



Google Cloud

Introduction to Building Batch Data Pipelines

Agenda

EL, ELT, ETL

Quality considerations

How to carry out operations in BigQuery

Shortcomings

ETL to solve data quality issues

The method you use to load data depends on how much transformation is needed



Extract and Load



Extract, Load, Transform



Extract, Transform, Load

When would you use EL?

Architecture	When you'd do it
Extract data from files on Google Cloud Storage	Batch load of historical data
Load it into BigQuery's native storage	Scheduled periodic loads of log files (e.g. once a day)
You can trigger this from Cloud Composer, Cloud Functions, or scheduled queries	But only if the data are already clean and correct!

When would you use ELT?

Architecture	When you'd do it
Extract data from files in Google Cloud Storage into BigQuery Transform the data on the fly using BigQuery views, or store into new tables	Experimental datasets where you are not yet sure what kinds of transformations are needed to make the data useable. Any production dataset where the transformation can be expressed in SQL

Agenda

EL, ELT, ETL

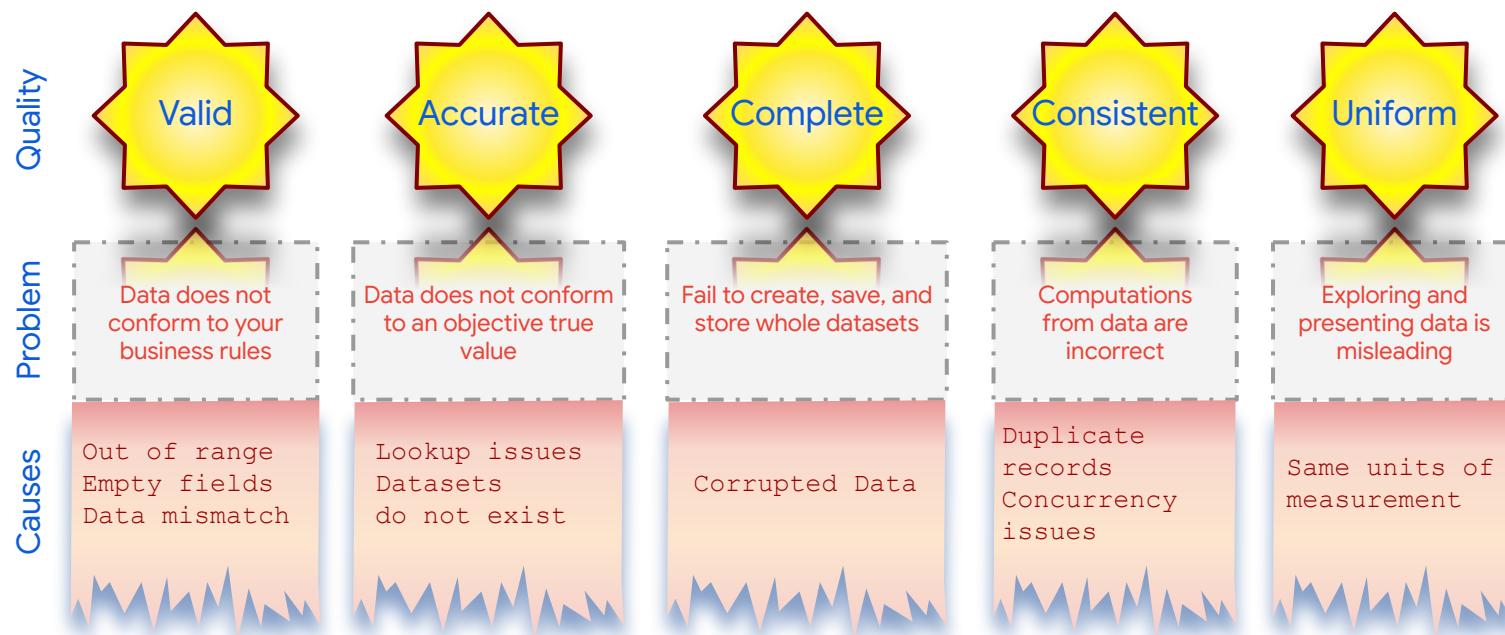
Quality considerations

How to carry out operations in BigQuery

Shortcomings

ETL to solve data quality issues

What are the purposes of Data Quality processing?



