# Scaling The Prototype

Capstone Step 6

# Process Updating ETL to Spark

This process was quite more involved than it had to be.  The main steps of the approach included:

1. Attempting and ultimately failing to debug **Azure Toolkit aztk**.

      This repository, prescribed by the project instructions, has become obsolete.

      The technical details of the issues are documented in **development_notes.txt**

2. Opted to use **Azure Databricks** instead of aztk.

3. Created **Azure Blob Storage Container**

      Uploaded data to blob container

      Configured entitlements of new **Databricks Cluster** to read from blob container.

4. Created notebook **pyspark_pipeline_nbook.ipynb** in Databricks with **Pyspark kernel** to load data from Azure blob storage.

      Translated code from Step 5 from Python to Pyspark.

      Observed noticeable speed improvement from locally deployed python script.

# Databricks Spark Notebook