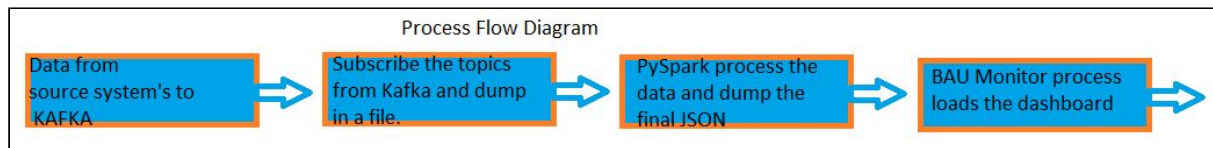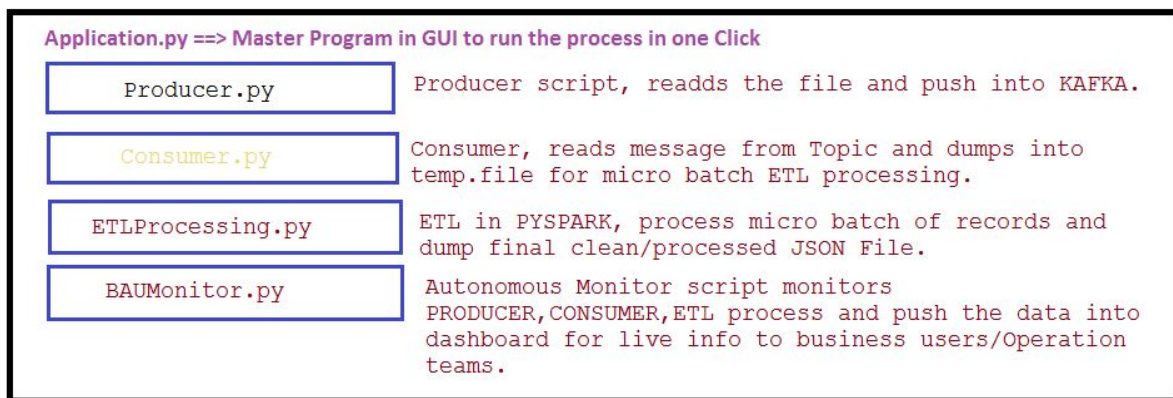# Data Streaming Application in kafka and Spark Processing

This is a standalone application for stream data processing from source, process the data in pyspark and load the data as JSON in the local directory..
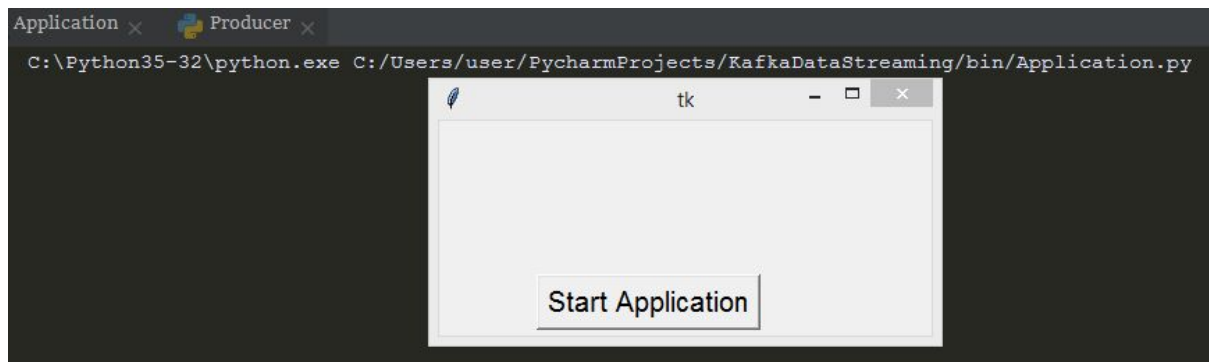
*The data pipeline is given below,*



