# Assignment 1
## Research Topic in AI

**Student Name(s)**      : Saurabh Shashikant Hebbalkar, Sampritha Hassan Manjunath
**Student IS(s)**        : 19232697, 19232922


## Introduction:

We have designed and implemented a Real-Time Feedback system for Taxi Drivers in New York City. The system streams the data, processes it and provides information to the Taxi Drivers regarding Peak Time and Busy Locations in the city. It also provides information regarding the accidents and generates a daily report of Peak Hour. We developed the application in Java and leveraged the Python Library (Pandas) for pre-processing the data.


## Pipeline Design:

Python is widely known for its time efficiency and Java is known for its structure and object orientation. We decided to use the best of both the worlds in designing our Microservice Architecture. There are approximately 87 Lakh unsorted records in the Yellow Taxi csv file, sorting the file manually in Excel causes data loss as Excel cannot handle large dataset (more than 1048576 rows). We first tried to sort the file in Java, but it was time consuming; therefore, we decided to use Python to sort the CSV file as it was more time efficient.

Once the file is sorted and saved, we stream the data from this file using Java Stream. The data is streamed sequentially. To generate the desired report, we need only the Pickup Timestamp and the Location ID; thus, we filter these two columns and ignore the other columns from the csv file. There were irrelevant records present in the csv file dating back to the year 2002; we also filtered them out and just fetched January 2018 data. This data is Streamed by the Publisher service.

We are using ActiveMQ, an open source message broker to communicate between the micro services. ActiveMQ connection is initiated at the Publisher service. As we have multiple micro services who rely on this processed data; our architecture implements Publish & Subscribe semantics. We initialize Topic at the Publisher and when the data is published, it is received by all the Subscribers who have subscribed to the Topic.
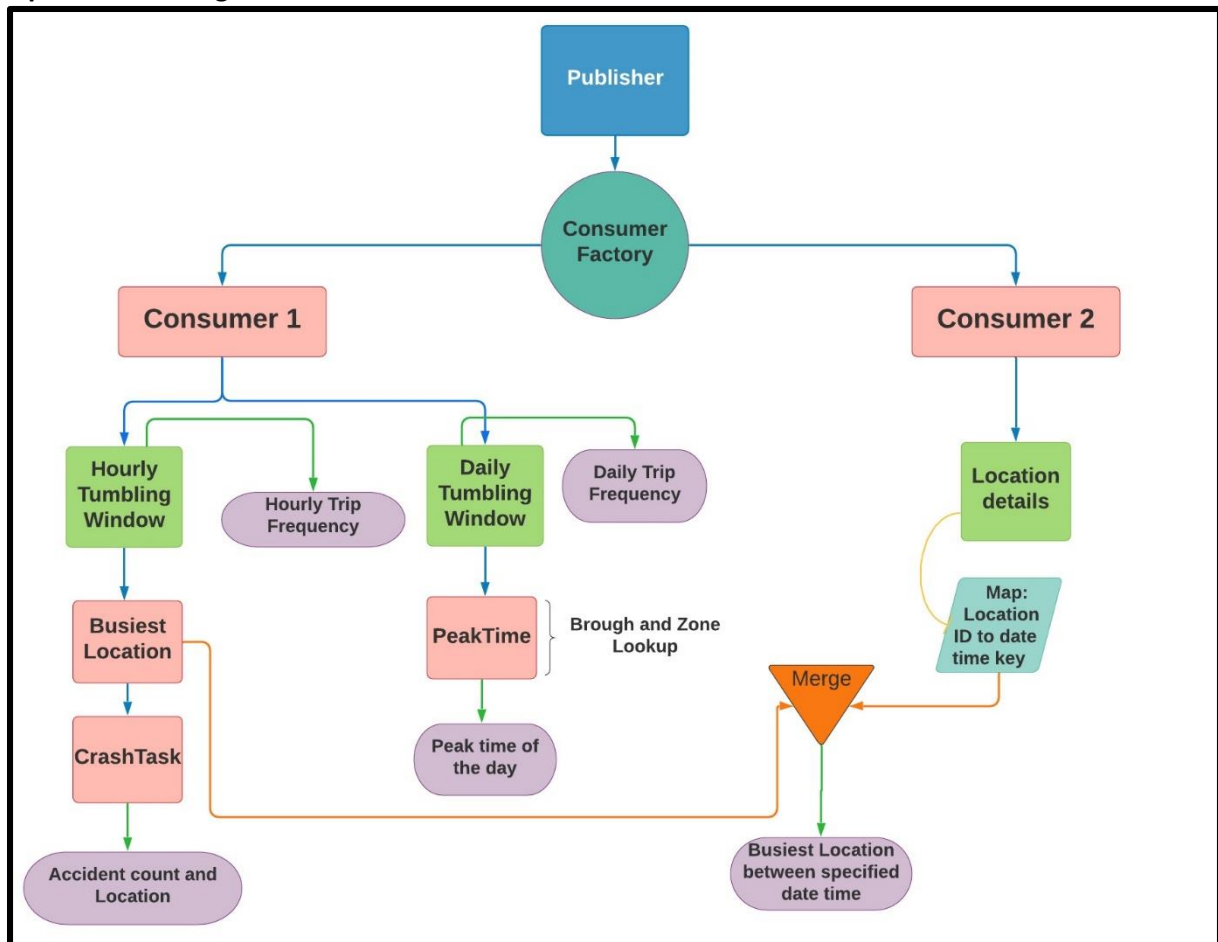
To have a reliable data communication between the Publisher & Subscriber, our inter-service communication checks for the data acknowledgement. When the message is received by the Subscriber, an acknowledgement is sent to the ActiveMQ and only then the next message is delivered to the subscriber. This process makes the system a bit slow, but avoiding it makes the system unreliable and affects the analytics. (Due to data and sequence loss)

