

## Big Data Application Architecture

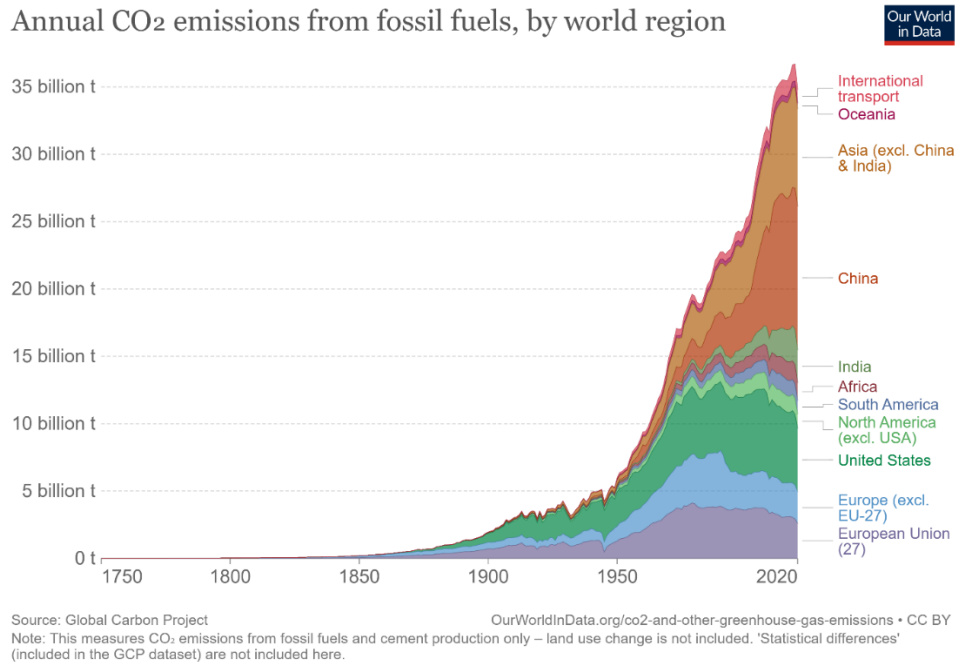
### Analyzing the impact of Greenhouse gas emissions on weather

Thanushri Rajmohan

#### Steps:

##### 1. Describing the datasets used:

Link to the greenhouse emissions dataset: [CO<sub>2</sub> and Greenhouse Gas Emissions - Our World in Data](#)



The greenhouse gas emissions in each country or region is represented in tonnes and the data is collected from 1971-2020 for each country. It has four columns – Country, three-letter country code, year, greenhouse gas emissions for that year. The header is removed from the dataset and uploaded to the Hadoop local file system using WinScp with “Hadoop” as the user and the private key putty file (.ppk). I had to choose WinScp because I am a windows user. Otherwise I could have tried uploading my data using the scp command from my local system to the Hadoop’s local file system.

FTP link to the weather dataset: [ftp://ftp.ncdc.noaa.gov/pub/data/ghcn/daily/by\\_year/](ftp://ftp.ncdc.noaa.gov/pub/data/ghcn/daily/by_year/)

The weather dataset that we used in the class was “Global summaries of the day” dataset for the states in the US from years 2000 to 2021. However, the dataset that I used is “Global summaries of the year” dataset for all countries/stations in the world from the years 1750 to 2021. Since the dataset is different and has different formats compared to the ones that we used in class (csv.gz whereas the ones used in class were tar zipped), I had to extract the weather dataset again from the ftp link using modified shell scripts ([getWeather.sh](#)) and put it on the local Hadoop file system.

2. Getting the weather data on HDFS:

The weather data is serialized and ingested into HDFS. If you unzip, “[weatherData\\_thanushrir\\_ingest.zip](#)”, you can find “[thanushrir](#)” folder which has the java code to serialize the inputs using Thrift. The uber-jar is generated in the “[/home/hadoop/thanushrir/src/target](#)” directory of Hadoop.

The yarn command to run the uber jar is: yarn jar <uber\_jar> <className> <local Hadoop directory> <HDFS output directory>

```
yarn jar uber-thanushrir-1.0-SNAPSHOT.jar  
edu.uchicago.mpcs53013.SerializeWeatherSummary  
/home/hadoop/thanushrir/project/weatherData_thanushrir /inputs/thanushrir_weatherData
```

Running the uber jar took 1.5 hours and the data was finally on HDFS in the “[/inputs/thanushrir\\_weatherData](#)” directory.

