

# Google Cloud UPDATES: 2023 年 Data Analytics / ML編

2023-02-20

Google Cloud

# Google Cloud Updates について

- このイベントでは、Google Cloud に関する四半期分のアップデートの振り返りを行っていきます
- 既存ユーザ様を主な対象としているため、基本的には、プロダクトの概要レベルの説明は行いません
- 出入りは自由ですが、退出時に[アンケート](#)にご協力下さい
- 質問は、Chat か、[こちらのフォーム](#)をご活用下さい
- 今回は、2022-10-01 ~ 2023-01-30 の DataAnalytics / AI,ML 系のプロダクトのアップデートの振り返りを行っていきます

# 本日のスピーカー



Yu Yamada  
BigQuery,  
ガバナンス



Yutaro Saito  
ETL & Workflow



Shingo Furuyama  
BQML, Vertex AI,  
その他

01

# Data Analytics - BigQuery ガバナンス

# Multi-statement transactions

## / 特徴

- 複数の SQL 文を1つのトランザクションとして実行
- DML、SELECT、一時表、一時関数の作成、削除をサポート
- INFORMATION\_SCHEMA.JOBS\_BY\_\* でトランザクション ID やトランザクションの成否が確認可能

## / ユースケース

- 複数の処理を束ねてアトミックに扱いたい場合
- マルチステートメントトランザクションを使っている他社 DB からのバッチ処理移行

```
BEGIN

BEGIN TRANSACTION;

CREATE TEMP TABLE tmp
  AS SELECT * FROM mydataset.NewArrivals WHERE warehouse =
'warehouse #1';

DELETE mydataset.NewArrivals WHERE warehouse = 'warehouse #1';

MERGE mydataset.Inventory AS I
USING tmp AS T
ON I.product = T.product
WHEN NOT MATCHED THEN
  INSERT(product, quantity, supply_constrained)
  VALUES(product, quantity, false)
WHEN MATCHED THEN
  UPDATE SET quantity = I.quantity + T.quantity;

DROP TABLE tmp;

COMMIT TRANSACTION;

EXCEPTION WHEN ERROR THEN
  -- Roll back the transaction inside the exception handler.
  SELECT @@error.message;
  ROLLBACK TRANSACTION;
END;
```

# Dataplex - Data Exploration Workbench

## Spark SQL

### / 特徴

- Cloud Storage / BigQuery にあるデータを Spark SQL / notebook で探索ができる
- スケジュールも可能
- 共有してコラボレーションが可能
- サーバレスなインフラで実行

The screenshot shows the Dataplex Spark SQL interface. On the left, there is a sidebar with a filter and a tree view of the workspace. The main area contains a script editor with a SQL query and a results pane showing a table of data.

```

1 select * from
2   `curateddata`.`nation`, `curateddata`.`region`
3   where n_regionkey=r_regionkey
  
```

Row	n_nationkey	n_name	n_regionkey	n_comment
0	16	MOZAMBIQUE	0	s. ironic, unusual asyntotes wake blithely r
1	15	MOROCCO	0	rms. blithely bold courts among the closely regular packages use fu
2	14	KENYA	0	pending excuses haggle furiously deposits. pending. express print
3	5	ETHIOPIA	0	ven packages wake quickly. regu
4	0	ALGERIA	0	haggle. carefully final deposits detect slyly agal
5	24	UNITED STATES	1	y final packages. slow fuses coggle quickly. quality silent statista
6	17	PERU	1	platelets. blithely pending dependencies use fluently across the eve
7	3	CANADA	1	ees hang ironic, silent packages. slyly regular packages are futurou
8	2	BRAZIL	1	y alongside of the pending deposits. carefully special packages and
9	1	ARGENTINA	1	al foves promise slyly according to the regular accounts. bold rep

## Notebook

### / ユースケース

- Dataplex 内で管理されているデータに対してすぐに分析や探索がしたい

The screenshot shows the Dataplex Notebook interface. The top menu includes File, Edit, View, Run, Kernel, Tabs, Settings, and Help. The main area contains a code cell with a BigQuery query and its execution results.

```

[5]: %bigquery
SELECT
  country_code,
  country_name,
  COUNT(DISTINCT region_code) AS num_regions
FROM
  `bigquery-public-data.google_trends.international_top_terms`
WHERE
  refresh_date = DATE_SUB(CURRENT_DATE, INTERVAL 2 DAY)
GROUP BY
  country_code,
  country_name
ORDER BY
  num_regions DESC;
  
```

Query complete after 0.00s: 100% [████████████████████] 4/4 [00:00:00:00, 1868.70query/s]  
Download[ing: 100% [████████████████████] 41/41 [00:01:00:00, 38.40rows/s]

[5]:	country_code	country_name	num_regions
0	TR	Turkey	81
1	TH	Thailand	77
2	VN	Vietnam	63
3	JP	Japan	47
4	RO	Romania	42
5	NG	Nigeria	37
6	IN	India	36

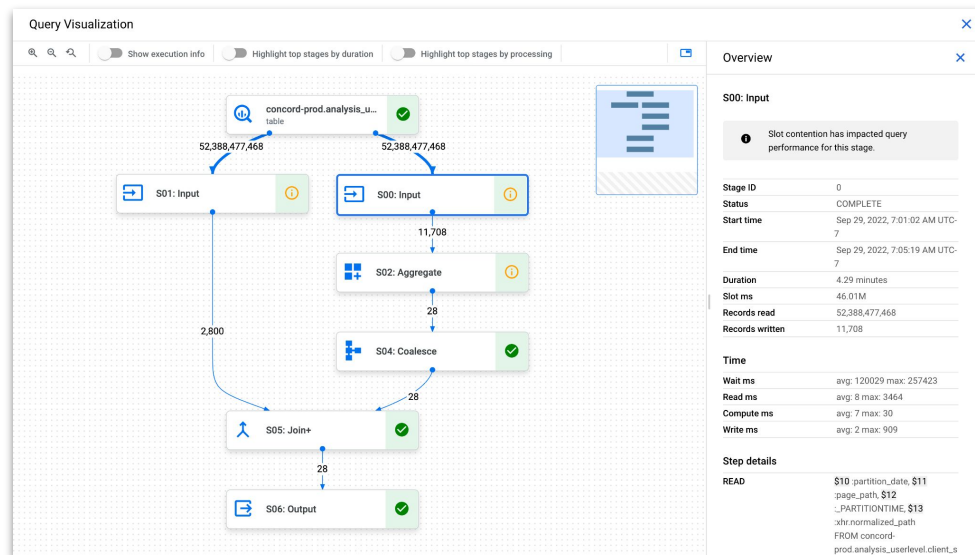
# BigQuery - Query Execution Graph

## / 特徴

- クエリ性能診断に必要な各種情報を可視化
  - クエリの実行計画
  - 各ステージの詳細の表示(スロット利用時間、読み取り件数、書き込み件数)
  - 処理量が多いステージのハイライト
  - 実行時間が長いステージのハイライト
  - パフォーマンス分析情報

## / ユースケース

- 遅いクエリの分析、原因の解析



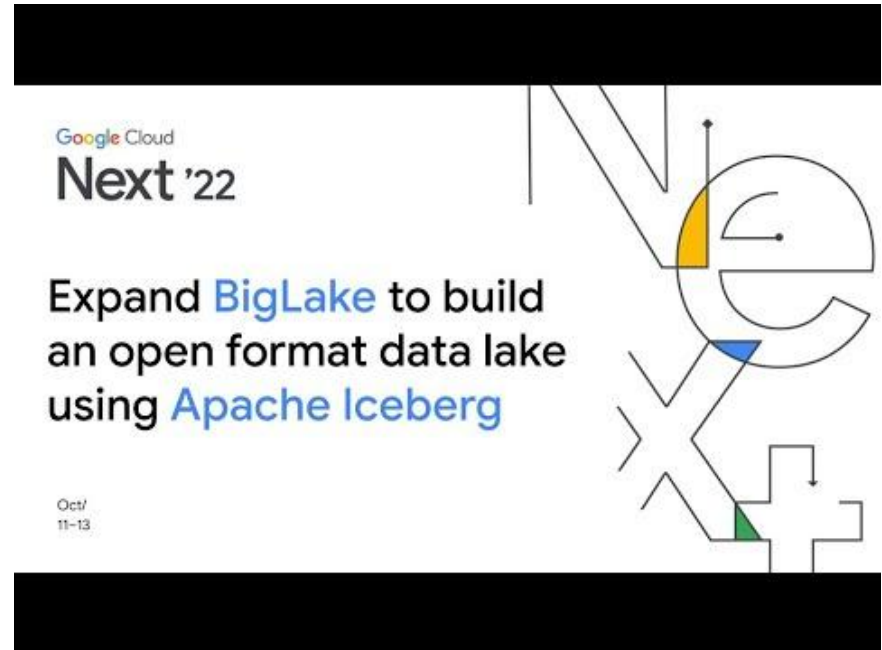
# BigQuery - BigLake for Iceberg tables

## / 特徴

- オープンソーステーブルフォーマットであるIceberg を BigLake でサポート
- Spark で作成した Iceberg table が自動的に BigQuery テーブルとして認識される
- Row, column level security, masking 等 BigLake のセキュリティ機能をIceberg テーブルに適用
- BigLake Metastore (Iceberg catalog) を参照することにより、object.list オペレーションが不要となり、Pruning により必要なファイルにのみアクセス。パフォーマンスメリットが得られる。