BigQuery can fix many data quality issues using SQL and Views



Agenda

EL, ELT, ETL

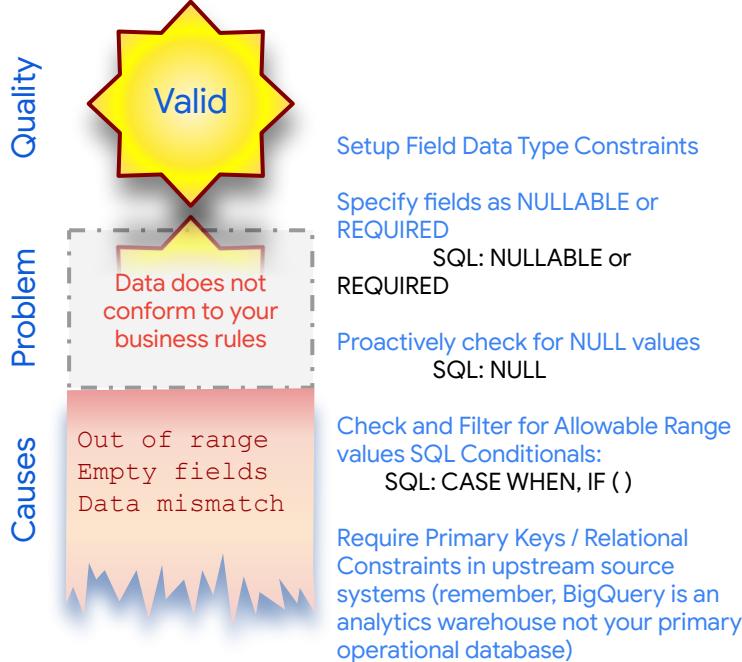
Quality considerations

How to carry out Transform operations in BigQuery

Shortcomings

ETL to solve data quality issues

Filter to identify and isolate invalid data



Setup Field Data Type Constraints

Specify fields as **NULLABLE** or **REQUIRED**

SQL: **NULLABLE** or **REQUIRED**

Proactively check for NULL values
SQL: **NULL**

Check and Filter for Allowable Range values
SQL Conditionals:
SQL: **CASE WHEN, IF ()**

Require Primary Keys / Relational Constraints in upstream source systems (remember, BigQuery is an analytics warehouse not your primary operational database)

Filter rows

WHERE (condition)

Filter aggregations

HAVING (condition)

