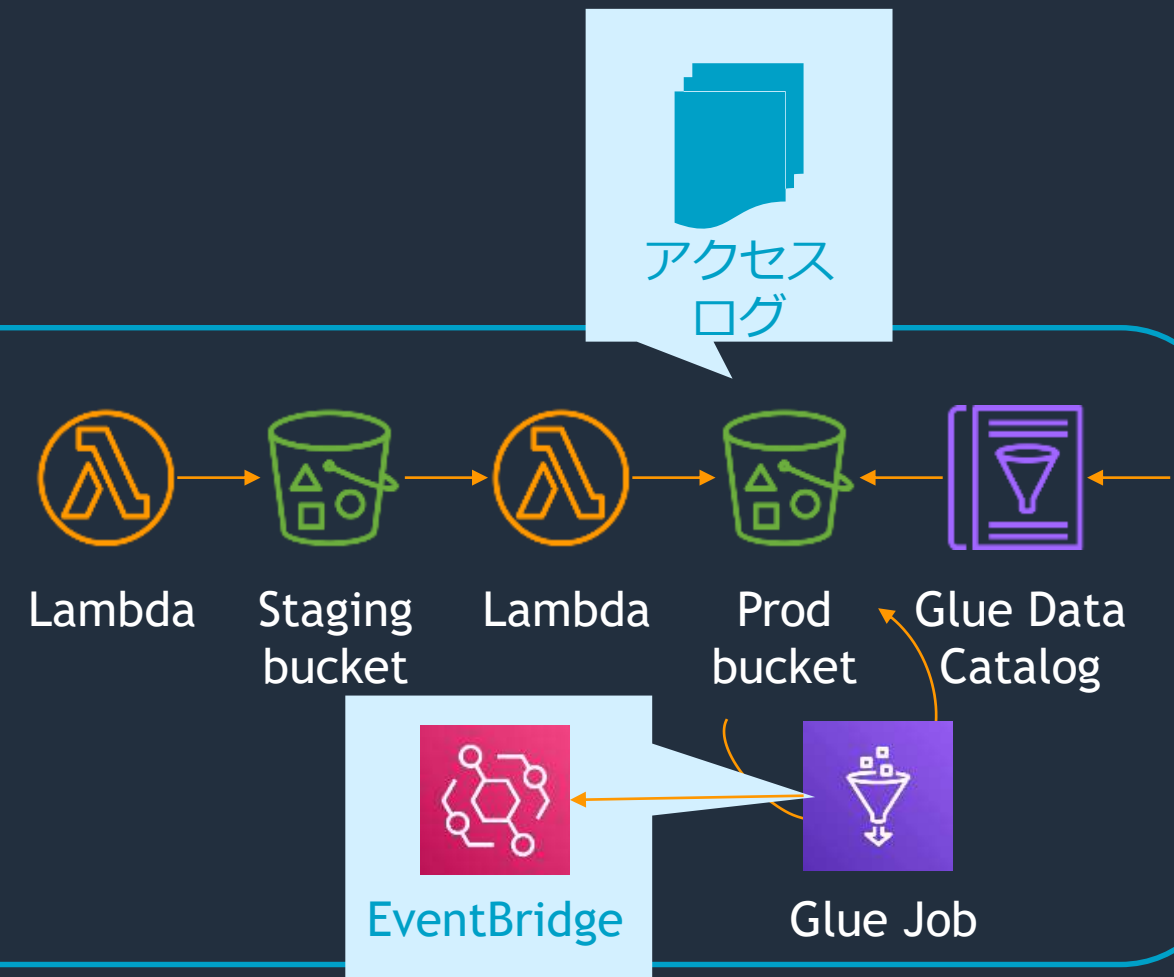
```
1 import boto3
2 from boto3.session import Session
3 import datetime
4 import time
5 import json
6
7 RETRY_COUNT = 20
8
9 client = boto3.client('athena')
10
11 allowed_latency_days = 2
12 dataset = "covid19-worldwide"
13 query = """
14     SELECT date FROM "covid-19"."world_cases_deaths_testing"
15     LEFT OUTER JOIN "covid-19"."country_codes"
16     ON "world_cases_deaths_testing"."iso_code"="country_codes"."alpha-3 code"
17     GROUP BY date
18     ORDER BY date DESC
19     LIMIT 1
20 """
21
22 print(query)
23 response = client.start_query_execution(
24     QueryString=query,
25     QueryExecutionContext={
26         'Database': 'covid-19'
27     },
28     ResultConfiguration={
29         'OutputLocation': 's3://aws-athena-query-results-us-east-2-647886179913/'
30     }
31 )
32
33 query_execution_id = response['QueryExecutionId']
34 print(query_execution_id)
35
36 for i in range(1, 1 + RETRY_COUNT):
37     query_status = client.get_query_execution(QueryExecutionId=query_execution_id)
38     query_execution_status = query_status['QueryExecution']['Status']['State']
39
```