## / ユースケース

- OSS フォーマットをベースとしたLakehouseアーキテクチャを BigQuery / BigLake で実現
- Spark ETL プロセスでIceberg テーブルを作成し、BigQuery からクエリ実行。





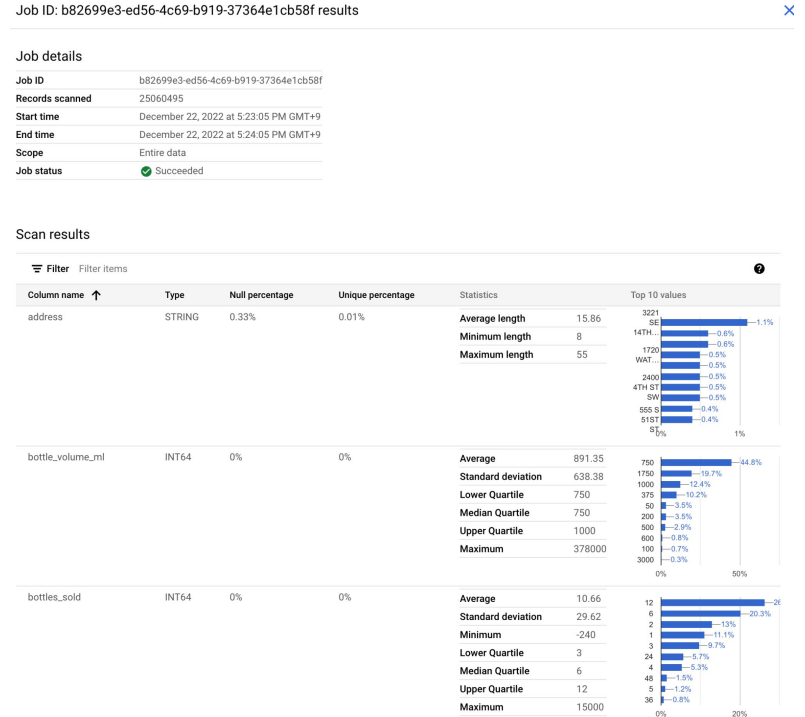
# Dataplex - Data Profiling

## / 特徴

- Dataplex の Asset として登録されている BigQuery テーブルに対して、自動プロファイリングにより以下の統計が取得できる
  - Null 率
  - ユニーク率
  - 統計情報 (Type により異なる)
  - Top 10 values
- Preview 期間は無料

## / Use case

- スキーマ情報が不明なデータの検査
- データ品質チェックの値確認



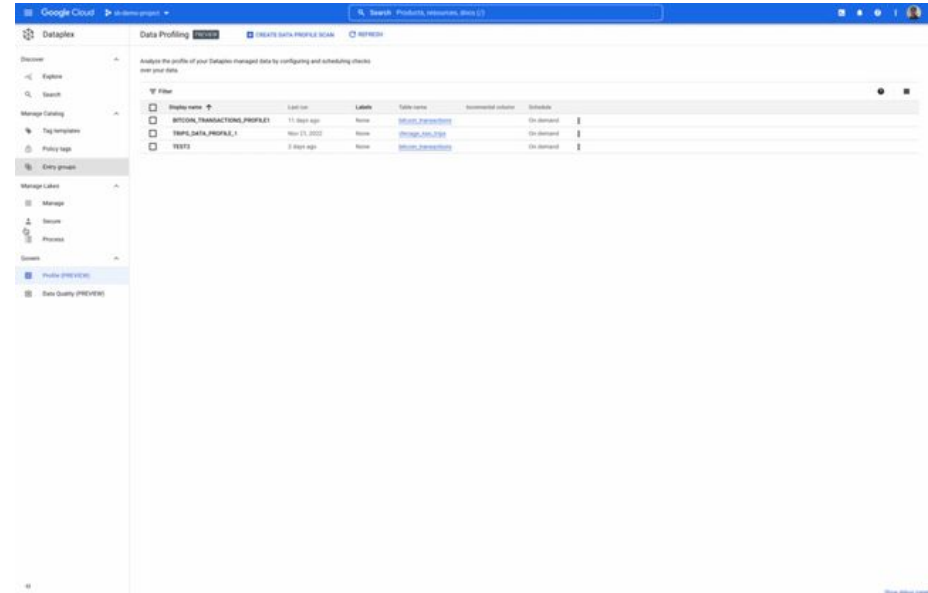
# Dataplex - Auto Data Quality (AutoDQ)

## / 特徴

- テンプレート化されたルールをもとに、データ品質チェックをクイックに行う
- Preview 期間は無料

## / Use case

- 自社のデータガバナンスが届かない相手先 (外部システム) からのデータ連携などで、データが要件を満たすか (or クレンジングされているか) チェック
- 通常の Dataplex DQ Task との使い分け
  - 簡易・クイックに => AutoDQ
  - しっかり・カスタム => 通常 DQ Task



# Dataplex - Auto Data Quality (AutoDQ)

## 事前定義ルール

- 範囲チェック(ex : min 1 - max 100)
- Not Null チェック
- 特定の値チェック(ex : USD, JPY, EUR のみ)
- 正規表現チェック(ex : 先頭 customer\* で始まる)
- ユニーク性チェック
- Statistic Range Expectation (ex : 金額の平均がxxx より高い)

Rule type (Name in Google Cloud console)	Row-level or aggregate rule	Description	Supported column types	Rule-specific parameters	Column types	Required parameters
RangeExpectation (Range check)	Row-level	Check if the value is between min and max.	All numeric, date, and timestamp type columns.	<b>Required:</b> <ul style="list-style-type: none"> <li>Passing threshold percentage</li> <li>mean, min, or max values: Specify at least one value.</li> </ul> <b>Optional:</b> <ul style="list-style-type: none"> <li>Enable strict min: If enabled, the rule check uses "&gt;" instead of "&gt;=".</li> <li>Enable strict max: If enabled, the rule check uses "&lt;" instead of "&lt;=".</li> <li>Enable ignore null: If enabled, null values are ignored in the rule check.</li> </ul>	All supported column types, except Record and Struct.	Required: Column and dimension from the supported parameters.
NonNullExpectation (Null check)	Row-level	Validate that column values are not NULL.	All supported column types.	<b>Required:</b> <ul style="list-style-type: none"> <li>Passing threshold percentage.</li> </ul>	All supported numeric column types.	<b>Required:</b> <ul style="list-style-type: none"> <li>mean, min, or max values: Specify at least one value.</li> </ul> <b>Optional:</b> <ul style="list-style-type: none"> <li>Enable strict min: If enabled, the rule check uses "&gt;" instead of "&gt;=".</li> <li>Enable strict max: If enabled, the rule check uses "&lt;" instead of "&lt;=".</li> <li>Enable ignore null: If enabled, null values are ignored in the rule check.</li> </ul>
SetExpectation (Set check)	Row-level	Check if the values in a column are one of the specified values in a set.	All supported column types, except Record and Struct.	<b>Required:</b> <ul style="list-style-type: none"> <li>Set of string values to check against.</li> <li>Passing threshold percentage.</li> </ul> <b>Optional:</b> <ul style="list-style-type: none"> <li>Enable ignore null: If enabled, null values are ignored in the rule check.</li> <li>Enable invert condition: If enabled, rule checks if the values are not from the set.</li> </ul>	All supported column types, except Record and Struct.	
RegexExpectation (Regex check)	Row-level	Check the values again a specified regular expression.	String	<b>Required:</b> <ul style="list-style-type: none"> <li>Regex pattern used to check.</li> <li>Passing threshold percentage.</li> </ul> <b>Note:</b> Google Standard SQL provides regular expression support using the <code>re2</code> library; See that documentation for its regular expression syntax.		

