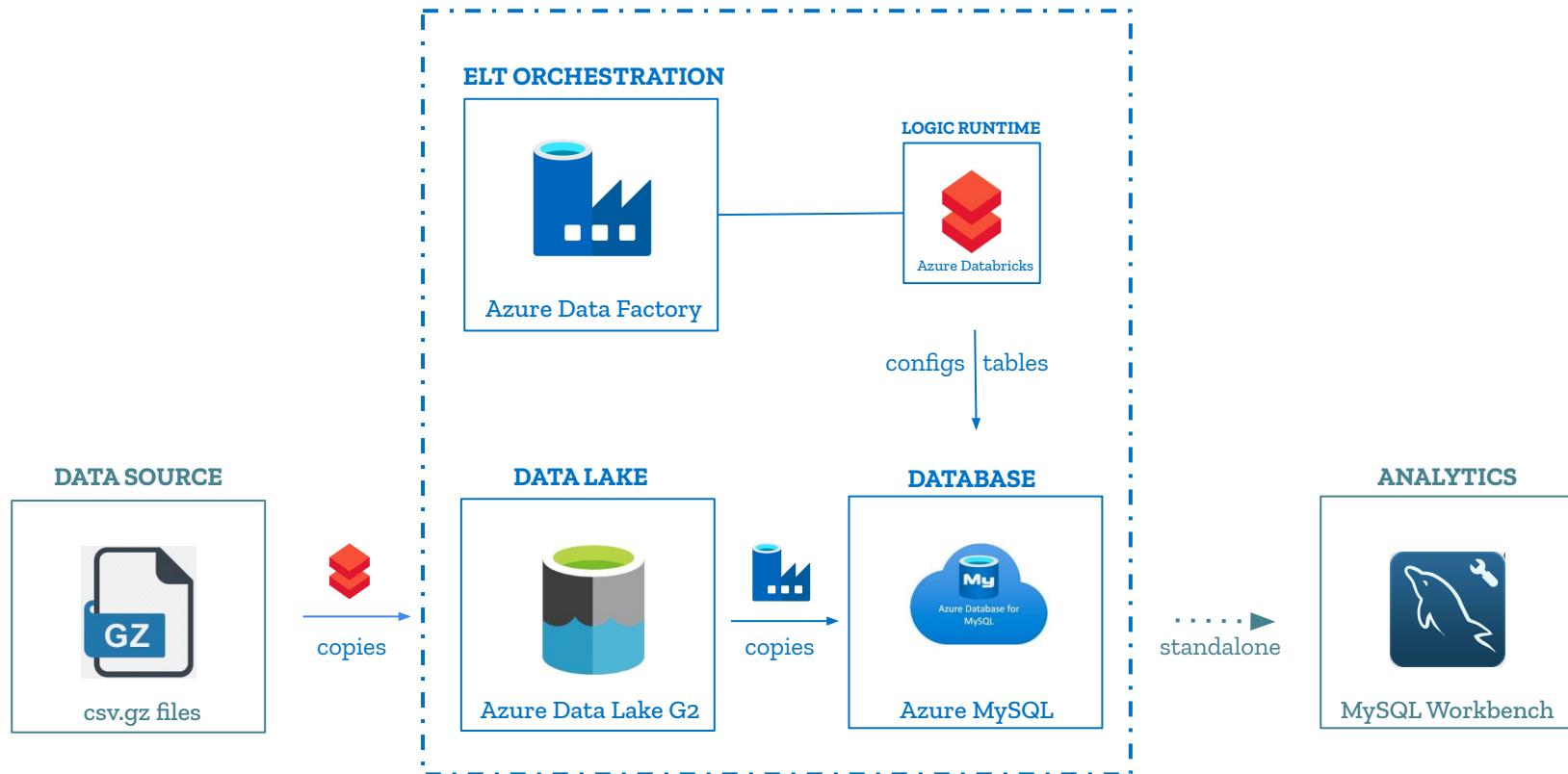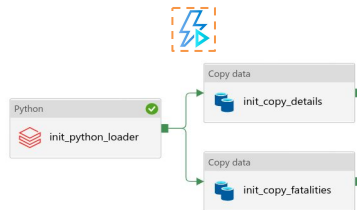# Final Submission