Filters NULLs but leave blanks

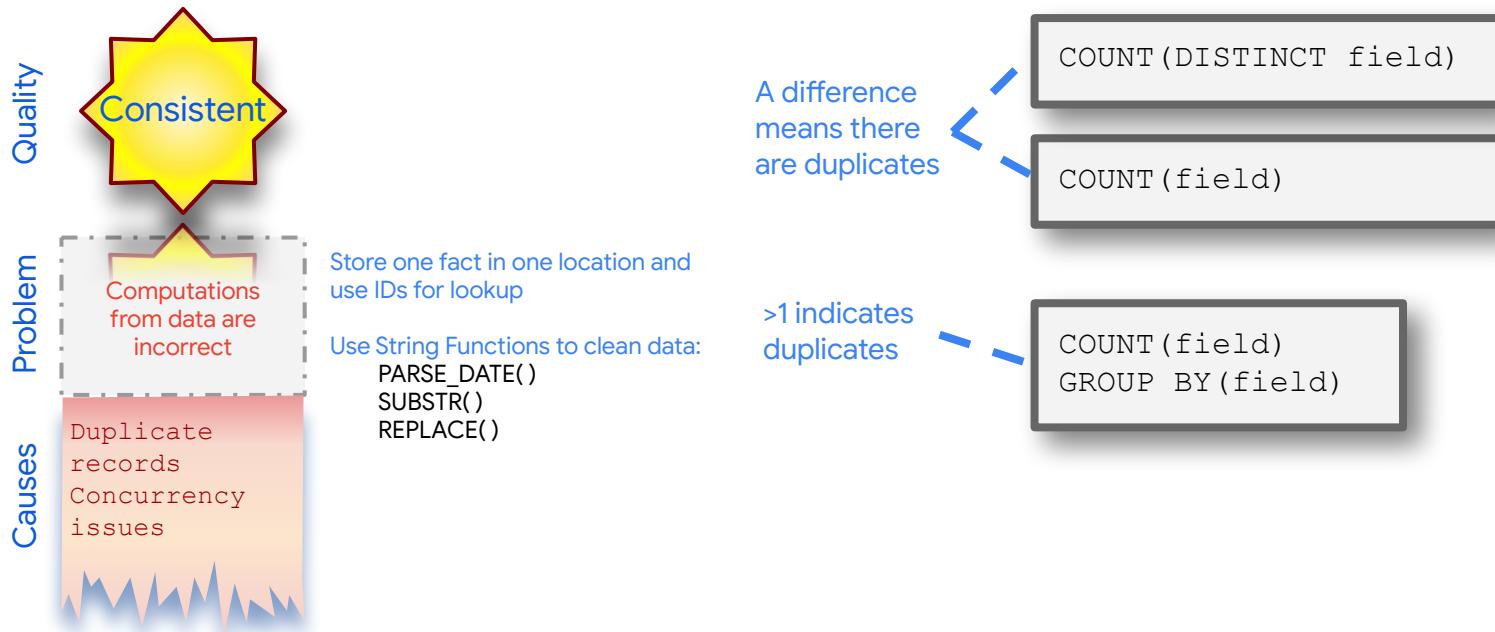
WHERE field IS NOT NULL

Filters NULLs and blanks

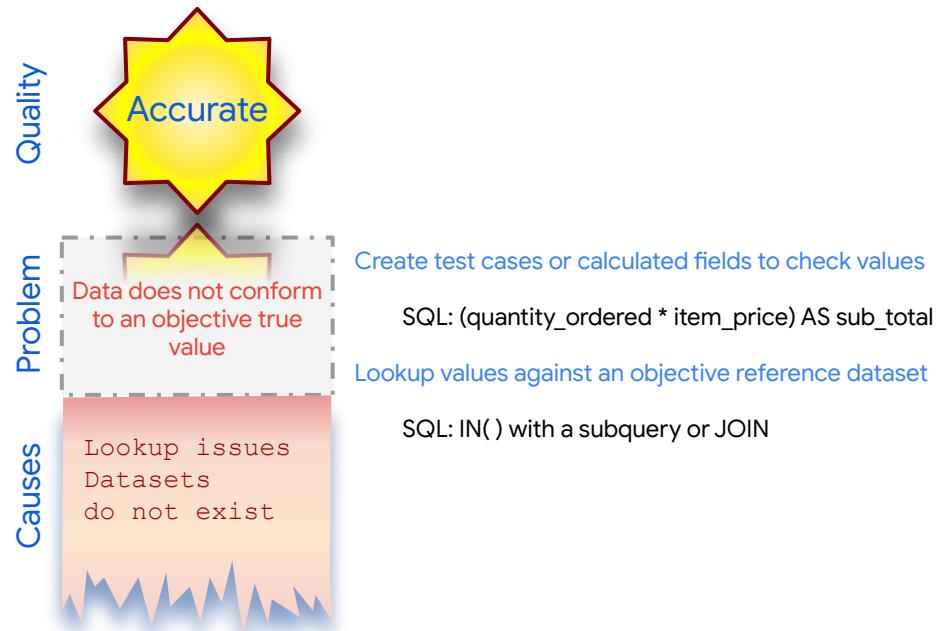
WHERE field IS NOT NULL AND field <> ""

A **NUL**L is the absence of data. A **BLANK** is a value of data.
Consider if you are trying to filter out both NULLS and BLANKS.