# BigQuery - Session の機能追加

## / 特徴

- Temp Function が Session 内で利用可能に
  - Session 内のみで有効
  - Session 終了後、自動的に削除
- Temp キーワードを含む SQL に OR REPLACE, IF NOT EXISTS が利用可能に
  - Temp function
  - Temp table

## / Use case

- Temp function / Temp table を使った処理のデバッグ

The screenshot shows the BigQuery console interface. At the top, there are buttons for RUN, SAVE, SHARE, SCHEDULE, and MORE. Below these is a SQL editor with the following code:

```

1 CREATE OR REPLACE TEMP FUNCTION AddFourAndDivide(x INT64, y INT64)
2 RETURNS FLOAT64
3 AS (
4   (x + 4) / y
5 );
6
7
8 SELECT
9   val, AddFourAndDivide(val, 2)
10 FROM
11 UNNEST([2,3,5,8]) AS val;

```

Below the editor, it shows "Processing location: US" and "Session Mode: ON". A red box highlights the SQL code, and a red arrow points from the text "Session ModeがONの場合、別々に実行することができる" to the "Session Mode: ON" status.

Query results

JOB INFORMATION		RESULTS	JSON	EXECUTION DETAILS	EXI
Row	val	f0_			
1	2	3.0			
2	3	3.5			
3	5	4.5			
4	8	6.0			

# BigQuery - Reservation の DDL の追加、改善

## / 特徴

- 購入済みのコミットメントの変更
  - プランの変更
  - リニューアルプランの変更
- 作成済みの予約の変更
  - アイドルスロットの利用の設定
  - スロット量
  - 同時実行数 (Query Queue 利用時)
- CREATE 時にオプションの指定が可能に
  - コミットメント、予約、割当

## / Use case

- プランの変更を SQL で実行
- BigQuery Script など SQL のみで BigQuery Reservation 関連の操作を実行する
  - ピーク時に Flex スロットを購入しオフピークにはもとに戻す

```
/* 追加でFlex Slotを100購入 */  
CREATE CAPACITY `admin_project.region-us.my-commitment`  
OPTIONS (  
  slot_count = 100,  
  plan = 'FLEX');  
  
/* 購入したスロットを既存の予約100スロット)に追加して200にする  
ALTER RESERVATION `admin_project.region-us.my-reservation`  
SET OPTIONS (  
  slot_capacity = 200  
);
```

# Dataplex - Business Glossary

## / 特徴

- ビジネス用語集を作成・管理
- 用語を BigQuery のカラムに関連付けて参照することが可能

## / Use case

- そのビジネスドメインに不案内な分析ユーザのデータ理解
- 用語が特殊であったり同じ用語でも組織内でコンテキストの違いで異なる意味になる場合

The screenshot displays the Google Cloud Dataplex Business Glossary interface. The top navigation bar includes the Google Cloud logo, a 'Knowledge repo' dropdown, a search bar, and utility icons. The main content area is divided into a left sidebar and a main panel.

**Left Sidebar (Glossaries):**

- Filter Search
- Customer Relationship Management
- Finance
- Marketing
- Operations
- Retail
  - Asset
  - Brand
  - Budget
  - Cash flow
  - Client
  - Cost
  - Cost of goods sold
  - Customer** (highlighted)
  - Customer relationship management
  - Distribution
  - Inventory
  - Point of sale
  - Price
  - Product

**Main Panel (Customer Entry):**

**Customer** [DELETE]

Term • Last modified: Jan 6, 2023, 6:10:21 PM

Steward: amanda.johnson@company.org

**Customer**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur.

Category	Description
Lorem ipsum dolor sit amet	Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur.
Ut enim ad minim veniam	Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat.

- Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.
- Lorem ipsum dolor sit amet, [consectetur adipiscing elit](#), sed do eiusmod tempor incididunt ut labore et dolore magna aliqua.

Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

**Related Terms**

- [Customer relationship management](#)

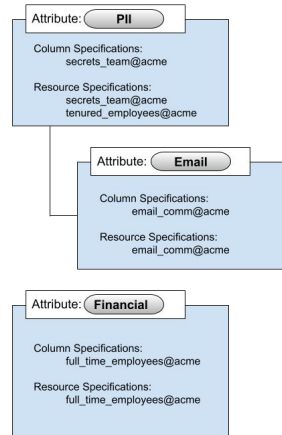
**Synonyms**

- [Client](#)

# Dataplex - Attribute store

## / 特徴

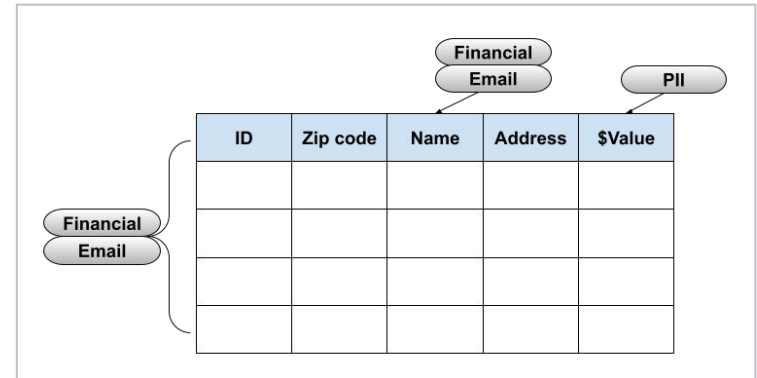
- Attribute を階層化してアクセスポリシー設定
- テーブルと列のACLを管理
- Dataplex の管理アセットに対して実施可能



## / Use case

- Dataplex でリソースを管理しているユーザで BigQuery のポリシータグのように階層的・一元的にテーブル・列レベルでのアクセス制御を管理したい

- Employees in `secrets_team@acme`, `tenured_employees@acme`, `full_time_employees@acme`, and `email_comm@acme` can access the table. This is because Dataplex merges the resource specifications of the attributes `Financial` and `Email`, and the attribute `Email` inherits the specifications from the attribute `PII`.
- Employees in `secrets_team@acme`, `email_comm@acme`, `full_time_employees@acme` can access the column `Name`. This is because Dataplex merges the column specifications of the attributes `Financial` and `Email`.
- Only employees in `secrets_team@acme` can query the column `$Value`.



02

# Data Analytics - ETL&Workflow



# BigQuery - Redshift Migration アセスメント

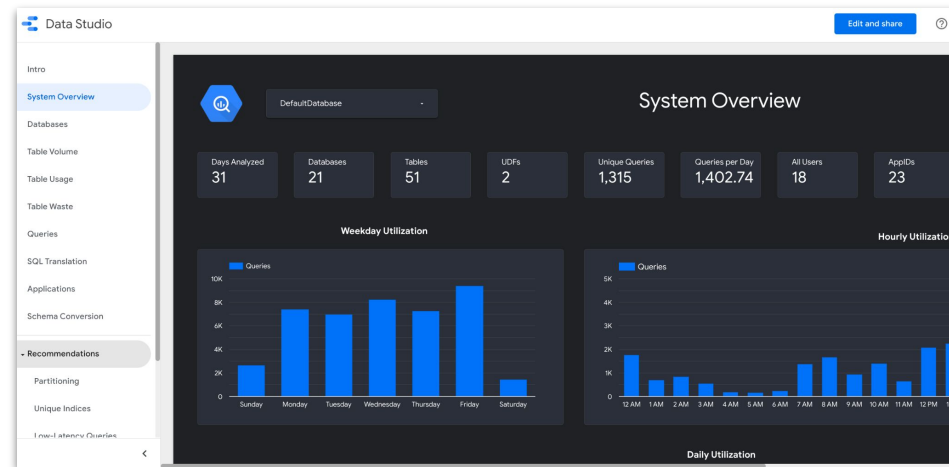
## / 特徴

Redshift に関する情報を分析し移行の複雑度、移行方法を評価

- 既存システム
  - データベース、スキーマ、テーブル等の情報
- BigQuery の安定状態
  - 自動変換可能なクエリ数、データ保存の費用、テーブルの最適化の提案
- 移行パス
  - 移行作業自体に関する情報