データレイクの副次的な正常性監視：異常傾向の検知

- ユーザーからのアクセスの異常を示すイベントを検知する
 - アクセス数の異常低下
 - データ収集や ETL ジョブの失敗

COVID-19 データレイクの場合

- S3 バケットへのアクセス数を S3 サーバーアクセスログで計測、異常検知
- ジョブの失敗を EventBridge で検知



異常状態から正常状態への復旧

- 正常状態を逸脱した場合、何らかの手段で正常状態に戻す
 - リトライ
 - キャパシティ追加
 - 管理者による操作 (データ収集、ETL パイプラインの修正など)

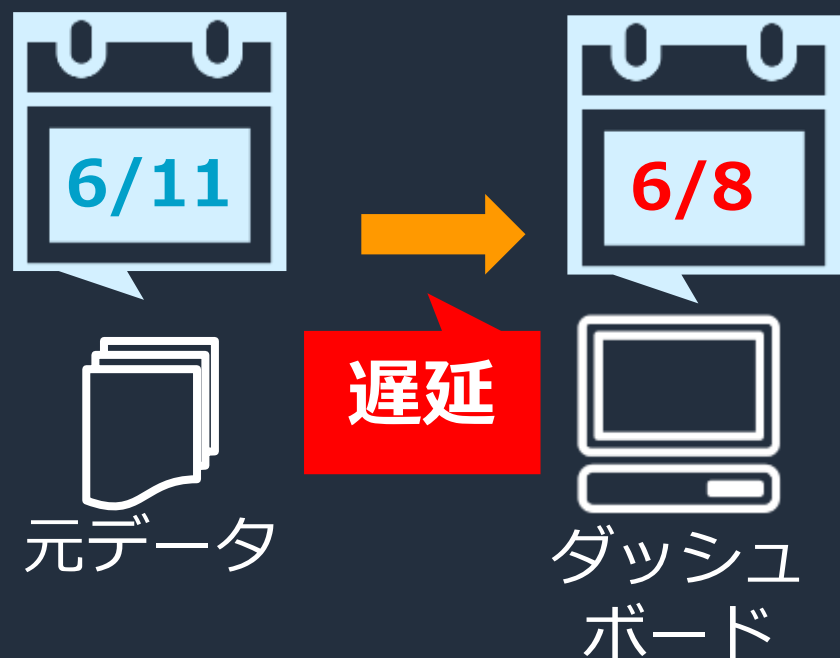
COVID-19 データレイクの場合

- データが利用可能になるまでの時間の SLA に抵触した場合にメール通知し管理者により対応

[covid19-lake-monitor]
Data latency increase detected



例1：データ鮮度の低下



原因

- データ収集・配信の遅延
- ETLパイプラインの遅延・失敗

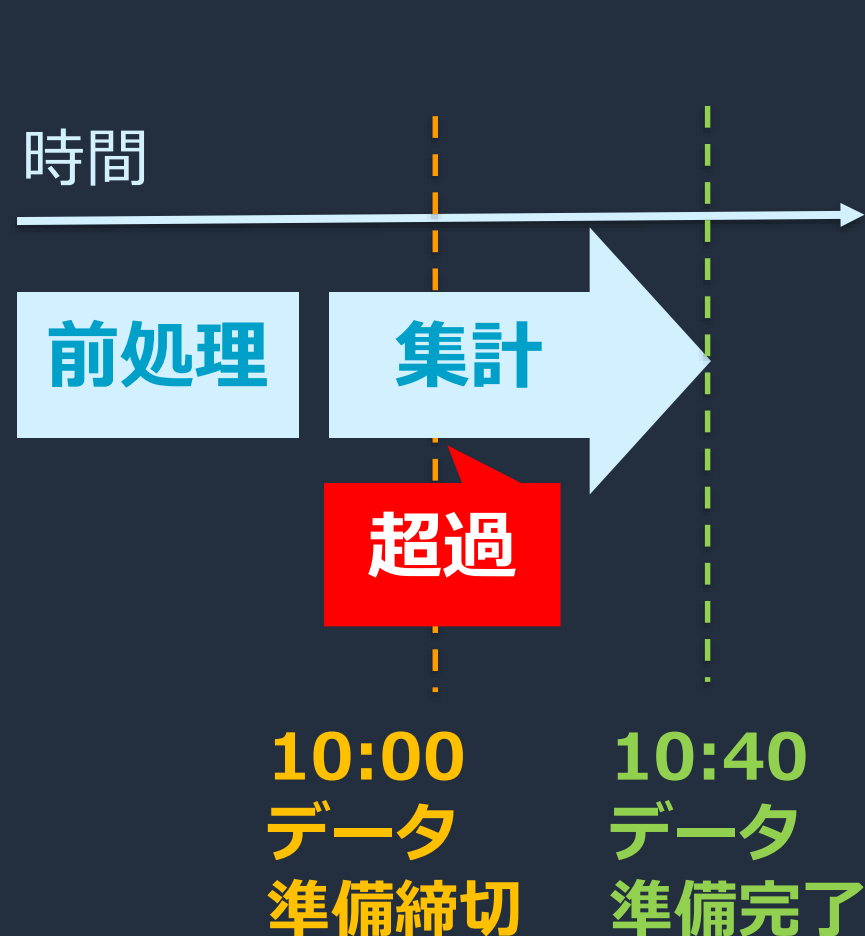
対応方法

- データ収集・配信の問題の修正
- ETLパイプラインの高速化・安定化

COVID-19 データレイクの場合

- 同時実行数超過による ETL ジョブ実行失敗 → ETLパイプラインの修正
- データ提供元からのファイル名の変更により、更新データが ETL から漏れる → ETLパイプラインの修正
