

Lakeflow Declarative Pipelines

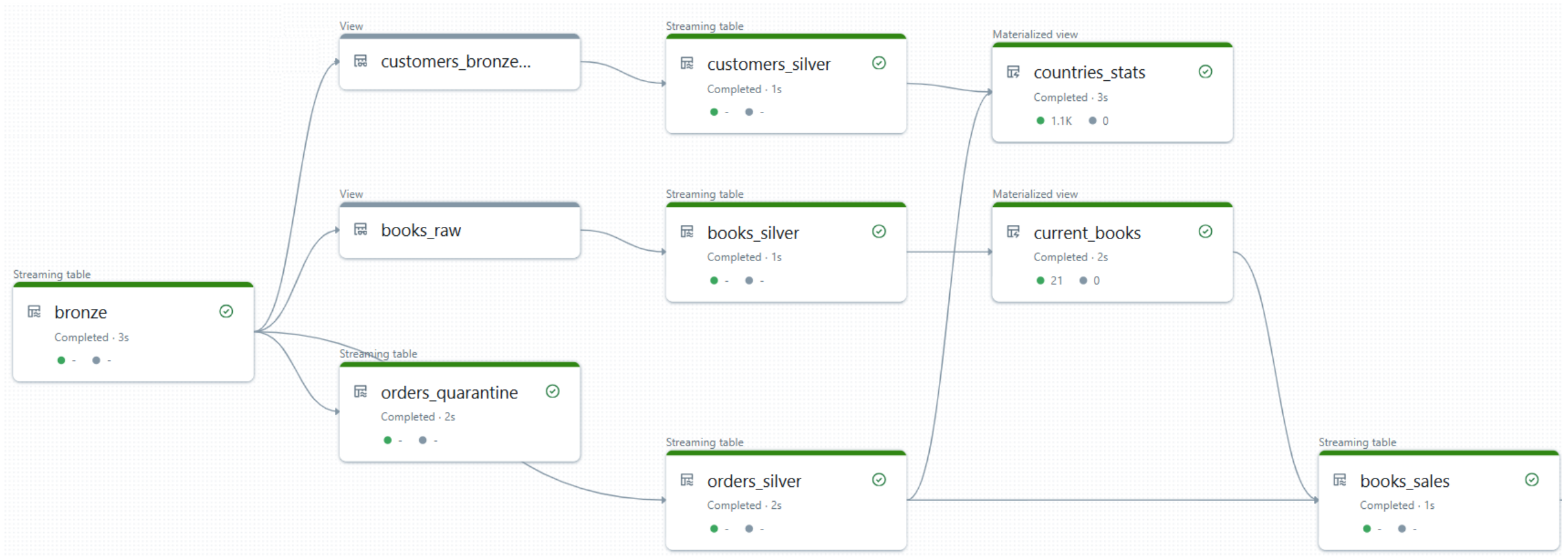
Learning Objectives

- ▶ Definition of Lakeflow Declarative Pipelines
- ▶ Differences with Apache Spark
- ▶ Object Types

Lakeflow Declarative Pipelines

- ▶ LDP is a declarative ETL framework powered by Apache Spark for building reliable and maintainable data pipelines.
- ▶ Previously known as Delta Live Tables (DLT)
- ▶ Open source: Spark Declarative Pipelines

Example



Benefits

- ▶ Automatic orchestration
- ▶ Handle checkpoints, retries, and optimization
- ▶ Easy to implement CDC, SCD Type 2, and data quality control

Spark vs. LDP

Spark

▶ `spark.readStream`

```
.format("cloudFiles")  
.option("cloudFiles.format", "json")  
.load('/some/path/')
```

`.writeStream`

```
.option("checkpointLocation", "/path")  
.table("orders_raw")
```

LDP

▶ `from pyspark import pipelines as dp`

`@dp.table`

`def orders_raw():`

`return (spark.readStream`

```
.format("cloudFiles")  
.option("cloudFiles.format", "json")  
.load("/some/path")  
)
```

Spark vs. LDP

Spark

- ▶ Can Not create streaming tables in Spark SQL syntax only. Need to pass by PySpark to register streaming tables

LDP

- ▶ Support creating streaming tables in SQL via **CREATE STREAMING TABLE** command

Creating LDP Objects (SQL)