## / ユースケース

- Redshift 移行時の移行支援
- PoC ワークロードの抽出



# BigQuery - ラップされた鍵を管理する SQL 関数

## / 特徴

- ラップされた鍵を管理する SQL 関数
  - KEYS.NEW\_WRAPPED\_KEYSET
  - KEYS.ROTATE\_WRAPPED\_KEYSET
  - KEYS.REWRAP\_KEYSET

## / ユースケース

- ラップされた鍵で列レベルの暗号化
- AEAD 暗号化をよりセキュアに実施したい

```
SET kms_resource_name =  
'gcp-kms://projects/my-project/locations/us/keyRings/my-key-ring/cryptoKeys/my-  
crypto-key';  
SET wrapped_keyset = (SELECT KEYS.NEW_WRAPPED_KEYSET(kms_resource_name,  
'AEAD_AES_GCM_256'));  
  
SELECT AEAD.ENCRYPT(  
  KEYS.KEYSET_CHAIN(  
    kms_resource_name,  
    wrapped_keyset),  
  'plaintext',  
  '') AS encrypted_content
```

# BigQuery - Load Data for S3 and Azure Blob

## / 特徴

- S3、Azure Blob 上にあるファイルを直接 BigQuery へロードが可能に(要 BigQuery Omni のスロット)
- US multi-region と US-EAST-4 region に対応

```
LOAD DATA INTO `mydataset.testparquet`  
FROM FILES (uris = ['s3://test-bucket/sample.parquet'],  
format = 'PARQUET')  
WITH CONNECTION `aws-us-east-1.test-connection`  
;
```

## / ユースケース

- S3 上にあるファイルをセキュアに BigQuery にロードして分析したい

# BigQuery - データの追加の改善

GA 2022-11-07

## / 特徴

- BigQuery にデータを追加する方法がよりわかりやすくなった。
- 一般的に利用されるローカルファイル、GCS、外部データソース
- Analytics Hub の Listing からの追加
- その他、様々なソースからの追加も可能に

The screenshot shows the 'データを追加' (Add Data) dialog box in BigQuery. It features a search bar at the top with the placeholder text 'データソースを検索' (Search for data source). Below the search bar, there are three main categories of data sources:

- 一般的なソース** (General sources): This section contains three cards: 'ローカルファイル' (Local files) for uploading local files, 'Google Cloud Storage' for object storage, and '外部データソースへの接続' (Connect to external data sources) for connecting to external sources.
- その他のソース** (Other sources): This section lists 24 results, including: 'プロジェクトを検索してスターを付ける' (Search for projects and star), '名前を指定してプロジェクトにスターを付ける' (Star projects by name), 'Analytics Hub' (Public, commercial, and limited public data sets), 'Google ドライブ' (Google Drive), 'Amazon S3 - Data Transfer' (Amazon S3 storage service), '公開データセット' (Public data sets), 'Cloud Bigtable' (Highly scalable NoSQL database), 'Informatica データローダ' (Informatica Data Loader), 'Fivetran Data Pipelines' (Automated data pipeline), 'Pub/Sub サブスクリプション' (Pub/Sub subscriptions), and 'Datastream' (MySQL, PostgreSQL, Oracle replication).

# Data Fusion - Private Instance 使用時における DNS での名前解決

## / 特徴

- プライベートインスタンスでソースとシンクにホスト名、ドメイン名を指定して設定可能
- IP アドレス直指定ではないので、ソースやシンクのフェイルオーバー等によるパイプラインの変更が不要
- Data Fusion 6.7.0 以降で利用可能

## / ユースケース

- オンプレや Private network 内のデータソースやシンクにホスト名でアクセスしたい

[Resolve domain names or hostnames](#)

# Composer - Airflow トリガーと Deferrable operators

## / 特徴

- ワーカーのアイドル期間中にワーカーズロットを解放して別のタスクで利用することが可能
- ワーカーを効率的に利用可能
- 簡単に利用可能
- Composer 2、且つ、Airflow 2.2.5、2.3.3以降で利用可能

## / ユースケース

- 処理を外部サービス(BigQuery / Dataproc)で実行していて且つ実行時間が長い

```
insert_query_job = BigQueryInsertJobOperator(  
    task_id="insert_query_job",  
    configuration={  
        "query": {  
            "query": INSERT_ROWS_QUERY,  
            "useLegacySql": False,  
        }  
    },  
    location=LOCATION,  
    deferrable=True,  
)
```

# Pub/Sub - BigQuery Subscriptions JSON Type のサポート

## / 特徴

- 全ての String の項目、data と attributes に 適用可能
- JSON Data Type の恩恵を受けられる
  - 格納効率(データサイズ)
  - JSON 関数の高速化

Field name	Type
<a href="#">subscription_name</a>	STRING
<a href="#">message_id</a>	STRING
<a href="#">publish_time</a>	TIMESTAMP
<a href="#">data</a>	JSON
<a href="#">attributes</a>	JSON

## / ユースケース

- JSON データを展開せずに Streaming で BigQuery に入れる

subscription_name	message_id	publish_time	data	attributes
projects/839210930825/subsc...	6335549661394078	2022-11-24 09:07:25.235000 U...	{'col1': 'hello', 'col2': 'world'}	{'googlient_schemaname': 'projects/tenishimdemo/schemas/sample', 'key1': 'value1', 'key2': 'value2', 'googlient_schemae
projects/839210930825/subsc...	6335583816702330	2022-11-24 09:17:55.651000 U...	{'col1': 'abc', 'col2': 'def'}	{'googlient_schemaname': 'projects/tenishimdemo/schemas/sample', 'googlient_schemaencoding': 'JSON'}
projects/839210930825/subsc...	6335710543489158	2022-11-24 09:33:27.009000 U...	{'col1': 'hello', 'col2': 'world'}	{'googlient_schemaencoding': 'JSON', 'key1': 'value1', 'googlient_schemaname': 'projects/tenishimdemo/sche

# BigQuery - Metadata Caching

## / 特徴

- メタデータをキャッシュしてクエリのパフォーマンスを向上
- BigLake テーブルと Object テーブルに対して有効
- Cloud Storage のリスティングを防ぐので大量のファイルをもつ BigLake テーブルに有効
- 期限を指定可能 (30 分～7 日間)
- リフレッシュは手動と自動が選択可能

## / ユースケース

- Hive Partitioned table など大量のファイルをもつテーブルを BigLake テーブルとして BigQuery から高速にアクセスしたい

```
CREATE EXTERNAL TABLE `my_dataset.object_table`  
WITH CONNECTION `us.my-connection`  
OPTIONS(  
  object_metadata = 'SIMPLE',  
  uris = ['gs://mybucket/*'],  
  max_staleness = INTERVAL 1 DAY,  
  metadata_cache_mode = 'AUTOMATIC'  
);
```



# Dataproц Metastore - Administrator Interface

## / 特徴

- (Dataproц Metastore をアタッチした) Dataproц や Hive のインスタンスがなくてもDataproц Metastore に保管されたメタデータの参照や管理ができるように
- Read-only operation ができること (roles/metastore.metadataQueryAdmin)
  - メタデータのクエリ
- Read and write operation (roles/metastore.metadataMutateAdmin)
  - データベース、テーブルやパーティションのロケーション変更
  - テーブルを別のデータベースへ移動

[Dataproц Metastore administrator interface](#)