By default, ActiveMQ creates a Message ID and Timestamp for each message which is being sent. As we have huge amount of data, creating ID and Timestamp for each message was consuming noticeable amount of time, hence we disabled the Message ID and Timestamp generation for each message. However, we are setting the delivery mode of the message to be Persistent which allows the messages to be written on the disk so that messages are not lost while broker restart.

Next, we have created a Consumer factory to initialize two consumers for the same Topic. In consumer factory we put the connection code which is commonly used to initiate both the consumers thus eliminating the redundancy. We have separated the processing of Location Mapping and Date & Time tumbling windows as two independent microservices as they both can be executed in parallel.

We are loading the crash data once independently and using it throughout the process to map the accidents, thus avoiding the need to stream the crash CSV every time to report the accident

**Pipeline Flow Diagram:**



## Pipeline Output:

Consumer 1 Microservice: Creates two tumbling windows, one for each hour of the day and one for the day. The Daily window counts the frequency of the trips that occurred each day and the hourly window counts the hourly trip frequency. As the data is received in sequence, we calculate the trip frequency for each hour, starting from $0^{th}$ Hour to $23^{rd}$ Hour. The Date and time is then passed to BusiestLocation Microservice to identify the busiest location. In the Daily window, we summarize the hourly data to generate the total number of trips happened in the day. The date is then passed to PeakTime microservice to identify the peak time of the day.

BusiestLocation Microservice: Looks up the Mapping information generated by Consumer 2. This data is, merged with the date and time information received from Consumer 1 to lookup Borough and Zones from the Location CSV and report the busiest location of specified date and time. The Busiest location also sends the borough name along with date and time to the CrashTask microservice to generate report on accidents for the specified date and time.

CrashTask Microservice: Counts the total number of accidents occurred in the given borough on specified date time and reports the same. It also reports on the total number of accidents occurred across the city for specified date time along with the location details.

PeakTime Microservice: Compares the trip frequency of each hour of the day and returns the time which has recorded the maximum trip frequency.

