# Hassle-free data ingestion into the Lakehouse

Make data ingestion simpler

**Burak Yavuz**
Engineering Manager, Databricks

**Benyue (Emma) Liu**
Staff Product Manager, Databricks
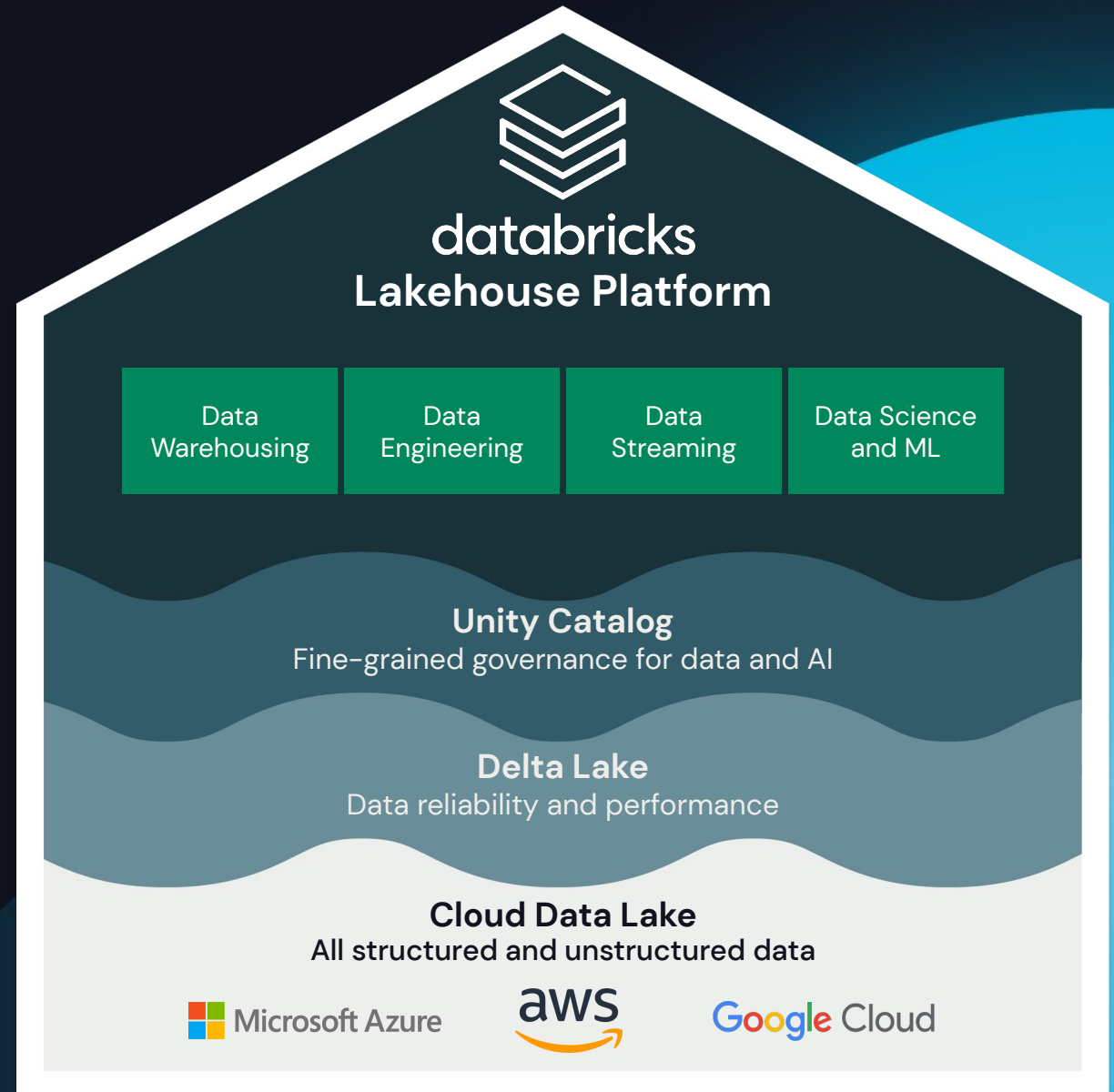
# Product safe harbor statement

This information is provided to outline Databricks' general product direction and is for informational purposes only. Customers who purchase Databricks services should make their purchase decisions relying solely upon services, features, and functions that are currently available. Unreleased features or functionality described in forward-looking statements are subject to change at Databricks discretion and may not be delivered as planned or at all.

