

eBook

# Low-Code Apache Spark™ and Delta Lake

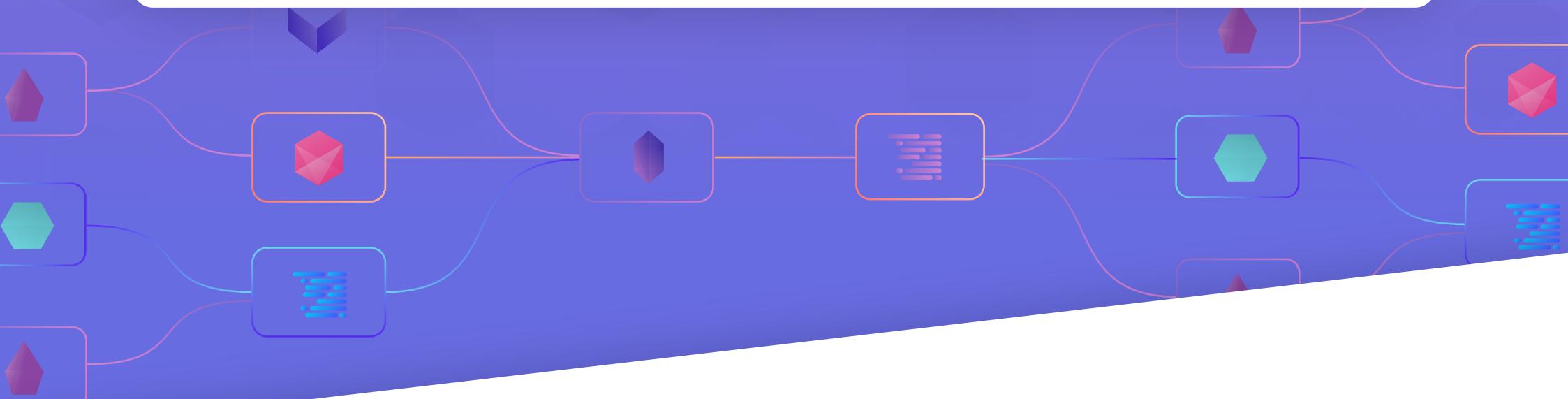
A Low-Code Lakehouse Guide



Prophecy



DELTA LAKE



Authors: Abhinav Jain, Scott Lyons, Maciej Szpakowski

## Data Lakehouse

Data Lakehouse is a relatively new term in the Big Data space that, like many things in the Data ecosystem, is a combination of two existing concepts: Data Warehouses and Data Lakes.

Data Warehouses (DWs) store large volumes of processed and highly-structured data while enabling efficient and performant querying on top of said data. Most DW implementations are based on [MPP architectures](#) and allow users to execute common SQL operations for standard databases such as inserts, updates, and even merges. Some are even [ACID compliant](#). However, with the ever-increasing quantities of data and the rise of unstructured data, warehouses become obsolete, expensive, and difficult to maintain.

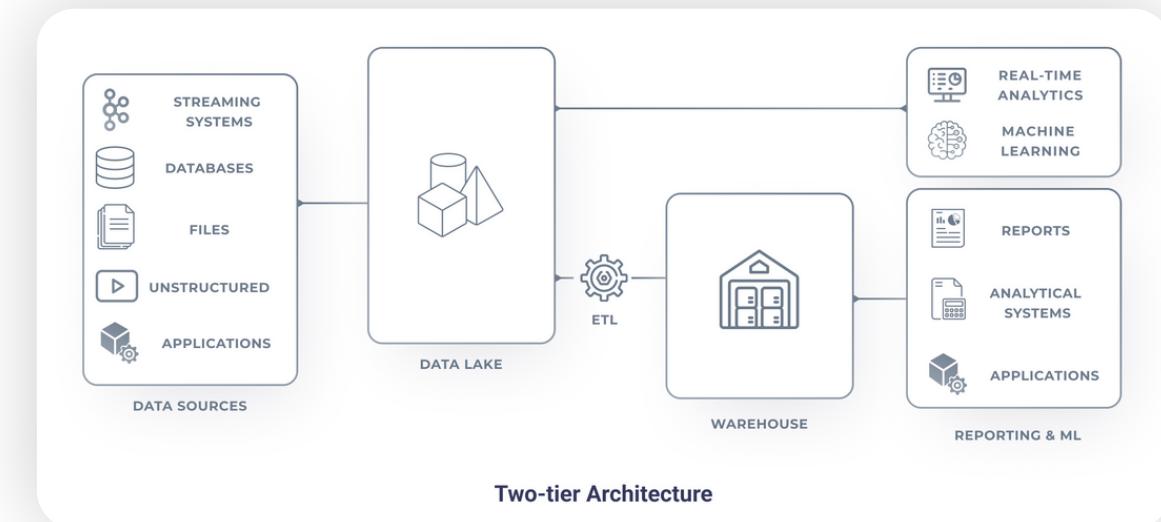
Data Lakes (DLs), while sharing similar underlying concepts, such as high-volume storage and efficient querying, generally tend to be more unstructured by nature. They're frequently used as the pre-processing area by data scientists and engineers to prepare data sourced from a variety of upstream sources for loading into a Data Warehouse.