- データ提供元からのデータ配信の停止 → データ提供元との交渉

例2：規定時間の超過



原因

- ETL パイプラインの遅延・失敗
- データ量、種類の増大

対応方法

- ETL パイプライン、データのチューニング
- データフォーマットの変換
- データ配置の最適化
- データ量の削減 (保存期間の調整, アーカイブ化)
- キャパシティ増強

COVID-19 データレイクの場合

→データの最新リビジョンのみを保持

(データ量がそこまで大きくないためこの種の問題は起きていない)

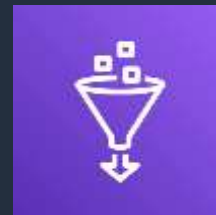
Tips : ETL パイプラインの最適化

小規模処理



AWS Lambda

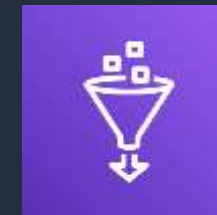
中規模処理



AWS Glue

Python Shell

大規模処理



AWS Glue

Spark

データを加工するコードのみを実装
Python であれば NumPy, Pandas などのライブラリも利用可能

- 実行時間に最大15分間の制限あり
- 豊富なトリガーを持ち、S3に配置されたタイミングで逐次処理することも可能
- Python以外の言語も選択可能

- 実行時間の制限なし
- Lambda に比べて利用できるメモリ量が多い（1GBまたは16GB）
- Athena、Redshift、EMR に対するSQLベースの処理も可能

