

ELEN 6889, Large-scale Streaming Processing

Project - Coronavirus Sub-hotspots Trend and News Fusion in United States

Professor Deepak S. Turaga

May 7th, 2020, presented By

Yanwen Jing (yj2556), Yanlin Liu (yl4238), Changxu Luo (cl3875), Weihan Chen (wc2681)

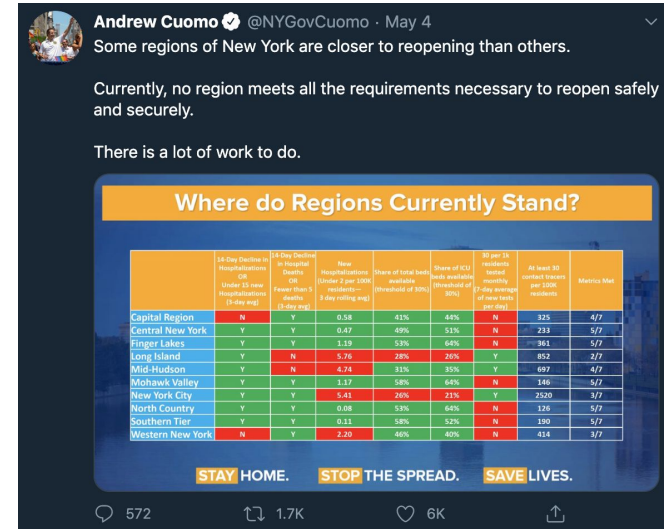
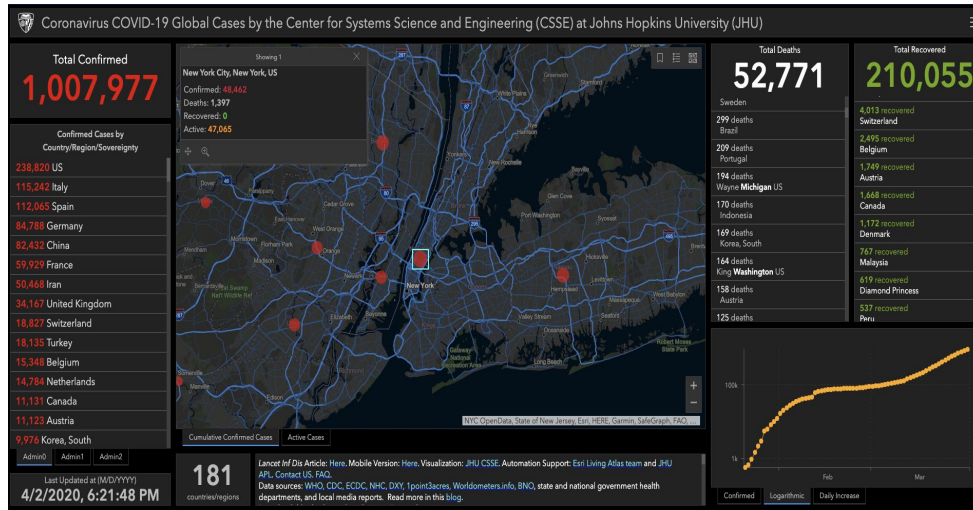


Problem

What we have now:

1. Coronavirus Cases Maps/Trending (<https://coronavirus.jhu.edu/map.html>)
2. Information from various News Organizations / US Government / Individual Users ... (twitter)

The map has only the numbers for the area, but no corresponding information/news (government actions, recommendations, progress, etc.)