Data Lakehouse aims to have the best of both worlds. DLs support both Batch and Streaming paradigms and any number of diverse workloads on top of them.

Apache Spark is the ETL/Computation engine of choice that has enabled many companies to perform efficient and scalable data engineering on the cloud without breaking the bank.

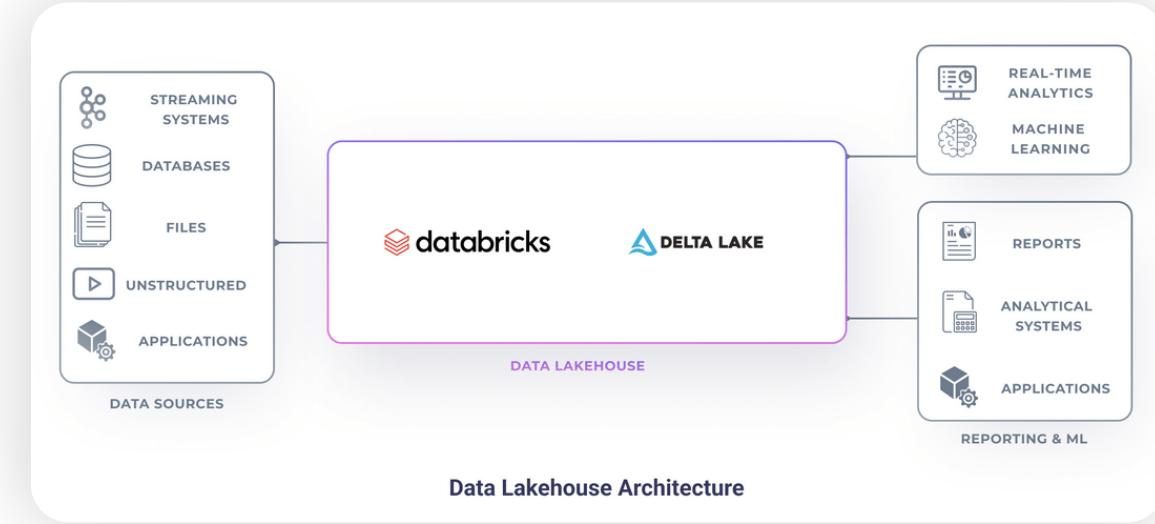
## The Architecture

Classically, the need for cost-effective storage, performant processing, and low-latency querying birthed a two-tier architecture. Data Lake for raw storage of unstructured data and Data Warehouse on top of it for high-performance reporting and querying. To integrate those layers ETL batch processes are used.



Even though this is one of the most commonly used architectures today, its use can result in a number of problems, such as:

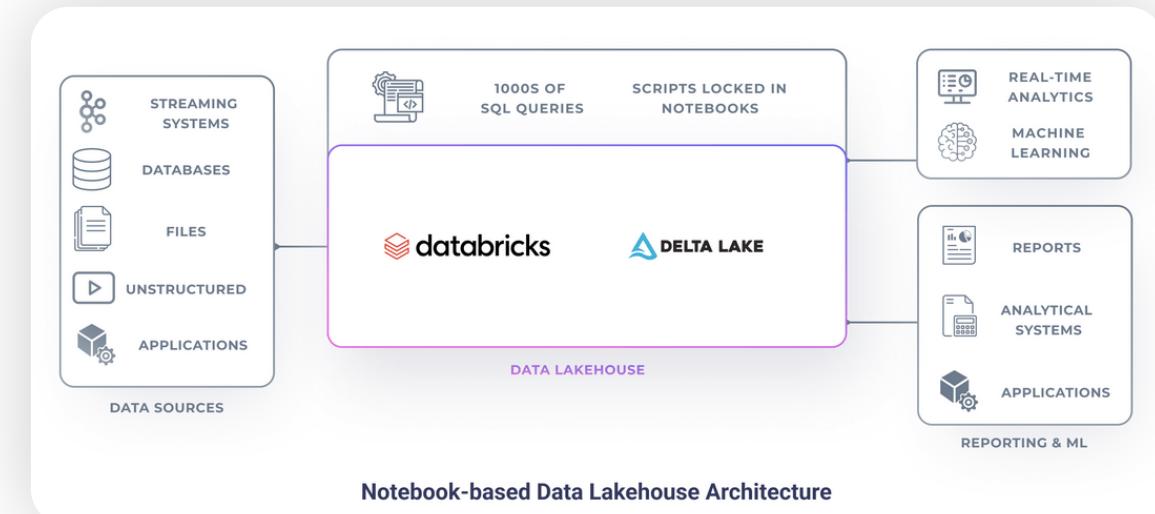
- Increased number of bugs due to a larger number of pipelines
- Stale data stored within the data warehouse
- Lack of ability to store unstructured information - SQL is often not enough.