Consumer 2 Microservice: Receives the pickup time and location data from the publisher. It creates a Multimap of the data with Date and Hour as the key and location ID as the Value. This mapping is used get the busiest location ID for the specified date time.

**Design Strength &Weakness:**
**Strengths:**
Pub-Sub Model (With Topic as Destination):
We are using a Pub-Sub model with One Publisher and Two Subscriber. The data is received by the subscribers at the same time and they can start processing independently. The output from these two microservices is merged and the final report is generated. This saves a lot of time as both the processes run on different cores parallelly.

Efficient Storage:
Since dataset is large, it contains data for the entire month. We only need to print the daily and hourly report once; therefore, we process the hourly and daily data and once the report is displayed, we clear the storage and is used for new day's information.

Object Oriented Approach:
- We are processing the data file using pure Java 8 streams. We have not used any framework and thus it keeps our code simple and easy to understand and modify if required.
- We are making use of static variable which helps to preserve the data which is stored in the variable. This data can be accessed by other microservices without the need of object instantiation
- We are using Multimap data structure provided by Google's guava jar. This helps in holding multiple values for a single key without the need of using map reduce approach.
- We are also running the microservices as independent thread. This makes the control from the calling function to return immediately and continue with processing next data without waiting for the result from the called function(thread).

**Weakness:**
Sorted Data:
We need the data to arrive in sequence; this system will work fine in real-time where the data arrives in sequence. If the data does not arrive in sequence, our system will still generate report but on the incomplete data.

Crash data Store:
We are currently loading the crash data only for January 2018 and storing it in a list. This list is used to generate a report on Accidents. We are not streaming the crash data and in future we can change the code to make it scalable.

Report stockpile:
We are currently not writing the report to any external source hence it limits us from referring to previous result. This can be avoided by implementing report logging.

**Scalability:**

- Whole taxi data is being streamed line by line without storing in memory. Even if the dataset is increased substantially, the system will be able the support it.
- We are using 1D array of size 24 to hold each day's trip count. On processing of each day's record, this array is nullified and used for next day. This approach will enable us to process any number of days without the need to provide extra storage.
  This approach is not only restricted to 1 month's data but, can also be used for multiple months.
- The system is reliable and generates expected results as each message is acknowledged and the message delivery is persistent. There is no issue of data loss or sequence loss in data communication.
- As we are using parallel processing and multi-threading, the resource allocation for the micro service is done efficiently.
- Try-Catch and Exception handling is done for all the classes there for as we scale, we will be able to handle the exceptions efficiently and solve new issues if they arise.
- We have created Consumer factory so that multiple subscribers can subscribe to specific Publisher topic at the same factory. If new subscribers want to subscribe to the same Publisher Topic, it just needs to be initiated at the factory and connection code can be reused.
- We can create many such Consumer Factories based on specific Topics and group them together eliminate redundancy.
- Free & Open source software are used to develop the system.
- Annotations available in Java provides the supplemental information about the program. For example, default run() method of the thread is overridden by the custom run method defined in each Thread class by using @Override
- Since the microservices are build using Java, a lot of libraries and resources are available and there are lot of cloud providers who can scale Java based microservices.

**Reference:**

1. Cugola, Gianpaolo, and Alessandro Margara. "Processing flows of information: From data stream to complex event processing." ACM Computing Surveys (CSUR) 44.3 (2012): 15
2. Jamshidi, Pooyan, et al. "Microservices: The journey so far and challenges ahead." IEEE Software, 2018.
3. Taibi, Davide, and Valentina Lenarduzzi. "On the definition of microservice bad smells." IEEE software, 2018
4. https://activemq.apache.org/components/artemis/documentation/1.5.0/perf-tuning.html
5. https://activemq.apache.org/slow-consumer-handling
6. https://activemq.apache.org/components/artemis/documentation/1.0.0/perf-tuning.html
7. https://activemq.apache.org/performance-tuning