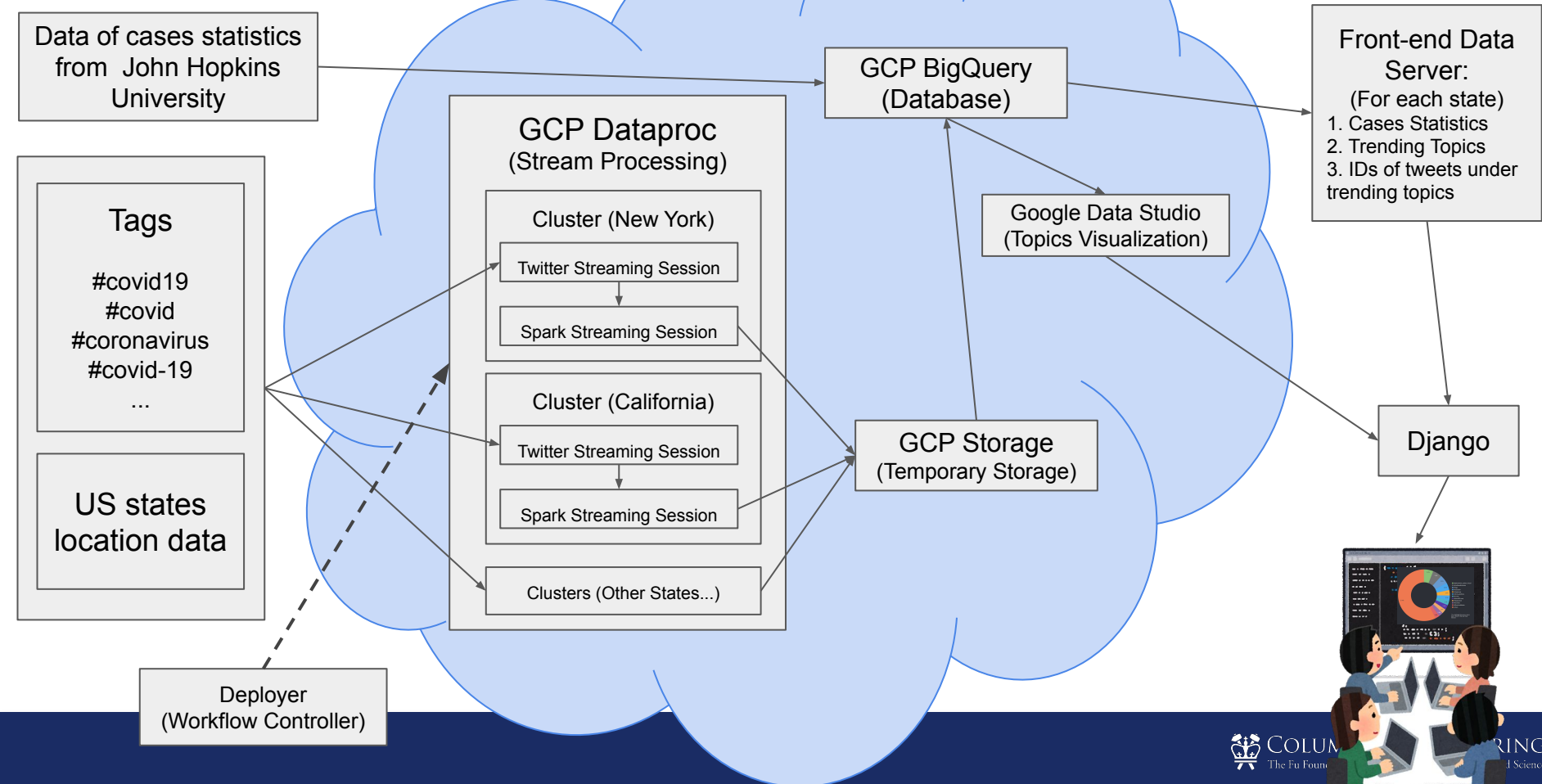
Problem

What we want now:

A US map with:


1. Daily updated coronavirus confirmed cases & deaths for each state
2. Plots about US confirmed cases & deaths trend over past few months
3. Corresponding information/news from various News Organizations/ US Government/ Individual Users ... (on twitter)
4. Geographically distribution for the sub-hotspots and trending of them by a streaming application.

Solution



Solution - confirmed cases & deaths

Dataset: on GCP BigQuery



COVID-19 Public Datasets

BigQuery Public Datasets Program

Datasets included in the COVID-19 public datasets program

[VIEW DATASET ↗](#)

**JHU Coronavirus COVID-19
Global Cases, by country**
Johns Hopkins University

Repository of aggregated
coronavirus COVID-19 cases by
JHU

Solution - confirmed cases & deaths

Confirmed cases & deaths
By state & date

```
self.query_string2 = """  
select state, date, sum(confirmed_cases) as confirmed_cases, sum(deaths) as deaths  
from bigquery-public-data.covid19_usafacts.summary  
group by state,date  
"""
```

```
def get_cases_deaths(self):  
    query_string = self.query_string2  
    dataframe = (  
        self.bqclient.query(query_string)  
        .result()  
        .to_dataframe()  
    )  
    self.cases_data = dataframe
```

Solution

```
def sendData(c_socket, tags):  
    """  
    send data to socket  
    """  
  
    testLocation = [-74, 40, -73, 41]  
    auth = OAuthHandler(CONSUMER_KEY, CONSUMER_SECRET)  
    auth.set_access_token(ACCESS_TOKEN, ACCESS_SECRET)  
    twitter_stream = Stream(auth, TweetsListener(c_socket)) # bind t.  
    twitter_stream.filter(track=tags,  
                          languages=['en'],  
                          # filter_level='medium'  
                          locations=testLocation  
                          )  
  
class twitter_client:  
    def __init__(self, TCP_IP, TCP_PORT):  
        self.s = s = socket.socket(socket.AF_INET, socket.SOCK_STREAM)  
        self.s.bind((TCP_IP, TCP_PORT))  
  
    def run_client(self, tags):  
        try:  
            self.s.listen(1)  
            while True:  
                print("Waiting for TCP connection...")  
                conn, addr = self.s.accept()  
                print("Connected... Starting getting tweets.")  
                sendData(conn, tags)  
                conn.close()  
        except KeyboardInterrupt:  
            exit
```

Waiting for TCP connection...