# Agenda

1.  Databricks Lakehouse overview

2.  Ingestion challenges

3.  Get started with ingestion in less than 10 minutes

4.  Demo

# Databricks
# Lakehouse overview

Your destination is the Lakehouse

databricks
**Lakehouse Platform**

| Data Warehousing | Data Engineering | Data Streaming | Data Science and ML |

**Unity Catalog**
Fine-grained governance for data and AI

**Delta Lake**
Data reliability and performance

**Cloud Data Lake**
All structured and unstructured data

Microsoft Azure     aws     Google Cloud

# Ingestion challenges

# Data ingestion presents challenges

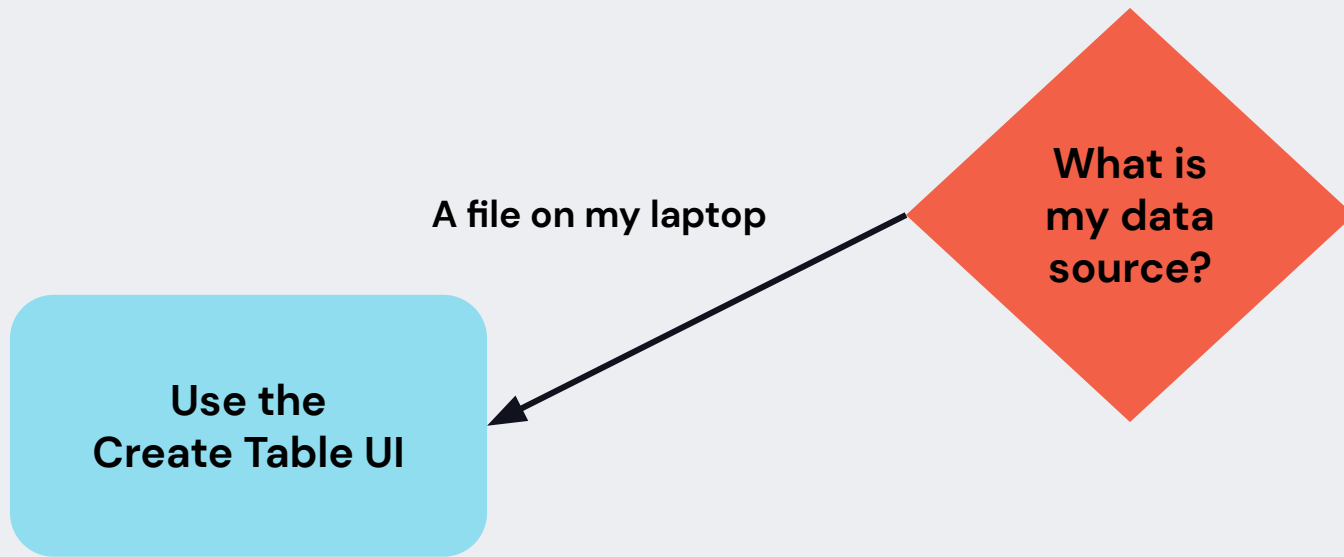| | |
|---|---|
| Too many data sources | Migrating existing tables |
| Figuring out what files to process | Different requirements on data freshness |
| Schema changes over time | Cloud configurations are complex |
| Fixing issues with bad data | Change-data-capture from OLTP databases |
| Scalability | Enable both ad hoc ingestion and production data pipeline |

**Cloud Data Lake**
Structured, semi-structured, and unstructured data

Get started with
Lakehouse ingestion
in less than
10 minutes!

What is
my data
source?

A file on my laptop

Use the
Create Table UI

# New ingestion features in Databricks SQL

## Create table from local files or cloud object stores

Quickly upload a local CSV & create a Delta table from it from DBSQL UI

**What's coming next?**

- Support for additional file types (JSON)

- UI to ingest from cloud storage like S3

- Unified Data Source UI to the Lakehouse

- SQL API to ingest from cloud storage

| Microsoft Azure | Databricks | | | | Portal  emma.liu@databricks.com |
|---|---|---|---|---|---|

**Create table in Databricks SQL**   Data McNugget 2.0 (XL)

airports.csv  9.84MB  ✕

hive_metastore | default | airports

☑ First row contains the header   Advanced attributes

| id | ident | type | name | latitude_deg | longitude_deg |
|---|---|---|---|---|---|
| 6523 | 00A | heliport | Total Rf Heliport | 40.07080078125 | -74.933601379394 |
| 323361 | 00AA | small_airport | Aero B Ranch Airport | 38.704022 | -101.473911 |
| 6524 | 00AK | small_airport | Lowell Field | 59.947733 | -151.692524 |
| 6525 | 00AL | small_airport | Epps Airpark | 34.86479949951172 | -86.7703018188470 |
| 6526 | 00AR | closed | Newport Hospital & Clinic Heliport | 35.6087 | -91.254898 |
| 322127 | 00AS | small_airport | Fulton Airport | 34.9428028 | -97.8180194 |
| 6527 | 00AZ | small_airport | Cordes Airport | 34.305599212646484 | -112.16500091552 |
| 6528 | 00CA | small_airport | Goldstone (GTS) Airport | 35.35474 | -116.885329 |
| 324424 | 00CL | small_airport | Williams Ag Airport | 39.427188 | -121.763424 |
| 322658 | 00CN | heliport | Kitchen Creek Helibase Heliport | 32.7273736 | -116.4597417 |
| 6529 | 00CO | closed | Cass Field | 40.622202 | -104.344002 |
| 6531 | 00FA | small_airport | Grass Patch Airport | 28.64550018310547 | -82.2190017700190 |

