

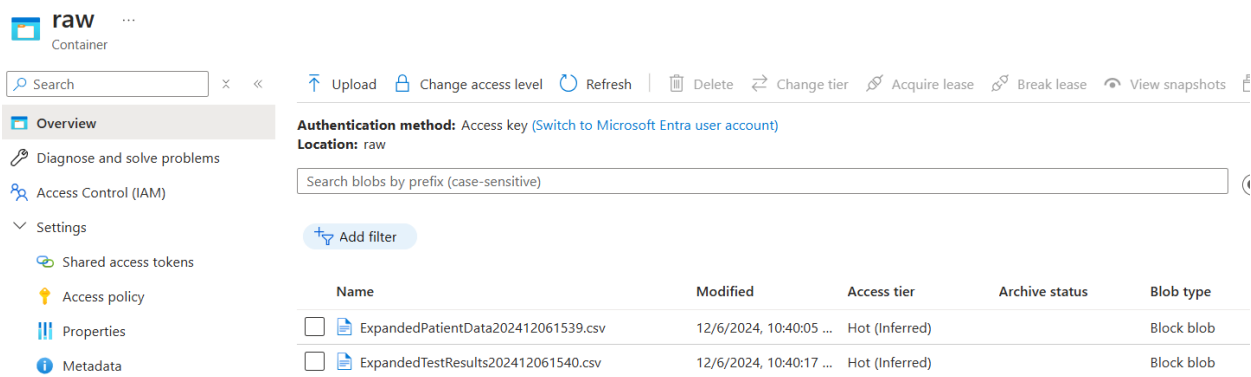
# Case Study: End-to-End Data Pipeline for COVID-19 Test Data Management

## Architecture

### Data Flow Steps:

1. \*\*Data Storage\*\*:

- Raw data was stored in Azure Blob Storage (HTTP endpoint).
- File types included CSVs containing patient and test details.



2. \*\*Data Cleaning and ETL\*\*:

- Databricks was used to clean and transform the data using PySpark.
- Key cleaning steps:
  - Removing rows with missing identifiers (e.g., `PatientID`, `TestID`).
  - Standardizing data types (e.g., casting `Age` to integer, `TestDate` to date).
  - Handling null values by replacing them with defaults like 'Unknown'.
- ETL steps included joining patient and test data, deriving new columns such as `ResultStatus`, and selecting relevant fields for reporting.

```
01:31 PM (2/25) 2

mount_point = "/mnt/httpforstorage/raw"

if any(mount.mountPoint == mount_point for mount in dbutils.fs.mounts()):
    # Unmount the existing mount point
    dbutils.fs.unmount(mount_point)
    print(f"Unmounted existing mount at {mount_point}")

try:
    dbutils.fs.mount(
        source="abfss://raw@httpforstorage.dfs.core.windows.net/",
        mount_point=mount_point,
        extra_configs=configs
    )
    print(f"Mounted successfully at {mount_point}")
except Exception as e:
    print(f"Error mounting: {e}")

/mnt/httpforstorage/raw has been unmounted.
Unmounted existing mount at /mnt/httpforstorage/raw
Mounted successfully at /mnt/httpforstorage/raw
```

### 3. \*\*Data Loading\*\*:

- Transformed data was saved in Delta format.
- Data was ingested into an Azure Synapse Analytics dedicated SQL pool for reporting and analytics.

---

```
01:31 PM (10/5) 6 Python

# Databricks Data Cleaning for COVID-19 Data
# Import required libraries
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, when, lit, countDistinct
from pyspark.sql.types import StringType, IntegerType, DateType

# Initialize Spark session
spark = SparkSession.builder.appName("COVIDDataCleaning").getOrCreate()

# Load Patient Data
patient_file_path = "/mnt/httpforstorage/raw/ExpandedPatientData202412061539.csv" # Update path if needed
patients_df = spark.read.format("csv").option("header", "true").load(patient_file_path)

# Load Test Data
test_file_path = "/mnt/httpforstorage/raw/ExpandedTestResults202412061540.csv" # Update path if needed
tests_df = spark.read.format("csv").option("header", "true").load(test_file_path)

# --- Data Cleaning on Patients Data ---
# Drop rows with missing PatientID
patients_df = patients_df.filter(col("PatientID").isNotNull())

# Cast PatientID and Age to Integer and TestDate to Date
def clean_patient_data(df):
    return (df.withColumn("PatientID", col("PatientID").cast(IntegerType()))
            .withColumn("Age", col("Age").cast(IntegerType()))
            .withColumn("TestDate", col("TestDate").cast(DateType())))
```