# Cloud Composer - Composer Local Development CLI tool

## / 特徴

- ローカルで実行されるApache Airflow 環境で DAG の開発ができるように

## / ユースケース

- DAG のテスト
- PyPI パッケージのテスト
- Airflow configuration option のテスト

# Cloud Composer - Environment snapshot / Scheduled snapshot

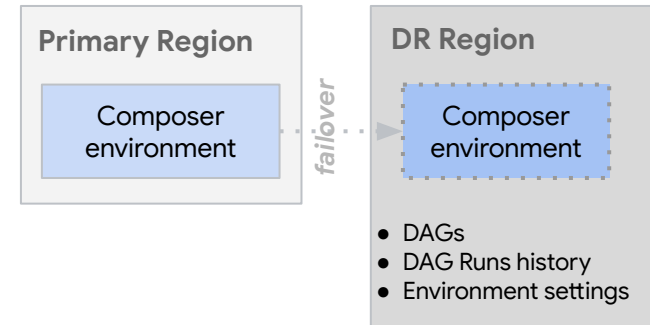
## / 特徴

- Composer 環境の設定と状態のバックアップを GCS バケットへ保管

## / ユースケース

- オペレーションエラーからの復旧
- 新しいバージョンで作成したComposer 環境にスナップショットを適用することで、side-by-side アップデートを実現
- DR 環境へのフェイルオーバー

### Disaster Recovery



# Cloud Data Fusion - Cloud Data Fusion 6.8.0

CDAP 6.8.0 リリース

## / 特徴

- Dataplex Batch Source と Dataplex Sink プラグインが GA
- Datastream を利用した Oracle から BigQuery への Replication が GA
  - GCS バケット経由の Replication の模様
- BigQuery batch source pushdown サポート
- AND trigger のサポート
- Replication job のバージョンアップグレードはできないため、利用している場合要注意

**BREAKING**

Upgrading the Cloud Data Fusion version for Replication jobs is broken. Upgrading Replication jobs to Cloud Data Fusion version 6.8.0 isn't recommended.

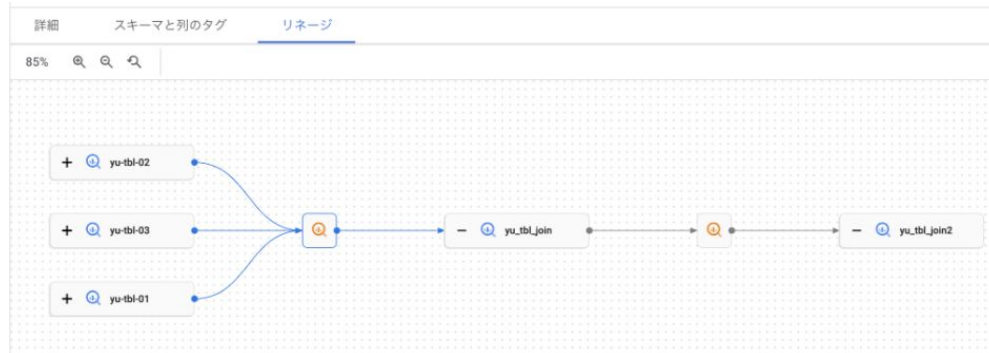
# Dataplex - Data Lineage

## / 特徴

- Table の依存関係を可視化
- Datamart 作成の SQL も確認可能
- BQ の UI から確認可能

## / Use case

- 生成元データの確認
- データガバナンス整備



# Cloud Composer - Data lineage

## / 特徴

- Composer 環境で Dataplex data lineage integration を有効にすることで、lineage 情報が Data Lineage API 経由で連携される
- 対応 version : Composer 2 (version 2.1.2+) (Airflow version 2.2.5+) のみ
- 対応 operators : BigQuery, GCS 関連 ([links](#))
- Airflow タスク実行の終了時に Lineage 情報が連携される (1-2 sec) 既存タスクのパフォーマンスには影響しない

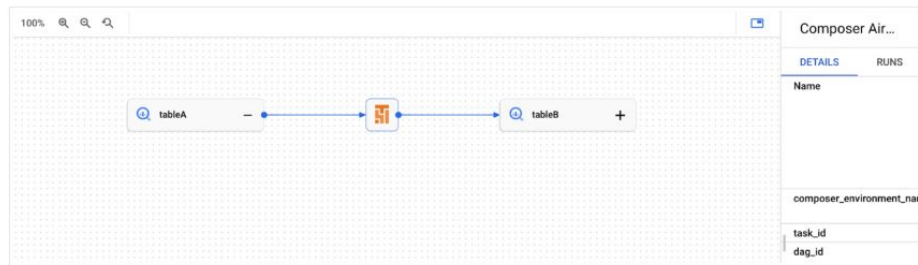
## / Use case

- Composer, Data Fusion, BigQuery, GCS, などを横断してパイプラインを構成しており、全体の Lineage 情報を統合管理したい

For example, running the following task:

```
task = BigQueryInsertJobOperator(
    task_id='snapshot_task',
    dag=dag,
    location='<dataset-location>',
    configuration={
        'query': {
            'query': 'SELECT * FROM dataset.tableA',
            'useLegacySql': False,
            'destinationTable': {
                'project_id': GCP_PROJECT,
                'dataset_id': 'dataset',
                'table_id': 'tableB',
            },
        },
    },
)
```

Results in creating the following lineage graph in Dataplex UI:



# Cloud Composer - Data lineage

/Tips

- Operator が実行する SQL がそのまま表示される  
INSERT などの部分が省略されているため注意が必要

```
# Task 3
t3 = BigQueryOperator(
    task_id='t3_bq_write_to_github_agg',
    use_legacy_sql=False,
    write_disposition='WRITE_TRUNCATE',
    allow_large_results=True,
    sql='''
#standardSQL
SELECT
    '{{ yesterday_ds_nodash }}' as date,
    repo,
    SUM(stars) as stars_last_7_days,
    SUM(IF(_PARTITIONTIME BETWEEN TIMESTAMP('{{ yesterday_ds }}')
        AND TIMESTAMP('{{ yesterday_ds }}'),
        stars, null)) as stars_last_1_day,
    SUM(forks) as forks_last_7_days,
    SUM(IF(_PARTITIONTIME BETWEEN TIMESTAMP('{{ yesterday_ds }}')
        AND TIMESTAMP('{{ yesterday_ds }}'),
        forks, null)) as forks_last_1_day
FROM
    `{{ PROJECT_ID }}_github_trends.github_daily_metrics`
WHERE _PARTITIONTIME BETWEEN TIMESTAMP('{{ macros.ds_add(ds, -6) }}')
AND TIMESTAMP('{{ yesterday_ds }}')
GROUP BY
    date,
    repo
'''.replace('{{ PROJECT_ID }}', PROJECT_ID),
    destination_dataset_table='{{ PROJECT_ID }}_github_trends.github_aggs{{ yesterday_ds_nodash }}.replace(
        '{{ PROJECT_ID }}', PROJECT_ID),
    dag=dag)

```

This is a partitioned table. [Learn more](#)

SCHEMA DETAILS PREVIEW LINEAGE **PREVIEW**

Changes made in other systems may not be reflected immediately on lineage.

100% 🔍 🔍 🔍

Query

DETAILS	RUNS
Name	projects/1062509424947/locations/us/processes/bfd0ef139dc4246c31eb5b9230a3e042
job_id	airflow_1671757819540913_069b88f68576ca241469852ae0811ead

```
#standardSQL
SELECT
    "20221214" as date,
    repo,
    SUM(stars) as stars_last_7_days,
    SUM(IF(_PARTITIONTIME BETWEEN TIMESTAMP("2022-12-14")
        AND TIMESTAMP("2022-12-14") ,
        stars, null)) as stars_last_1_day,
    SUM(forks) as forks_last_7_days,
    SUM(IF(_PARTITIONTIME BETWEEN TIMESTAMP("2022-12-14")
        AND TIMESTAMP("2022-12-14") ,
        forks, null)) as forks_last_1_day
FROM
    `my-project-test01-argolis.github_trends.github_daily_metrics`
WHERE _PARTITIONTIME BETWEEN TIMESTAMP("2022-12-09")
AND TIMESTAMP("2022-12-14")
GROUP BY
    date,
    repo

```

# BigQuery - query Cloud SQL with Private connection

## / 特徴

- 従来は Public connection のみだった Cloud SQL に対する Federated Query が Private Connection にも対応

## / Use case

- Cloud SQL を 閉域網内のみで利用 (Private Connection モード) されているお客様で、BigQuery を用いてデータ分析を行いたい

