# COVID-19 CLOUD DATA WAREHOUSE AND ANALYTICS

#### Team 1:

Kunal Jaiswal, Sharvari Karnik, Keval Vankudre

Project Dashboard URL: <a href="https://bit.ly/CovidDataIntegration">https://bit.ly/CovidDataIntegration</a>

#### Objectives:

Create a Data Warehouse by integrating COVID public datasets from different sources and design dashboards to get insights.

- Perform source systems analysis on data sources (& compare)
- Ingest & Integrate from/to various sources to generate BI Dashboards
- Gain experience with cloud databases & cloud DW
- Gain experience with data integration & data pipeline tools using data from both cloud and on
- premise (notebook) sources

#### **Deliverables:**

- Dashboards visualizing COVID data tracking US states and counties
  - Microsoft PowerBI, Tableau or Google Data Studio
- Data should be loaded into:
  - Google BigQuery
  - Azure SQL (note: likely limitations on student credits)
  - SQL Server or MySQL (as a fallback if team exhausts their free student credits)
- Review & compare data integration tools

#### **Tools Used:**

- Talend Data Quality Open Studio For Data Profiling
- Stitch Used to load COVID public dataset and flat files to Big Query
- Talend Pipeline Designer For creating staging area and for Data Integration
- Talend Big Data For Data Integration

#### **Cloud Technologies Used:**

- Google Big Query
- Azure SQL

# BI tool used:

Power BI

# List of COVID dataset sources used:

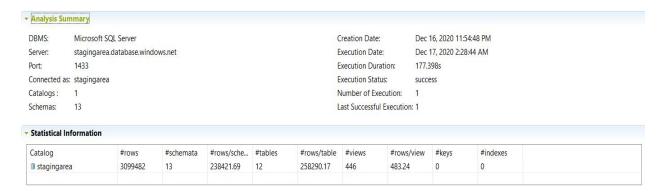
- Open Data By state and county
- New York Times Data By state and county
- The COVID Tracking Project By state
- Oxford Policy Tracker By State
- American Heart Association (AHA) By county
- Flat files (from the open source github) Demographics

# **Source System Analysis:**

Source system analysis has been performed in Talend Data Profiling. Below snippets give an overview about the type of data that we have on our hands.

# Connection Overview Analysis:

This analysis returns an overview of the content of your database. It computes the number of tables and the number of rows per table for each catalog and/or schema, etc. It also counts the number of indexes and primary keys, etc.



From the above analysis we get a gist of data that we have our hands.

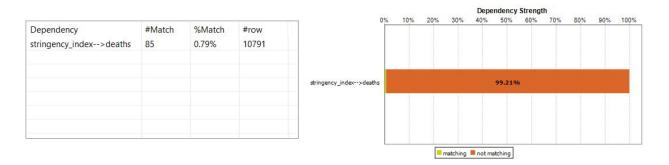
# Column Analysis:

This type helps to analyze and manually assign the indicators on each column, such as number of nulls, frequency table, summary statistics, pattern matching indicators, etc.



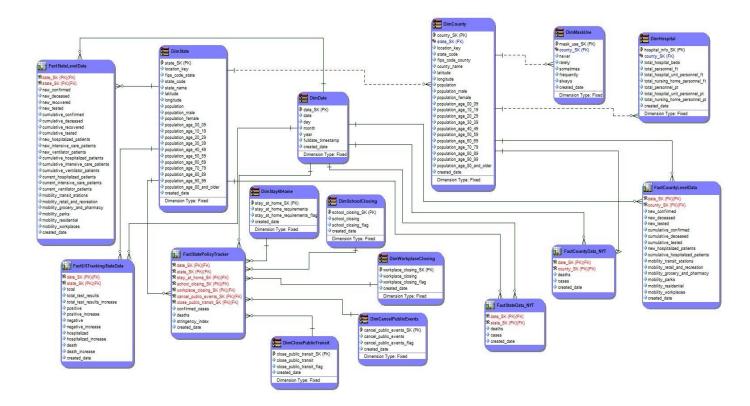
# Functional Dependency Analysis:

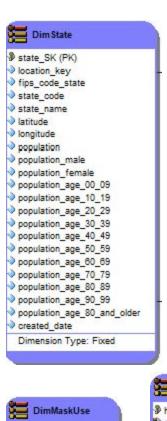
This analysis enables you to detect anomalies in your column dependencies. It determines to which extent the values of a column A determines the values of a column B. This relationship is noted as A->B.