**Appendix:**

**BusiestLocation.java**

```java
1.  package yellowcab;
2.
3.  import java.io.IOException;
4.  import java.nio.file.Files;
5.  import java.nio.file.Paths;
6.  import java.util.AbstractMap;
7.  import java.util.Collection;
8.  import java.util.Map;
9.  import java.util.function.Function;
10. import static java.util.stream.Collectors.counting;
11. import static java.util.stream.Collectors.groupingBy;
12. import java.util.stream.Stream;
13. import static yellowcab.CrashTask.getCrashCountOnBusiestLoc;
14. import static yellowcab.Consumer2.locationID;
15.
16. public class BusiestLocation implements Runnable {
17.
18.     String key;
19.     String fileName = "taxi+_zone_lookup.csv";
20.
21.     // Set Key to the input key
22.     public BusiestLocation(String key) {
23.         this.key = key;
24.     }
25.
26.     @Override
27.     public void run() {
28.         try {
29.
30.             // Collection to store all the location ID for secified key
31.             Collection<Integer> busiestLoc = locationID.get(key);
32.
33.             // Count the frequency of trips started at Pickup location
34.             Map<Integer, Long> collect = busiestLoc.stream().collect(groupingBy(Fun
    ction.identity(), counting()));
35.
36.             // get location ID where maximum trip booking are initiated
37.             int maxLoc = collect.entrySet().stream().max((entry1, entry2) -
    > entry1.getValue() > entry2.getValue() ? 1 : -1).get().getKey();
38.
39.             // Borough and Zone Lookup for the identified busy location ID
40.             try ( Stream<String> data = Files.lines(Paths.get(fileName)).filter(s -
    > !s.isEmpty())) {
41.                 data.skip(1).map(s -
    > new AbstractMap.SimpleEntry<>(s, s.split(",")))
42.                         .filter(s -
    > s.getValue()[0].equals(String.valueOf(maxLoc)))
43.                         .forEach(s -> {
44.                             System.out.println("Buisest Location : " + s.getValue()
    [1] + "," + s.getValue()[2]);
45.                             try {
46.                                 // Get Accident count on the busiest Location
47.                                 getCrashCountOnBusiestLoc(Integer.parseInt(key.spli
    t(",")[0]),
48.                                         Integer.parseInt(key.split(",")[1]), s.getV
    alue()[1].replaceAll("^\"|\"$", ""));
49.                             } catch (IOException ex) {
50.                             }
51.                         });
52.             }
```

```
53.        } catch (IOException e) {
54.        }
55.
56.    }
57. }
```