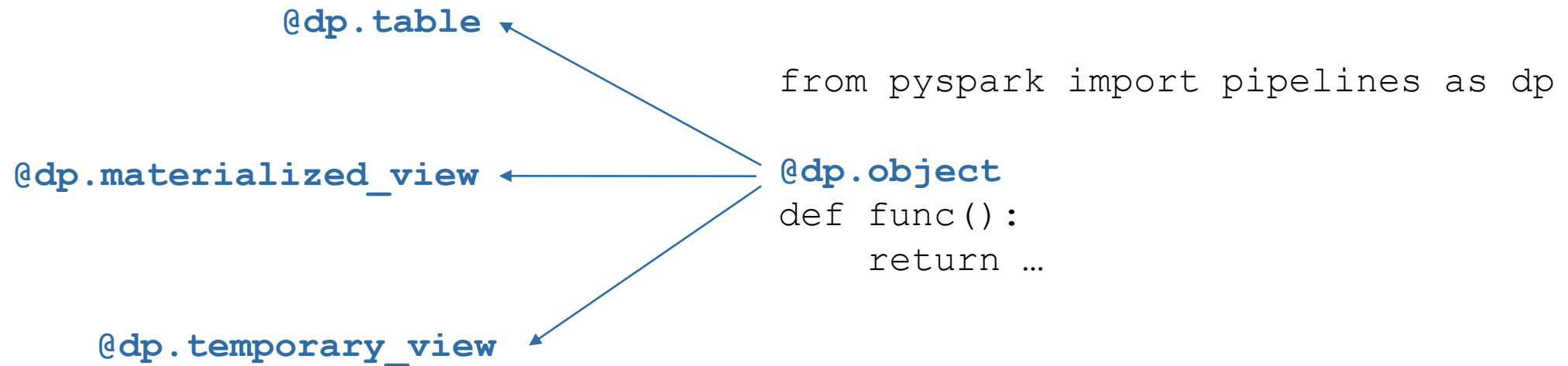
CREATE OR REFRESH **OBJECT** <object_name>

STREAMING TABLE

MATERIALIZED VIEW

TEMPORARY VIEW

Creating LDP Objects (Python)



LDP Object Types

Streaming Tables

- ▶ Permanent objects
- ▶ Handle incremental refresh
- ▶ Used for data ingestion from streaming sources (append-only):
 - ▶ `spark.readStream`
 - ▶ **STREAM()** SQL function
- ▶ Support near real-time data ingestion

Materialized Views

- ▶ Permanent objects
- ▶ Handle full or incremental refresh (Serverless-only)
- ▶ Used for precomputing complex BI queries, or for data ingestion from non-streamable sources:
 - ▶ `spark.read`
- ▶ Not designed for low-latency use cases

Temporary Views

- ▶ Temporary objects
- ▶ Handle temporarily processed data
- ▶ Used for intermediate transformations and data quality checks

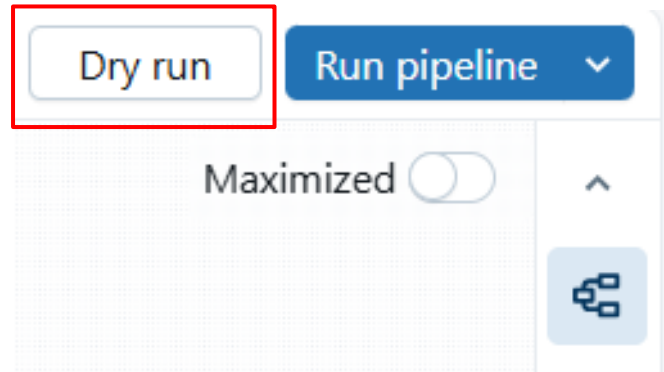
LDP vs DLT

	LDP	DLT (Old Syntax)
Module import	<pre>from pyspark import pipelines as dp</pre>	<pre>import dlt</pre>
Streaming Table	<pre>@dp.table def func() return spark.readStream ...</pre>	<pre>@dlt.table def func() return spark.readStream ...</pre>
Materialized View	<pre>@dp.materialized_view def func() return spark.read ...</pre>	<pre>@dlt.table def func() return spark.read ...</pre>
Temporary View	<pre>@dptemporary_view def func () return ...</pre>	<pre>@dlt.view def func() return ...</pre>

Validating Code

LDP

- Files (.py or .sql)



DLT

- Notebooks