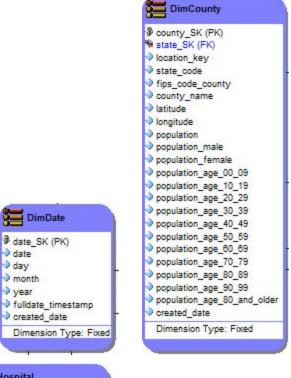


From this we can say that, stringency\_index determines only 0.79% of deaths. From this we can also infer that to determine the deaths more detailed analysis apart from stringency\_index is needed.

#### **Data Model**









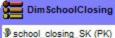


date

day

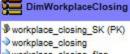
year





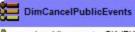
school\_closing school\_closing\_flag created date

Dimension Type: Fixed



workplace\_closing\_flag created date

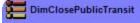
Dimension Type: Fixed



cancel\_public\_events\_SK (PK) cancel\_public\_events

cancel\_public\_events\_flag created date

Dimension Type: Fixed



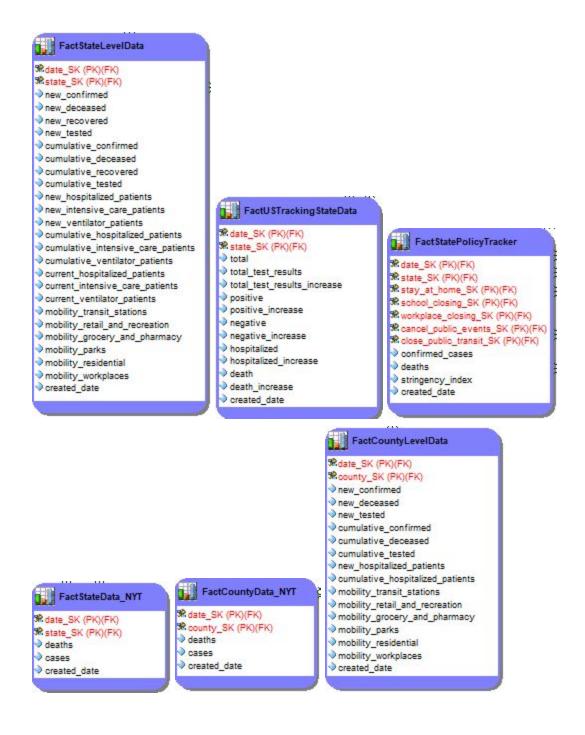
Close\_public\_transit\_SK (PK)

close public transit

close\_public\_transit\_flag

created date

Dimension Type: Fixed



# Review and compare data integration tools:

#### **Talend Big Data**

- Significantly easier to create data pipelines in Talend Big Data than Talend pipeline designer
- Easy to access database and create connections

- Can change and compare data types easily
- Can perform join operations easily
- Easy to understand and interpret the errors
- Takes a bit longer to process, integrate and ingest the data in the data warehouse than Talend Pipeline Designer

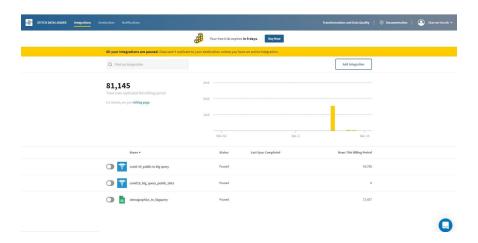
#### Stitch

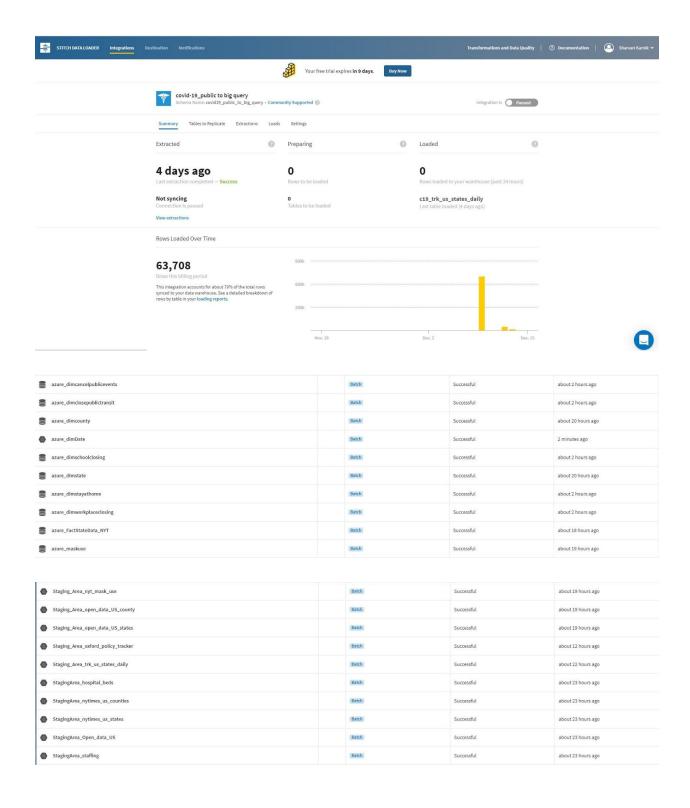
- Easy to use
- Comes with plethora of connectors for building connections with different sources which helped in seamless data ETL
- Fast and efficient