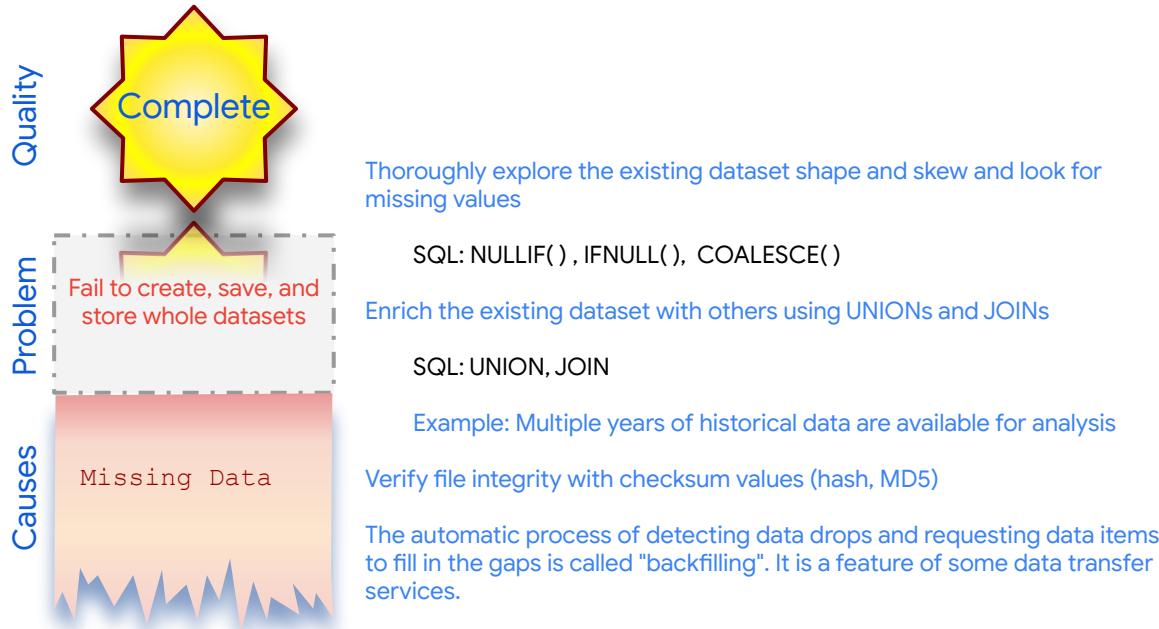
Detect duplication, enforce uniqueness for consistency



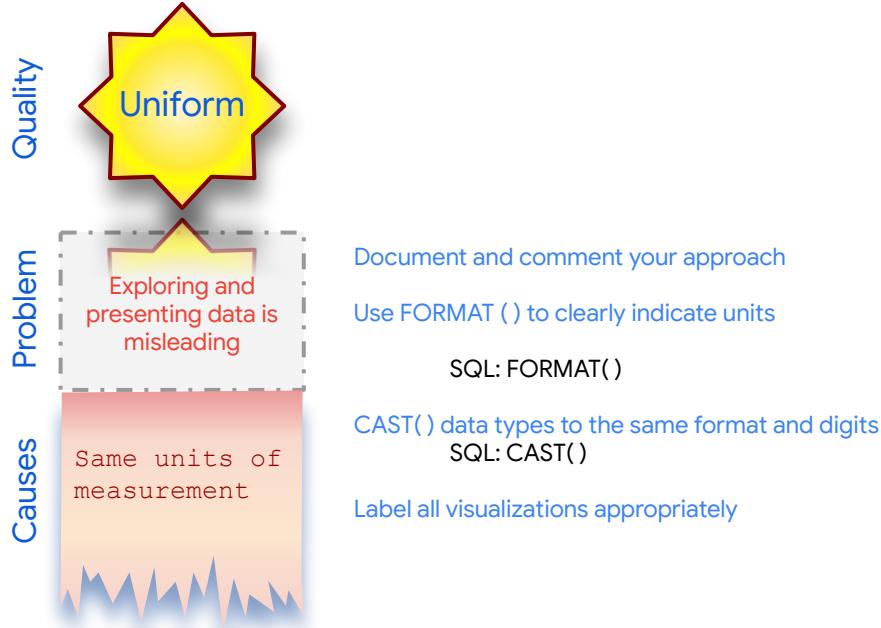
Test data against known good values for accuracy



Identify and fill in missing values for completeness.