## Consumer1.java

```java
1.  package yellowcab;
2.
3.  import java.util.logging.Level;
4.  import java.util.logging.Logger;
5.  import javax.jms.JMSException;
6.  import javax.jms.Message;
7.  import javax.jms.MessageListener;
8.  import javax.jms.TextMessage;
9.  import static yellowcab.YellowTaxi.thread;
10.
11. public class Consumer1 implements MessageListener {
12.
13.     static String key;
14.
15.     // Initialize 1D array to hold counter for each record received
16.     static int[] dateTime = new int[24];
17.
18.     @Override
19.     public void onMessage(Message message) {
20.         try {
21.             int dayFreq = 0, hourFreq = 0;
22.
23.             if (message instanceof TextMessage) {
24.                 TextMessage textMessage = (TextMessage) message;
25.
26.                 /*
27.                     When the message is not reached eof.
28.                         1. Enable Hourly Tumbling Window
29.                         2. Enable Daily Tumbling Window
30.                  */
31.                 if (!(textMessage.getText().equals("eof"))) {
32.
33.                     // Split the message into date and time
34.                     // Also split the date into dd, mm and yyyy to track the counte
    r and array writing
35.                     // split the time to hh mm and ss
36.                     String day = textMessage.getText().split(",")[0].split(" ")[0];

37.                     int dd = Integer.parseInt(day.split("-")[2]);
38.                     int mm = Integer.parseInt(day.split("-")[1]);
39.                     int yyyy = Integer.parseInt(day.split("-")[0]);
40.                     String time = textMessage.getText().split(",")[0].split(" ")[1]
    ;
41.                     int hh = Integer.parseInt(time.split(":")[0]);
42.
43.                     // Hourly Tumbling window
44.                     if (hh == 0 && dateTime[hh] == 0) {
45.                         dateTime[hh] += 1;
46.                     } else if (hh != 0 && dateTime[hh] == 0) {
47.                         dateTime[hh] += 1;
48.                         hourFreq = dateTime[hh - 1];
49.                         key = dd + "," + (hh - 1);
50.                         System.out.println("Total Trips between " + (hh - 1) + " to
     " + hh + " : " + hourFreq);
51.                         thread(new BusiestLocation(key), false);
```

```java
52.                         }
53.
54.                         // Daily Tumbling window
55.                         else if (hh == 0 && dateTime[23] != 0) {
56.                             hourFreq = dateTime[23];
57.                             key = (dd - 1) + "," + 23;
58.                             System.out.println("Total Trips between " + 23 + " to " + 2
    4 + " : " + hourFreq);
59.                             // Identify and print busiest location
60.                             thread(new BusiestLocation(key), false);
61.
62.                             // Summarise whole day's trip frequency
63.                             for (int i = 0; i < 24; i++) {
64.                                 dayFreq += dateTime[i];
65.                             }
66.                             Thread.sleep(500);
67.                             System.out.println("Total Trips made on " + (dd - 1) + "-
    " + mm + "-" + yyyy + " : " + dayFreq);
68.                             // Identify and print peak time
69.                             thread(new PeakTime(), false);
70.                             Thread.sleep(500);
71.                             dateTime = new int[24];
72.                             dateTime[hh] += 1;
73.                         } else {
74.                             dateTime[hh] += 1;
75.                         }
76.                     }
77.                     /*
78.                         If the stream reaches the eof, print:
79.                             1. 23rd to 24th hourly statistics
80.                             2. Last day's peak time and total trip frequency
81.                         This data is hardcoded as we need to process only for January 2
    018
82.                     */
83.                     else if (textMessage.getText().equals("eof")) {
84.
85.                         hourFreq = dateTime[23];
86.                         key = 31 + "," + 23;
87.                         System.out.println("Total Trips between " + 23 + " to " + 24 +
    " : " + hourFreq);
88.                         // Identify and print busiest location
89.                         thread(new BusiestLocation(key), false);
90.
91.                         // Summarise whole day's trip frequency
92.                         for (int i = 0; i < 24; i++) {
93.                             dayFreq += dateTime[i];
94.                         }
95.                         Thread.sleep(500);
96.                         System.out.println("Total Trips made on " + 31 + "-" + 01 + "-
    " + 2018 + " : " + dayFreq);
97.                         // Identify and print peak time
98.                         thread(new PeakTime(), false);
99.                     }
100.                 }
101.             } catch (NumberFormatException | JMSException e) {
102.             } catch (InterruptedException ex) {
103.                 Logger.getLogger(Consumer1.class.getName()).log(Level.SEVERE, nu
    ll, ex);
104.             }
105.         }
106.     }
```