Create   Cancel

GA NOW

What is my data source?

A file on my laptop → Use the Create Table UI

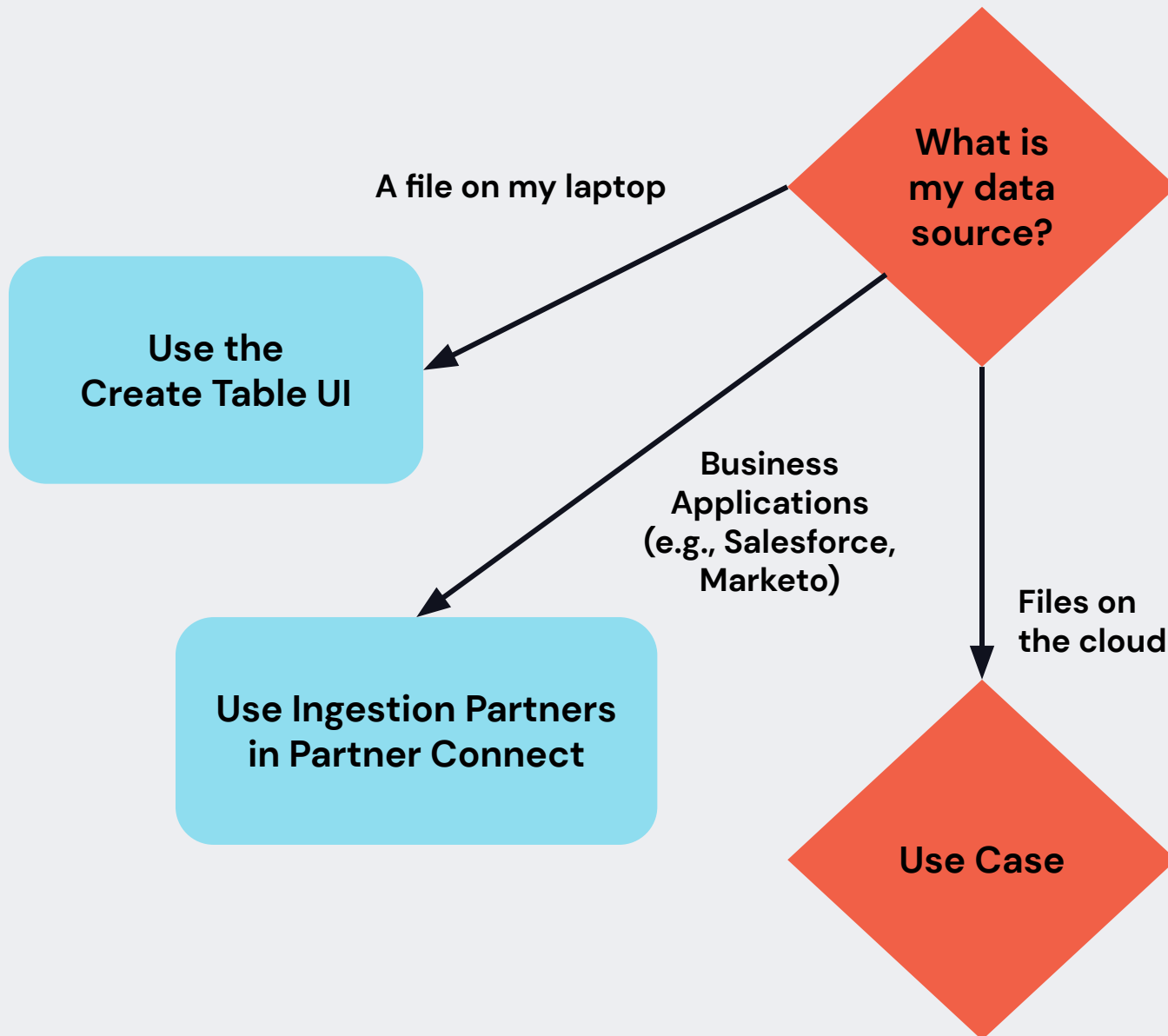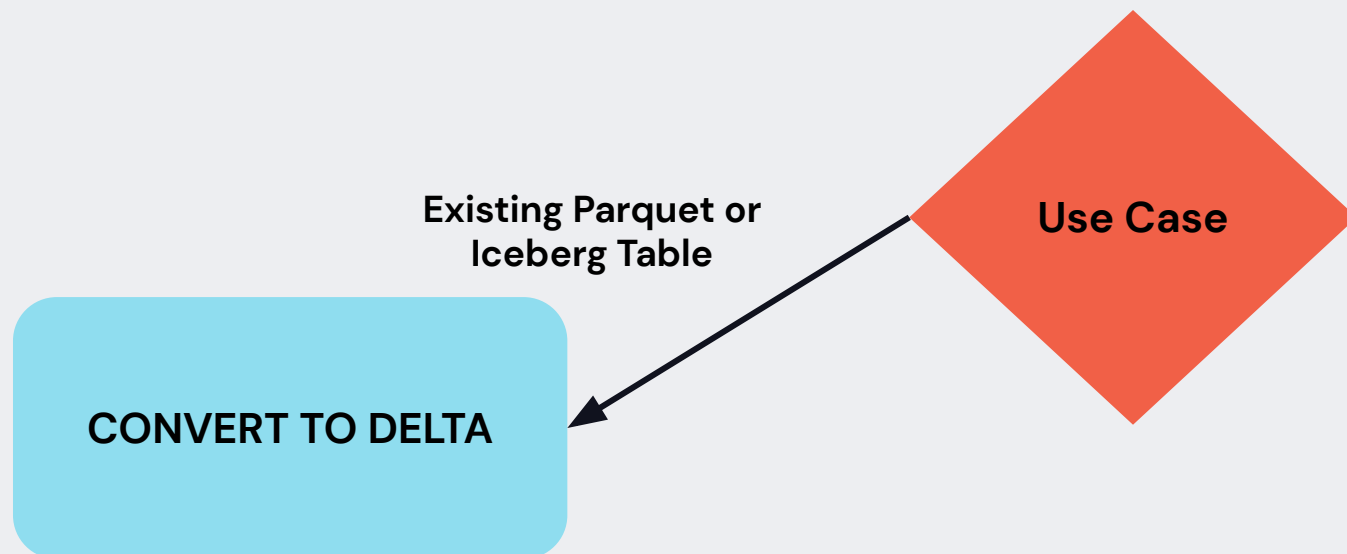Business Applications (e.g., Salesforce, Marketo) → Use Ingestion Partners in Partner Connect