Make data types and formats explicit for uniformity.



Demo

ELT to improve data quality in
BigQuery

Agenda

EL, ELT, ETL

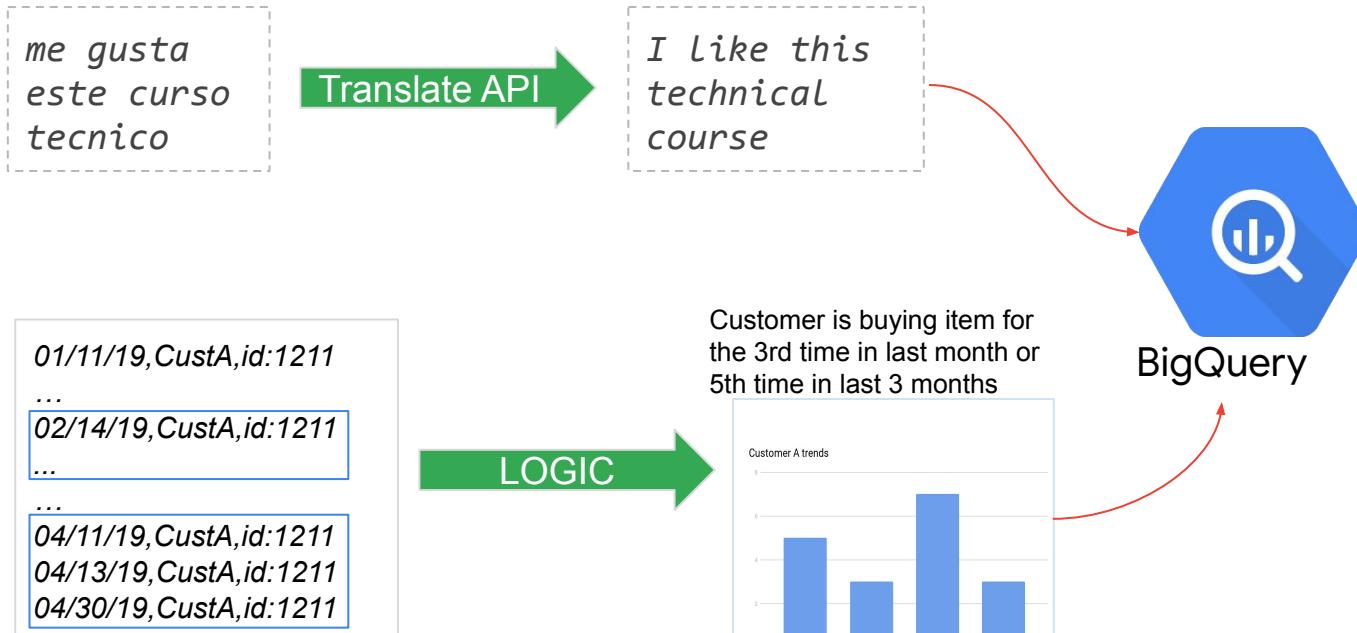
Quality considerations

How to carry out operations in BigQuery

Shortcomings

ETL to solve data quality issues

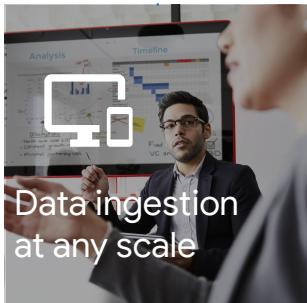
What if the transformations cannot be expressed in SQL? Or are too complex to do in SQL?