## Consumer2.java

```java
1.  package yellowcab;
2.
3.  import com.google.common.collect.ArrayListMultimap;
4.  import com.google.common.collect.Multimap;
5.  import javax.jms.Connection;
6.  import javax.jms.JMSException;
7.  import javax.jms.Message;
8.  import javax.jms.MessageConsumer;
9.  import javax.jms.MessageListener;
10. import javax.jms.Session;
11. import javax.jms.TextMessage;
12.
13. public class Consumer2 implements MessageListener {
14.
15.     // Multimap to store all the Pickup Location at specified date time
16.     static Multimap<String, Integer> locationID = ArrayListMultimap.create();
17.
18.     Connection connection = null;
19.     Session session = null;
20.     MessageConsumer consumer2 = null;
21.
22.     // Constrcutor to set Connetion and session details
23.     public Consumer2(Connection connection, Session session, MessageConsumer consum
    er2) {
24.
25.         this.connection = connection;
26.         this.session = session;
27.         this.consumer2 = consumer2;
28.     }
29.
30.     @Override
31.     public void onMessage(Message message) {
32.         try {
33.             if (message instanceof TextMessage) {
34.                 TextMessage textMessage = (TextMessage) message;
35.
36.                 // If eof receieved, end the consumer, session and connection
37.                 if (textMessage.getText().equals("eof")) {
38.                     consumer2.close();
39.                     session.close();
40.                     connection.close();
41.
42.                 } else {
43.
44.                     // Get Date Time and Pickup Location ID from the Message
45.                     int day = Integer.parseInt(textMessage.getText().split(",")[0].
    split(" ")[0].split("-")[2]);
46.                     int time = Integer.parseInt(textMessage.getText().split(",")[0]
    .split(" ")[1].split(":")[0]);
47.                     int locID = Integer.parseInt(textMessage.getText().split(",")[1
    ]);
48.
49.                     // Generate Key as Date,Time
50.                     String key = day + "," + time;
51.
52.                     // Add multiple Pickup location IDs to date,time key
53.                     locationID.put(key, locID);
54.                 }
55.             }
56.         } catch (NumberFormatException | JMSException e) {
57.         }
58.     }
59. }
```

## ConsumerFactory.java

```java
1.  package yellowcab;
2.
3.  import java.io.IOException;
4.  import java.util.logging.Level;
5.  import java.util.logging.Logger;
6.  import javax.jms.Connection;
7.  import javax.jms.Destination;
8.  import javax.jms.JMSException;
9.  import javax.jms.MessageConsumer;
10. import javax.jms.Session;
11. import org.apache.activemq.ActiveMQConnectionFactory;
12.
13. public class ConsumerFactory implements Runnable {
14.
15.     @Override
16.     public void run() {
17.
18.         try {
19.             // Administrative object
20.             ActiveMQConnectionFactory connectionFactory = new ActiveMQConnectionFac
    tory("tcp://localhost:61616");
21.
22.             Connection connection = connectionFactory.createConnection();
23.
24.             // Create session object
25.             Session session = connection.createSession(false, Session.AUTO_ACKNOWLE
    DGE);
26.
27.             // Administrative object
28.             Destination dest = session.createTopic("myTopic");
29.
30.             // Consumer1 subscribes to customerTopic
31.             MessageConsumer consumer1 = session.createConsumer(dest);
32.             consumer1.setMessageListener(new Consumer1());
33.
34.             // Consumer2 subscribes to customerTopic
35.             MessageConsumer consumer2 = session.createConsumer(dest);
36.             consumer2.setMessageListener(new Consumer2(connection, session, consume
    r2));
37.
38.             connection.start();
39.
40.         } catch (JMSException e) {
41.         }
42.     }
43. }
```

