

A Kafka-Based Centralized Platform for Smart Vehicle Supervising

Filipe Pires [85122], João Alegria [85048]

Software Architecture

Department of Electronics, Telecommunications and Informatics

University of Aveiro

April 18, 2020

Introduction

This report aims to describe the work developed for the second assignment of the course of 'Software Architecture', focused on a platform that collects and processes information from simulated vehicles.

The aim of the assignment was to explore the capabilities of a technology like Kafka in standard IoT solutions. The platform itself isn't meant to provide a field-tested solution to a problem or a set of problems, rather it is supposed to show how communications via Kafka can empower developers in a time that modularity is more important than ever and system components must be prepared to easily transfer data streams between each other.

So in this report we present the architecture of our solution and the Kafka-related configurations, justifying them according to what we learned and found to be most suitable for each specific use case. We also mention how the work was distributed amongst the authors.

All code developed is publicly accessible in our GitHub repository:

<https://github.com/FilipePires98/AS/>

1 IoT, Kafka and Connected Devices

Internet of Things (IoT) is becoming an increasing topic of interest among technology giants and business communities. IoT Components are interconnected devices over a network, which are embedded with sensors, software and smart applications so they can collect and exchange data with each other or with cloud / data centers.

One of the areas in which IoT is paving its way is the connected vehicles. According to Gartner predictions (*1*), by this year there should be about a quarter-billion connected vehicles on the road, which are more automated, providing new in-vehicle services such as enhanced navigation system, real-time traffic updates, weather alerts and integration with monitoring dashboards. In order to process the data generated by IoT connected vehicles, data is streamed to central processors usually located in the cloud. The collected information can be analysed and data can be extracted and transformed to the final result, which can be sent back to the vehicle or to a monitoring dashboard.

In this project we explore a hypothetical use case where communication between devices (or entities) is done through Kafka. Apache Kafka (*2*) is high-throughput distributed messaging system in which multiple producers send data to Kafka cluster and which in turn serves them to consumers. It is a distributed, partitioned, replicated commit log service. The Java application we developed and that is described in this report is a simplified version of an IoT data processing and monitoring application for connected vehicles, aimed to to explore Kafka capabilities for messaging between entities.

1.1 The Data

As having actual connected vehicles with processing units capable of collecting car data and transmitting it to other entities was out of the scope of the project, this is simulated through a simple text file containing 1 transmission (or message) per line. This data file with the name of `CAR.TXT` is placed under a specific directory and is read by our platform for processing.

In order to dynamically generate such source data, we developed a small script in Python that receives as parameters the number of cars to be simulated and the total number of messages to be generated from those cars and stored in the text file. The script is called `generateCAR.py` and is placed in the scripts package inside our project. By default, it creates 10 different cars and writes 100 messages, but these numbers can be easily modified inside the script.

We used a random-based approach to generate each aspect of a message. Unique register codes are created to represent vehicles, as well as their status and speeds. Message types are not generated with a specific pattern, but we made it far more likely to generate a message of type HEARTBEAT (see section 1.2) than of any other type.

1.2 The Messages

Messages can be of 1 of 3 types, each with a specific purpose and format:

- HEARTBEAT - the simplest type of message meant to notify the platform that the car is still connected.

Format: | car_reg | timestamp | 00 |

- SPEED - message meant to inform the platform about the current speed of the car. This allows the system to determine whether the car is going under the speed limit or not and trigger an alarm if necessary.

Format: | car_reg | timestamp | 01 | speed |

- STATUS - messages meant to detect whether the smart sensor detects any malfunction. This in theory could help the platform provide or suggest a solution for the malfunction to the car driver considering the entire network of connected vehicles.