The screenshot shows the Google Cloud BigQuery interface. On the left, the Explorer pane displays a list of external connections under the project 'my-project-tetsuki-argolis'. The connection 'us-central1.test-connection' is selected. The main pane shows a query editor with the following SQL:

```
1 SELECT * FROM EXTERNAL_QUERY("us-central1.test-connection",
2 "select * from information_schema.tables;");
```

The 'Query results' pane shows the following table:

Row	TABLE_CATALOG	TABLE_SCHEMA	TABLE_NAME	TABLE_TYPE	ENGIN
1	def	mysql	audit_log_rules	BASE TABLE	Innoi
2	def	mysql	audit_log_rules_expanded	BASE TABLE	Innoi
3	def	mysql	audit_log_supported_ops	BASE TABLE	Innoi
4	def	mysql	cloudsql_replica_index	BASE TABLE	Innoi
5	def	mysql	columns_priv	BASE TABLE	Innoi
6	def	mysql	component	BASE TABLE	Innoi
7	def	mysql	db	BASE TABLE	Innoi
8	def	mysql	default_roles	BASE TABLE	Innoi
9	def	mysql	engine_cost	BASE TABLE	Innoi
10	def	mysql	func	BASE TABLE	Innoi
11	def	mysql	general_log	BASE TABLE	CSV
12	def	mysql	global_grants	BASE TABLE	Innoi
13	def	mysql	gtid_executed	BASE TABLE	Innoi
14	def	mysql	heartbeat	BASE TABLE	Innoi



# 03

## AI, ML - BQML, Vertex AI, その他

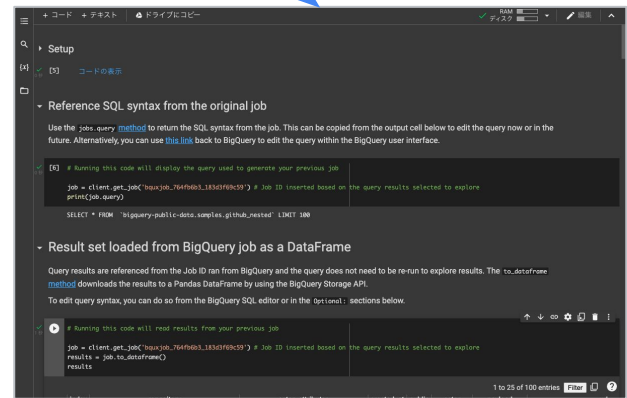
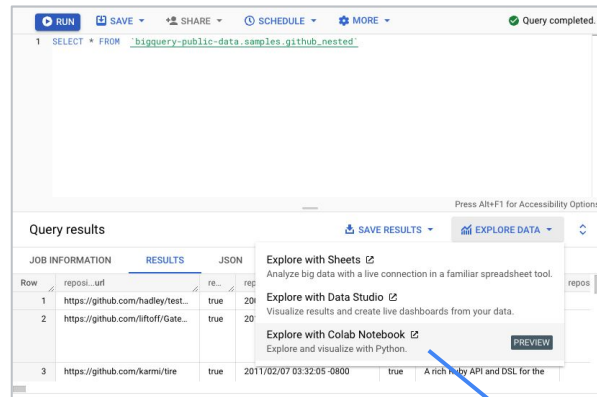
# Explore query results in Colab

## / 特徴

- BigQuery の SQL 実行結果を Colab Notebooks で分析できる

## / ユースケース

- BigQuery にデータがあり、Colab でデータサイエンスのタスクを実行している



[Explore data in Colab](#)

# Stored procedure for Apache Spark

## / 特徴

- BigQuery から SQL で Apache Spark の pyspark のコードがストアードプロシージャとして実行可能
- 実行環境には Dataproc を指定できる

## / ユースケース

- PySpark で実装されたワークロードがあり、BigQuery のワークロードと統合して容易に運用

## / 注意点

- Scala は未対応

[Apache Spark のストアード プロシージャを操作する](#)

[Connect to Apache Spark](#)

## PySpark コードを直接書いた場合のストアードの作成

```
CREATE OR REPLACE PROCEDURE my_bq_project.my_dataset.spark_proc()
WITH CONNECTION `my-project-id.us.my-connection`
OPTIONS(engine="SPARK")
LANGUAGE PYTHON AS R"""
from pyspark.sql import SparkSession

spark = SparkSession.builder.appName("spark-bigquery-demo").getOrCreate()

# Load data from BigQuery.
words = spark.read.format("bigquery") \
    .option("table", "bigquery-public-data:samples.shakespeare") \
    .load()
words.createOrReplaceTempView("words")

# Perform word count.
word_count = words.select('word', 'word_count').groupBy('word').sum('word_count')
word_count.show()
word_count.printSchema()

# Saving the data to BigQuery
word_count.write.format("bigquery") \
    .option("writeMethod", "direct") \
    .save("wordcount_dataset.wordcount_output")
"""
```

## GCS にある PySpark コードを参照するパターン

```
CREATE PROCEDURE my_bq_project.my_dataset.spark_proc()
WITH CONNECTION `my-project-id.us.my-connection`
OPTIONS(engine="SPARK", main_file_uri="gs://my-bucket/my-pyspark-main.py")
LANGUAGE PYTHON
```

# Analytics Hub

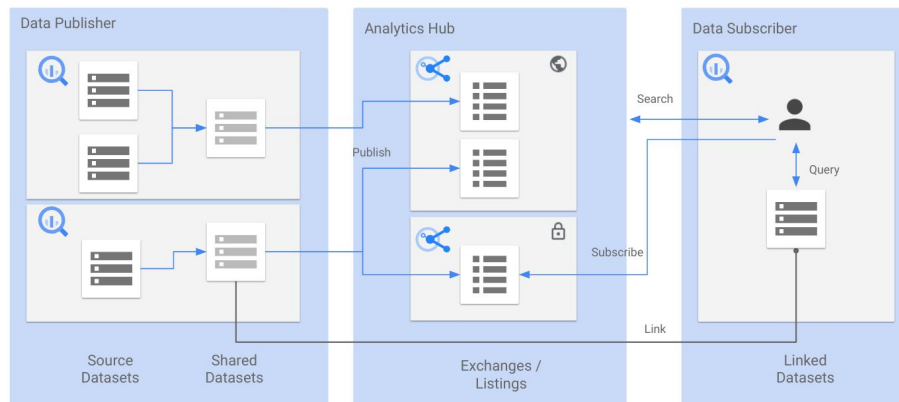
## / 特徴

- 組織の境界を越えてデータと分析情報を共有するデータ交換プラットフォーム
- セルフサービスで、データを複製せず、セキュアに共有可能
- Google の Public Dataset とかけあわせて分析もできる
- Dataplex(Data Catalog) 統合で検索可能に
- 公開されたデータのサブスクリプションの一覧と削除が可能に

## / ユースケース

- 自社のデータを他社に提供し収益化
- 組織内のデータ共有のセルフサービス化

## Analytics Hub



# BigQuery - BI Engine の新しいダッシュボード メトリックス

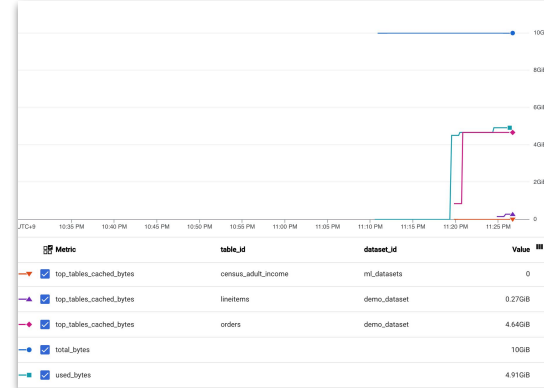
## / 特徴

- BI Engine Top Tables Cached Bytes
  - 上位のテーブルのキャッシュ上のバイト数
- BI Engine Fallback Count
  - BI Engine を利用しなかったクエリの割合(クエリ数/秒)
  - 原因ごとに表示も可能