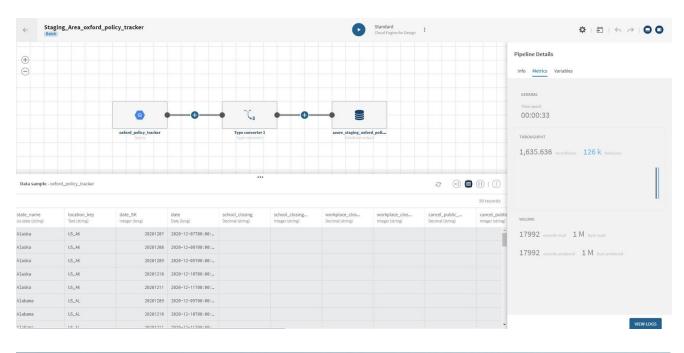
### **Talend Pipeline Designer**

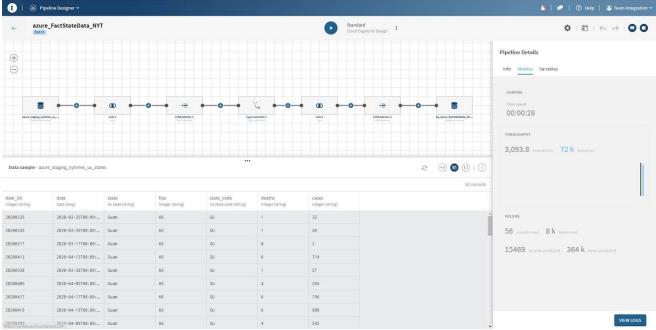
- Easy drag and drop ETL tool
- Can use different processor to play with the data
- Can preview the input data and output data after using a processor
- Takes considerably less time to process, integrate and ingest the data in the data warehouse than Talend Big Data
- Difficult to change the data types of the attributes and perform join operations with other source
- It would be better to access a particular data source if it had a proper dataset organization
- No truncation option before running a pipeline
- A bit difficult to understand the errors
- Might cause timeout errors if the data is too large to process

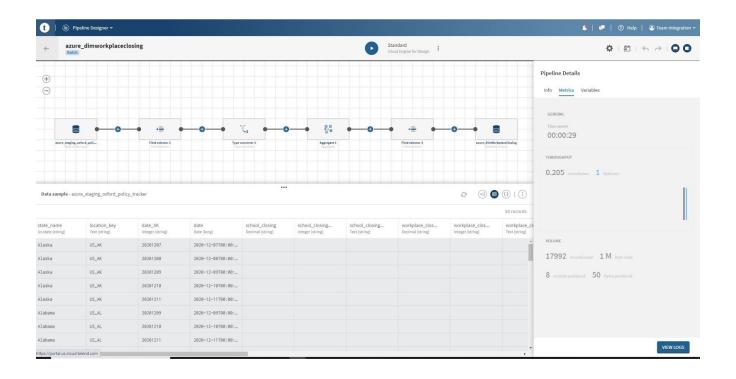
#### Screenshots:

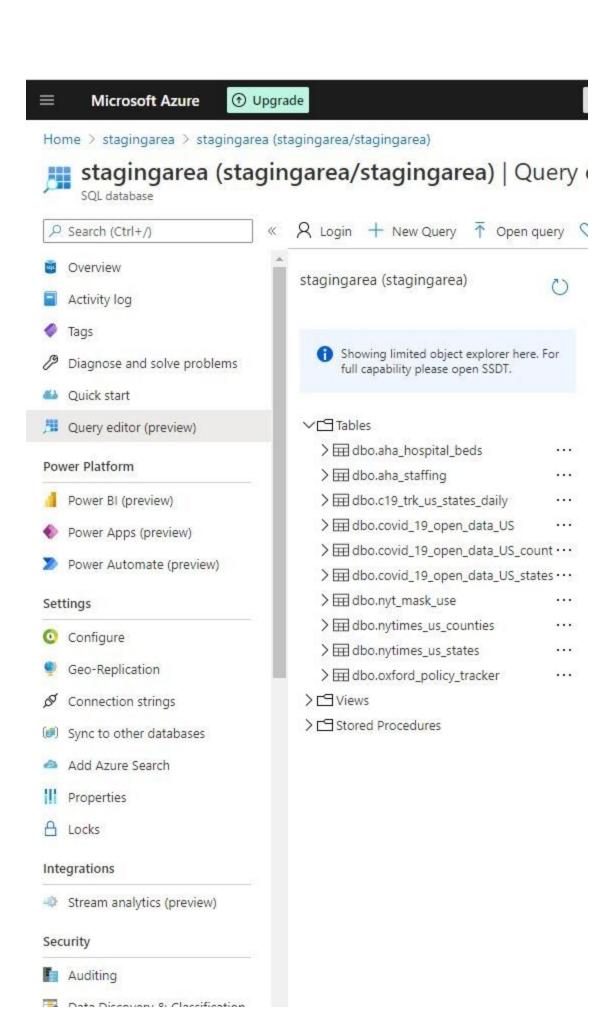


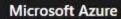








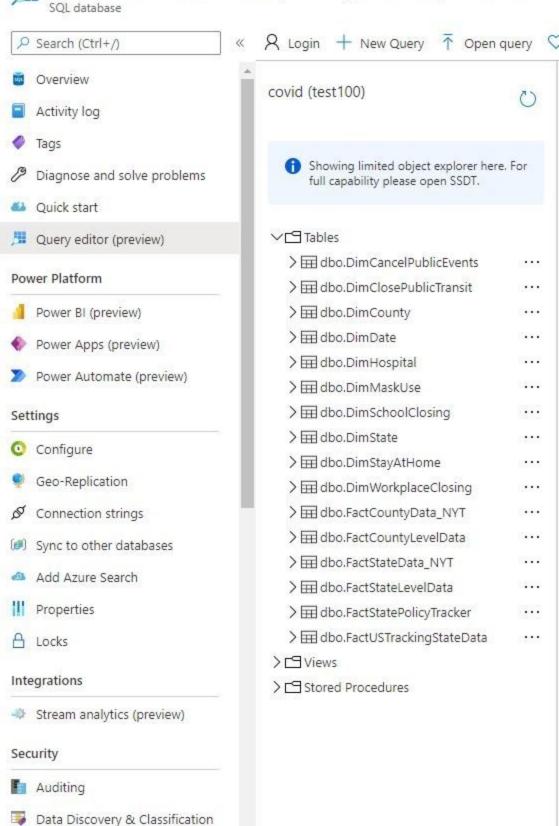






# Home > covid (test100/covid)

# covid (test100/covid) | Query editor (preview)



1 SELECT TOP (1000) \* FROM [dbo].[FactCountyLevelData]

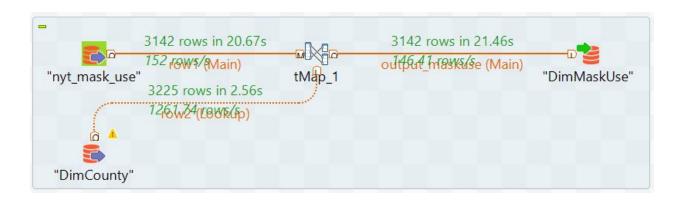
# Results Messages

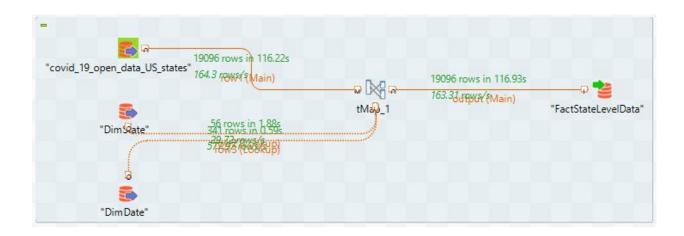
O Search to filter items								
date_SK	county_SK	new_confirmed	new_deceased	new_tested	cumulative_confirmed	cumulative_deceased	cumulative_tested	
20200208	1621	0	0	0	0	0	0	
20200208	1622	0	0	0	0	0	0	
20200208	1623	0	0	0	Ö	0	0	
20200208	1624	0	0	0	0	0	0	
20200208	1625	0	0	0	0	0	0	
20200208	1626	0	0	0	0	0	0	
20200208	1627	0	0	0	0	0	0	
20200208	1628	0	0	0	0	0	0	