# Databricks Partner Connect



Easily discover and connect data, analytics, and AI tools to your lakehouse

# Ingesting files on the Cloud

**Use Case**

**Existing Parquet or Iceberg Table**

**CONVERT TO DELTA**

# Ingesting files on the Cloud

**Existing Parquet or Iceberg Table**

**Use Case**

**CONVERT TO DELTA**

**Scale/Schema Evolution**

**Auto Loader**

# Auto Loader

- Available in Python & Scala (and SQL in Delta Live Tables!)
- Incremental loading
- Exactly once ingestion
- Scalable for large amounts of data
- Designed for structured, semi-structured and unstructured data

```python
df = spark.readStream.format("cloudFiles")
    .option("cloudFiles.format", "json")
    .load("/path/to/table")
```

"We ingest more than
**5 petabytes per day**
with Auto Loader"
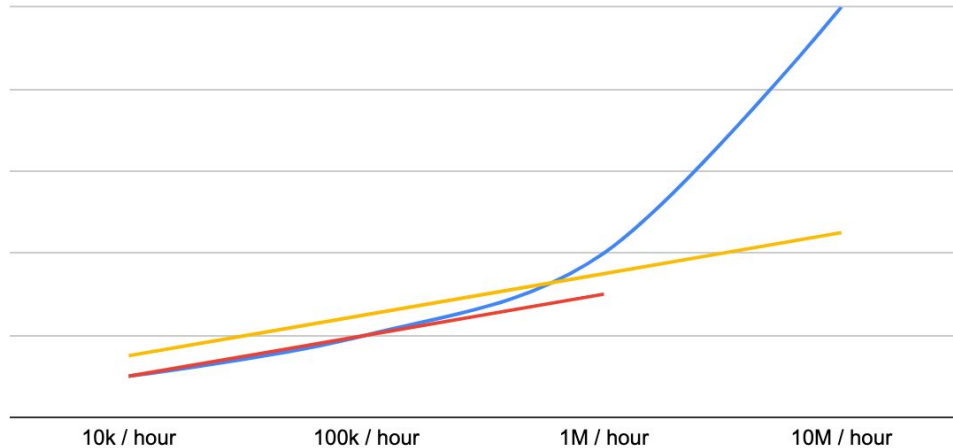
# Auto Loader performance—
# File Notifications & Directory Listing
## Comparison