PySpark や Scalaで実装
必要なものはコードのみ

- 実行時間の制限なし
- 複数のワーカーで並列分散処理
- 数100GB以上の大量データ処理も可能
- DPU 設定によりスケーラビリティを調整可能

Tips : データのパフォーマンス最適化

- ファイルフォーマット変換・圧縮
 - テキストベースのファイル (JSON, CSV等) を、分析に適したカラムナフォーマット (Parquet, ORC) に変換
 - 非圧縮ファイルを、分析に適した形式で圧縮
- 集約 (コンパクション)
 - ファイル処理のオーバーヘッドを小さくするために、大量の小さいファイルを 128MB程度のまとまりに集約
- パーティショニング
 - 単一ディレクトリにフラットに配置したファイルを、日付で区切ったディレクトリに保存

Tips : 保存期間の調整、アーカイブ化

- S3 ライフサイクル設定例
 - **1年経過**
→ストレージクラスを Standard-IA に変更
 - **5年経過**
→削除

The screenshot shows the 'Lifecycle Rule' configuration page in the AWS console. The page title is 'ライフサイクルルール' (Lifecycle Rule). The breadcrumb navigation shows the steps: '名前とスコープ' (Name and Scope), '移行' (Transition), '有効期限' (Expiration), and '確認' (Confirm), with '確認' being the current step (indicated by a '4' in a circle). The configuration details are as follows:

- 名前とスコープ** (Name and Scope):
 - 名前: ManageRetention
 - スコープ: バケット全体 (All buckets)
- Transitions** (Transitions):
 - オブジェクトの現行バージョン (Current version of object)
 - 移行先: Standard-IA 次の後: 365 日間 (Transition to Standard-IA after 365 days)
- 有効期限** (Expiration):
 - 失効期限: 1825 日間 (Expiration period: 1825 days)

At the bottom right, there are two buttons: '戻る' (Back) and '保存' (Save).

例3：意図せぬクエリ結果

NO DATA



ダッシュ
ボード

原因

- データ収集・配信・ETL パイプラインの遅延・失敗
- データ内容の変化、意図せぬ削除
- アクセスパターンの変化

対応方法

- データ収集・配信の問題の修正
- ETL パイプラインの修正
- データのリストア

COVID-19 データレイクの場合

- 収集したデータ内容の変更により想定したデータが含まれなくなった
→**データ収集方法の修正 (API, データソース等)**
- S3 バケットに保存される過去バージョンのファイルからリストア

Tips : S3 バケットのデータリストア

aws サービス リソースグループ

Amazon S3 > covid19-lake > rearc-covid-19-world-cases-deaths-testing > csv > covid-19-world-cases-deaths-testing.csv

covid-19-world-cases-deaths-testing.csv **最新バージョン**

概要 プロパティ アクセス権限 S3 Select

開ける ダウンロード 名前をつけてダウンロード 公開する コピーパス

所有者
68862915a0b1ecf44a6ecb256b2eb4df3c4b169bbf5d73bf638bdd717ba1dcf0

最終更新日時
6月 5, 2020 1:42:11 午前 GMT+0900

Etag
f102fb75c5f2202365e7ca5b3081bf5b

ストレージクラス
スタンダード

サーバー側の暗号化
なし

サイズ
3.2 MB

キー
rearc-covid-19-world-cases-deaths-testing/csv/covid-19-world-cases-deaths-testing.csv

