# Analytics Use Cases on Databricks

databricks

**SECTION 5.1**

# How to Build a Marketing Analytics Solution Using Fivetran and dbt on Databricks

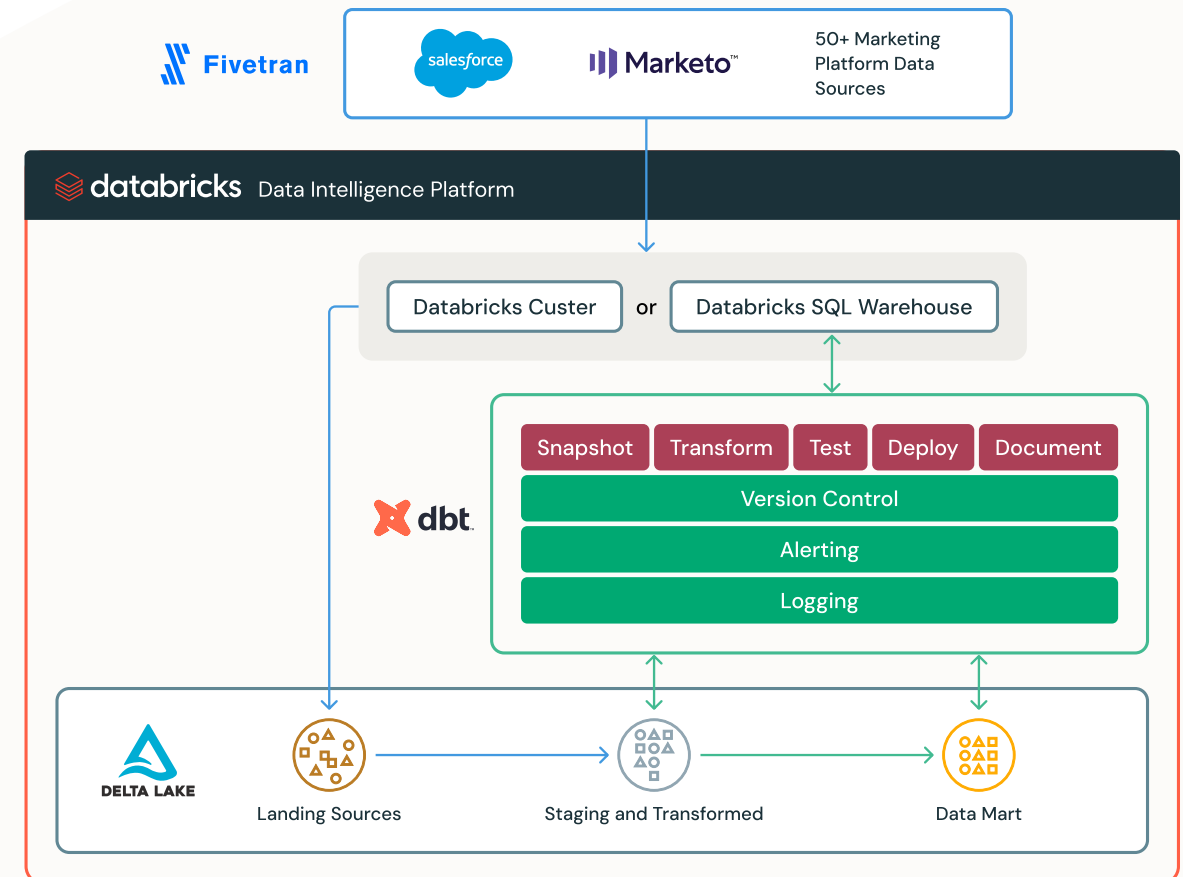by Tahir Fayyaz, Bilal Aslam and Robert Saxby

Marketing teams use many different platforms to drive marketing and sales campaigns, which can generate a significant volume of valuable but disconnected data. Bringing all of this data together can help to drive a large return on investment, as shown by Publicis Groupe who were able to increase campaign revenue by as much as 50%.

Databricks, which unifies data warehousing and AI use cases on a single platform, is the ideal place to build a marketing analytics solution: we maintain a single source of truth, and unlock AI/ML use cases. We also leverage two Databricks partner solutions, Fivetran and dbt, to unlock a wide range of marketing analytics use cases including churn and lifetime value analysis, customer segmentation, and ad effectiveness.

Fivetran allows you to easily ingest data from 50+ marketing platforms into Delta Lake without the need for building and maintaining complex pipelines. If any of the marketing platforms' APIs change or break, Fivetran will take care of updating and fixing the integrations so your marketing data keeps flowing in.

dbt is a popular open source framework that lets lakehouse users build data pipelines using simple SQL. Everything is organized within directories, as plain text, making version control, deployment, and testability simple. Once the data is ingested into Delta Lake, we use dbt to transform, test and document the data. The transformed marketing analytics data mart built on top of the ingested data is then ready to be used to help drive new marketing campaigns and initiatives.
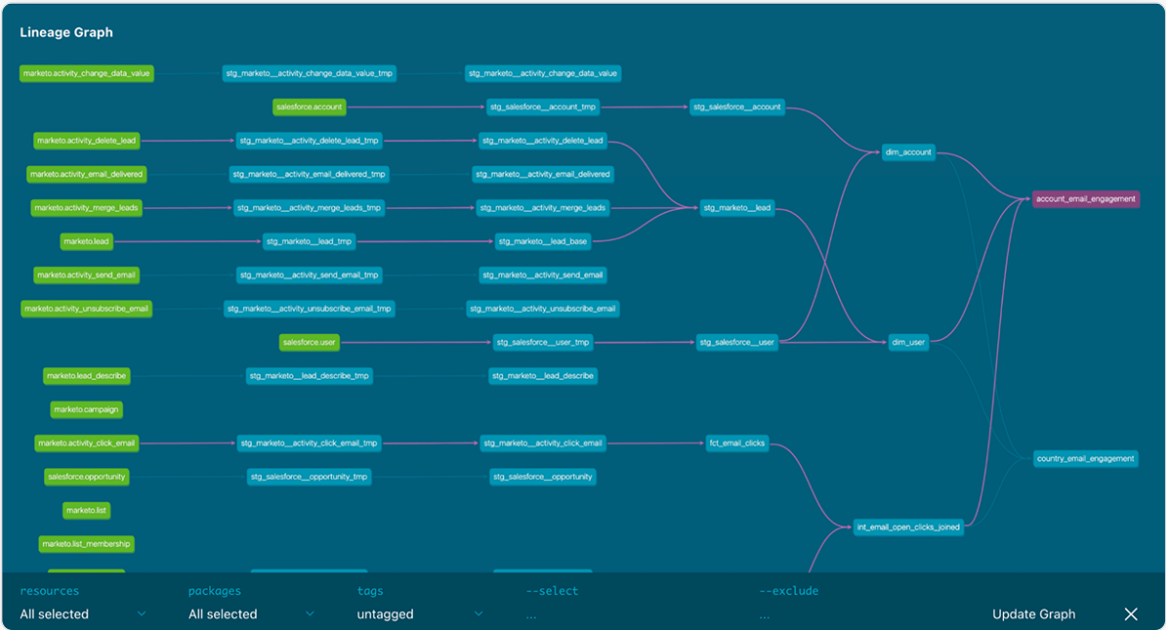


Fivetran and dbt can read and write to Delta Lake using a Databricks cluster or Databricks SQL warehouse

Both Fivetran and dbt are a part of Databricks Partner Connect, a one-stop portal to discover and securely connect data, analytics and AI tools directly within the Databricks Platform. In just a few clicks you can configure and connect these tools (and many more) directly from within your Databricks workspace.