3. Getting the greenhouse emissions data on HDFS:

I uploaded the greenhouse emissions data in “[thanushrir/project/greenhouse\\_emissions/](#)” directory. You might not find the file there because I removed it due to storage issues on the local file system as mentioned by Prof. Mike. So, I used a hdfs -put command to put the file on HDFS.

```
hdfs dfs -put thanushrir_greenhouse_emissions.csv /inputs/greenhouse_emissions
```

4. Create tables in Hive:

Run the commands in “[load\\_weatherData.hql](#)” to load the weather data in hive using the beeline command.

From the next image, you can see that the table “[thanushrir\\_project\\_weathersummary](#)” has the following columns – Country Code, Station, Year, Month, Day, Value Type and Mean Temperature.

- Country code – FIPS two letter code
- Station – Station ID
- Value Type – TMAX for Maximum Temperature, TMIN for Minimum Temperature, PRCP for precipitation and TAVG for Average Temperature. In our case, we need only the average temperature, so we later filter out and take into consideration the TAVG values only.
- Mean Temperature – This column has the mean temperature of that station over the entire year.

```

0: jdbc:hive2://localhost:10000/default> SELECT * FROM thanushrir_project_WeatherSummary LIMIT 5;
INFO : Compiling command(queryId=hive_20221111223915_cce0eb1f-dd0a-424f-b408-fd45fb1570fa): SELECT
* FROM thanushrir_project_WeatherSummary LIMIT 5
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Semantic Analysis Completed (retrial = false)
INFO : Returning Hive schema: Schema(fieldSchemas:[FieldSchema(name:thanushrir_project_weathersummary.countrycode, type:string, comment:null), FieldSchema(name:thanushrir_project_weathersummary.station, type:string, comment:null), FieldSchema(name:thanushrir_project_weathersummary.year, type:smallint, comment:null), FieldSchema(name:thanushrir_project_weathersummary.month, type:tinyint, comment:null), FieldSchema(name:thanushrir_project_weathersummary.day, type:tinyint, comment:null), FieldSchema(name:thanushrir_project_weathersummary.valuetype, type:string, comment:null), FieldSchema(name:thanushrir_project_weathersummary.meantemperature, type:double, comment:null)], properties:null)
INFO : EXPLAIN output for queryid hive_20221111223915_cce0eb1f-dd0a-424f-b408-fd45fb1570fa : STAGE
DEPENDENCIES:
  Stage-0 is a root stage [FETCH]

STAGE PLANS:
  Stage: Stage-0
    Fetch Operator
      limit: 5
    Processor Tree:
      TableScan
        alias: thanushrir_project_weathersummary
        GatherStats: false
      Select Operator
        expressions: countrycode (type: string), station (type: string), year (type: smallint), month (type: tinyint), day (type: tinyint), valuetype (type: string), meantemperature (type: double)
        outputColumnNames: _col0, _col1, _col2, _col3, _col4, _col5, _col6
        Limit
          Number of rows: 5
          ListSink

INFO : Completed compiling command(queryId=hive_20221111223915_cce0eb1f-dd0a-424f-b408-fd45fb1570fa); Time taken: 0.027 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20221111223915_cce0eb1f-dd0a-424f-b408-fd45fb1570fa): SELECT
* FROM thanushrir_project_WeatherSummary LIMIT 5
INFO : Completed executing command(queryId=hive_20221111223915_cce0eb1f-dd0a-424f-b408-fd45fb1570fa); Time taken: 0.001 seconds
INFO : OK
INFO : Concurrency mode is disabled, not creating a lock manager
+-----+-----+-----+-----+-----+-----+-----+
| thanushrir_project_weathersummary.countrycode | thanushrir_project_weathersummary.station | thanushrir_project_weathersummary.year | thanushrir_project_weathersummary.month | thanushrir_project_weathersummary.day | thanushrir_project_weathersummary.valuetype | thanushrir_project_weathersummary.meantemperature |
+-----+-----+-----+-----+-----+-----+-----+
| AG | | AG000060390 | | 1971 |
| | TMAX | | 96.0 |
| AG | | AG000060390 | | 1971 |
| | TMIN | | 51.0 |
| AG | | AG000060390 | | 1971 |
| | PRCP | | 442.0 |
| AG | | AG000060590 | | 1971 |
| | TMAX | | 193.0 |
| AG | | AG000060590 | | 1971 |
| | TMIN | | 7.0 |
+-----+-----+-----+-----+-----+-----+-----+
5 rows selected (0.251 seconds)

```

Run the commands in “[load\\_greenhouse\\_emissions.hql](#)” to load greenhouse emissions data in hive using the beeline command. The table “[thanushrir\\_greenhouse\\_emissions](#)” stored as ORC form has the greenhouse emissions data.