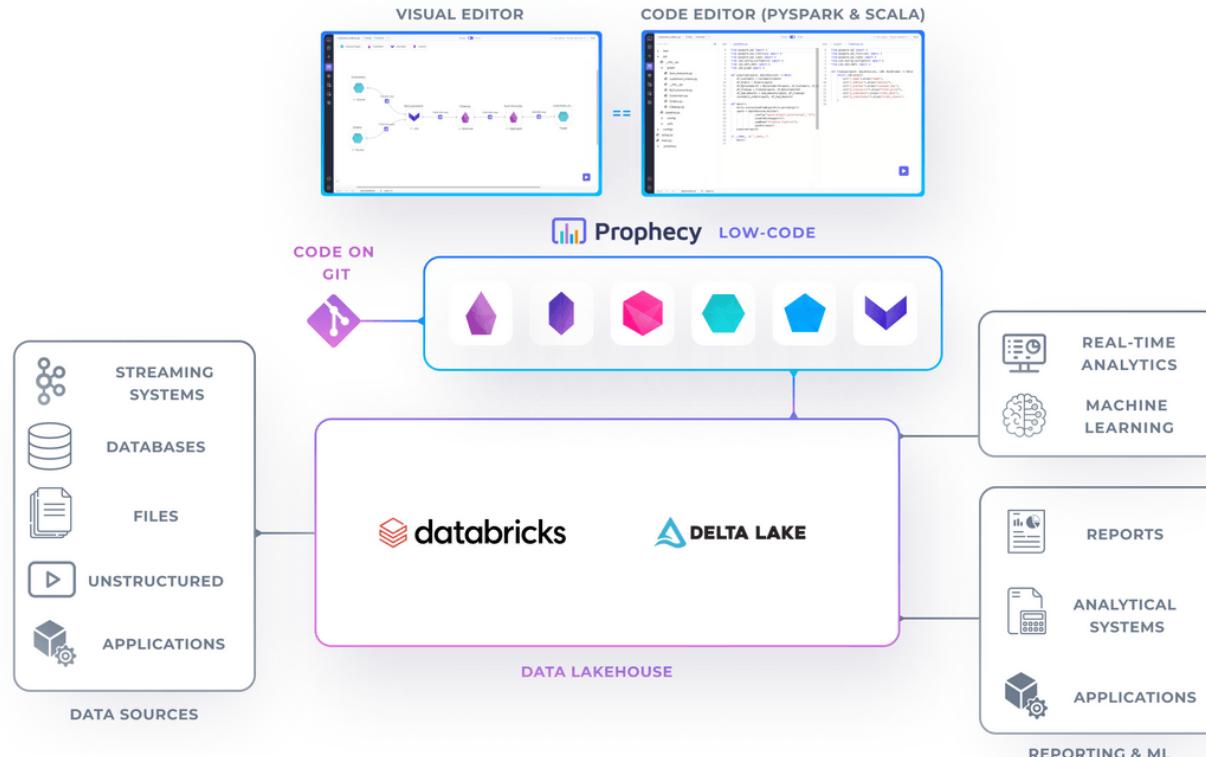
The introduction of a Lakehouse architecture, using [Delta](#) as the underlying storage format and Spark as the querying engine, attempts to solve the two-tiered architecture by unifying it into a single layer. Delta enables the simplicity of a data warehouse with advanced update operations and ACID guarantees while supporting the large majority of Cloud vendor storage offerings. This makes data more accessible, cheaper to query, and reduces the number of buggy data pipelines.

## Making Data Lakehouse even easier

Even though Spark and Delta are the perfect basis for your future data pipeline, and the Data Lakehouse architecture allows us to enable a tremendous amount of new use cases, the usability and productivity still remain a challenge. Delta is yet another technology that your team has to learn, requires you to know how to write code in an optimal way, and common use cases like performing slowly changing dimensions require quite a few lines of repeatable code.



Along comes Prophecy. Gone are the chaotic collection of disconnected ETL notebooks and haphazard orchestration scripts, replaced with a visually built data flow that is turned into well-engineered and efficient Spark or PySpark code stored directly into Git. You can even start leveraging all the best practices that have been classically available only for traditional software engineers, including Unit Testing and CI/CD. Once you're confident that the pipeline will produce the Data you're looking for, you can then use Prophecy to schedule the jobs.