Capstone

# Background

# Architecture

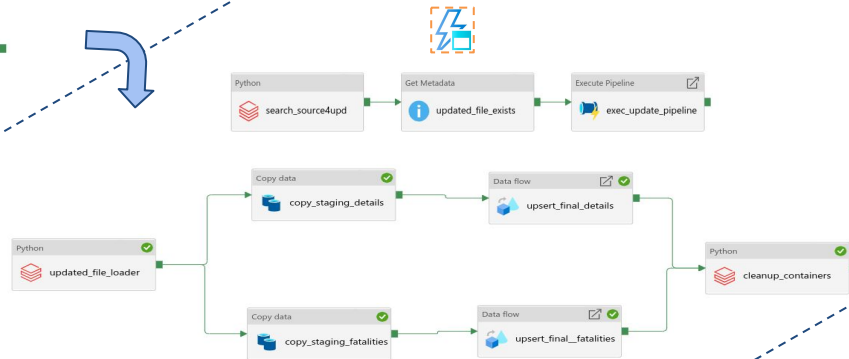# ETL Overview

Three Azure Data Factory pipelines tailored for different use case based on source data release cadence
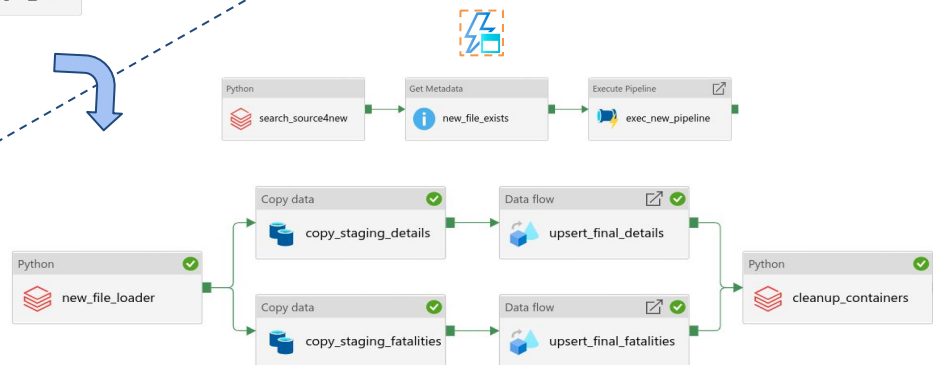
## Initial Load

Python
init_python_loader

Copy data
init_copy_details

Copy data
init_copy_fatalities

## Monthly Update

Python
search_source4upd

Get Metadata
updated_file_exists

Execute Pipeline
exec_update_pipeline

Python
updated_file_loader

Copy data
copy_staging_details

Data flow
upsert_final_details

Copy data
copy_staging_fatalities

Data flow
upsert_final_fatalities

Python
cleanup_containers

## Yearly New

Python
search_source4new

Get Metadata
new_file_exists

Execute Pipeline
exec_new_pipeline

Python
new_file_loader

Copy data
copy_staging_details

Data flow
upsert_final_details

Copy data
copy_staging_fatalities

Data flow
upsert_final_fatalities

Python
cleanup_containers

# Deep Dive: "Yearly New" ETL

1) **New File Trigger**

   ● Annually on April 17th Triggers "Extract New File" Pipeline

2) **"Extract New File" Pipeline**

   ● check source url for new file drop
     ○ pull if available
       ■ **success:** triggers "Load+Transform New File" Pipeline
       ■ **failure:** triggers retry in 1 day