オブジェクト URL
<https://covid19-lake.s3.us-east-2.amazonaws.com/rearc-covid-19-world-cases-deaths-testing/>

aws サービス リソースグループ

Amazon S3 > covid19-lake > rearc-covid-19-world-cases-deaths-testing > csv > covid-19-world-cases-deaths-testing.csv

covid-19-world-cases-deaths-testing.csv **最新バージョン**

6月 5, 2020 1:42:11 午前 GMT+0900 (最新バージョン)	スタンダード	ダウンロード	削除
6月 4, 2020 1:42:07 午前 GMT+0900	スタンダード	ダウンロード	削除
6月 3, 2020 1:42:07 午前 GMT+0900	スタンダード	ダウンロード	削除
6月 2, 2020 1:42:09 午前 GMT+0900	スタンダード	ダウンロード	削除
6月 1, 2020 1:42:10 午前 GMT+0900	スタンダード	ダウンロード	削除
5月 31, 2020 1:42:03 午前 GMT+0900	スタンダード	ダウンロード	削除

サーバー側の暗号化
なし

サイズ
3.2 MB

キー
rearc-covid-19-world-cases-deaths-testing/csv/covid-19-world-cases-deaths-testing.csv

オブジェクト URL
<https://covid19-lake.s3.us-east-2.amazonaws.com/rearc-covid-19-world-cases-deaths-testing/csv/covid-19-world-cases-deaths-testing.csv>

利用状況の計測と データレイクの継続的改善

データレイク、誰がどのくらい使ってる？

直近の1週間で
利用量がトップ20のユニークユーザーと
テーブルごとのアクセス数を
分析しといてくれる？



上司



データレイク管理者
● (初心者)