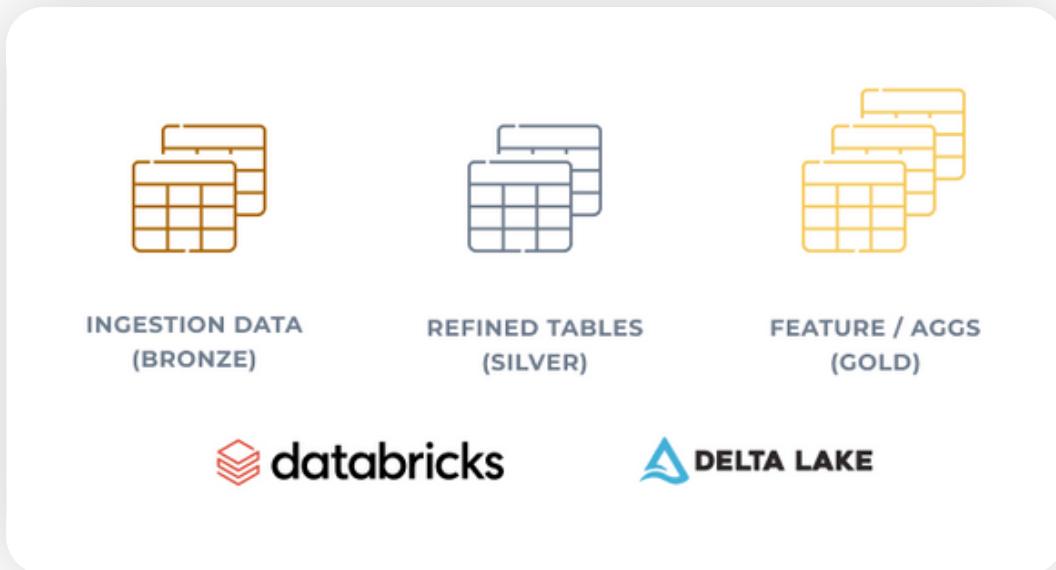




## Bronze, Silver, and Gold layers

With the Data Lakehouse architecture, all your data flows from the operational systems directly to your [object store](#). However, this often results in poor data management and governance, which leads to the creation of a [data swamp](#). To prevent this it's paramount to organize your Lakehouse well. Proper governance, lineage, and taxonomies are also essential (but we will leave that for a future blog).

A common architecture is based on the organization pattern that uses tables that correspond to different quality levels of data:



- **Bronze tables** - raw data directly available after data ingestion, often used as historical store data and a single source of truth
- **Silver tables** - clean transformed data (e.g. clean training data for machine learning - the feature store)
- **Gold tables** - aggregated datasets, ready for consumption by downstream customers (e.g. reports or model predictions)

Building such architecture allows anyone to directly leverage data at different levels, depending on their needs. However, building it often requires multiple data pipelines at each stage that have to be often collaboratively worked upon. Using a low-code tool can make this job significantly easier.

## Pipeline example

To demonstrate how this architecture can be built, we're going to use two datasets:

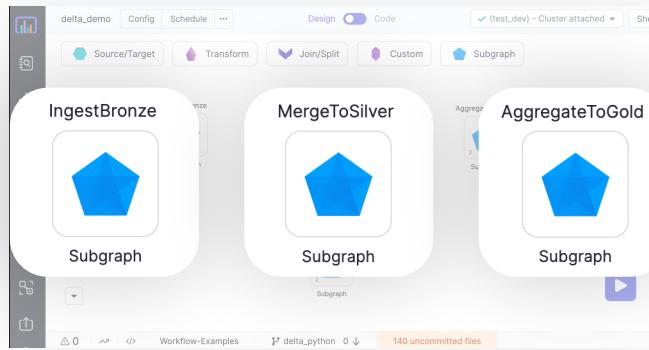
**1. Orders** - Table that contains a new order for every row. This table is queried directly from an operational system, so every time we query this table it might contain different information for the order\_id. As an example, once the order is delivered the order\_delivered\_dt changes from null to a valid date.

order_id	customer_id	order_dt	order_amount	order_status	order_approved_dt	order_delivered_dt
1	A	2022-01-12	500	delivered	2022-01-12	2022-01-14
2	A	2022-01-14	1000	approved	2022-01-14	null
3	B	2022-01-12	200	delivered	2022-01-12	2022-01-14
4	B	2022-01-14	200	approved	2022-01-14	null
5	C	2022-01-14	500	approved	2022-01-14	null
.						
.						
.						
.						

**2. Customers** - Table that contains a customer per row, with only occasional inserts and updates (e.g. when a customer changes their address).

customer_id	customer_zip_code	customer_city	customer_state	updated_dt
A	10000	City1	State1	2022-01-01
B	10010	City2	State1	2022-01-02
C	20000	City3	State2	2022-01-11
D	20010	City4	State2	2022-01-12

Our reporting pipeline is going to consist of three separate pieces:



#### a. IngestBronze - ingestion of source data

The screenshot shows the configuration of the 'IngestBronze' subgraph. It is connected to a 'raw\_orders\_bronze' dataset. The dataset configuration window is open, showing the 'PROPERTIES' tab. The 'Dataset' dropdown is set to 'delta\_demo\_raw\_orders'. The 'Schema' section displays the following schema:

Name	Type	Option...
order_id	string	<input checked="" type="checkbox"/>
customer_id	string	<input checked="" type="checkbox"/>
order_dt	date	<input checked="" type="checkbox"/>
order_amount	string	<input checked="" type="checkbox"/>
order_status	string	<input checked="" type="checkbox"/>
order_approved_dt	date	<input checked="" type="checkbox"/>
order_delivered_dt	date	<input checked="" type="checkbox"/>

Below the schema, there are fields for 'Write Mode' (set to 'overwrite') and 'Column delimiter' (set to ','), and a checkbox for 'First row is header' which is checked.

b. **MergeToSilver** - cleansing the source data and creation of delta tables

The screenshot displays a data pipeline interface with two main sections:

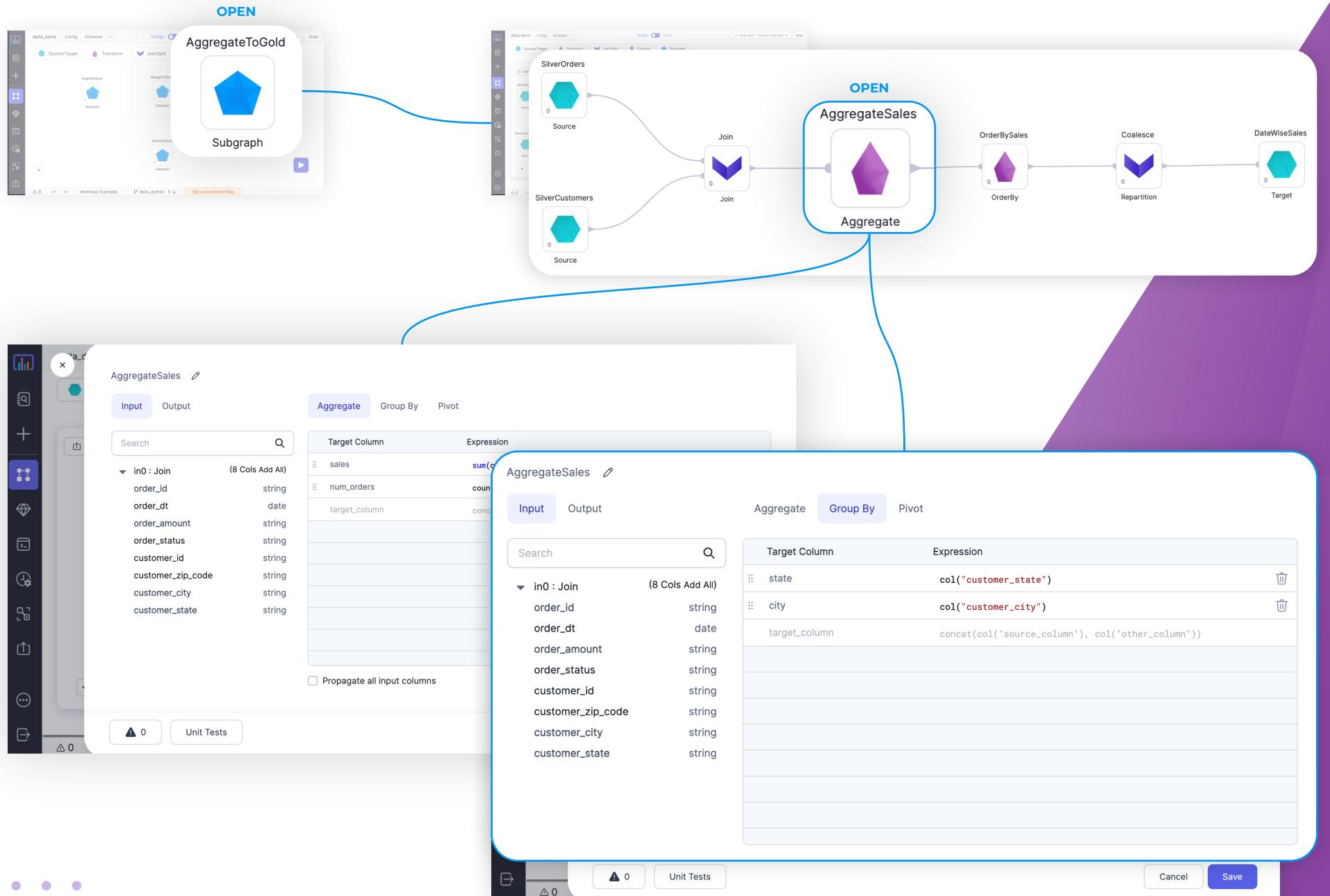
- Workflow Overview:** On the left, a complex workflow is shown with various nodes and connections. A specific node, "MergeToSilver", is highlighted with a blue box and labeled "OPEN". This node has a "Subgraph" icon below it. Other nodes include "bronze\_orders" (Source), "orders\_delta" (Target), "bronze\_customers" (Source), "Clean\_scd2" (Reformat), "customers\_scd2..." (Target), "Clean\_scd3" (Reformat), and "customers\_scd3..." (Target). Arrows indicate data flow between these nodes.
- Configuration Details:** On the right, a detailed configuration panel for the "orders\_delta" dataset is shown. The "Dataset" dropdown is set to "delta\_demo\_orders\_merge". The "PROPERTIES" tab is selected, showing the "Write Mode" as "merge". The "Merge Condition" section contains the following code:
 

```
1 col("source.order_id")==col("target.order_id")
```

 A callout box highlights this condition with the text: "Merge condition: checks if you need to merge this row".
- Data Samples:** At the bottom right, a table titled "Data samples for bronze\_orders (out)" shows the first 5 rows of data from the "bronze\_orders" source. The columns are: order\_id, customer\_id, order\_dt, order\_amount, order\_status, and order\_approved\_dt. The data is as follows:
 

order_id	customer_id	order_dt	order_amount	order_status	order_approved_dt
2	A	2022-01-14	1000	delivered	2022-01-14
4	B	2022-01-14	200	delivered	2022-01-14
6	D	2022-01-16	700	approved	2022-01-16
7	A	2022-01-16	400	approved	2022-01-16
8	C	2022-01-16	800	approved	2022-01-16

### c. AggregateToGold - the creation of a business report



After every change the code for our entire pipeline is automatically generated in the backend in PySpark.