Connected... Starting getting tweets.

TEXT:RT @Thumbel28421928: Deathbed FaceTime calls
iPad visitors only
Nurses stifle sobs

#coronavirus #CoronaHaiku

#QuarantineHaiku #CovidHaiku...

TEXT:RT @Shem_Infinite: Don't forget, tomorrow at 12 PM EST
questions...

Twitter Streaming Session

Get the latest streaming data from twitter

Solution

Time: 2020-05-01 05:08:25

```
(' #covid19', 98)
(' #MayDay2020', 3)
(' #lockdown', 3)
(' #Odisha', 3)
(' #India', 2)
(' #bestgames', 2)
(' #BREAKING:', 2)
(' #Ohio', 2)
(' #caceroleatelachota', 2)
(' #HospitalPlaylist', 2)
...
```

The screenshot shows the Tableau interface. At the top, there is a search bar with the text "Search for your tables and datasets" and a magnifying glass icon on the left and a question mark icon on the right. Below the search bar, the left pane displays a list of tables under the heading "hashtag_dataset". The tables listed are: "hashtagsAK", "hashtagsAL", "hashtagsAR", "hashtagsAZ", "hashtagsCA" (which is highlighted with a grey background), "hashtagsCO", "hashtagsCT", "hashtagsDC", "hashtagsDE", "hashtagsFL", "hashtagsGA", "hashtagsHI", "hashtagsIA", and "hashtagsID". Each table name is preceded by a small icon representing a table grid.

Spark Streaming Session

Process the Data and Store them in Google Bigquery

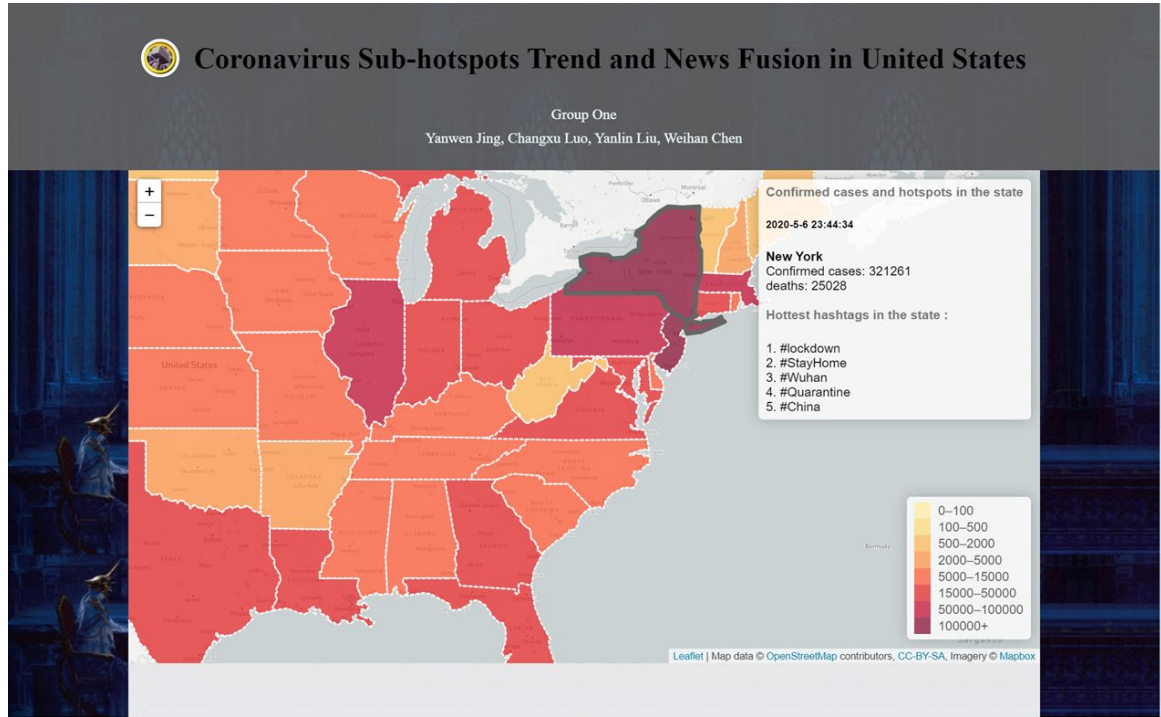
Solution

```
QUERY = (  
    'SELECT * FROM `hashtag_dataset.` + '%s'%i + '~`',  
    'ORDER BY count DESC ',  
    'LIMIT 5')  
query_job = client.query(QUERY) # API request  
rows = query_job.result() # Waits for query to finish
```

```
{'created_at': 'Thu May 07 03:31:25 +0000 2020', 'id': 1258237976255115264},  
{'created_at': 'Thu May 07 03:31:25 +0000 2020', 'id': 1258237975919562753},  
{'created_at': 'Thu May 07 03:31:24 +0000 2020', 'id': 1258237974556311553},  
{'created_at': 'Thu May 07 03:31:24 +0000 2020', 'id': 1258237973235187712},  
{'created_at': 'Thu May 07 03:31:23 +0000 2020', 'id': 1258237968847863808},  
{'created_at': 'Thu May 07 03:31:21 +0000 2020', 'id': 1258237960878731269},
```