databricks

## How to build a marketing analytics solution

In this hands-on demo, we will show how to ingest Marketo and Salesforce data into Databricks using Fivetran and then use dbt to transform, test, and document your marketing analytics data model.
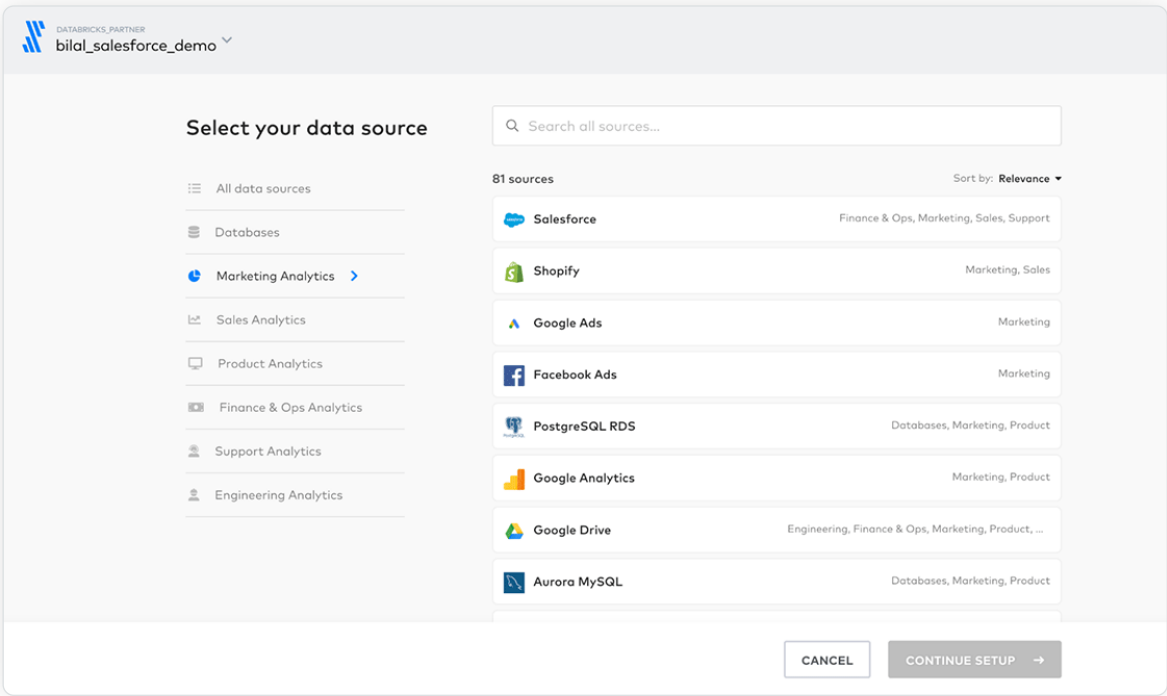
All the code for the demo is available on Github in the workflows-examples repository.



dbt lineage graph showing data sources and models

The final dbt model lineage graph will look like this. The Fivetran source tables are in green on the left and the final marketing analytics models are on the right. By selecting a model, you can see the corresponding dependencies with the different models highlighted in purple.
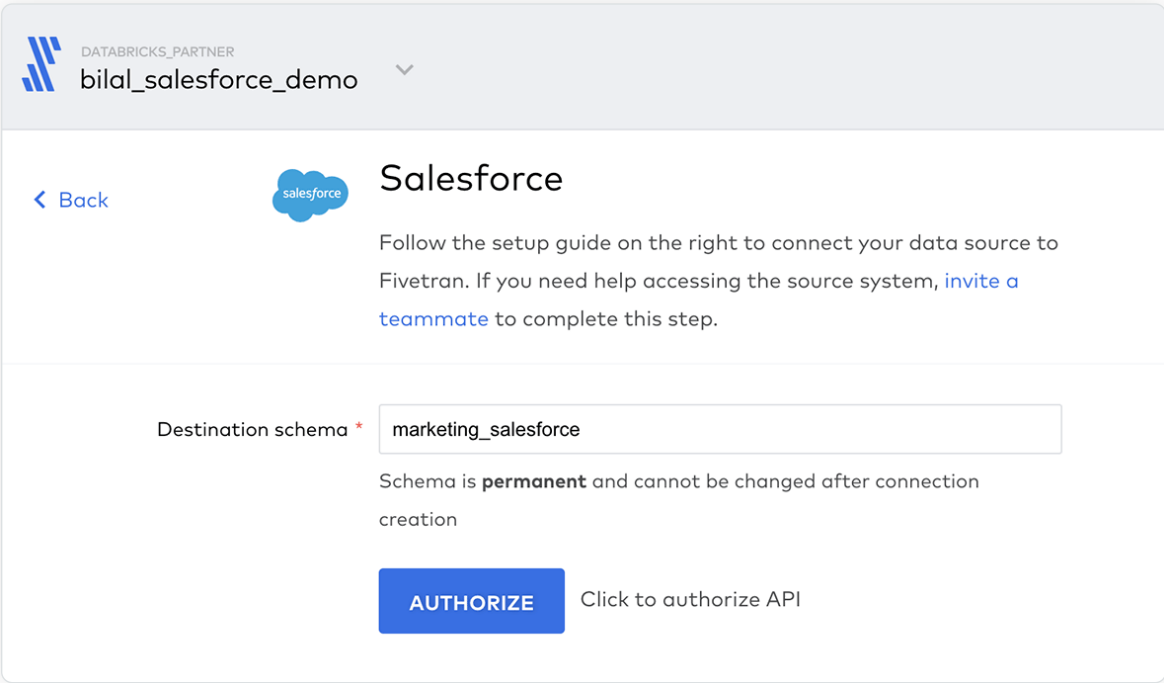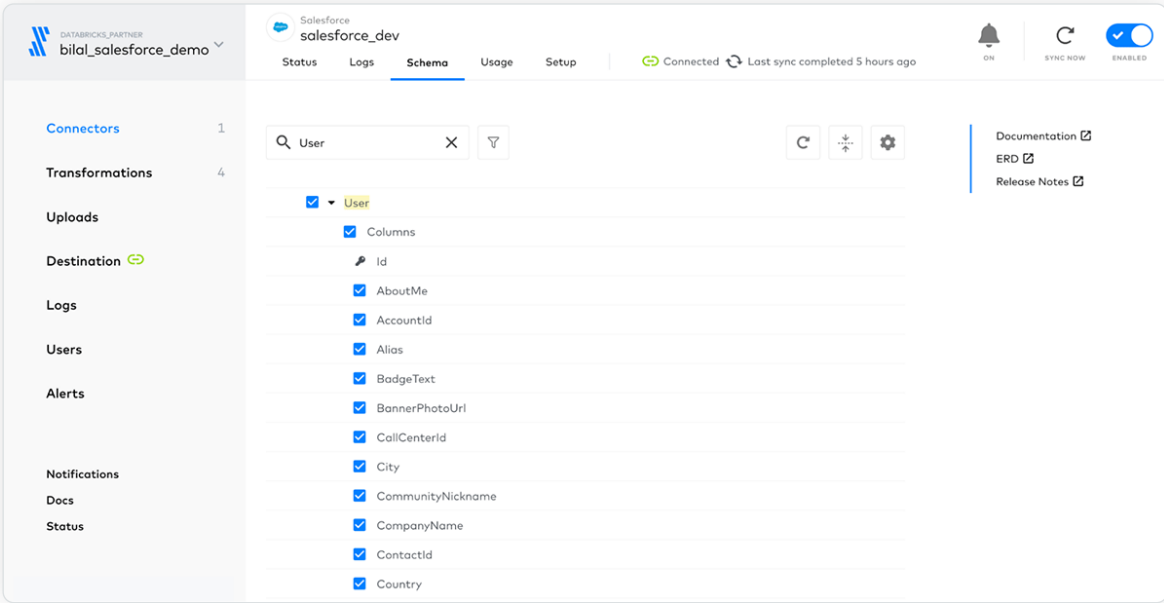
## Data ingestion using Fivetran



Fivetran has many marketing analytics data source connectors

Create new Salesforce and Marketo connections in Fivetran to start ingesting the marketing data into Delta Lake. When creating the connections Fivetran will also automatically create and manage a schema for each data source in Delta Lake. We will later use dbt to transform, clean and aggregate this data.