In just a few minutes we've built our report while following the best practices in terms of structuring the code. Now it's time to schedule it.

The screenshot shows the Databricks code editor with the file `pipeline.py` open. The code is a PySpark pipeline for processing customer data. It includes imports for `pyspark.sql`, `pyspark.sql.functions`, `pyspark.sql.types`, `job.config.ConfigStore`, `job.udfs.UDFs`, `prophecy.utils`, and `job.graph`. The code defines two main functions: `pipeline` and `main`. The `pipeline` function calls `A_IngestBronze`, `B_MergeToSilver`, and `C_AggregateToGold` subfunctions. The `main` function creates a SparkSession builder, configures parallelism, and initializes a pipeline. It then sets the prophecy metadata pipeline URI, starts the MetricsCollector, and runs the pipeline. Finally, it checks if the script is run directly and calls the `main` function. The code ends with a `delta` save operation.

```
from pyspark.sql import *
from pyspark.sql.functions import *
from pyspark.sql.types import *
from job.config.ConfigStore import *
from job.udfs.UDFs import *
from prophecy.utils import *
from job.graph import *

def pipeline(spark: SparkSession) -> None:
    A_IngestBronze(spark)
    B_MergeToSilver(spark)
    C_AggregateToGold(spark)
    ViewOutputs(spark)

def main():
    spark = SparkSession.builder\
        .config("spark.default.parallelism", "4")\
        .config("spark.sql.legacy.allowUntypedFunctions", "true")\
        .enableHiveSupport()\
        .appName("Prophecy Pipeline")\
        .getOrCreate()\
        .newSession()
    Utils.initializeFromArgs(spark, parse_args())
    spark.conf.set("prophecy.metadata.pipeline.uri", "dbfs:/FileStore/data_eng/delta_demo")
    MetricsCollector.start(spark = spark, pipelineId = 1)
    pipeline(spark)
    MetricsCollector.end(spark)

if __name__ == "__main__":
    main()

if __name__ == "__main__":
    delta.write\
        .format("delta")\
        .option("overwriteSchema", True)\
        .mode("overwrite")\
        .save("dbfs:/FileStore/data_eng/delta_demo")
```

The screenshot shows the Databricks pipeline builder interface. On the left, a pipeline graph is visible with nodes like `IngestBronze`, `MergeToSilver`, `AggregateToGold`, and `ViewOutputs`. A blue box highlights the `Schedule` button on the top navigation bar. A large blue callout box points from this button to a modal dialog titled "Schedule this Pipeline". The dialog contains instructions: "This screen helps you schedule this pipeline. You can edit and extend the job in the job builder." It shows an "Existing Jobs" table with one entry: `delta_deploy` (Enabled, Schedule: 0 0 6 \* \* ?). On the right, there is a "Create New" form with fields for Name (`delta_deploy_daily`), Fabric / Scheduler (`test_dev / Databricks`), Cluster (`small`), Schedule Interval (`0 0 0 1/1 * ?`), Alerts on the full job (checkbox), Alerts Email (`ops_team@acme.com`), and When should we send alerts? checkboxes for "on start" and "on success". A "Create New" button is at the bottom right of the dialog.

## Slowly Changing Dimensions (SCD)

So far we've built a set of pipelines that produce our report without taking into account any historical data. Historical information, however, can contribute a tremendous value to the business, as it allows us to draw patterns over time.

A commonly followed pattern in the data warehouse space has been called "Slowly changing dimensions". It allows us to store and manage both current and historical data over time directly within our tables. This used to be very difficult to do in a traditional Data Lake due to the lack of table merge and updates. However, Delta makes it possible.

There are three most popular types of SCD:

- **Type 1** - Overwriting
- **Type 2** - Creating dimension records
- **Type 3** - Creating historical columns

Consider a simple example where our initial data contains a Customer table in the operational database. After a Customer changes their address, the row with their corresponding customer\_id is updated.

If we're ingesting the data from our operational database using tools like [Fivetran](#) or [AWS DMS](#), we'd likely see two files created for us.

customer_id	zip_code	state
123	94105	San Francisco
<i>customers/2022-01-22T00:00:00.parquet</i>		

And, after the customer changes their address:

customer_id	zip_code	state
123	10001	New York
<i>customers/2022-05-22T13:23:23.parquet</i>		

## SCD Type 1

The standard operational databases store data using the SCD1 method to limit the amount of data present in the database at any point in time and speed up the querying.