## CrashTask.java

```java
1.  package yellowcab;
2.
3.  import java.io.IOException;
4.  import java.nio.file.Files;
5.  import java.nio.file.Paths;
6.  import java.util.AbstractMap;
7.  import java.util.ArrayList;
8.  import java.util.Collections;
9.  import java.util.HashSet;
10. import java.util.List;
11. import java.util.Set;
12. import java.util.stream.Stream;
```

```
13.
14.  public class CrashTask {
15.
16.      // Data Strcuture to store Crash data
17.      static List<String> crashData = new ArrayList<>();
18.      static String fileName2 = "Crash.csv";
19.      static long count;
20.
21.      public static void run() {
22.          // Stream Crash.csv
23.          try ( Stream<String> data = Files.lines(Paths.get(fileName2)).filter(s -
     > !s.isEmpty())) {
24.              data.map(s -> new AbstractMap.SimpleEntry<>(s, s.split(",")))
25.                      // Filter Jan 2018 from the csv
26.                      .filter(s -> (s.getValue()[0].split("/")[0].equals("01")
27.                      && s.getValue()[0].split("/")[2].equals("2018")))
28.                      .forEach(s -> {
29.                          // Handle empty values in the Borough field
30.                          if (s.getValue()[2].isEmpty()) {
31.                              s.getValue()[2] = "Location Not Available";
32.                          }
33.                          String crash = s.getValue()[0] + "," + s.getValue()[1] + ",
     " + s.getValue()[2];
34.                          // Add Date, Time and Borough deatils to the list
35.                          crashData.add(crash);
36.
37.                      });
38.          } catch (IOException ex) {
39.          }
40.      }
41.
42.      public static void getCrashCountOnBusiestLoc(int date, int time, String borough
     ) throws IOException {
43.          // List to store crash details of specified date time
44.          List<String> temp = new ArrayList<>();
45.          crashData.parallelStream()
46.                  .map(s -> new AbstractMap.SimpleEntry<>(s, s.split(",")))
47.                  // Filter crashData on specified date time
48.                  .filter(s -
     > Integer.parseInt(s.getValue()[0].split("/")[1]) == date
49.                  && Integer.parseInt(s.getValue()[1].split(":")[0]) == time)
50.                  .forEach(s -> {
51.                      String crash = s.getValue()[0] + "," + s.getValue()[1] + "," +
     s.getValue()[2];
52.                      // Add crash details of specified date time to a temperory list

53.                      temp.add(crash);
54.                  });
55.
56.          // Count the number of accidents occured in specified bororugh
57.          count = temp.parallelStream()
58.                  .map(s -> new AbstractMap.SimpleEntry<>(s, s.split(",")))
59.                  .filter(s -
     > s.getValue()[2].toLowerCase().contains(borough.toLowerCase()))
60.                  .count();
61.          System.out.println("Number of accidents at " + borough + " : " + count);
62.
63.          // Count the number of accidents occured overall city
64.          count = temp.parallelStream().count();
65.          System.out.println("Number of accidents across city : " + count);
66.
67.          // Get the Location details of the accidents occured overall city
68.          List<String> t = new ArrayList<>();
69.          temp.stream()
70.                  .map(s -> new AbstractMap.SimpleEntry<>(s, s.split(",")))
71.                  .forEach(s -> t.add(s.getValue()[2]));
```

```
72.
73.            Set<String> distinct = new HashSet<>(t);
74.            // Print the frequency in each location
75.            distinct.forEach((s) -> {
76.                System.out.println(s.toUpperCase() + ": " + Collections.frequency(t, s)
      );
77.            });
78.            System.out.println();
79.        }
80. }
```

## PeakTime.java

```
1.  package yellowcab;
2.
3.  import java.util.ArrayList;
4.  import java.util.Collections;
5.  import java.util.List;
6.
7.  // Thread to calculate peak time
8.  public class PeakTime implements Runnable {
9.
10.     // Define start and end time of a day
11.     int startTime = 00;
12.     int endTime = 23;
13.
14.     // List to hold frequency of each hour
15.     List<Integer> tripFreq = new ArrayList<>();
16.
17.     // List to identify the time frame
18.     List<Integer> timeFrame = new ArrayList<>();
19.
20.     @Override
21.     public void run() {
22.
23.         // for each time, add the frequency to the list
24.         for (int i = startTime; i <= endTime; i++) {
25.             timeFrame.add(i);
26.             // fetch frequency from Consumer 1
27.             tripFreq.add(Consumer1.dateTime[i]);
28.         }
29.
30.         // Calculate max frequency
31.         int maxVal = Collections.max(tripFreq);
32.
33.         // Get the index of the max frequency
34.         // This gives the time frame of peak hour
35.         int maxIdx = tripFreq.indexOf(maxVal);
36.         System.out.println("Peak time : " + maxIdx + " to " + (maxIdx + 1));
37.         System.out.println("------------------------------------------------
      \n");
38.         tripFreq.clear();
39.         timeFrame.clear();
40.     }
41. }
```