For the demo name the schemas that will be created in Delta Lake marketing_salesforce and marketing_marketo. If the schemas do not exist Fivetran will create them as part of the initial ingestion load.
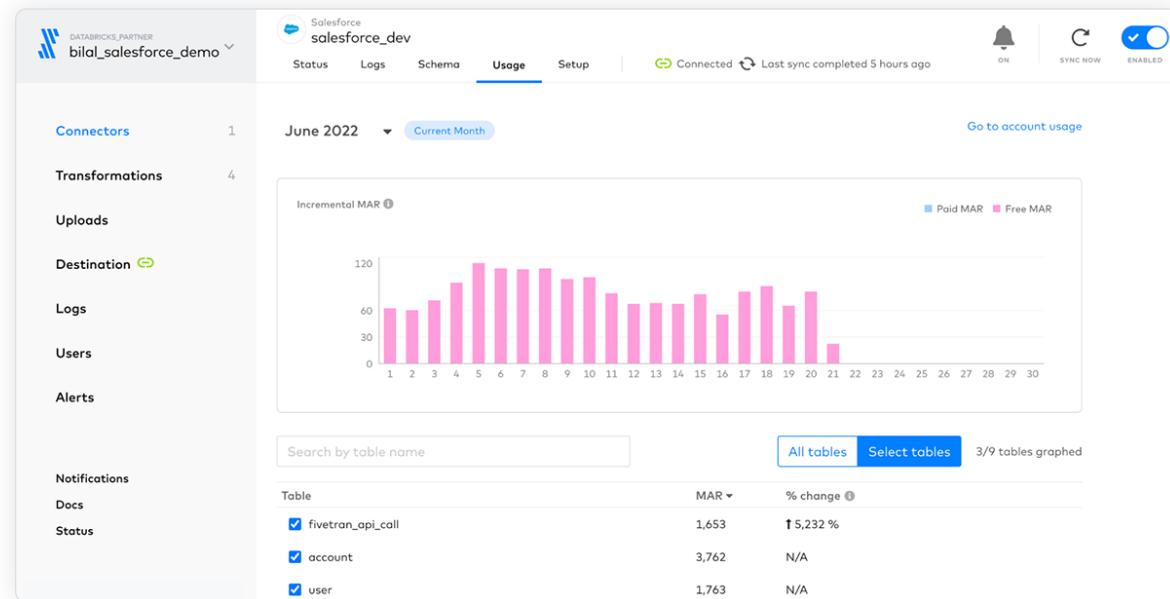


Define a destination schema in Delta Lake for the Salesforce data source



Select which data source objects to synchronize as Delta Lake tables

You can then choose which objects to sync to Delta Lake, where each object will be saved as individual tables. Fivetran also makes it simple to manage and view what columns are being synchronized for each table:



Fivetran monitoring dashboard to monitor monthly active rows synchronized

Additionally, Fivetran provides a monitoring dashboard to analyze how many monthly active rows of data are synchronized daily and monthly for each table, among other useful statistics and logs.

## Data modeling using dbt

Now that all the marketing data is in Delta Lake, you can use dbt to create your data model by following these steps

### Setup dbt project locally and connect to Databricks SQL

Set up your local dbt development environment in your chosen IDE by following the set-up instructions for dbt Core and dbt-databricks.

Scaffold a new dbt project and connect to a Databricks SQL Warehouse using dbt init, which will ask for following information.

```
$ dbt init
Enter a name for your project (letters, digits, underscore):
Which database would you like to use?
[1] databricks
[2] spark

Enter a number: 1
host (yourorg.databricks.com):
http_path (HTTP Path):
token (dapiXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX):
schema (default schema that dbt will build objects in):
threads (1 or more) [1]:
```

Once you have configured the profile you can test the connection using:

```
$ dbt debug
```

databricks

## Install Fivetran dbt model packages for staging

The first step in using the Marketo and Salesforce data is to create the tables as sources for our model. Luckily, Fivetran has made this easy to get up and running with their pre-built Fivetran dbt model packages. For this demo, let's make use of the `marketo_source` and `salesforce_source` packages.

To install the packages just add a `packages.yml` file to the root of your dbt project and add the `marketo-source`, `salesforce-source` and the `fivetran-utils` packages:

```
1   packages:
2     - package: dbt-labs/spark_utils
3       version: 0.3.0
4     - package: fivetran/marketo_source
5       version: [">=0.7.0", "=0.4.0", "
6
7   To download and use the packages run
```

```
1   $ dbt deps
```

You should now see the Fivetran packages installed in the packages folder.

## Update dbt_project.yml for Fivetran dbt models

There are a few configs in the `dbt_project.yml` file that you need to modify to make sure the Fivetran packages work correctly for Databricks.

The `dbt_project.yml` file can be found in the root folder of your dbt project.

**spark_utils overriding dbt_utils macros**

The Fivetran dbt models make use of macros in the `dbt_utils` package but some of these macros need to be modified to work with Databricks which is easily done using the `spark_utils` package.

It works by providing shims for certain `dbt_utils` macros which you can set using the dispatch config in the `dbt_project.yml` file and with this dbt will first search for macros in the `spark_utils` package when resolving macros from the `dbt_utils` namespace.

```
1   dispatch:
2     - macro_namespace: dbt_utils
3       search_order: ['spark_utils', 'dbt_utils']
```

**Variables for the marketo_source and salesforce_source schemas**

The Fivetran packages require you to define the catalog (referred to as database in dbt) and schema of where the data lands when being ingested by Fivetran.

Add these variables to the `dbt_project.yml` file with the correct catalog and schema names. The default catalog is `hive_metastore` which will be used if `_database` is left blank. The schema names will be what you defined when creating the connections in Fivetran.

```
1   vars:
2     marketo_source:
3       marketo_database: # leave blank to use the default hive_metastore catalog
4       marketo_schema: marketing_marketo
5     salesforce_source:
6       salesforce_database: # leave blank to use the default hive_metastore catalog
7       salesforce_schema: marketing_salesforce
```

databricks

**Target schema for Fivetran staging models**

To avoid all the staging tables that are created by the Fivetran source models being created in the default target schema it can be useful to define a separate staging schema.

In the `dbt_project.yml` file add the staging schema name and this will then be suffixed to the default schema name.

```
1  models:
2    marketo_source:
3      +schema: your_staging_name # leave blank to use the default target_schema
4    salesforce_source:
5      +schema: your_staging_name # leave blank to use the default target_schema
```

Based on the above, if your target schema defined in `profiles.yml` is `mkt_analytics`, the schema used for `marketo_source` and `salesforce_source` tables will be `mkt_analytics_your_staging_name`.

**Disable missing tables**

At this stage you can run the Fivetran model packages to test that they work correctly.

```
1  dbt run -select marketo_source
```
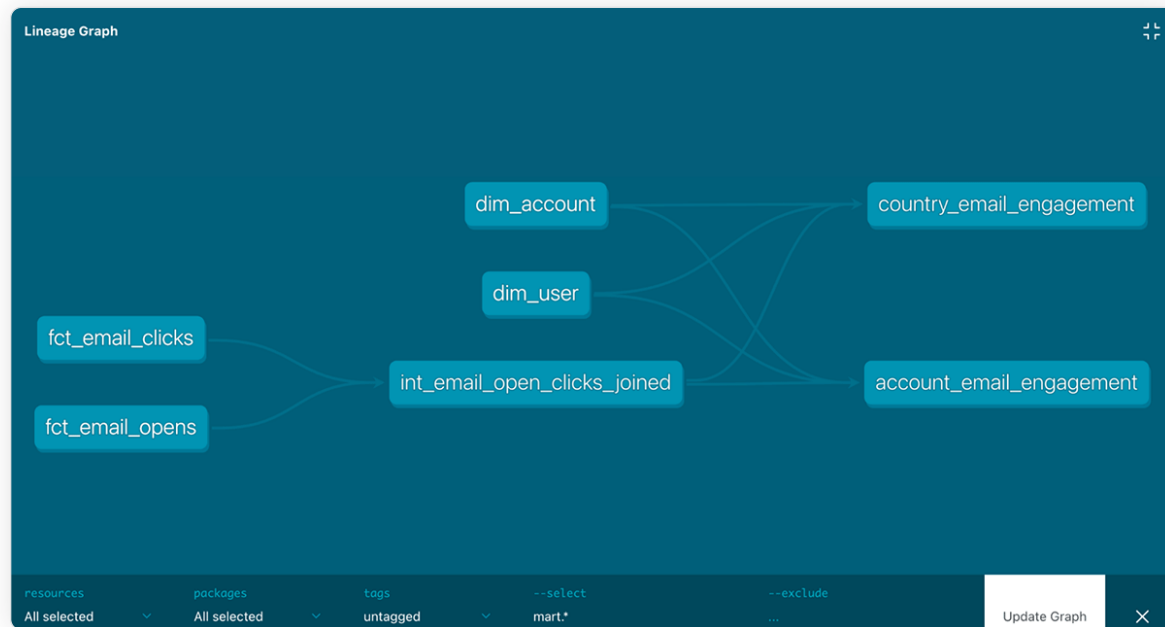
```
1  dbt run -select marketo_source
```

If any of the models fail due to missing tables, because you chose not to sync those tables in Fivetran, then in your source schema you can disable those models by updating the `dbt_project.yml` file.