3) **"Load+Transform New File" Pipeline**

# "Extract New File" Pipeline

**Databricks Python** Activity:
**spark cluster:** initialized with environment variables for auth
**runs:** `python new_file_to_blob_only.py`
  **if:** source url has dropped a new filename for new year
  **then:** extract new csv.gz file
    **from:** source url
    **to:** data lake "newfiles" container

**Get Metadata** Activity:
**if:** file exists in "newfiles" container
  **then**: continue

**Execute Pipeline** Activity:
**trigger:** "Load+Transform New File" Pipeline

| Python ✅ | Get Metadata ✅ | Execute Pipeline ↗ ✅ |
|---|---|---|
| 🗇 search_source4new | ⓘ new_file_exists | 🛢 exec_new_pipeline |

---

# "Load+Transform New File" Pipeline

**Copy Data** Activity
**copy:** all data from new year csvgz file
  **from:** data lake "newfiles" container
  **to:** mysql staging table

**Data Flow** Activity
**insert:** new rows based on primary key (equates to all)
  **from:** mysql staging table
  **to:** mysql final table

**Databricks Python** Activity
**run:** `python datalake_housekeeping_new.py`
  **copy** new 2022 files
    **from:** "newfiles" container
    **to:** "allfiles" container
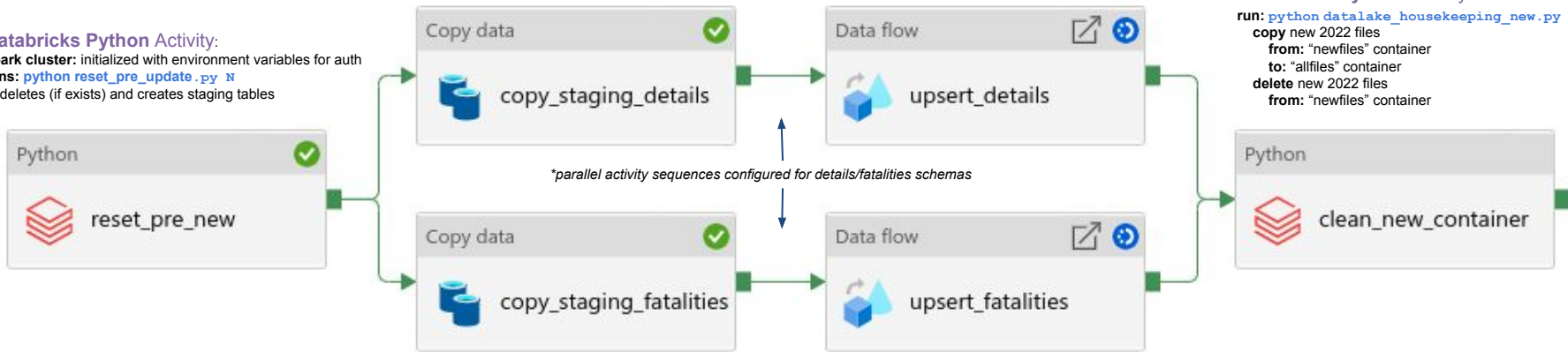  **delete** new 2022 files
    **from:** "newfiles" container

**Databricks Python** Activity:
**spark cluster:** initialized with environment variables for auth
**runs:** `python reset_pre_update.py N`
  deletes (if exists) and creates staging tables

| Python ✅ | Copy data ✅ | Data flow ↗ 🔵 | Python |
|---|---|---|---|
| 🗇 reset_pre_new | 🛢 copy_staging_details | 🔼 upsert_details | 🗇 clean_new_container |

*parallel activity sequences configured for details/fatalities schemas*

| Copy data ✅ | Data flow ↗ 🔵 |
|---|---|
| 🛢 copy_staging_fatalities | 🔼 upsert_fatalities |

# Development Process

# Testing

# Deployment

# Monitoring

# Lessons Learned