5. Adding a separate country table to Hive:

The weather data has only the FIPS code with station ID. We need to join the weather data to the greenhouse emissions data using the FIPS code from weather data and the country name in the greenhouse emissions data. We get rid of the three letter country code in the greenhouse emissions data as it is not necessary. “[ghcnm-countries.csv](#)” is the csv file that maps FIPS country code to the country names. This is loaded to HDFS in the same way by loading to local HDFS using WinScp and using hdfs put command to put it on HDFS. This file is in the “[/inputs/thanushrir\\_country](#)” on HDFS. Run the commands in “[load\\_countryData.hql](#)” to create a “[thanushrir\\_country\\_data](#)” table in Hive.

6. Joining Hive tables to create one large final table to store it in HBase:

Join the greenhouse emissions data with the countries data based on the FIPS code using hive commands in “[join\\_greenhouse\\_country.hql](#)”. The table created in Hive is “[thanushrir\\_project\\_greenhouse](#)”.

Now, the weather data is grouped by FIPS country code and year by running the commands in “[group\\_weatherData.hql](#)”.

Now this resultant table is in turn joined with the “[thanushrir\\_project\\_greenhouse](#)” hive table based on the FIPS code and the year. The commands are present in “[join\\_weather\\_greenhouse.hql](#)”. The Hive table that is created in this step “[thanushrir\\_project\\_weather\\_greenhouse](#)” is the final Hive table which we will be using for processing everywhere.

7. Generating HBase table for this Hive table:

Create table “[thanushrir\\_project\\_hbase\\_HiveStyle](#)” in hbase shell with column key as “values”.

```
hbase:001:0> create 'thanushrir_project_hbase_HiveStyle','values';
```

```
hbase:002:0> scan 'thanushrir_project_hbase_HiveStyle';
```

```
ROW          COLUMN+CELL
```

```
0 row(s)
```

```
Took 0.0881 seconds
```

Now, run the commands in “[hbase\\_storage\\_hivestyle.hql](#)” in Hive to store the hive table.

Run the following commands in HBase to see if the tables values are successfully loaded.

```
hbase:005:0> get 'thanushrir_project_hbase_HiveStyle','US2006';
```

```
COLUMN
```

```
CELL
```

```
values:avg_temp      timestamp=2022-11-18T00:12:45.585, value=78.70165
```

```
values:country       timestamp=2022-11-18T00:12:45.585, value=United States
```

```
values:emissions     timestamp=2022-11-18T00:12:45.585, value=6057163300
```

```
1 row(s)
```

```
Took 0.0043 seconds
```

```
hbase:007:0> get 'thanushrir_project_hbase_HiveStyle','UK2000';
```

COLUMN	CELL
values:avg_temp	timestamp=2022-11-18T00:12:45.569, value=77.47375
values:country	timestamp=2022-11-18T00:12:45.569, value=United Kingdom
values:emissions	timestamp=2022-11-18T00:12:45.569, value=569033660

1 row(s)  
Took 0.0037 seconds

#### 8. Generating Hive style HBase tables from Spark Scala:

Execute all the commands in “[Spark\\_HiveStyle\\_tables.scala](#)”. Create a HBase table using hbase shell with the following command:

```
hbase:001:0> create 'thanushrir_project_spark_hivestyle','values';
```

Created table thanushrir\_project\_spark\_hivestyle

Took 0.9098 seconds

Export the Spark Hive Style table to HBase by running the commands in

“[Spark\\_export\\_HiveStyle\\_table\\_to\\_hbase.hql](#)” in hive using beeline command. If you type “get 'thanushrir\_project\_spark\_hivestyle','US2006';” in hbase shell, it should return the values stored in the hbase tables.

#### 9. Generating Object style HBase tables from Spark Scala:

Execute all the commands in “[Spark\\_ObjectStyle\\_tables.scala](#)”. Create a HBase table using hbase shell with the following command:

```
hbase:001:0> create 'thanushrir_project_spark_objectstyle','values';
```

Created table thanushrir\_project\_spark\_objectstyle

Took 0.9098 seconds

Export the Spark Object Style table to HBase by running the commands in

“[Spark\\_export\\_ObjectStyle\\_table\\_to\\_hbase.hql](#)” in hive using beeline command. If you type “get 'thanushrir\_project\_spark\_objectstyle','US2006';” in hbase shell, it should return the values stored in the hbase tables.