For example if the email bounced and email template tables are missing from the Marketo source schema you can disable the models for those tables by adding the following under the models config:

```
1  models:
2    marketo_source:
3      +schema: your_staging_name
4      tmp:
5        stg_marketo__activity_email_bounced_tmp:
6          +enabled: false
7        stg_marketo__email_template_history_tmp:
8          +enabled: false
9      stg_marketo__activity_email_bounced:
10         +enabled: false
12     stg_marketo__email_template_history:
13         +enabled: false
```

## Developing the marketing analytics models



dbt lineage graph showing the star schema and aggregate tables data model

Now that the Fivetran packages have taken care of creating and testing the staging models you can begin to develop the data models for your marketing analytics use cases which will be a star schema data model along with materialized aggregate tables.

For example, for the first marketing analytics dashboard, you may want to see how engaged certain companies and sales regions are by the number of email campaigns they have opened and clicked.

To do so, you can join Salesforce and Marketo tables using the Salesforce user email, Salesforce `account_id` and Marketo `lead_id`.

The models will be structured under the mart folder in the following way.

```
1   marketing_analytics_demo
2   |-- dbt_project.yml
3   |-- packages.yml
4   |-- models
5       |-- mart
6           |-- core
7           |-- intermediate
8           |-- marketing_analytics
```

You can view the code for all the models on Github in the /models/mart directory and below describes what is in each folder along with an example.

**Core models**

The core models are the facts and dimensions tables that will be used by all downstream models to build upon.

The dbt SQL code for the `dim_user` model:

```
1   with salesforce_users as (
2       select
3           account_id,
4           email
5       from {{ ref('stg_salesforce__user') }}
6       where email is not null and account_id is not null
7   ),
8   marketo_users as (
9       select
10          lead_id,
12          email
13      from {{ ref('stg_marketo__lead') }}
14  ),
15  joined as (
16      select
17          lead_id,
18          account_id
19      from salesforce_users
20          left join marketo_users
21          on salesforce_users.email = marketo_users.email
22  )
23
24  select * from joined
```

You can also add documentation and tests for the models using a yaml file in the folder.

There are 2 simple tests in the `core.yml` file that have been added

```
 1   version: 2
 2
 3   models:
 4    - name: dim_account
 5      description: "The Account Dimension Table"
 6      columns:
 7        - name: account_id
 8          description: "Primary key"
 9          tests:
10            - not_null
12    - name: dim_user
13      description: "The User Dimension Table"
14      columns:
15        - name: lead_id
16          description: "Primary key"
17          tests:
18            - not_null
```

**Intermediate models**

Some of the final downstream models may rely on the same calculated metrics and so to avoid repeating SQL you can create intermediate models that can be reused.

The dbt SQL code for `int_email_open_clicks_joined` model:

```
 1   with opens as (
 2          select *
 3          from {{ ref('fct_email_opens') }}
 4   ), clicks as (
 5          select *
 6          from {{ ref('fct_email_clicks') }}
 7   ), opens_clicks_joined as (
 8
 9      select
10        o.lead_id as lead_id,
12        o.campaign_id as campaign_id,
13        o.email_send_id as email_send_id,
14        o.activity_timestamp as open_ts,
15        c.activity_timestamp as click_ts
16      from opens as o
17        left join clicks as c
18        on o.email_send_id = c.email_send_id
19        and o.lead_id = c.lead_id
20
21   )
22
23   select * from opens_clicks_joined
```

databricks

**Marketing Analytics models**

These are the final marketing analytics models that will be used to power the dashboards and reports used by marketing and sales teams.

The dbt SQL code for `country_email_engagement` model:

```sql
with accounts as (
        select
        account_id,
        billing_country
        from {{ ref('dim_account') }}
), users as (
        select
        lead_id,
        account_id
        from {{ ref('dim_user') }}
), opens_clicks_joined as (

    select * from {{ ref('int_email_open_clicks_joined') }}

), joined as (

        select *
        from users as u
        left join accounts as a
        on u.account_id = a.account_id
        left join opens_clicks_joined as oc
        on u.lead_id = oc.lead_id

)

select
        billing_country as country,
        count(open_ts) as opens,
        count(click_ts) as clicks,
        count(click_ts) / count(open_ts) as click_ratio
from joined
group by country
```

## Run and test dbt models

Now that your model is ready you can run all the models using
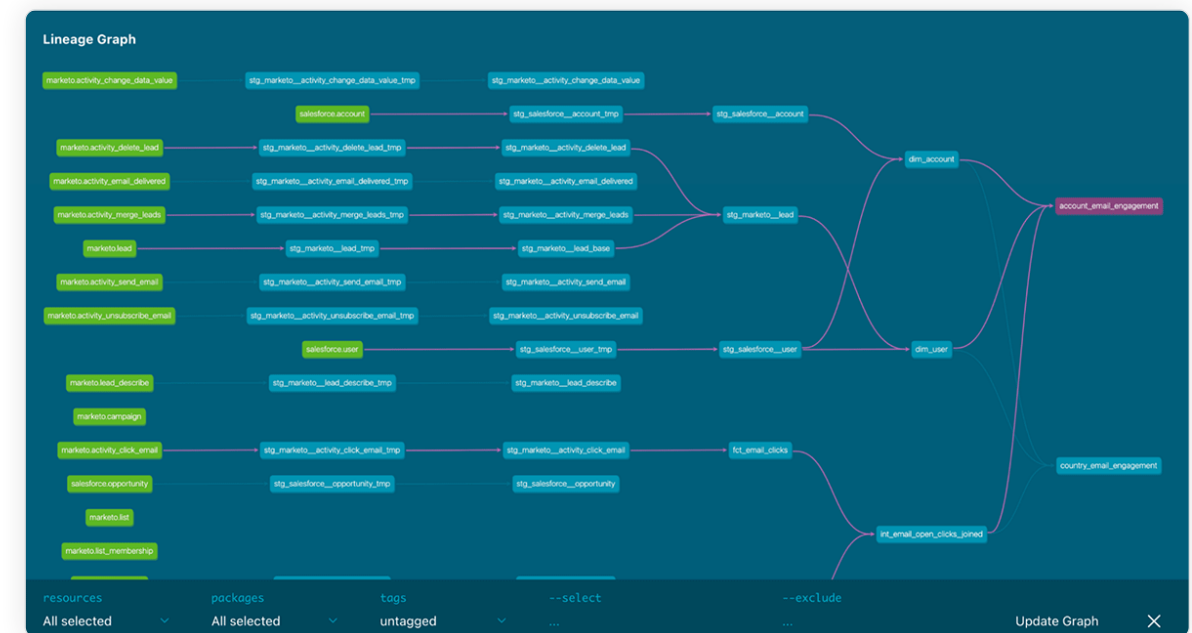
```
dbt run
```