### ## Implementation

### ### 1. Data Storage in Azure Blob Storage

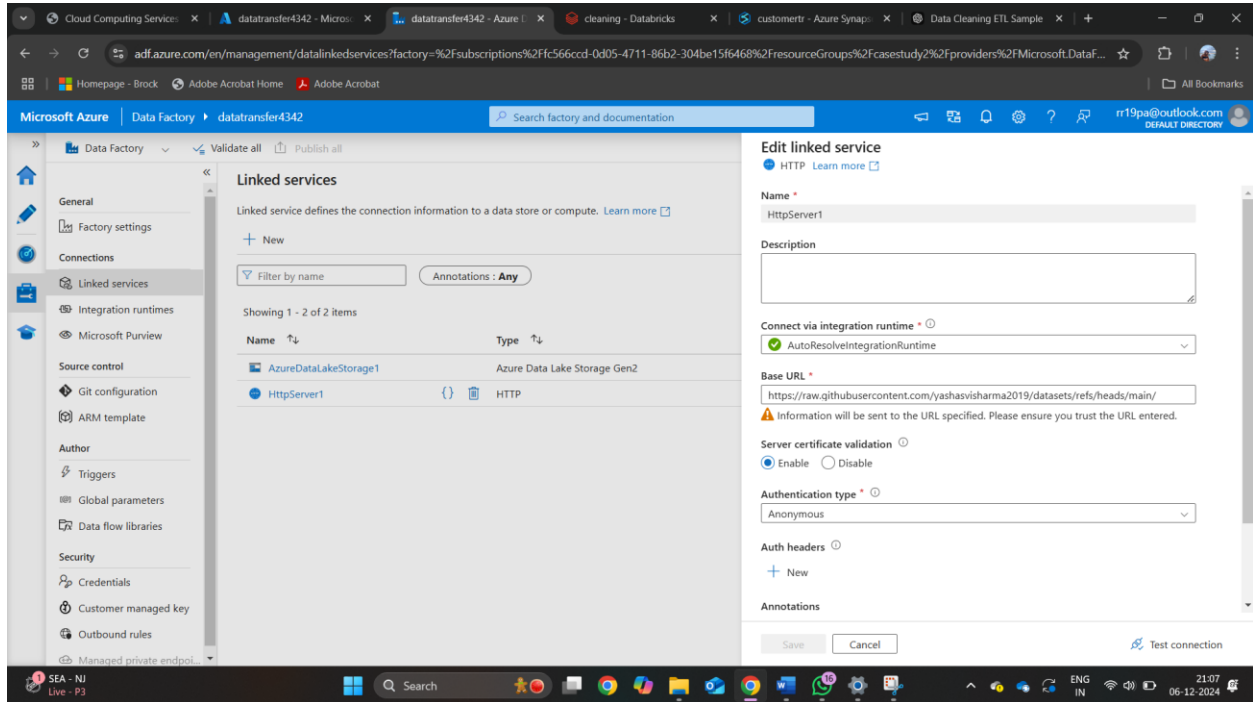
#### - \*\*Source Data:\*\*

- Patient data: `ExpandedPatientData.csv`

- Test results: `ExpandedTestResults.csv`

#### - \*\*Storage Location:\*\*

- Data was uploaded to an Azure Blob Storage container via HTTP.



### ### 2. Data Cleaning and Transformation with Databricks

#### #### Data Cleaning

#### - \*\*Patients Data\*\*:

- Dropped rows with null `PatientID`.

- Cast `PatientID` and `Age` to `IntegerType` and `TestDate` to `DateType`.

- Replaced null `COVIDStatus` values with 'Unknown'.

#### - \*\*Tests Data\*\*:

- Removed rows with null `TestID` or `PatientID`.

- Cast `TestID`, `PatientID` to `IntegerType` and `ResultDate` to `DateType`.

- Replaced null `LabName` values with 'Unknown'.

#### #### Data Transformation

- Joined patient and test data on `PatientID`.
- Added a derived column `ResultStatus` based on test results ('Positive' or 'Negative').
- Selected and renamed columns for analytics.

#### #### Error Handling

- Missing columns or unresolved variables were handled with validation steps and exceptions.
- Write operations to storage were encapsulated in try-except blocks to capture and log errors.

### ### 3. Data Loading to Synapse Analytics

- **Output Storage**: Transformed data was saved in Delta format in Azure Data Lake.
- **Synapse Pipeline**:
  - Data from Delta tables was ingested into an Azure Synapse Analytics dedicated SQL pool.

---

## ## Results

- Improved data quality by resolving inconsistencies and handling missing values.
- Achieved a centralized, optimized SQL-based data store for reporting and analytics.
- Enabled real-time insights for healthcare decision-making.

---

## ## Conclusion

This end-to-end pipeline effectively addressed the challenges of managing raw COVID-19 test data. The integration of Azure Blob Storage, Databricks, and Synapse Analytics streamlined the process from data ingestion to reporting. The scalable architecture supports future enhancements for additional datasets or analytics requirements.

## ## Future Work

- Automate the entire pipeline using Azure Data Factory.
- Implement advanced analytics on Synapse using machine learning models.