## Script Flow Diagram:



## Execution Steps:

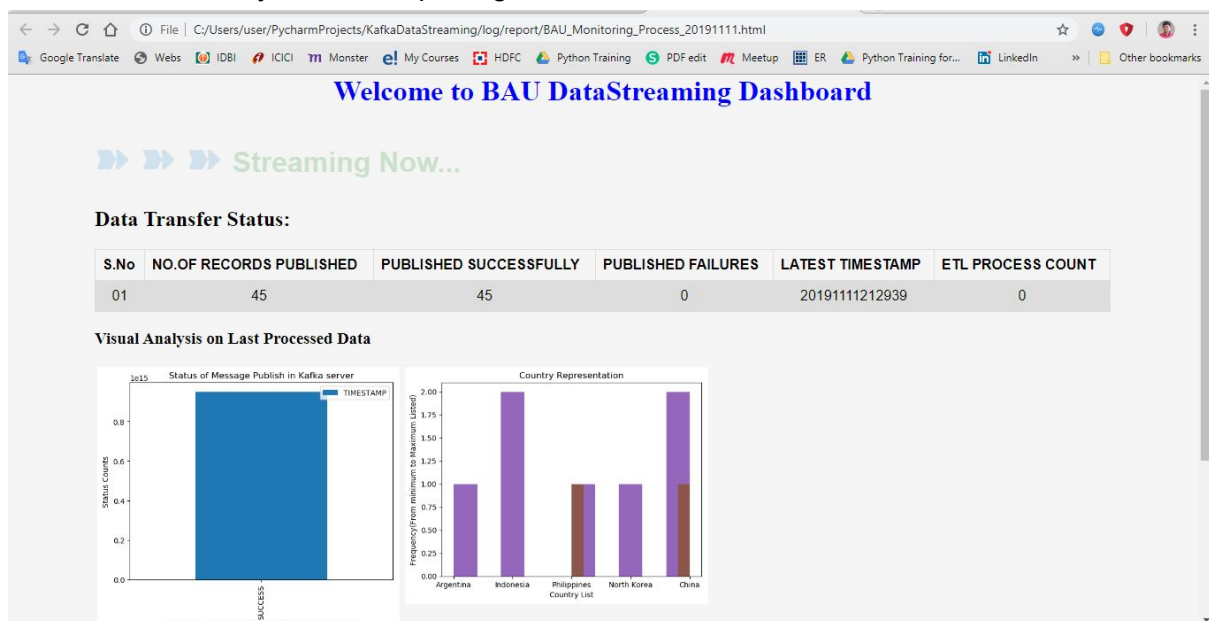User just need to execute the Applciation.py a GUI opens and should Start application as below,

Upon clicking the Start application,process starts as,



Dashboard View:

Since this is just a temporary project parameters needs to be tuned fined for excellent performance and better GUI/UI/UX performance. For now,the Dashboard(HTML) should be manually refreshed over browser due to the limitation of my personal laptop configurations, haven't linked to any realtime reporting.



## SNAPSHOT 2:

**ETL Process Status:**

| S.No | UNIQ ID's | UNIQ COUNTRY | MAXIMUM LISTED COUNTRY | MINIMUM LISTED COUNTRY | MALE COUNTS | FEMALE COUNTS |
|------|-----------|--------------|------------------------|------------------------|-------------|---------------|
| 01   | 2         | 2            | Philippines            | Philippines            | 2           | 0             |

The code is designed so robust and versatile in extending to other applications/further integrations with new softwares/framework.Almost the code is kind of plug-play!

NOTE:

ETL process has been segregated so that complex logic can be handled separately and actual data streaming cycle remains un disturbed incase of time dependent streaming from Multiple data sources.

Data process happens in multiple micro batches to process the data efficiently using the best of the available resources/timekeeping in the fact that downstream receives good amount of data for their incremental data process.