And then run the tests using

```
dbt test
```

## View the dbt docs and lineage graph



dbt lineage graph for the marketing analytics model

Once your models have run successfully you can generate the docs and lineage graph using

```
$ dbt docs generate
```
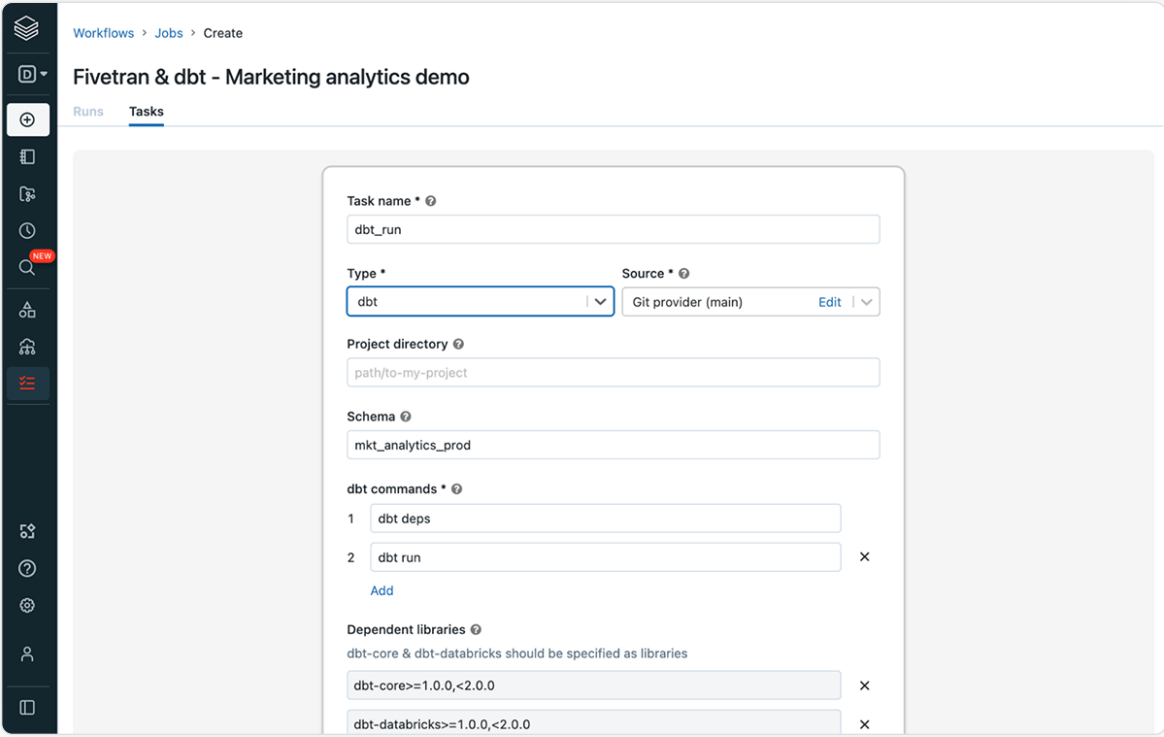
To then view them locally run

```
$ dbt docs serve
```

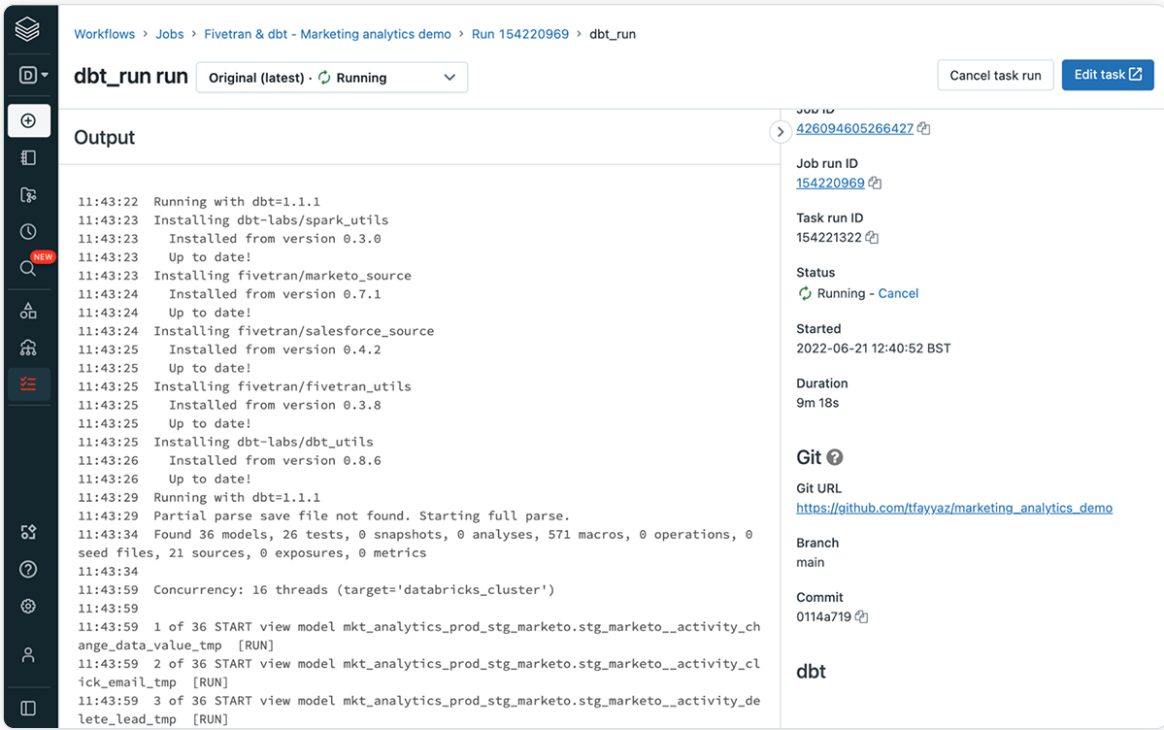databricks

# Deploying dbt models to production

Once you have developed and tested your dbt model locally you have multiple options for deploying into production one of which is the new dbt task type in Databricks Workflows (private preview).

Your dbt project should be managed and version controlled in a Git repository. You can create a dbt task in your Databricks Workflows job pointing to the Git repository.



Using a dbt task type in Databricks Workflows to orchestrate dbt

As you are using packages in your dbt project the first command should be dbt deps followed by dbt run for the first task and then dbt test for the next task.



Viewing the dbt logs for each dbt run

For each run you can see the logs for each dbt command helping you debug and fix any issues.

## Powering your marketing analytics with Fivetran and dbt

As shown here, using Fivetran and dbt alongside the Databricks Data Intelligence Platform allows you to easily build a powerful marketing analytics solution that is easy to set up, manage and flexible enough to suit any of your data modeling requirements.

To get started with building your own solution, visit the documentation for integrating Fivetran and dbt with Databricks and re-use the marketing_ analytics_demo project example to quickly get started.

The dbt task type in Databricks Workflows is in private preview. To try the dbt task type, please reach out to your Databricks account executive.

databricks

SECTION 5.2

# Claims Automation on Databricks

Smart claims increases efficiency by automating all aspects of claims processing from ingestion, analysis, and decision-making.

by Anindita Mahapatra and Marzi Rasooli

According to the latest reports from global consultancy EY, the future of insurance will become increasingly **data-driven**, and **analytics enabled**. The recent focus on the cloud has improved access to advanced technological infrastructure, but most organizations still need help implementing and leveraging these capabilities. It's time to shift the focus on operationalizing services to realize value.
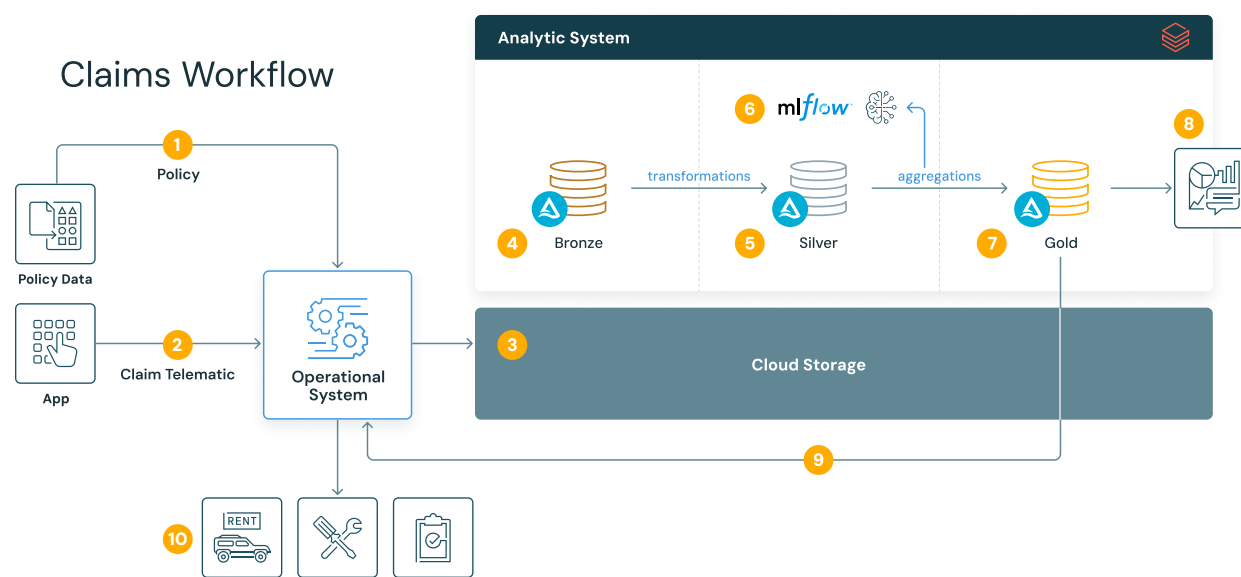
In today's economic circumstances, insurance companies face an ever-increasing number of challenges. Insurers are being forced to leverage data to their advantage and **innovate** at an accelerated pace. For personal P&C insurers, this means an increased focus on **personalization** and customer retention. Brand loyalty is at an all-time low, with customers continuously shopping for more competitive rates and better overall experiences, which increases the risk of churn. An increase in fraudulent claims further erodes profit margins. Insurers need to find additional ways to reduce costs and better manage risks.