困った・・・
何も記録していないぞ・・・

利用状況の計測

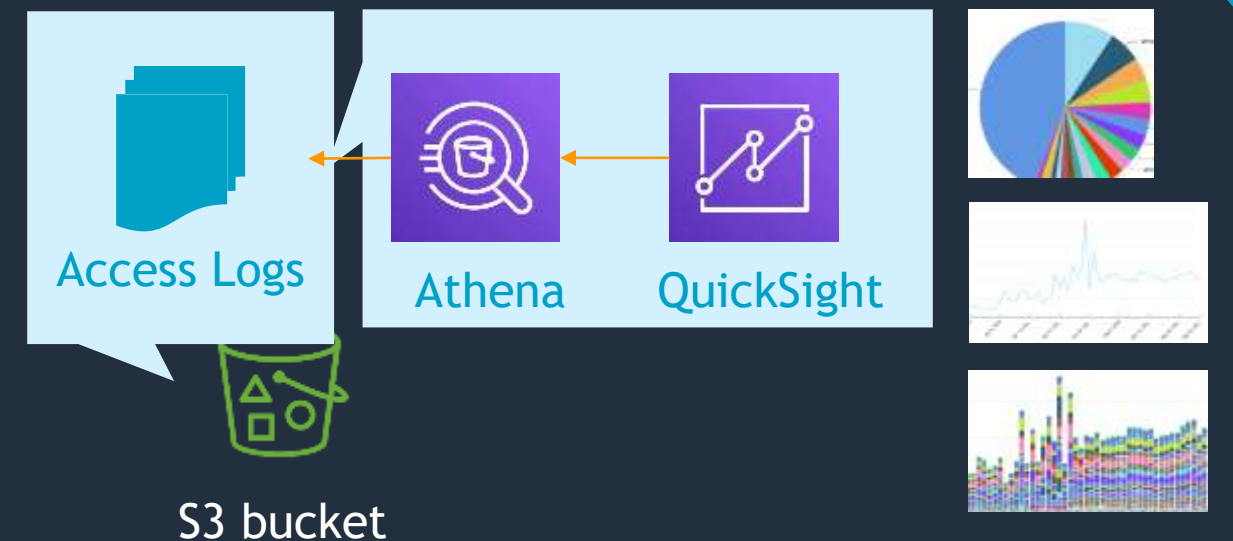
継続的改善にはユーザー利用状況の計測が必須

アクティビティの種類と計測方法

- データ参照: S3 サーバーアクセスログ、CloudTrail 等
- クエリ実行: クエリログ

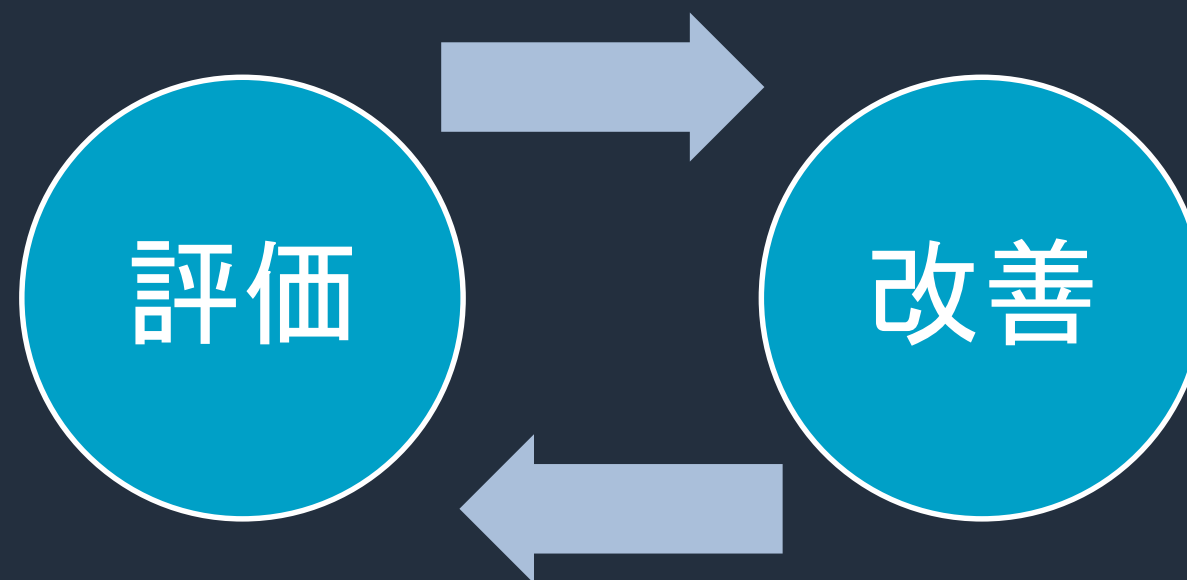
COVID-19 データレイクの場合

- S3 サーバーアクセスログを整形
- ユーザーアクセスログ分析用の QuickSight ダッシュボードを作成



データレイクの継続的改善

データレイクは塩漬けにせず継続的に進化させることが大切



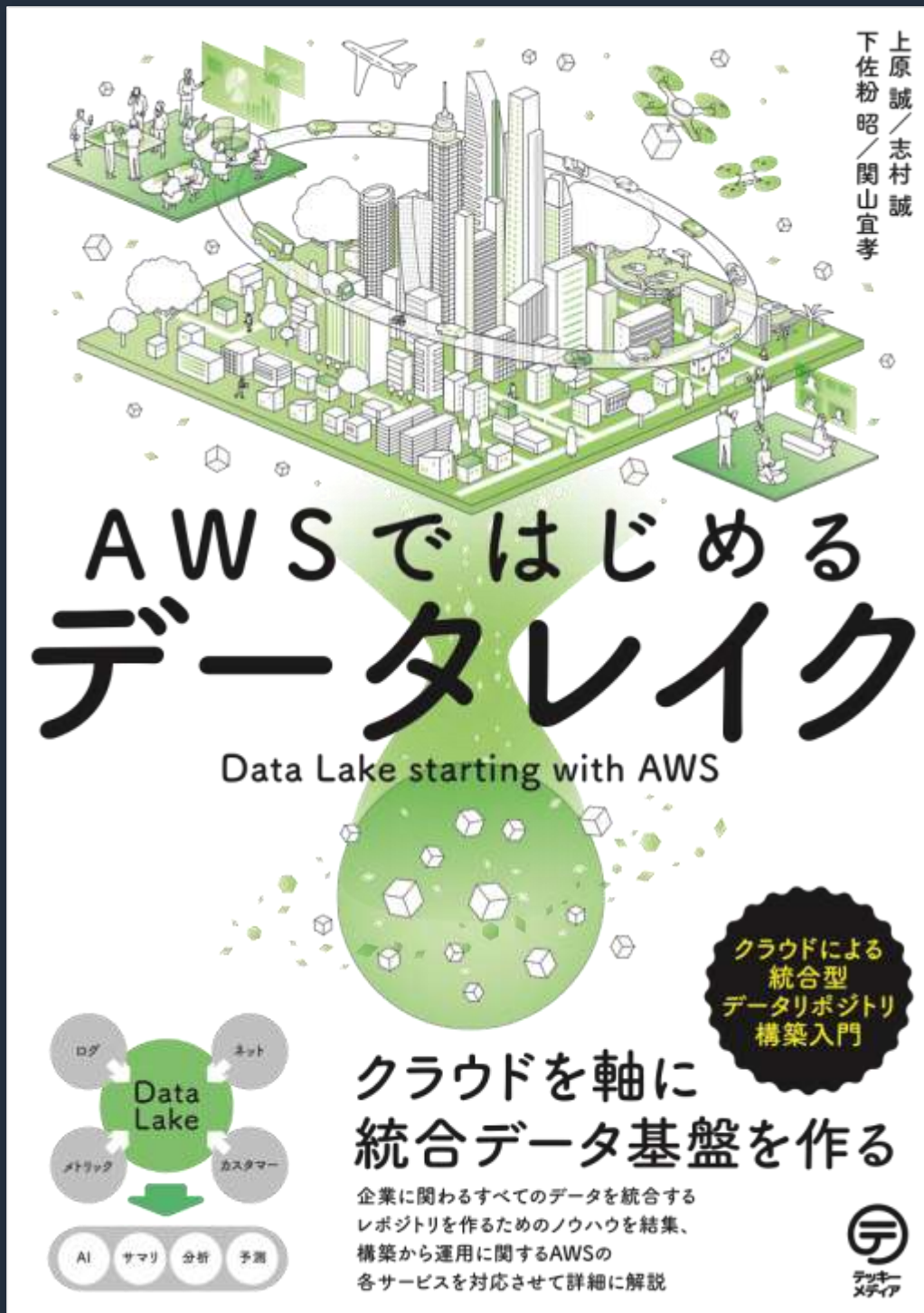
- 利用状況の分析
- フィードバックの収集

- 収集対象データの拡充
- 利用エンジンの拡充
- ドキュメント整備
- 利用ポリシーの改善

まとめ

- データレイクの運用
 - 正常状態を定義し、SLA を遵守する
 - ユーザーからの End-to-End のアクセスパターンに対するデータレイクの動作の正常性をチェックする
- 利用状況の計測とデータレイクの継続的改善
 - 継続的な改善には利用状況の計測が必須
 - データレイクは塩漬けにせず継続的に進化させることが大切

「AWSではじめるデータレイク」 2020年7月発売予定



内容についての注意点

- 本資料では2020年6月11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おります。日本居住者のお客様には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.

Q&A

関連リンク

AWSではじめるデータレイク 出版記念 オンラインセミナー

- 第1回 データレイクははじめの一步
 - 資料: <https://go.aws/3df1Ebo>
- 第2回 データレイクの構築と蓄積されたデータの活用方法
 - 資料: <https://go.aws/3cxGvbe>

関連セミナーのご案内

AWSではじめるデータレイク 出版記念 オンラインセミナー

- 第4回 Glue, Lake Formation, Athena, EMR 最新アップデート
 - 2020年6月18日(木) 11:00 ~ 12:00
 - <https://aws-seminar.smktg.jp/public/seminar/view/2674>