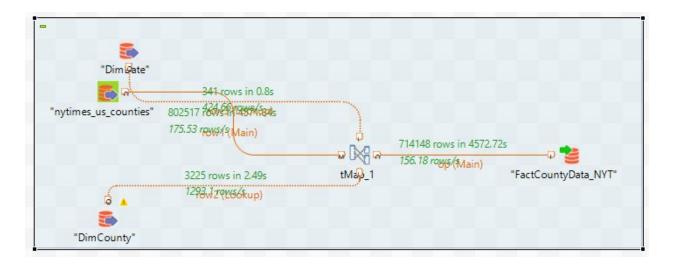
1 SELECT TOP (1000) \* FROM [dbo].[DimMaskUse]

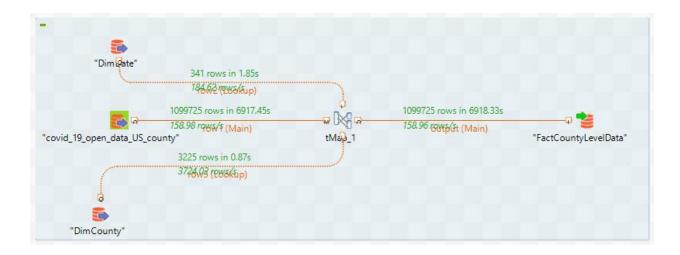
#### Results Messages

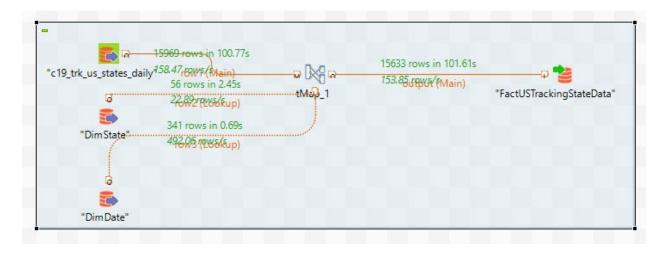
© Search to filter items									
mask_use_SK	county_SK	never	rarely	sometimes	frequently	always	created_date		
1	1	0.053	0.074	0.134	0.295	0.444	2020-12-16T21:04:48.913000		
2	2	0.083	0.059	0.098	0.323	0.436	2020-12-16T21;04:48.9130000		
3	3	0.067	0.121	0.12	0.201	0.491	2020-12-16T21:04:48.9130000		
4	4	0.02	0.034	0.096	0.278	0.572	2020-12-16T21:04:48.913000		
5	5	0.053	0.114	0.18	0.194	0.459	2020-12-16T21:04:48.9130000		
6	6	0.031	0.04	0.144	0.286	0.5	2020-12-16T21:04:48.913000		
7	7	0.102	0.053	0.257	0.137	0.451	2020-12-16T21:04:48.9130000		
8	8	0.152	0.108	0.13	0.167	0.442	2020-12-16T21:04:48.9130000		

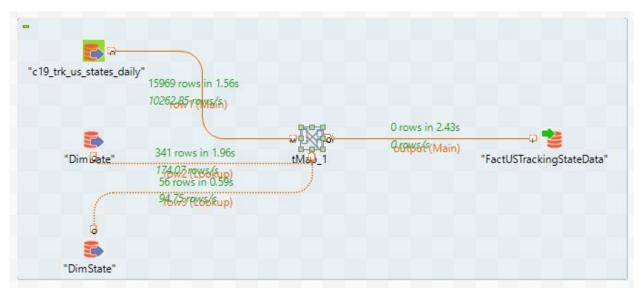


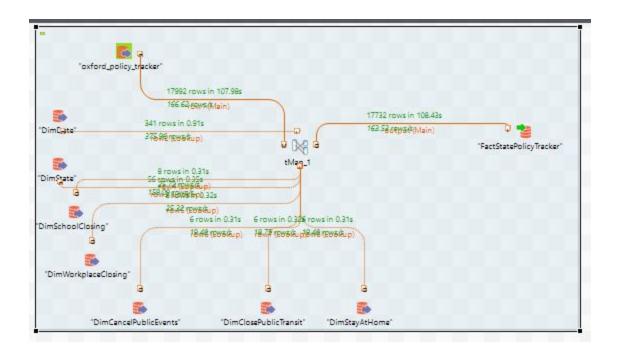












# PowerBI Dashboard:

URL: <a href="https://bit.ly/CovidDataIntegration">https://bit.ly/CovidDataIntegration</a>