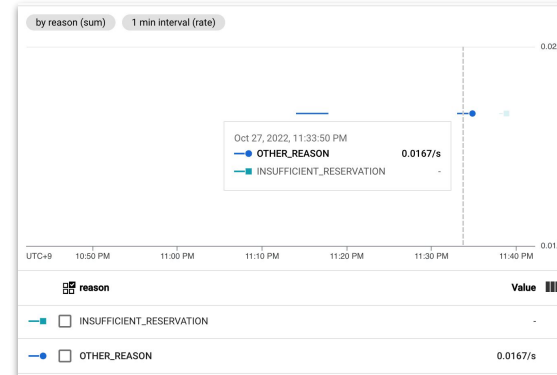
## / ユースケース

- BI Engine が有効に利用されているかの確認
- 目的のテーブルがキャッシュされているか
- BI Engine の効果を測る

### BI Engine Top Tables Cached Bytes



### BI Engine Fallback Count



# BigQuery - Search Index

## / 特徴

- テキストインデックスによりテーブル内から高速に特定の文字列を検索
- 列名が不明でもテーブル全体から高速に検索
- STRING と Native JSON データに対応
- 専用の SEARCH 関数を提供
- 索引のメンテナンスは BigQuery が自動的に実施
- スロットの利用量とクエリコストを削減
- **Tokyo リージョンでも使用可能に**

## / ユースケース

- ログ分析
- セキュリティ監査
- 選択性の高い文字列検索

## / 制限

- あいまい検索やタイポ等のドキュメント系の検索には未対応

## Search Indexの作成

```
CREATE SEARCH INDEX my_index ON Logs(ALL COLUMNS);
```

## Search Indexを利用した検索

```
SELECT * FROM Logs WHERE SEARCH(Logs, 'bar');
```

# BigQuery - リモート関数

## / 特徴

- SQL から Cloud Functions / Cloud Run を実行可能
- BQ のユーザ定義関数で実装できない処理が実装可能 (外部APIの呼び出し)
- 複数の言語で実装可能 (Node.js, Python, Go, Java, .Net, Ruby, PHP)
- [VPC-SC のサポート](#)

## / ユースケース

- SQLで
  - Cloud DLP で匿名化、再識別化
  - 画像、動画、自然言語 API の利用
- ELT ができる範囲が広くなりアーキテクチャのシンプル化

## Cloud Functionsで関数の定義

```
import json
_MAX_LOSSLESS=9007199254740992

def remote_add(request):
    try:
        return_value = []
        request_json = request.get_json()
        calls = request_json['calls']
        for call in calls:
            return_value.append(sum([int(x) if isinstance(x, str) else x for x in call if x
is not None]))
        replies = [str(x) if x > _MAX_LOSSLESS or x < -_MAX_LOSSLESS else x for x in
return_value]
        return_json = json.dumps( { "replies" : replies } )
        return return_json
    except Exception as inst:
        return json.dumps( { "errorMessage": 'something unexpected in input' } ), 400
```

## BigQueryでリモート関数の作成

```
CREATE OR REPLACE FUNCTION demo_dataset.remote_add(x INT64, y INT64)
RETURNS INT64
REMOTE WITH CONNECTION `tetsu-project01.us.remote_add_connection`
OPTIONS (endpoint =
'https://us-central1-tetsu-project01.cloudfunctions.net/remote_add');
```

## SQLから実行

```
SELECT val, demo_dataset.remote_add(val, 2)
FROM UNNEST([NULL,2,3,5,8]) AS val;
```

# BigQuery - Object Tables

## / 特徴

- BigQuery で Cloud Storage 上の非構造化データを SQL で分析可能に
- ファイルのメタデータ情報の保持
- Data という疑似列も参照可能でBQML と一緒に利用可能
- Signed Object URL を生成し、Remote functionへ渡すことも可能
- BigQuery のセキュリティも適用可能

## / ユースケース

- Cloud Storage 上の画像、動画等のファイルをSQLを使って機械学習にかけて分析する
- 画像、動画等のファイルをBigQuery 上のその他のデータとかけあわせて分析する

Field name	Type	Mode
<a href="#">uri</a>	STRING	NULLABLE
<a href="#">generation</a>	INTEGER	NULLABLE
<a href="#">content_type</a>	STRING	NULLABLE
<a href="#">size</a>	INTEGER	NULLABLE
<a href="#">md5_hash</a>	STRING	NULLABLE
<a href="#">updated</a>	TIMESTAMP	NULLABLE
▼ <a href="#">metadata</a>	RECORD	REPEATED
<a href="#">name</a>	STRING	NULLABLE
<a href="#">value</a>	STRING	NULLABLE

```
CREATE TABLE my_dataset.my_inference_results AS
SELECT uri, content_type, vision_feature
FROM ML.PREDICT(
  MODEL my_dataset.vision_model,
  SELECT ML.DECODE_IMAGE(data) AS vision_input
  FROM my_dataset.object_table
);
```

[Introduction to object tables](#)



# BigQuery ML - Vertex AI Model Registry との統合

## / 特徴

- Vertex AI Model Registry に BigQuery ML モデルを登録して管理、バージョンニング、モニタリングできる
- 登録したモデルを Vertex AI Endpoints へデプロイ
- 評価指標の比較、追跡ができる

```
CREATE MODEL bqmlga4.churn_logreg
OPTIONS(
  MODEL_REGISTRY="vertex_ai",
  vertex_ai_model_id="bqmlga4_churn_logreg",
  MODEL_TYPE="LOGISTIC_REG",
  INPUT_LABEL_COLS=["churned"]
) AS
SELECT (省略)
```

```
CREATE MODEL bqmlga4.churn_logreg2
OPTIONS(
  MODEL_REGISTRY="vertex_ai",
  vertex_ai_model_id="bqmlga4_churn_logreg",
  MODEL_TYPE="LOGISTIC_REG",
  INPUT_LABEL_COLS=["churned"]
) AS
SELECT (省略)
```

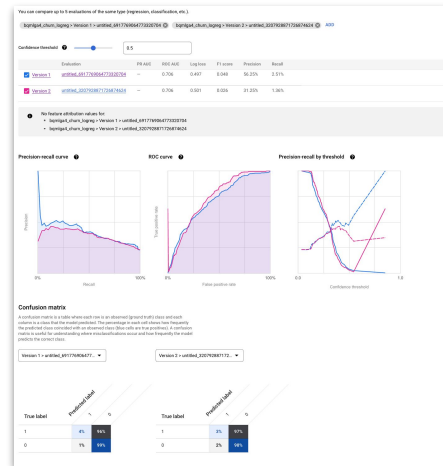
Version ID	Alias	Status
2	-	Ready
1	default	Ready

Model Registry に登録

## / Use case

- BigQuery で作成したモデルを使い Vertex AI でオンライン予測をする
- BigQuery ML を含めた ML Ops
- Dataform (preview) でデータとモデルをビルドしてステージングする

## 評価指標の比較



# BigQuery ML - Vision model の画像分析のサポート

## / 特徴

- Vision model を import として BigQuery ML でオブジェクトの検出やOCR 等が可能
- BigQuery の Object テーブルで GCS 上のファイルにアクセス
- Tensorflow モデル
  - ImageNet
  - ResNet 50

## / Use case

- BigQuery ML で画像の推論や文字の認識をする

```
CREATE MODEL my_dataset.my_vision_model
OPTIONS(
  model_type = 'TENSORFLOW',
  model_path =
  'gs://my_bucket/my_model_folder/*');
```

```
SELECT * FROM
ML.PREDICT(
  MODEL my_dataset.vision_model,
  (SELECT uri, ML.DECODE_IMAGE(data) AS input
  FROM my_dataset.object_table)
);
```

# BigQuery ML - Sparse Input のサポート

## / 特徴

- 値がほとんどゼロまたは空であるデータに対するトレーニングにおける空間効率が改善される(1 hot-encodingなどによって発生する多数の0の無駄を省ける)
- Array [STRUCT < int, numerical > ] で Sparse なカテゴリ変数を効率的に保存

## / Use case

- 多数のカテゴリ変数をもつデータでモデルをトレーニングする

	bigquery	features	hello	now	sparse	supports	world
Hello World	0	0	1	0	0	0	1
BigQuery supports sparse features now	1	1	0	1	1	1	0

ARRAY<STRUCT<k INT64, v INT64>> で表すと

Hello World は、[(2, 1), (6, 1)]

BigQuery supports .... は、[[0,1],[1,1],[3,1],[4,1],[5,1]]

[Support sparse input for feature columns](#)  
[Sparse Features Support in BigQuery](#)