Automating and optimizing the claims-handling process is one area that can significantly reduce costs through time saved and lesser reliance on human capital. Furthermore, effectively leveraging insights from data and advanced analytics can substantially reduce the overall exposure to risk.

The motivation behind the "Smart Claims" solution accelerator is simple — improve the claims handling process to enable faster settlement, lower processing costs, and deliver quicker insights about potentially fraudulent claims, with the lakehouse. Implementing the lakehouse paradigm simplifies the current architectural landscape and sets the scene for future expansion across the organization. The accompanying assets can be found here.

databricks

# Reference Architecture and Workflow

A typical claims workflow involves some level of orchestration between operational systems such as Guidewire and analytical systems like Databricks. The diagram below shows an example of such a workflow for an automotive insurer.



Smart Claims Reference Architecture and Workflow

Automating and optimizing the claims handling process requires a deep understanding of customer interaction with operational systems and the various sources of information available for analysis.

In this example, we assume that customers primarily interact through a mobile application, and from there, they can submit claims and monitor the status of existing cases. This touch point offers vital information about customer behavior. Another important source of information is IoT devices installed in customer vehicles. Telematics data can be streamed to the operational and analytical systems, providing valuable insights into customer-driving behavior and patterns. Other external data sources may include weather and road conditions data that supplement the traditional data categories such as vehicle characteristics (make, model, year), driver characteristics and exposure/coverage (limits, deductibles)

Access to additional data sources can become increasingly important, especially in the absence of data from traditional sources such as credit bureaus. Credit scores from bureaus usually form the basis for risk modeling, assessing the exposure for drivers, which ultimately impacts their premiums. On the other hand, data from mobile applications and IoT devices provide a more personalized view of customer behavior, which could be used to create a more accurate indicator of the risk associated with a given party. This alternative, **behavioral-based** approach to risk modeling and pricing is essential for delivering a hyper-personalized customer experience.

The lakehouse architecture powered by the Databricks Data Intelligence Platform is the only platform that combines all the required features and services to support a future-proof claims-handling process. From streaming to machine learning and reporting, Databricks offers the best platform to build an end-to-end solution for the insurance industry of tomorrow.

The following steps capture the overall flow:

- Policy data is ingested.
- Telematics data is continuously ingested from IoT sensors. A claimant submits claims data via a mobile app.
- All the operational data is ingested into cloud storage.
- This is incrementally loaded as 'Raw data' into Delta Bronze tables
- The data is wrangled and refined via various data transformations
- Data is scored using the trained model
- The predictions are loaded to a Gold table
- The Claims Dashboard is refreshed for visualization
- The resulting insights are fed back to the operational system. This provides the feedback loop of pulling data from Guidewire and passing 'Next Best Action' back to Guidewire in real time to know which claims should be prioritized.
- The Claims Decisioning workflows use these generated insights to route the case appropriately. (Eg. approve repair expenses, rental reimbursement, or alert authorities)

## How the lakehouse paradigm aids Smart Claims

The lakehouse architecture enables all data personas (data engineers, data scientists, analytic engineers, and BI analysts) to work collaboratively on a single platform. Supporting all big-data workloads and paradigms (e.g., batch processing, streaming, DataOps, ML, MLOps, and BI) in a single, collaborative platform greatly simplifies the overall architecture, improves stability, and reduces cost significantly.
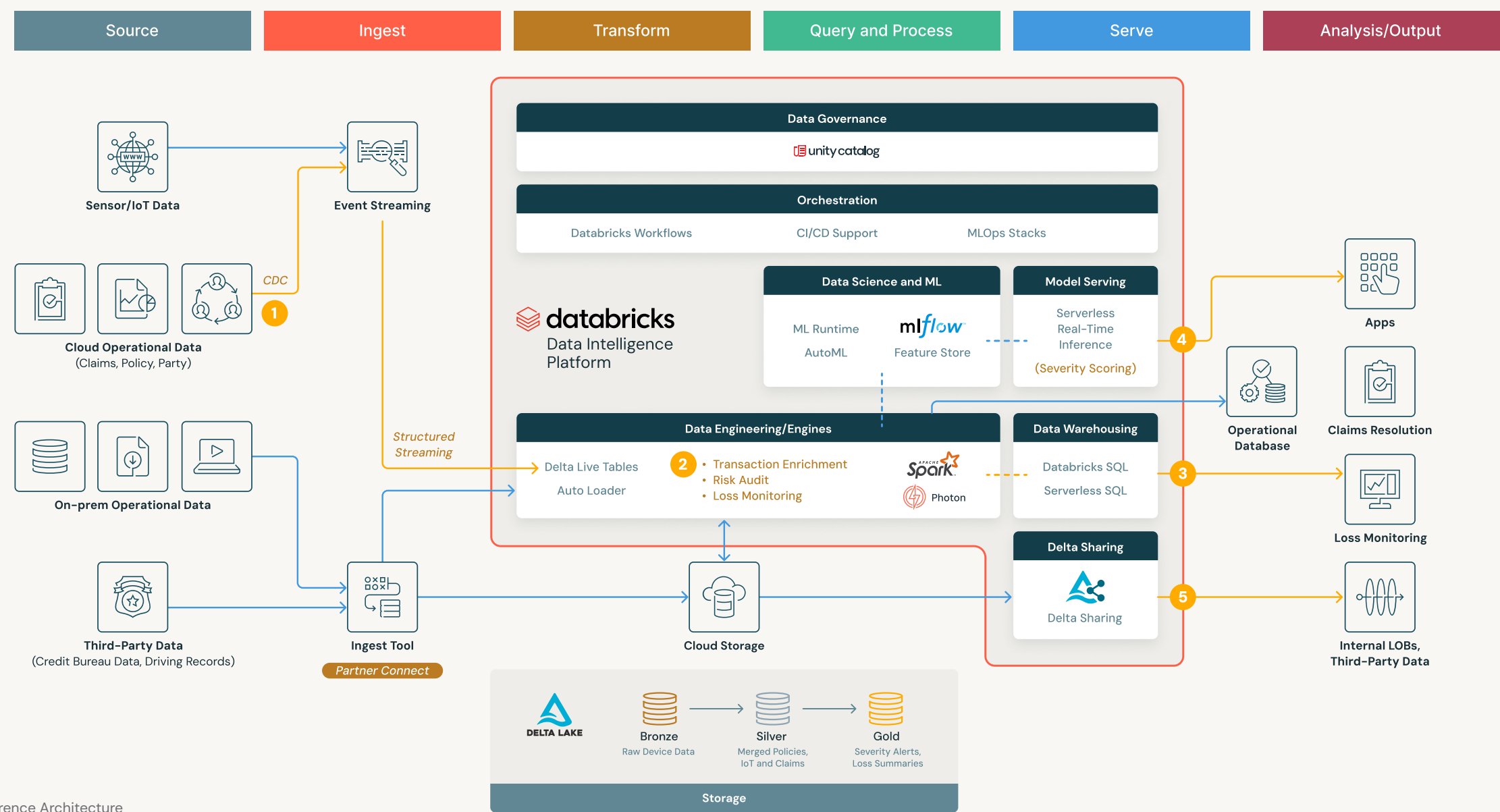
Databricks Delta Live Tables (DLT) pipelines offer a simple, declarative framework to develop and implement workloads quickly. It also provides native support for data quality management with granular constraints to guarantee the integrity of outputs.