Build ETL pipelines in Dataflow and land the data in BigQuery

Architecture	When you'd do it
<p>Extract data from Pub/Sub, Google Cloud Storage, Cloud Spanner, Cloud SQL, etc.</p> <p>Transform the data using Cloud Dataflow</p> <p>Have Dataflow pipeline write to BigQuery</p>	<p>When the raw data needs to be quality-controlled, transformed, or enriched before being loaded into BigQuery.</p> <p>When the data loading has to happen continuously, i.e. if the use case requires streaming.</p> <p>When you want to integrate with continuous integration / continuous delivery (CI/CD) systems and perform unit testing on all components.</p>

Google Cloud offers a range of ETL tools



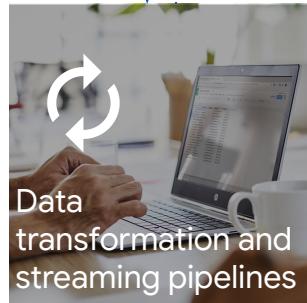
Cloud Pub/Sub



Data Transfer Services



Cloud IoT Core



Cloud Dataflow



Cloud Dataproc



Cloud Data Fusion



BigQuery



Cloud Storage



Cloud BigTable



AI Platform



Google Data Studio



AI Platform Notebooks



BI tools from Technology Partners



Google Sheets

Orchestration



Cloud Composer

Agenda

EL, ELT, ETL

Quality considerations

How to carry out operations in BigQuery

Shortcomings

ETL to solve data quality issues

Cases when you look beyond Dataflow and BigQuery

Issue	Solution
Latency, throughput	Dataflow to Bigtable
Reusing Spark pipelines	Cloud Dataproc
Need for visual pipeline building	Cloud Data Fusion

< Action Safe

Title Safe >

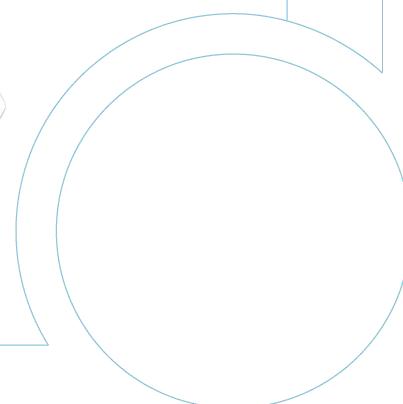
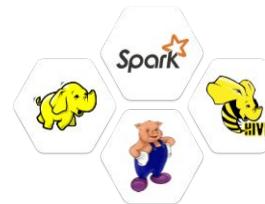
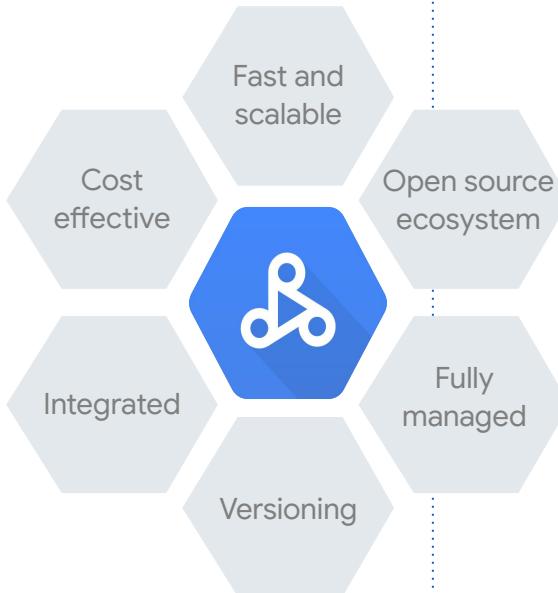
Cloud Dataproc is a managed service for batch processing, querying, streaming, and ML