# BigQuery ML - ARIMA\_PLUS\_XREG モデル

## / 特徴

- 多変量時系列予測が可能に
- 数値、カテゴリ、配列の自動特徴量エンジニアリング
- 季節の傾向や休日の検出などARIMA\_PLUSモデルの機能も使用できる
- Reservation のみで利用可能

## / Use case

- 多変量解析による時系列予測

[シアトルの大気質データからの多変量時系列予測](#)  
[The CREATE MODEL statement for the ARIMA\\_PLUS\\_XREG model](#)

```
CREATE OR REPLACE MODEL
  `bqml_tutorial.seattle_pm25_xreg_model`
OPTIONS (
  MODEL_TYPE = 'ARIMA_PLUS_XREG',
  time_series_timestamp_col = 'date',
  time_series_data_col = 'pm25')
AS
SELECT
  date,
  pm25,
  temperature,
  wind_speed
FROM
  `bqml_tutorial.seattle_air_quality_daily`
WHERE
  date
  BETWEEN DATE('2012-01-01')
  AND DATE('2020-12-31')
```

# BigQuery ML - Vertex AI Pipelines での BigQuery ML 用のコンポーネント

## / 特徴

- 20 以上のコンポーネントが利用可能
- パイプラインのメタデータのトラッキングが可能
- Vertex AI のオンライン予測などシームレスな連携が可能に

## / Use case

- Vertex AI Pipeline で BigQuery ML を含めたML のパイプラインを作成する

## 主要なコンポーネント

- BigqueryQueryJobOp
- BigqueryCreateModelJobOp
- BigqueryExportModelJobOp
- BigqueryPredictModelJobOp
- BigqueryEvaluateModelJobOp
- BigqueryDropModelJobOp
- BigqueryEvaluateModelJobOp
- BigqueryExplainForecastModelJobOp
- BigqueryExplainPredictModelJobOp
- BigqueryForecastModelJobOp

## Appendix

# その他のアップデート

# その他

## Cloud Composer

October 06, 2022

**ANNOUNCEMENT**

Starting from January 2023, the default version for new Cloud Composer environments changes from Cloud Composer 1 to Cloud Composer 2. New environments will use the latest Cloud Composer 2 version with the default Airflow 2 version ( `composer-2-airflow-2` ). Currently, the default version is `composer-1-airflow-1.10.15` .

## Dataproc

October 03, 2022

**ANNOUNCEMENT**

Preemptible [SPOT VMs](#) can be used as [secondary workers](#) in a Dataproc cluster. Unlike [legacy preemptible VMs](#) with a 24-hour maximum lifetime, Spot VMs have no maximum lifetime.

# その他

## BigQuery

2022/11/14

FEATURE

The [slot estimator](#) helps you manage slot capacity based on historical performance metrics. This feature is now [generally available \(GA\)](#).

## Dataproc

2022/11/14

FEATURE

If a Dataproc Metastore service uses the [gRPC endpoint protocol](#), a Dataproc or self-managed cluster located in any region can [attach to the service](#).

2022/11/17

FEATURE

Dataproc Serverless for Spark supports [Spark and System metrics](#). These metrics are enabled by default. Spark driver and executor metrics can be customised using overrides.



# その他

Pub/Sub

2022/12/1

ANNOUNCEMENT

Exactly once delivery is now GA.

# その他

## BigQuery

2022/12/8

FEATURE

The [demo query guide](#) helps you query a public dataset from Google Trends and is now in [preview](#).

## Data Fusion

2022/11/30

DEPRECATED

Cloud Data Fusion version 6.4 is no longer [supported](#) as of December 1, 2022. You should upgrade your instances to run in a supported version. For instructions, see [Upgrading your Cloud Data Fusion environment](#).

2022/12/01

FEATURE

Google Cloud Platform Plugins version 0.20.4 is generally available ([GA](#)) in Cloud Data Fusion versions 6.7.1 and 6.7.2. This version includes Dataplex [Source](#) and [Sink](#) plugins in GA. For more information, see the [CDAP Hub release log](#).

FEATURE

Google Cloud Platform Plugins version 0.19.3 is generally available ([GA](#)) in Cloud Data Fusion version 6.6.0. This version includes Dataplex [Source](#) and [Sink](#) plugins in GA. For more information, see the [CDAP Hub release log](#).

# その他

## BigQuery

December 22, 2022

FEATURE

The [Lineage tab](#) in the table properties page lets you track how your data moves and transforms through BigQuery. This feature is now in [preview](#).

CHANGED

BigQuery now blocks [saving query results to Google Drive](#) from projects inside a [VPC Service Controls protected perimeter](#).

## PubSub Lite

December 19, 2022

FEATURE

Pub/Sub Lite now supports [export subscriptions](#). You can use an export subscription to export Pub/Sub Lite messages to a destination Pub/Sub topic. This feature is [generally available \(GA\)](#).

# その他

## BigQuery

### FEATURE

The following geography functions are now [generally available](#) (GA):

- `ST_ISCLOSED`: Returns `TRUE` for a non-empty geography, where each element in the geography has an empty boundary.
- `ST_ISRING`: Checks if a geography is a linestring and if the linestring is both closed and simple.

## Dataproc

### FEATURE

[Dataproc Serverless for Spark](#) now supports `spark.dataproc.diagnostics.enabled` property that enables auto diagnostics on Batch failure. Note that enabling auto diagnostics will hold compute and storage quota after Batch is complete and until diagnostics is finished.

# その他

## BigQuery ML Jan 03

CHANGED

Customers can use BigQuery ML to train and run models on BigLake in Cloud Storage. See [Data Cloud Blog](#) and [End to end unstructured data use cases demo](#).

## Dataflow Jan 03

FEATURE

Starting in version 2023-01-03-RC00, the [Google-provided Dataflow templates](#) support ES6 syntax for JavaScript user-defined functions (UDFs). This change is backwards-compatible. ES5 syntax and existing user-defined functions are still supported.

When you run Google-provided templates using the latest version, your jobs are upgraded automatically on restart. If you want to keep running an earlier version of a template, when you run the template, specify version `2022-12-15-00_RC00` or earlier.

# その他

## Cloud Data Fusion Jan05

FEATURE

The [SAP SuccessFactors Batch Source](#) plugin is available in [Preview](#). You can connect your data pipeline to an SAP SuccessFactors Source and a BigQuery Sink with this plugin in Cloud Data Fusion versions 6.5.1 and later.

# その他

## Cloud Composer

January 25, 2023

FEATURE

**Airflow 2.4.3** is available in Cloud Composer images.

# その他

February 01, 2023

FEATURE

The BigQuery Data Transfer Service can now [transfer data from Azure Blob Storage](#) into BigQuery. This feature is now in [preview](#).

January 31, 2023 

FEATURE

[Azure workload identity federation](#) is now [generally available \(GA\)](#) for BigQuery Omni connections. You can now [create a connection for federated identity](#) using Google Cloud console.

CHANGED

**Cloud console updates:** When you create datasets, select locations to run specific queries, or create exchanges in [Analytics Hub](#), you now see separate options for multi-region and specific regions. Based on your selection, you see a list with more options.

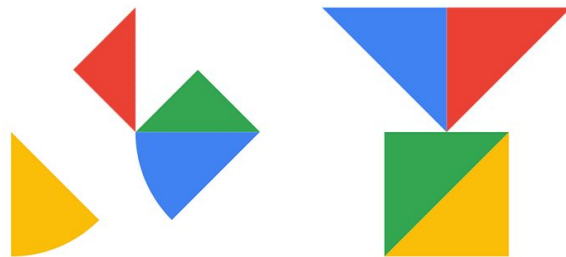
January 30, 2023

FEATURE

You can search for BigQuery partners in the [BigQuery Partner Center](#). This feature is in [Preview](#).



# Google Cloud Day '23 Tour



- 📍 東京 (オンライン) 5月23日(火)～25日(木)
- 📍 大阪 (ハイブリッド) 6月2日(金)
- 📍 名古屋 (ハイブリッド) 6月22日(木)
- 📍 福岡 (ハイブリッド) 6月30日(金)

企業のDXを加速する、そのヒントを4都市からお届けします。

今すぐ登録 [goo.gle/gcd23\\_1p](https://goo.gle/gcd23_1p)