#### 10. Running the applications on local machines:

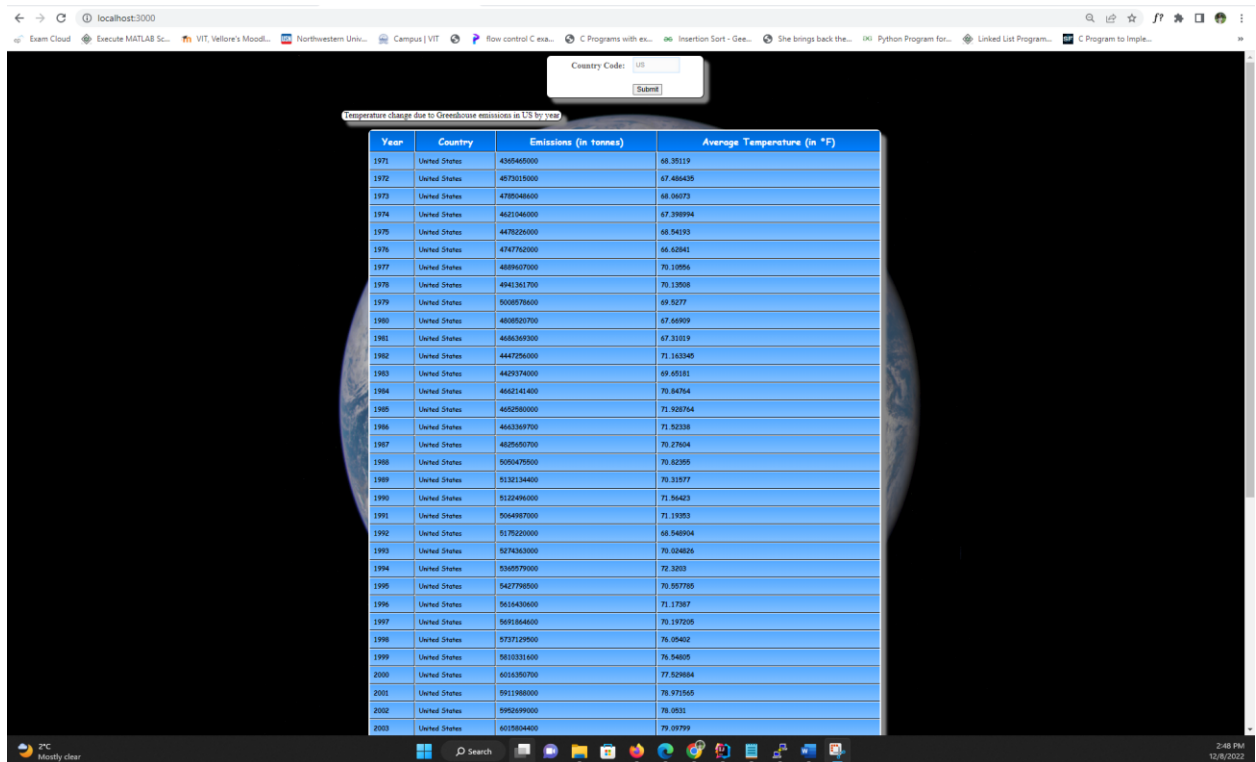
There are three different ways that these apps use the tables.

- “[Thanushri\\_Rajmohan\\_Project\\_HBase\\_Hive](#)” uses the HBase table “thanushrir\_project\_hbase\_HiveStyle” stored initially using plain hive commands.
- “[Thanushri\\_Rajmohan\\_Project\\_Spark\\_HiveStyleTable](#)” uses the spark hive style HBase table “”
- “[Thanushri\\_Rajmohan\\_Project\\_Spark\\_ObjectStyleTable](#)” uses the spark object style HBase table “”

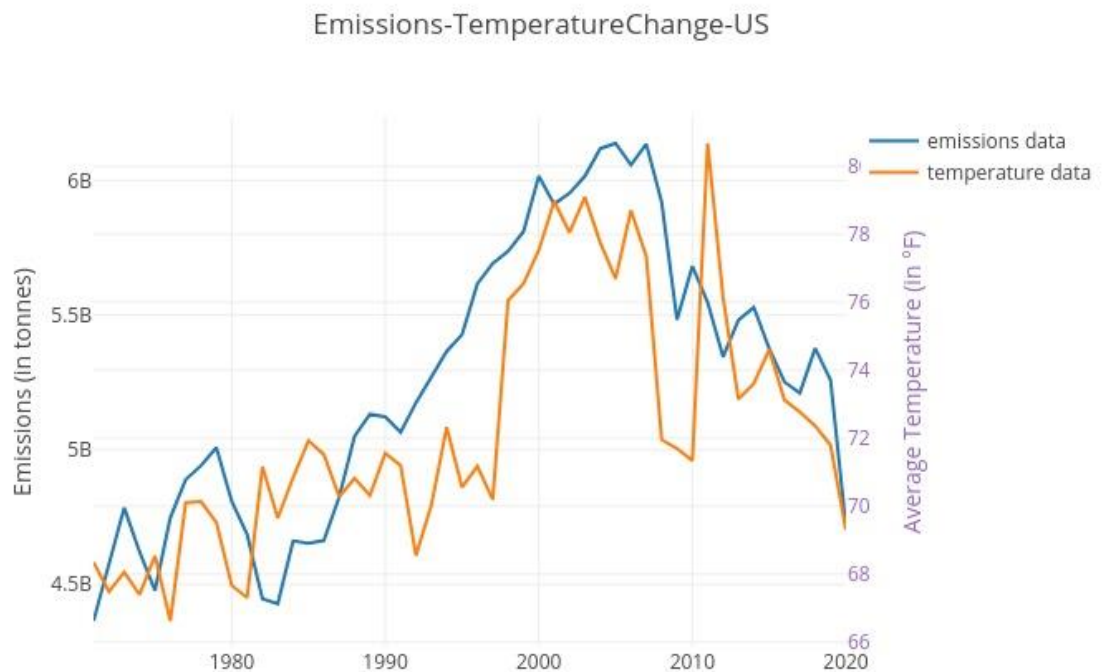
All these applications have the same code. Just the tables that have been used in each of the applications are different.