ML and AI workloads can easily be created and managed with MLflow for reproducibility and auditability. MLFlow simplifies the entire model lifecycle, from experimenting through model deployment, serving, and archiving. ML can be run on all types of data including unstructured data beyond text (images, audio, video, etc). In this solution, we will use computer vision capabilities to assess damage to the vehicle.

Finally, Databricks SQL provides a fast and efficient engine to query curated and aggregated data. These insights can then be packaged and served through interactive Dashboards within minutes.

Unity Catalog provides a multi-cloud, centralized governance solution for all data and AI assets including files, tables, machine learning models and dashboards with built-in search, discovery, automated workload lineage.
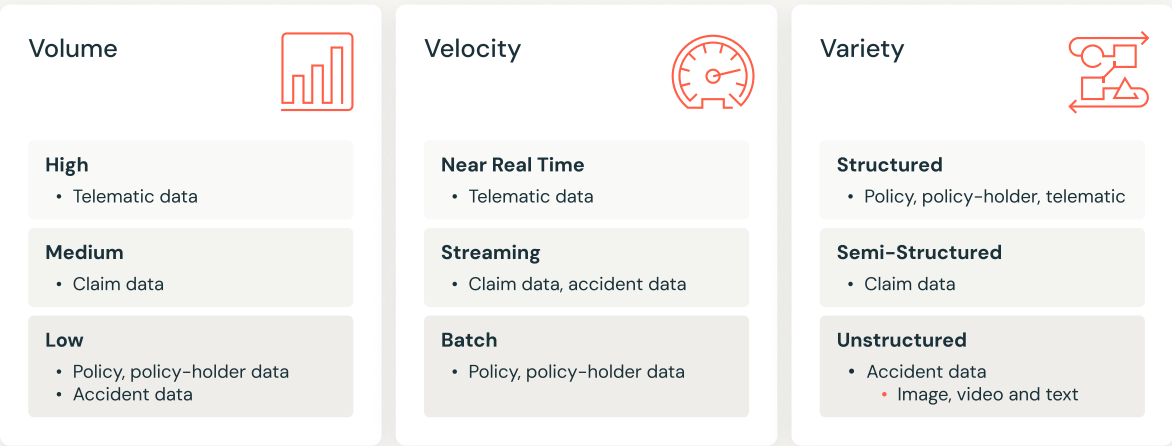
databricks

The diagram below shows a reference architecture for the lakehouse in the context of typical insurance use cases:



Insurance Reference Architecture

## Data Ingestion using DLT and Multitask Workflows

Automating the claims-handling process starts with optimizing the ingestion and data engineering workflow. The figure below offers a summary of the typical data sources encountered including structured, semi-structured and unstructured. Some sources are slower-moving, while others update more rapidly. Additionally, some sources might be additive, requiring appending, while others offer incremental updates and must be treated as slow-changing dimensions.

**Volume**

**High**
- Telematic data

**Medium**
- Claim data

**Low**
- Policy, policy-holder data
- Accident data

**Velocity**

**Near Real Time**
- Telematic data

**Streaming**
- Claim data, accident data

**Batch**
- Policy, policy-holder data

**Variety**

**Structured**
- Policy, policy-holder, telematic

**Semi-Structured**
- Claim data

**Unstructured**
- Accident data
  - Image, video and text

Sample Datasets used in Claims Processing

DLT can simplify and operationalize the data processing pipeline. The framework offers support for Auto Loader to facilitate ingestion from streaming sources, efficient auto-scaling to handle sudden changes in data volumes, and resiliency via a restart of task failure.
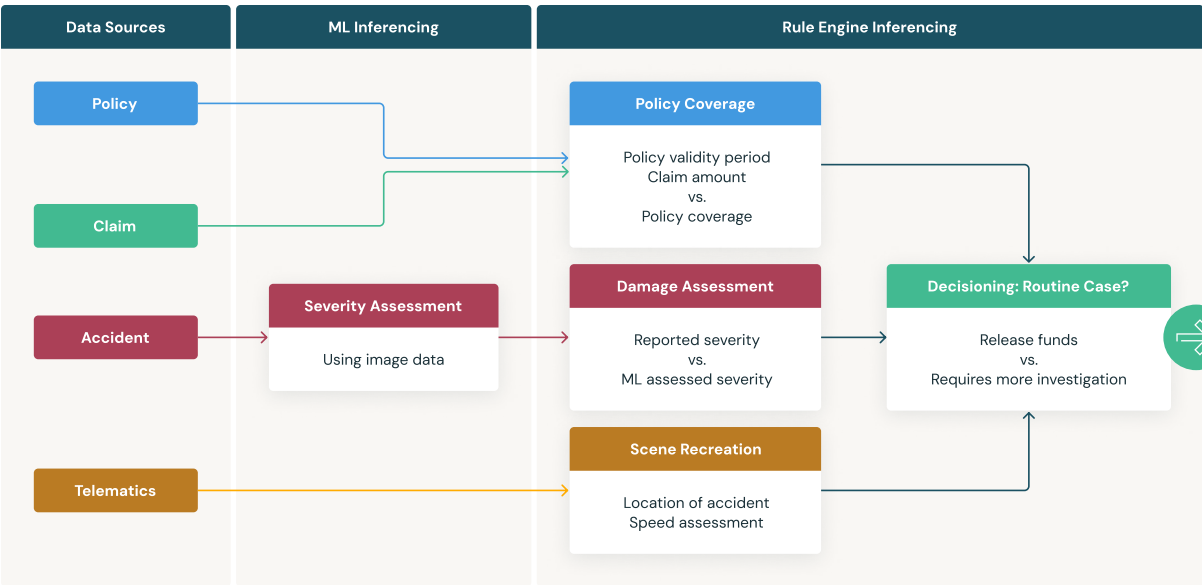
Databricks Workflows can accommodate multiple tasks and workloads (e.g., notebooks, DLT, ML, SQL). Workflows support repair-and-run and compute sharing across tasks – enabling robust, scalable, cost-effective workloads. Additionally, Workflow can easily be automated through schedules or programmatic invoking via REST APIs.

## Insight Generation using ML and Dynamic Rules Engine

Leveraging ML is essential to uncovering previously unknown patterns, highlighting new insights, and flagging suspicious activity. However, combining ML and traditional rules-based approaches can be even more powerful.

Within the claims-handling process, ML can be used for several use cases. One example would be using computer vision and ML for assessing and scoring images submitted with vehicle insurance claims. Models can be trained to focus on the validity and severity of damages. Here, MLFlow can be crucial in simplifying the model training and serving process with its end-to-end MLOps capabilities. MLFlow offers a serverless model serving through REST APIs. Trained models can be operationalized and put into production with the click of a button.