### Relative Performance



**Why two modes?**

- Directory Listing is **simple** to set up and works well for small throughputs
  - Optimized Cloud APIs reduce RPC costs
- *Incremental Listing* is a **special case** of Directory Listing where files are incrementally ordered. **Automagically** detected.
  - 2022-05-01-0001.csv
  - 2022-05-01-0002.csv
  - Etc.
- File Notifications allow **scaling** to millions of files/hour, but require permissions

# Auto Loader schema management

| Schema Location | Schema Inference | Data Rescue | Schema Evolution Mode |
|---|---|---|---|
| (cloudFiles.schemaLocation) | (cloudFiles.inferColumnTypes) | (cloudFiles.rescuedDataColumn) | (cloudFiles.schemaEvolutionMode) |

- Stores changes to the inferred schema over time

- Infer column data types or treat everything as a String

- Rescues data that does not match your schema expectation

- Add new columns
- Fail on new columns
- Rescue new columns

# Ingesting files on the Cloud



Use Case

Existing Parquet or Iceberg Table

CONVERT TO DELTA

Ad–hoc/Simplicity/ Reloading files

COPY INTO

Scale/Schema Evolution
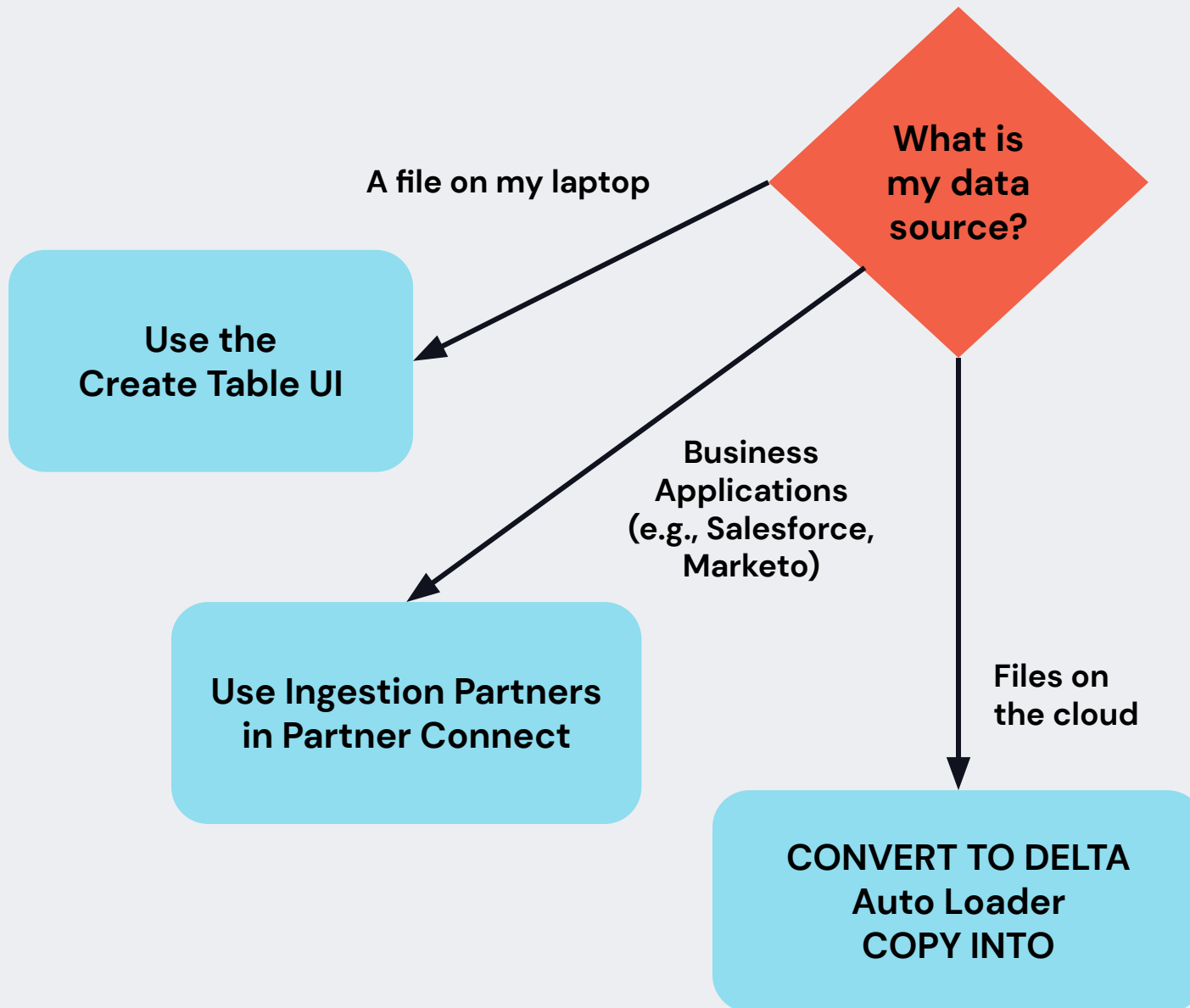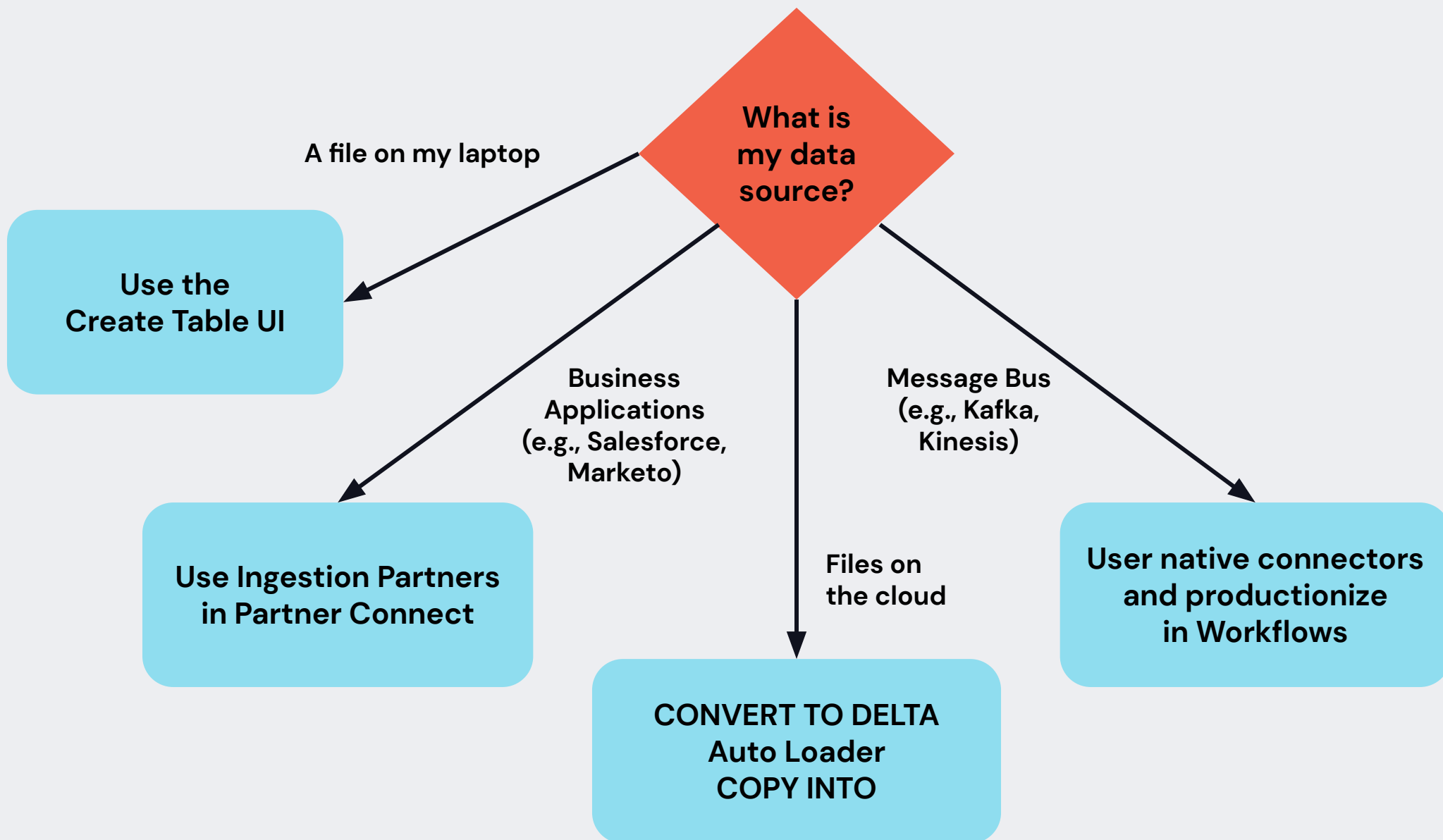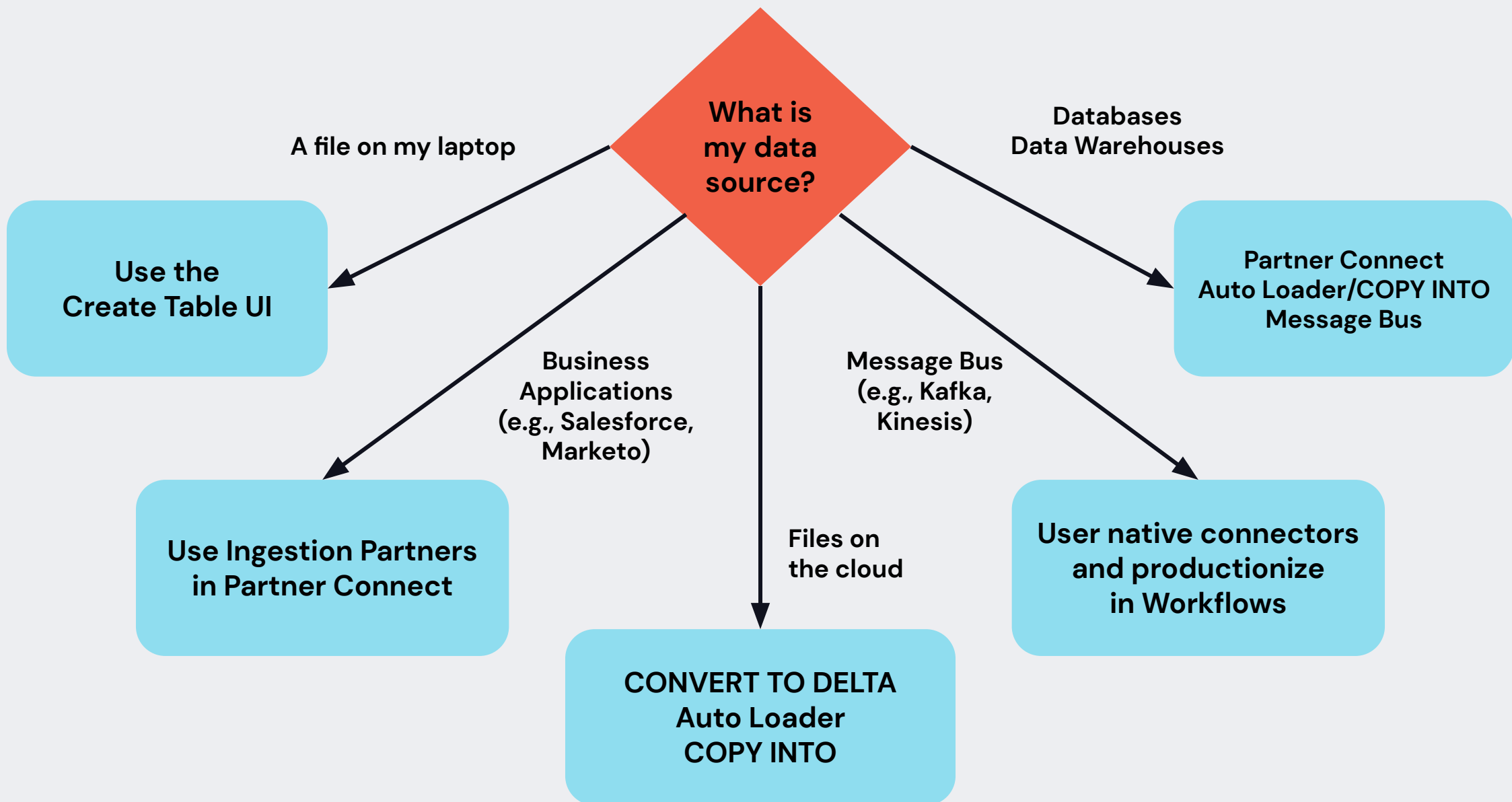
Auto Loader