One additional functionality that I incorporated in this was to autogenerate a graph for each country code that is submitted to view the greenhouse gas emissions and the weather changes in a single graph as a better form of representation. I used the plotly library which can be installed using `npm install plotly`. Plotly requires a api token which I generated and the api key along with the code to autogenerate the graph using plotly is in app.js file.

Once each of the zip files of codes are unzipped, each of the application with the app.js code must be built and the node.js configurations called “RunApp” have to be set up to run it locally on port 3000 with “3000 localhost 8070”. Then run the project with localhost:3000/ and submit FIPS code of any country (UK for United Kingdom, US for United States, CA for Canada, IN for India, etc.) and you should get your results in a table format below. In addition to that, if you go to your IntelliJ console, you should see a json object with an url to plotly (something similar to this - <https://chart-studio.plotly.com/~thanushrir/8> ). Click on this url and then you can see the graph that is generated for the country code that you entered as input. You can save it for use at a later point of time.



```
{
  streamstatus: undefined,
  url: 'https://chart-studio.plotly.com/~thanushrir/8',
  message: '',
  warning: '',
  filename: 'Emissions-TemperatureChange-US',
  error: ''
}
```



#### 11. Web app which sends real-time data to kafka consumer:

I first created a kafka topic.

```
kafka-topics.sh --create --zookeeper z-2.mpcs530142022.7vr20l.c19.kafka.us-east-1.amazonaws.com:2181,z-1.mpcs530142022.7vr20l.c19.kafka.us-east-1.amazonaws.com:2181,z-3.mpcs530142022.7vr20l.c19.kafka.us-east-1.amazonaws.com:2181 --replication-factor 1 --partitions 1 --topic thanushrir_project_greenhouseWeather
```



When you unzip

“[Thanushri\\_Rajmohan\\_Project\\_app\\_with\\_weather\\_\\_greenhouse\\_form.zip](#)”, you should find the project that takes as input weather and greenhouse emissions data for a year from the user using the “[submit-weather-emissions.html](#)” form. The app.js file takes this input and puts it in the kafka consumer.

View the kafka-consumer using the following command.

```
kafka-console-consumer.sh --bootstrap-server b-2.mpcs530142022.7vr20l.c19.kafka.us-east-1.amazonaws.com:9092,b-3.mpcs530142022.7vr20l.c19.kafka.us-east-1.amazonaws.com:9092,b-1.mpcs530142022.7vr20l.c19.kafka.us-east-1.amazonaws.com:9092 --topic thanushrir_project_greenhouseWeather --from-beginning
```

Every app from now on uses “[thanushrir\\_project\\_spark\\_hivestyle](#)” hbase table to run. To run this project, once you unzip, set SFTP with name Webserver to ec2-server in the deployment options. Exclude node modules and set mappings deployment path to “[/thanushrir/weather\\_greenhouse\\_form\\_app](#)”. Change the RunApp configuration of the project to “3000 localhost 8070 localhost:9092”. Upload your entire project directly to “[/home/ec2-user/thanushrir/weather\\_greenhouse\\_form\\_app](#)” by right clicking on the project directory -> Deployment -> Upload to Webserver. Go to the above directory in the ssh shell and run the following. (3004 is my assigned port)

```
npm install
```

```
npm install kafka-node --no-optional
```

```
node app.js 3004 ec2-54-166-56-39.compute-1.amazonaws.com 8070 b-2.mpcs530142022.7vr20l.c19.kafka.us-east-1.amazonaws.com:9092,b-3.mpcs530142022.7vr20l.c19.kafka.us-east-1.amazonaws.com:9092,b-1.mpcs530142022.7vr20l.c19.kafka.us-east-1.amazonaws.com:9092
```

Because port 3004 is occupied by the load balanced servers currently, running this command might give an error.