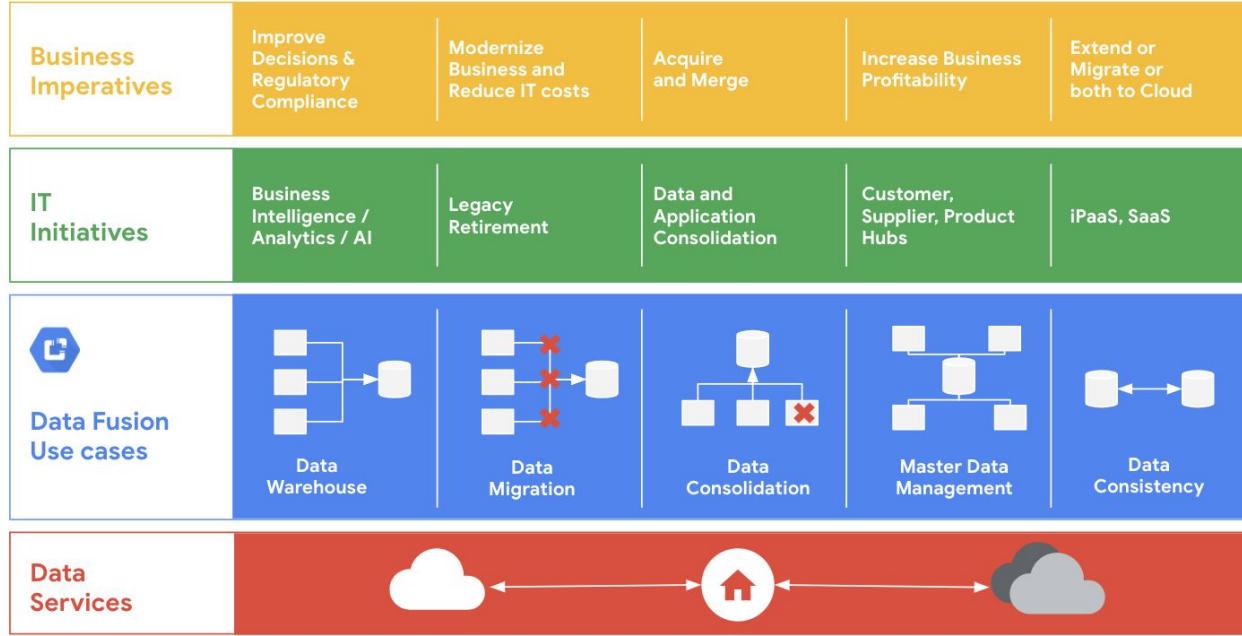
Google Cloud Platform

< Action Safe

Title Safe >



Cloud Data Fusion is a fully-managed, cloud native, enterprise data integration service for quickly building and managing data pipelines



Tracking lineage in ETL pipelines can be important



Where it came from



The processes it has
been through



Its present location and condition

Lineage: Metadata about the data

- What format is it in?
- What qualities does it have?
- Is it fit for the intended use?
- Can you transform or process it to make it fit for the intended use?

Labels on datasets, tables, and views can help track lineage

Label

Key:	Value
------	-------

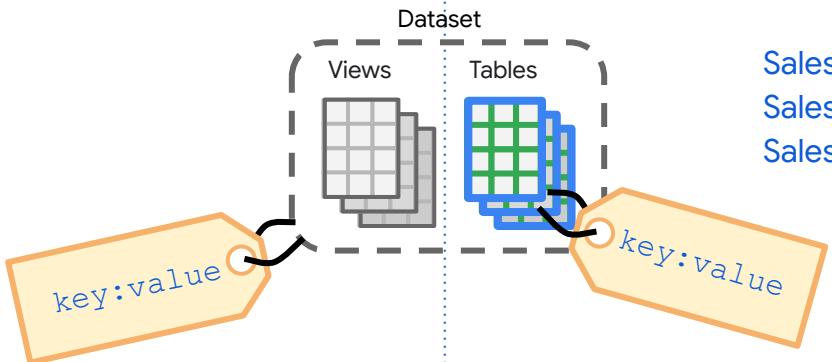
Example:

A series of similar tables:

Salesdata: Europe

Salesdata: March

Salesdata: Repeat customers



< Action Safe

Title Safe >

View your datasets and labels in Data Catalog