## Publisher.java

```
1.  package yellowcab;
2.
3.  import java.io.IOException;
```

```java
4.  import java.nio.file.Files;
5.  import java.nio.file.Paths;
6.  import java.util.AbstractMap;
7.  import java.util.stream.Stream;
8.
9.  import javax.jms.Connection;
10. import javax.jms.DeliveryMode;
11. import javax.jms.Destination;
12. import javax.jms.JMSException;
13. import javax.jms.Message;
14. import javax.jms.MessageProducer;
15. import javax.jms.Session;
16.
17. import org.apache.activemq.ActiveMQConnectionFactory;
18.
19. public class Publisher implements Runnable {
20.
21.     @Override
22.     public void run() {
23.
24.         try {
25.
26.             // Administrative object
27.             ActiveMQConnectionFactory connectionFactory = new ActiveMQConnectionFac
    tory("tcp://localhost:61616");
28.
29.             Connection connection = connectionFactory.createConnection();
30.             connection.start();
31.
32.             // Create session object
33.             Session session = connection.createSession(false, Session.AUTO_ACKNOWLE
    DGE);
34.
35.             // Administrative object
36.             Destination dest = session.createTopic("myTopic");
37.
38.             // Producer object
39.             MessageProducer producer = session.createProducer(dest);
40.
41.             // Avoid disk writing by setting Delivery Mode to NON_PERSISTENT
42.             producer.setDeliveryMode(DeliveryMode.PERSISTENT);
43.
44.             /* Disable Message ID and Time Stamp creation for each message
45.                 at producer end */
46.             producer.setDisableMessageID(true);
47.             producer.setDisableMessageTimestamp(true);
48.
49.             String fileName = "SortedTaxiData.csv";
50.
51.             // Stream Sorted CSV row by row
52.             try (Stream<String> data = Files.lines(Paths.get(fileName)).filter(s -
    > !s.isEmpty())) {
53.                 data.skip(1).map(s -
    > new AbstractMap.SimpleEntry<>(s, s.split(",")))
54.                         // Filter by only Jan 2018 records
55.                         .filter(s -> (s.getValue()[1].split(" ")[0].split("-
    ")[1].equals("01")
56.                         && s.getValue()[1].split(" ")[0].split("-
    ")[0].equals("2018")))
57.                         .forEach(s -> {
58.                             try {
59.                                 // Fetch Pickup DateTime and PickUp Location and cr
    eate message
60.                                 Message message = session.createTextMessage(s.getVa
    lue()[1] + "," + s.getValue()[7]);
61.                                 // Send Message to ActiveMQ
```

```
62.                            producer.send(message);
63.                        } catch (JMSException e) {
64.                        }
65.                    });
66.            } catch (IOException ex) {
67.            }
68.            // Send End OF File Message to ActiveMQ
69.            // This will help to close the session once the file is processed
70.            producer.send(session.createTextMessage("eof"));
71.            connection.close();
72.        } catch (JMSException e) {
73.        }
74.    }
75. }
```

## YellowTaxi.java

```java
1.  package yellowcab;
2.
3.  public class YellowTaxi {
4.
5.      public static void main(String[] args) throws Exception {
6.          // load crash data
7.          CrashTask.run();
8.
9.          // initialize publisher
10.         thread(new Publisher(), false);
11.
12.         // Initialize gateway
13.         thread(new ConsumerFactory(), false);
14.     }
15.
16.     // Run Threads
17.     public static void thread(Runnable runnable, boolean daemon) {
18.         Thread brokerThread = new Thread(runnable);
19.         brokerThread.setDaemon(daemon);
20.         brokerThread.start();
21.     }
22. }
```