Instead run this command with a random port 3060 to check if the app is working. If this also doesn't work and gives port in use error, try another port below 3100.

```
node app.js 3060 ec2-54-166-56-39.compute-1.amazonaws.com 8070 b-2.mpcs530142022.7vr20l.c19.kafka.us-east-1.amazonaws.com:9092,b-3.mpcs530142022.7vr20l.c19.kafka.us-east-1.amazonaws.com:9092,b-1.mpcs530142022.7vr20l.c19.kafka.us-east-1.amazonaws.com:9092
```

Now if you go to the following link, you should see your app running.

<http://ec2-52-73-126-153.compute-1.amazonaws.com:3060/>



<http://ec2-52-73-126-153.compute-1.amazonaws.com:3060/submit-weather-emissions.html>

Type input and enter Submit. Once you do that, your kafka consumer should have your input values as a JSON object report like below. You can see Canada is in the Kafka Consumer as the last value.

But if you go to your home page and try to view the newly entered value for that country, the new input won't be there. This is because the values from kafka consumer are not updated in the HBase table. So the next step is going to be how to put the values in kafka consumer into HBase table.

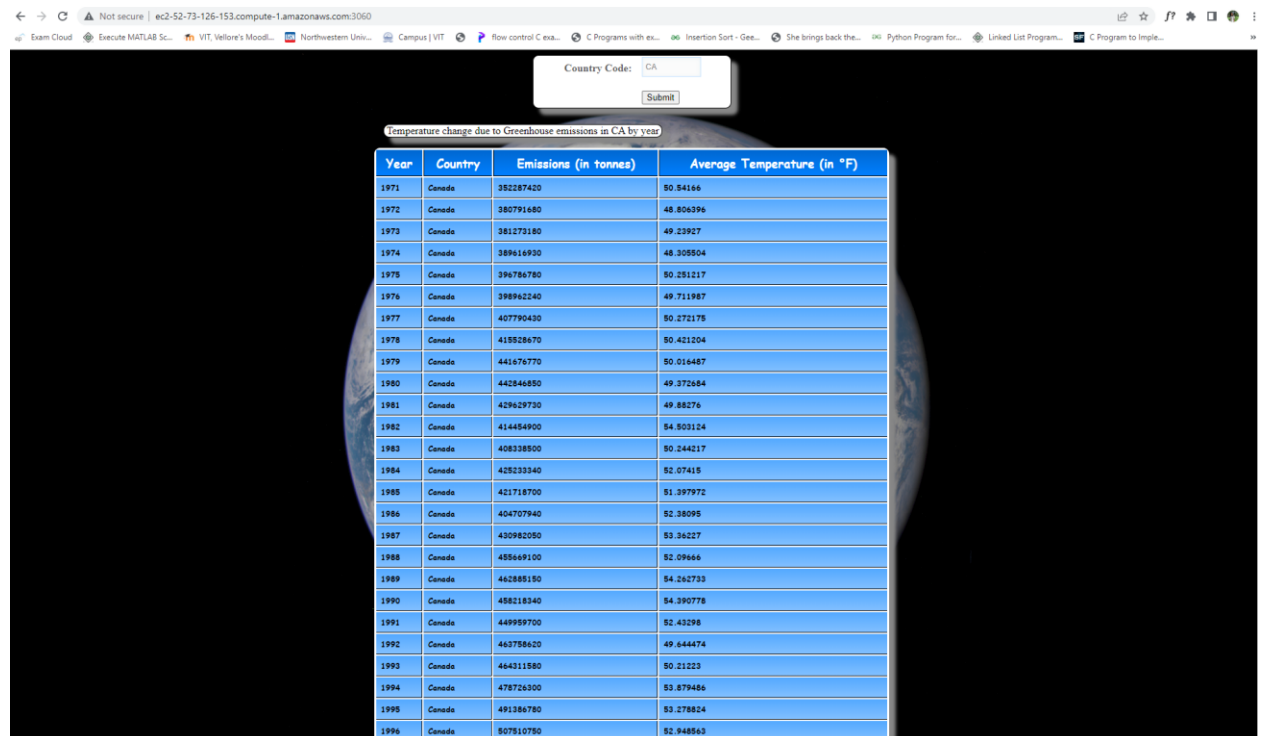
My streaming application writes to the HBase table directly from the speed layer data. This was the best way to insert data into my application that takes as input float, string values. This is because HBase Increment is not the best way because increment is

allowed only for "long" datatype values and I have string and float values for my application. So I used the HBase put command to directly update my HBase table.