databricks

On the other hand, rules engines offer flexible ways of defining known operational characteristics and statistical checks, which can be automated and applied without requiring human interaction. Flags are raised whenever data does not comply with preset expectations and are sent for human review and investigation. Incorporating such an approach with ML-based workflows offers additional oversight and significantly reduces the time claims investigators require to dissect and review flagged cases.
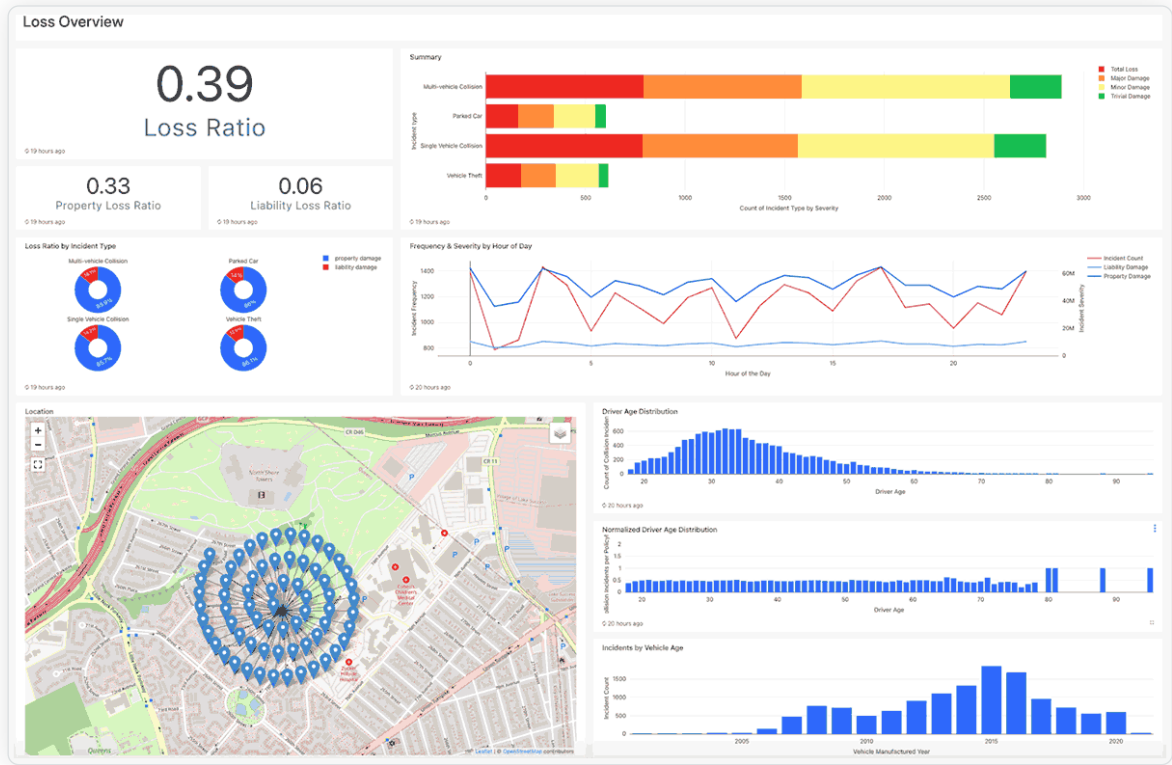


ML and Rule Engine Inferencing

## Insight visualization using Dashboards

In this example, we created two dashboards to capture critical business insights. The dashboards include the following:
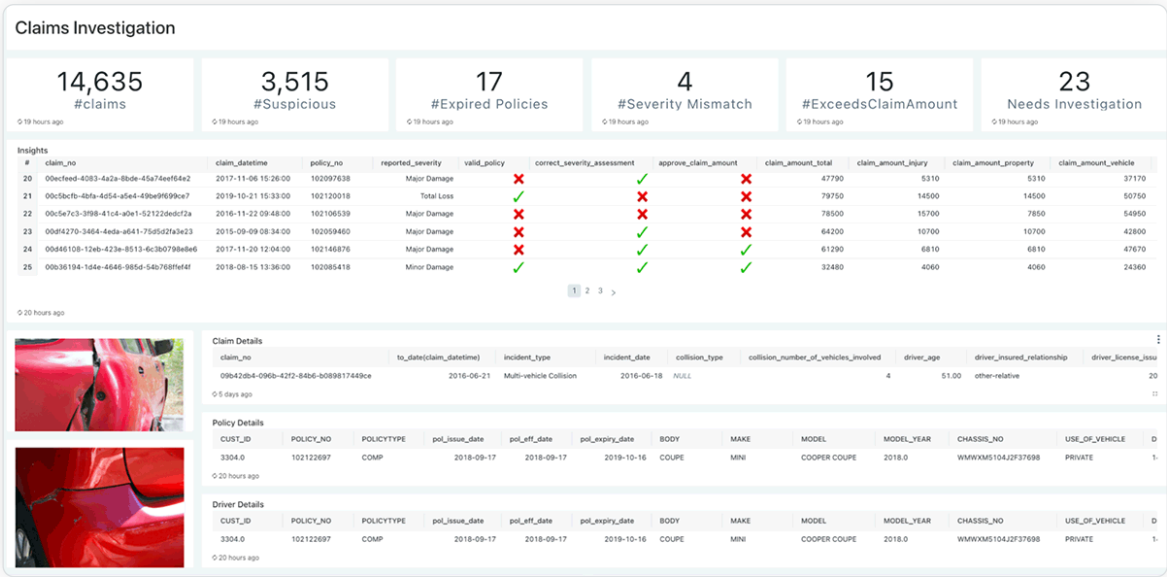
- A **Loss Summary** dashboard for a high-level view of the overall business operations; and

- A **Claims Investigation** dashboard with a granular view of claims details to understand the specifics of a given case.



databricks

Analyzing recent trends can further aid in reviewing similar cases such as :

- Loss Ratio is computed by insurance claims paid plus adjustment expenses divided by total earned premiums. E.g. typical average Loss Ratio (all coverages combined, Bodily Injury, and Physical Damage) for personal auto should be around 65%

- Summary visualization captures count of incident type by damage severity

- Trend lines over various features/dimensions

- Geographic distribution of policies

The Claims Investigation dashboard facilitates faster investigation by providing all relevant information around a claim allowing the human investigator to drill down to a specific claim to see details such as images of the damaged vehicle, Claim, Policy and Driver details, Telematic data draws the path taken by the vehicle, Reported data is contrasted with assessed data insights.



Claim Investigation Dashboard

Provides recent claims that are auto-scored in the pipeline using ML inferencing and rule engine

- A green tick is used to denote auto-assessment matches claims description

- A red cross indicates a mismatch that warrants further manual investigation

databricks

## Summary

Innovation and personalization are essential for insurance firms to differentiate themselves from the competition. Databricks provides a data intelligence platform for insurers to enable and accelerate innovation with an open, secure, extensible architecture that easily integrates with third–party tools and services. This Solution Accelerator demonstrates how the paradigm can be applied to claims handling. Further, the Databricks ecosystem offers a range of capabilities to enable data teams and business stakeholders to collaborate and generate and share insights that support business decisions and drive tangible value to the bottom line.

The technical assets, including pipeline configurations, models, and sample data used in this example, can be accessed here or directly on Git.

databricks

# Design Patterns for Batch Processing in Financial Services

Laying the foundation for automating workflows

by Eon Retief

## Introduction

Financial services institutions (FSIs) around the world are facing unprecedented challenges ranging from market volatility and political uncertainty to changing legislation and regulations. Businesses are forced to accelerate digital transformation programs; automating critical processes to reduce operating costs and improve response times. However, with data typically scattered across multiple systems, accessing the information required to execute on these initiatives tends to be easier said than done.

Architecting an ecosystem of services able to support the plethora of data-driven use cases in this digitally transformed business can, however, seem to be an impossible task. This blog will focus on one crucial aspect of the modern data stack: batch processing. A seemingly outdated paradigm, we'll see why batch processing remains a vital and highly viable component of the data architecture. And we'll see how Databricks can help FSIs navigate some of the crucial challenges faced when building infrastructure to support these scheduled or periodic workflows.

## Why batch ingestion matters

Over the last two decades, the global shift towards an instant society has forced organizations to rethink the operating and engagement model. A digital-first strategy is no longer optional but vital for survival. Customer needs and demands are changing and evolving faster than ever. This desire for instant gratification is driving an increased focus on building capabilities that support real-time processing and decisioning. One might ask whether batch processing is still relevant in this new dynamic world.

While real-time systems and streaming services can help FSIs remain agile in addressing the volatile market conditions at the edge, they do not typically meet the requirements of back-office functions. Most business decisions are not reactive but rather, require considered, strategic reasoning. By definition, this approach requires a systematic review of aggregate data collected over a period of time. Batch processing in this context still provides the most efficient and cost-effective method for processing large, aggregate volumes of data. Additionally, batch processing can be done offline, reducing operating costs and providing greater control over the end-to-end process.

The world of finance is changing, but across the board incumbents and startups continue to rely heavily on batch processing to power core business functions. Whether for reporting and risk management or anomaly detection and surveillance, FSIs require batch processing to reduce human error, increase the speed of delivery, and reduce operating costs.

databricks