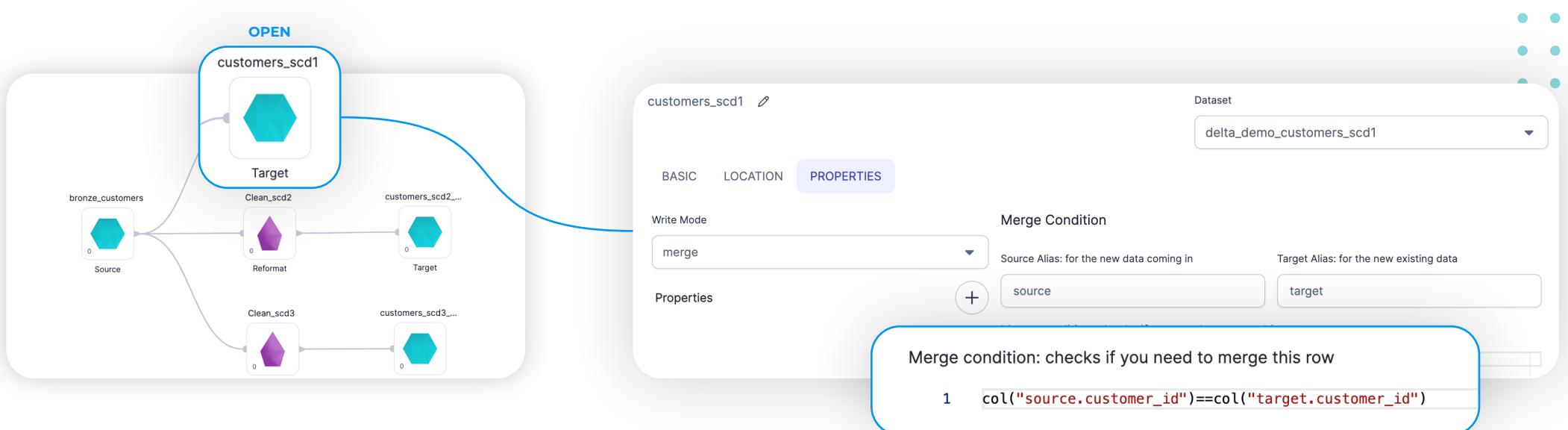
In our Lakehouse, we'd like to store the information using the Type 1 SCD which would simply mean updating and overwriting the row in the existing silver customers table.

This would result in the following table

customer_id string	customer_zip_code string	customer_city string	customer_state string	updated_dt date
A	10001	City1	State1	2022-01-15

*customers\_silver\_scd1*

SCD1 doesn't maintain the historical information, thus it's only useful when the change of information over time is not important. Building a pipeline that updates the existing customers, based on the latest record is very simple:

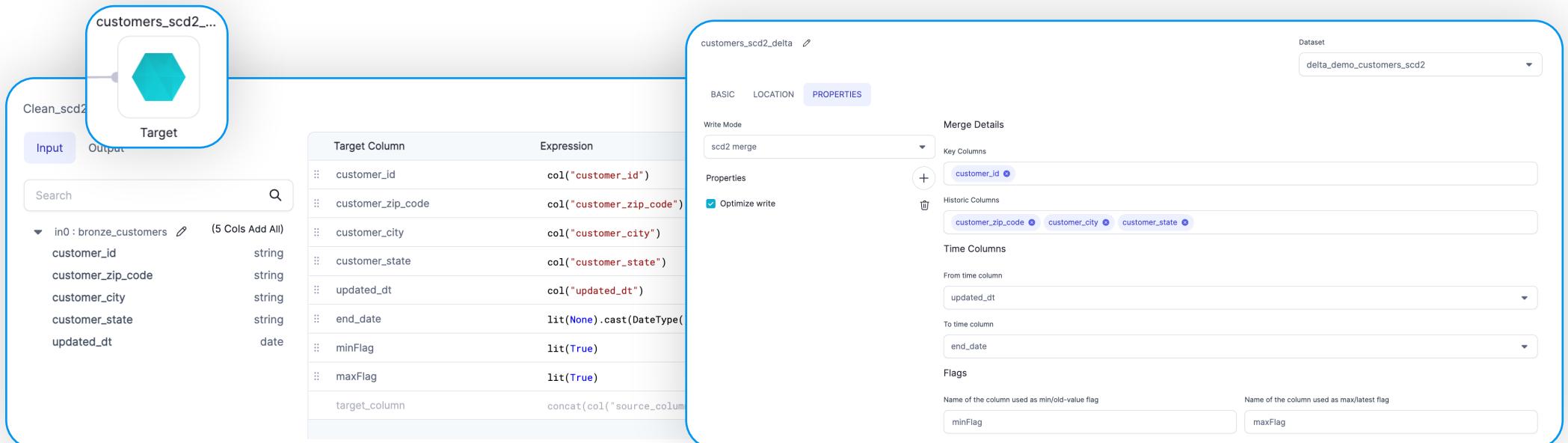


## SCD Type 2

SCD2 is another approach to handle slowly changing dimensions. In this case, current and historical data is stored in a data warehouse. The purpose of an SCD2 is to preserve the history of changes. If a customer changes their address, for example, or any other attribute, SCD2 allows analysts to link facts back to the customer and their attributes in the state they were at the time of the fact event.