Unzip “[Thanushri\\_Rajmohan\\_Project\\_app\\_speed\\_layer.zip](#)” and build the project. Your uber jar files should be generated (in my case its already there in the Hadoop file system). Cd to “[/home/hadoop/thanushrir/src/target](#)” after starting a Hadoop user ssh session. Run the uber jar using the following spark-submit command.

```
spark-submit --master local[2] --driver-java-options "-  
Dlog4j.configuration=file:///home/hadoop/ss.log4j.properties" --class StreamWeather  
uber-Thanushri_Rajmohan_Project_app_speed_layer-1.0-SNAPSHOT.jar b-  
2.mpcs530142022.7vr20l.c19.kafka.us-east-1.amazonaws.com:9092,b-  
3.mpcs530142022.7vr20l.c19.kafka.us-east-1.amazonaws.com:9092,b-  
1.mpcs530142022.7vr20l.c19.kafka.us-east-1.amazonaws.com:9092
```

Keep the application in the 11<sup>th</sup> step running and go to the same link to submit a new value. Once you submit the result and go to the home page and enter “CA” as input, you should see the new value for 2021 in your result table.



1993	Canada	464311880	80.21223
1994	Canada	478726300	83.879486
1995	Canada	491286780	83.278824
1996	Canada	507810780	82.948563
1997	Canada	521839970	84.737354
1998	Canada	529649540	82.54774
1999	Canada	544063800	82.203068
2000	Canada	566690400	80.215893
2001	Canada	588777840	84.992302
2002	Canada	564050900	81.194965
2003	Canada	581308300	85.01588
2004	Canada	579592040	83.550377
2005	Canada	574653630	83.705946
2006	Canada	568450940	84.56187
2007	Canada	593515800	86.47071
2008	Canada	576586340	82.395956
2009	Canada	543967600	81.812172
2010	Canada	558560400	80.263412
2011	Canada	567054140	83.860394
2012	Canada	568223170	81.078197
2013	Canada	572613300	84.475006
2014	Canada	569839700	84.94489
2015	Canada	574298240	76.89507
2016	Canada	569525300	77.67343
2017	Canada	571544640	80.074066
2018	Canada	584369180	80.67596
2019	Canada	564714200	83.75817
2020	Canada	834863840	85.84958
2021	Canada	812415271	82.12497

```
[hadoop@ip-172-31-95-18 project]$ kafka-console-consumer.sh --bootstrap-server b-2.mpc530142022.7vr201.ci9.kafka.us-east-1.amazonaws.com:9092,b-3.mpc530142022.7vr201.ci9.kafka.us-east-1.amazonaws.com:9092,b-1.mpc530142022.7vr201.ci9.kafka.us-east-1.amazonaws.com:9092 --topic thanushrir_project_greenhouseWeather --from-beginning
{"country_code":"UK","country":"United Kingdom","year":"2021","emissions":"341500000","avg_temp":"73.7"}
{"country_code":"UK","country":"United Kingdom","year":"2021","emissions":"341500000","avg_temp":"73.7"}
{"country_code":"UK","country":"United Kingdom","year":"2021","emissions":"341500000","avg_temp":"73.7"}
{"country_code":"UK","country":"United Kingdom","year":"2021","emissions":"341500000","avg_temp":"73.4"}
{"country_code":"UK","country":"United Kingdom","year":"2021","emissions":"341500000","avg_temp":"73.7"}
{"country_code":"UK","country":"United Kingdom","year":"2021","emissions":"341500000","avg_temp":"73.7"}
{"country_code":"UK","country":"United Kingdom","year":"2021","emissions":"341500000","avg_temp":"0"}
{"country_code":"UK","country":"United Kingdom","year":"2021","emissions":"341500000","avg_temp":"73.7"}
{"country_code":"UK","country":"United Kingdom","year":"2021","emissions":"326","avg_temp":"75.78"}
{"country_code":"UK","country":"United Kingdom","year":"2021","emissions":"341500000","avg_temp":"73.702411"}
{"country_code":"UK","country":"United Kingdom","year":"2022","emissions":"341500000","avg_temp":"73.789124"}
{"country_code":"UK","country":"United Kingdom","year":"2021","emissions":"341500000","avg_temp":"73.782"}
{"country_code":"UK","country":"United Kingdom","year":"2022","emissions":"341500000","avg_temp":"73.789124"}
{"country_code":"US","country":"United States","year":"2021","emissions":"142875141","avg_temp":"78.9214124"}
{"country_code":"UK","country":"United Kingdom","year":"2020","emissions":"341500000","avg_temp":"73.561"}
{"country_code":"UK","country":"United Kingdom","year":"2021","emissions":"341500000","avg_temp":"73.7"}
{"country_code":"UK","country":"United Kingdom","year":"2020","emissions":"341500000","avg_temp":"73.75"}
{"country_code":"UK","country":"United Kingdom","year":"2022","emissions":"356812111","avg_temp":"74.89124"}
{"country_code":"UK","country":"United Kingdom","year":"2021","emissions":"345671102","avg_temp":"74.2"}
{"country_code":"CA","country":"Canada","year":"2021","emissions":"812415271","avg_temp":"82.12497"}
{"country_code":"CA","country":"Canada","year":"2021","emissions":"812415271","avg_temp":"82.12497"}
```