Use the data stored in Google Bigquery to obtain the tweets and tweets_id about the corresponding topic.

Solution



Use Django to integrate all the data and finish the visualization part.

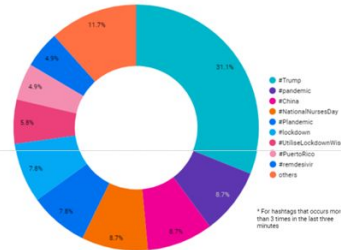
New York

COVID-19 Situation:

Confirmed cases: 321261

Deaths: 25028

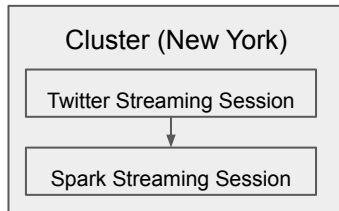
Hottest Topics in State:



Detail information of each state will be showed if click on the map. These information includes a more detailed map, summary of confirm cases, deaths and top 10 hottest twitter hashtags in the state, and some selected tweets in this area.



Streaming Part - Algorithm

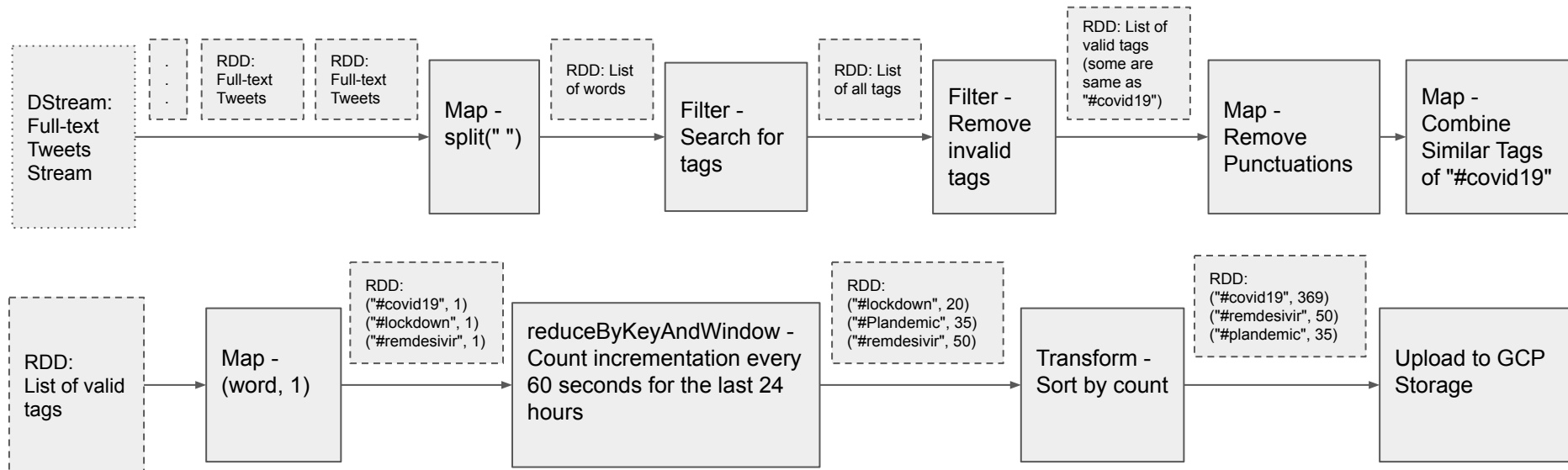


For each state, we maintain two cluster jobs on the cloud:

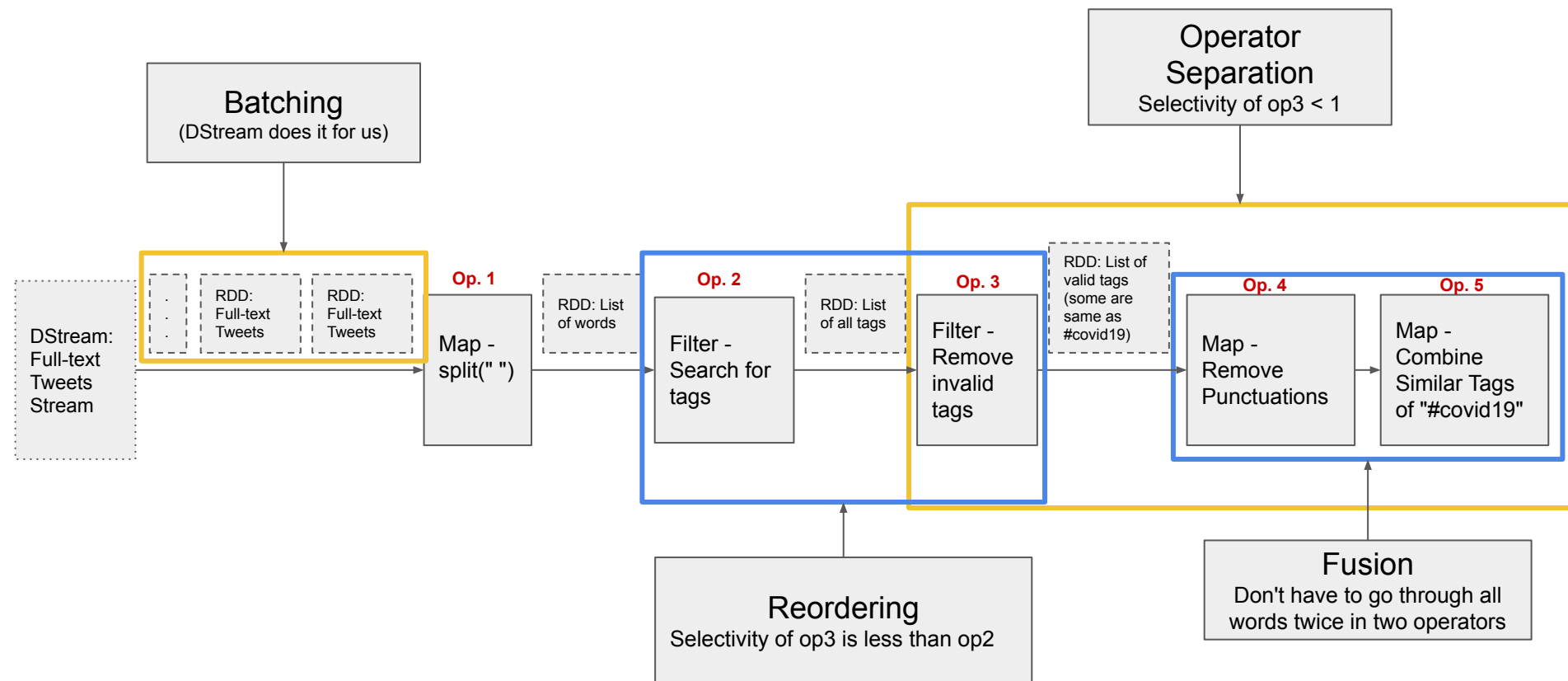
Job 1: getting real-time tweets, filter with states' location and tags pool, send to job 2

Job 2: analyze the tweets stream, get trending topics, send to GCP storage

Here we focus on details and optimization for job 2.



Streaming Part - Optimization





Demo / Results

TRANSCENDING DISCIPLINES, TRANSFORMING LIVES

What we want to improve in the future:

1. Apply semantic analysis / topic modeling to the trending sub-topics with the corresponding tweets, so that we can learn the topics in real-time and extend the presenting of topics-related information to the news literature without #hashtags
2. Plots about US confirmed cases & deaths trend over past few months
3. More detailed information like what are people focusing on in different counties
4. Optimization on the frontend, especially try to directly “connect” it to the stream



Thanks for listening

TRANSCENDING DISCIPLINES, TRANSFORMING LIVES