# COPY INTO

- SQL command
- Idempotent and incremental
- Great when source directory contains ~ thousands of files
- Schema automatically inferred

```
COPY INTO my_delta_table
FROM 's3://my-bucket/path/to/csv_files'
FILEFORMAT = CSV
FORMAT_OPTIONS ('header'='true','inferSchema'='true')
```
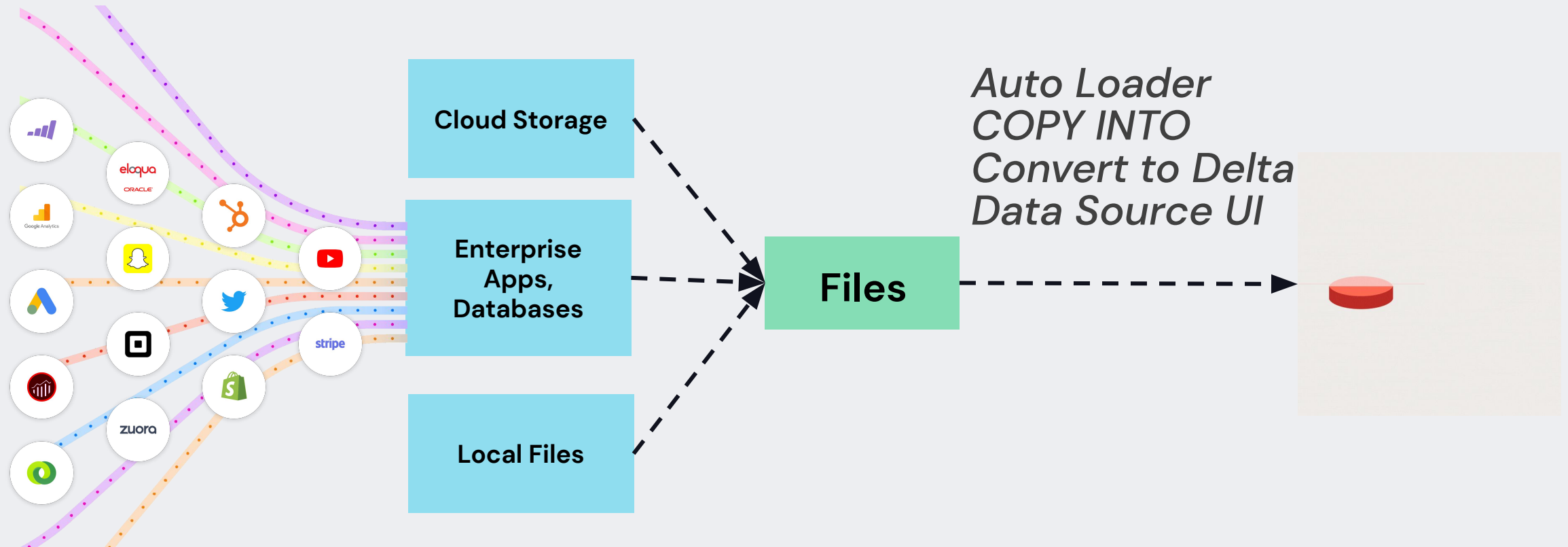
What is my data source?

A file on my laptop → Use the Create Table UI

Business Applications (e.g., Salesforce, Marketo) → Use Ingestion Partners in Partner Connect

Files on the cloud → CONVERT TO DELTA Auto Loader COPY INTO

**A file on my laptop**

**What is my data source?**

**Databases Data Warehouses**

**Use the Create Table UI**

**Partner Connect Auto Loader/COPY INTO Message Bus**

**Business Applications (e.g., Salesforce, Marketo)**

**Message Bus (e.g., Kafka, Kinesis)**

**Use Ingestion Partners in Partner Connect**

**Files on the cloud**

**User native connectors and productionize in Workflows**

**CONVERT TO DELTA Auto Loader COPY INTO**

# Demo

# Common ingestion pattern



Cloud Storage

Enterprise Apps, Databases

Local Files

Files

*Auto Loader*
*COPY INTO*
*Convert to Delta*
*Data Source UI*

# Related ingestion sessions

**Moving to the Lakehouse:
Fast & Efficient Ingestion with
Auto Loader**

- Virtual session

- Speakers: Emma Liu & Eric Maynard

**Ingest Data into Lakehouse with
COPY INTO**

- Virtual session

- Speakers: Yaohua Zhao & Jackie Zhang

# Learn more

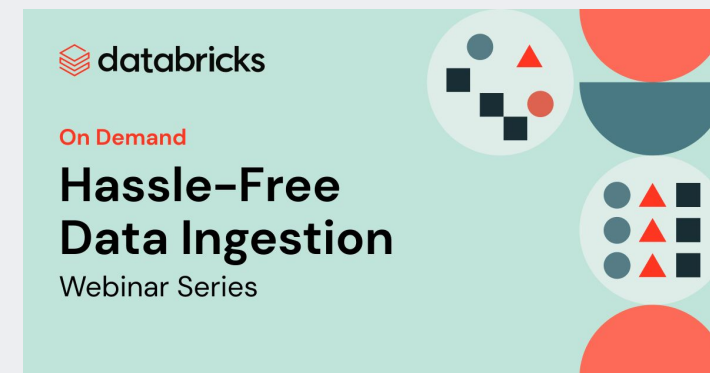## Read the documentation



**Ingest data into Delta Lake**

## Read an ebook



**All Roads Lead to the Lakehouse**

A deep dive into data ingestion with the lakehouse

## Watch Ingestion webinar series



**Easily load data into Delta Lake to power analytics, data science and machine learning**

DATA+AI
SUMMIT 2022

# Thank you!

## Q&A