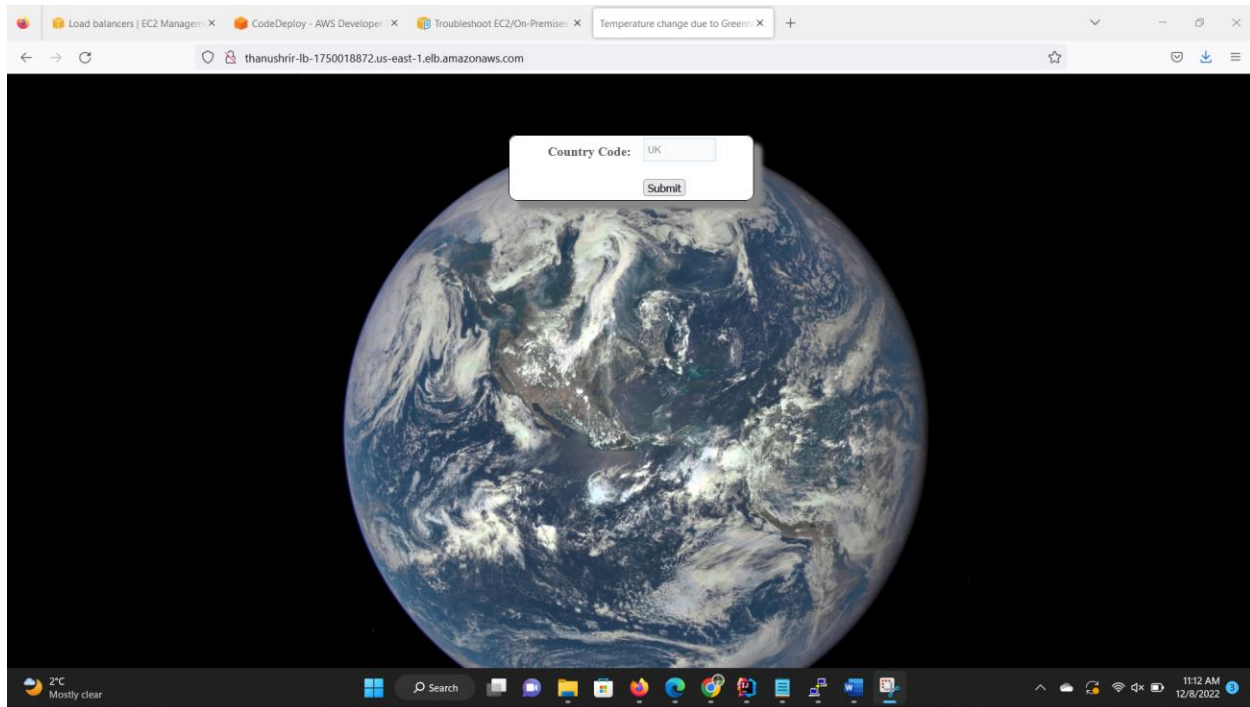
You can see the new Canada value that you entered in step 12 in your Kafka Console Consumer.

### 13. Deploying the application:

The project is deployed to CodeDeploy and the following link can be used to run the CodeDeploy application.

<http://thanushrir-lb-1750018872.us-east-1.elb.amazonaws.com/>

When you enter an input in <http://thanushrir-lb-1750018872.us-east-1.elb.amazonaws.com/submit-weather-emissions.html>, the application sends the values to the kafka consumer. Kafka consumer gets messages from the code deploy application but these new values are not updated in hbase because streaming app is not uploaded on codedeploy. Following are some screenshots of the app using loadbalancer DNS name.



Country Code:

Temperature change due to Greenhouse emissions in UK by year

Year	Country	Emissions (in tonnes)	Average Temperature (in °F)
1971	United Kingdom	660388200	69.250946
1972	United Kingdom	648026300	65.582596
1973	United Kingdom	659577100	69.48291
1974	United Kingdom	617183600	69.54029
1975	United Kingdom	603247040	70.76023
1976	United Kingdom	598526000	71.59135
1977	United Kingdom	604361700	68.98773
1978	United Kingdom	604713500	69.78485
1979	United Kingdom	644312800	67.202965
1980	United Kingdom	579035400	69.578156
1981	United Kingdom	560554500	69.89161
1982	United Kingdom	548240400	73.897865
1983	United Kingdom	545485100	73.39319
1984	United Kingdom	529108160	71.870766
1985	United Kingdom	559627300	68.94406
1986	United Kingdom	568554000	67.70338
1987	United Kingdom	571667600	69.36798
1988	United Kingdom	570293760	72.6376

2°C  
Mostly clear

11:13 AM  
12/8/2022