customer_id	customer_zip_code	customer_city	customer_state	updated_dt	end_date	minFlag	maxFlag
string	string	string	string	date	date	boolean	boolean
A	10000	City1	State1	2022-01-01	2022-01-15	true	false
A	10001	City1	State1	2022-01-15	null	false	true

customers\_silver\_scd2



Here we have created 3 extra columns which would make it easy to consume this data in several downstream pipelines:

- **end\_date** - the date until the row was active
- **minFlag** - if true, the oldest record for that *customer\_id*
- **maxFlag** - if true, the latest record for that *customer\_id*

## SCD Type 3

A Type 3 SCD stores two versions of values for certain selected level attributes. Each record stores the previous value and the current value of the selected attribute. When the value of any of the selected attributes changes, the current value is stored as the old value and the new value becomes the current value.

customer_id	previous_customer_zip_code	customer_zip_code	previous_customer_city	customer_city	previous_customer_state	customer_state	updated
string	string	string	string	string	string	string	date
A	10000	10001	City1	City1	State1	State1	2022-0

customers\_silver\_scd3

The screenshot shows the configuration interface for a Type 3 SCD transformation. On the left, there is a preview of the source dataset, 'customers\_scd3...', which contains a single row with columns: customer\_id, previous\_customer\_zip\_code, customer\_zip\_code, previous\_customer\_city, customer\_city, previous\_customer\_state, customer\_state, and updated. The 'customer\_id' column has a value 'A', and the 'updated' column has a value '2022-0'. A blue box highlights this preview area.

The main configuration area is titled 'customers\_scd3\_delta'. It includes tabs for BASIC, LOCATION, and PROPERTIES. The PROPERTIES tab is active, showing:

- Write Mode:** merge
- Merge Condition:** col("source.customer\_id")==col("target.customer\_id")
- Properties:** Overwrite table schema (checkbox checked)

Below the merge condition, there is a section for **Custom Clauses**, which includes:

- Action:** update
- Only when the additional condition is true:** 1
- Replace default update with these expressions (optional):**

Target Column	Expression
previous_customer_	col("target.customer_zip_code")
customer_zip_code	col("source.customer_zip_code")
previous_customer_	col("target.customer_city")
customer_city	col("source.customer_city")
previous_customer_	col("target.customer_state")

## I've configured my slowly changing dimension gems, but how do I know if they are working?

Let's compare SCD1, SCD2, and SCD3 silver tables.

Using SCD1 merge, we don't retain any history. The customer records represent the latest changes, so we know that Customer A lives in zip code 10001.

customer_id string	customer_zip_code string	customer_city string	customer_state string	updated_dt date
A	10001	City1	State1	2022-01-15
B	10010	City2	State1	2022-01-02
C	30000	City3	State2	2022-01-15
D	20010	City4	State2	2022-01-12
E	30001	City5	State3	2022-01-15

customers\_silver\_scd1

Using SCD2 merge, we retain some interesting history. Customers A and C have moved across zip codes while customers B, D, and E have not. We can see the current zip code for each customer by checking the end\_date field. Customer A has moved from 10000 to 10001.

customer_id string	customer_zip_code string	customer_city string	customer_state string	updated_dt date	end_date date	minFlag boolean	maxFlag boolean
A	10000	City1	State1	2022-01-01	2022-01-15	true	false
A	10001	City1	State1	2022-01-15	null	false	true
B	10010	City2	State1	2022-01-02	null	true	true
C	20000	City3	State2	2022-01-11	2022-01-15	true	false
C	30000	City3	State2	2022-01-15	null	false	true
D	20010	City4	State2	2022-01-12	null	true	true
E	30001	City5	State3	2022-01-15	null	true	true

customers\_silver\_scd2

Using SCD3 merge, we can see Customer A's history in a slightly different format. We retain just one record for each customer, and we can see the current (eg 10001) and previous (eg 10000) zip codes as fields for each record.

customer_id string	previous_customer_zip... string	customer_zip_code string	previous_customer_city string	customer_city string	previous_customer_state string	customer_state string	updated_date
A	10000	10001	City1	City1	State1	State1	2022-0
B	null	10010	null	City2	null	State1	2022-0
C	20000	30000	City3	City3	State2	State2	2022-0
D	null	20010	null	City4	null	State2	2022-0
E	null	30001	null	City5	null	State3	2022-0

customers\_silver\_scd3

Phew that was a lot! We saw how to organize data in layers our Delta Lake with an example Pipeline in Prophecy. Then we got into the details of merging slowly changing data using SCD1, SCD2, and SCD3 options. Prophecy's low-code environment makes it easy for anyone to build this use case and contribute the resulting 100% open source Spark code to their git repo.

## How can I try Prophecy?

Prophecy is available as a SaaS product where you can add your Databricks credentials and start using it with Databricks. You can use an Enterprise Trial with Prophecy's Databricks account for a couple of weeks to kick the tires with examples. We also support installing Prophecy in your network (VPC or on-prem) on Kubernetes.

[SIGN UP FOR A TRIAL](#)

**Wherever you are in your cloud journey,  
we're the product partner for your success!**