Format: | car_reg | timestamp | 02 | status |

The car_reg corresponds to the register code of each vehicle and the timestamp corresponds to the time instance when the message was created, the remaining elements are self explanatory. As we will see, each message type is treated differently both in their purpose and in the care with which their transmission is done.

```
| 45-SH-72 | 1586183268975 | 02 | OK |  
| 28-MC-82 | 1586183269976 | 02 | OK |  
| 42-UW-71 | 1586183272978 | 00 |  
| 73-FD-20 | 1586183273979 | 00 |  
| 28-MC-82 | 1586183274980 | 01 | 20 |  
| 55-LZ-42 | 1586183276981 | 00 |  
| 64-IY-98 | 1586183281986 | 00 |  
| 45-SH-72 | 1586183286989 | 00 |  
| 42-UW-71 | 1586183311005 | 00 |  
| 80-DE-01 | 1586183315006 | 00 |  
| 30-UU-59 | 1586183319009 | 00 |  
| 78-ST-77 | 1586183324009 | 00 |  
| 28-MC-82 | 1586183325011 | 02 | KO |  
| 30-UU-59 | 1586183327012 | 01 | 0 |  
| 73-FD-20 | 1586183328012 | 01 | 130 |
```

Example of a portion of the CAR.TXT file.

2 System Architecture

Modular architecture is very appealing for IoT solutions, as it offers a way to manage the complexity of a problem by breaking it down to smaller and more easily manageable modules. Plus, IoT components are constantly changing, so limiting the effects of such organic-like transformations to individual modules allows developers to keep in control of the system's growth.

Software applications are embracing distributed, real time and on-the-cloud as the new norm. They are focused on providing real service rather than just being obsessed on lists of features. This is what drives the increase in awareness for the importance of modularity in software development. And this is partly what was meant to be explored in this assignment. So in this chapter we focus on presenting the interacting entities of our vehicle supervising platform and the components that constitute their functionality.

2.1 Entities

There are 4 entities, each with its own responsibilities. These entities simulate the management of the car network, with representative features of what could be accomplished by a fully implemented version. Following is their description and a simple block diagram in Figure 1.

- **CollectEntity** - its role is to collect data from connected vehicles (represented by `CAR.TXT`) and produce messages to the other entities with the gathered information (through communication channels established *a priori* and explained in chapter 3).
- **ReportEntity** - as the name states, this entity is responsible for reporting all that is transmitted by **CollectEntity**, writing it to `REPORT.TXT`.
- **BatchEntity** - **Batch's** role is to make possible the computation of any relevant metrics and calculations related to the information collected; in our case it merely accomplishes the same as **ReportEntity**, writing results to `BATCH.TXT`.
- **AlarmEntity** - its role is to trigger alarms and present them to the platform user when some relevant event occurs; in our case alarms are triggered when a car surpasses the predefined speed limit of 120 Km/h; all alarms are written to `ALARM.TXT`.

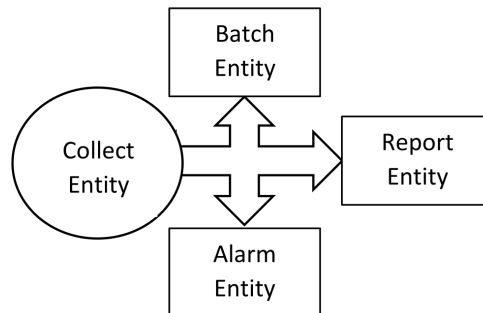


Figure 1: Basic block diagram for the platform centralized in the Collect Entity, taken from (3).

2.2 Components

The Java application is organized with packages. Each package contains a set of files with common natures so that managing source code is made easy. The "data" package contains both the input (CAR.TXT) and the output data (REPORT.TXT, BATCH.TXT and ALARM.TXT). The "entities" package contains every entity class. The "scripts" package holds project scripts, including the Kafka initialization and deletion scripts and the CAR.TXT file generation script.

The "message" package contains all classes related to the messages sent through Kafka; here we have the Message class, with all that was described in section 1.2, and our implementations of the serializer and deserializer classes of Message instances. Kafka has serialization classes available for non-structured data such as regular integers or strings, but we intended to send a structured message for greater usability, so implementing a serialization process for such structure was required in order to transmit data through Kafka.

Finally, the "kafkaUtils" package provides Kafka utilities such as EntityAction. This interface, implemented by all consumer entities, ensures that such entities have means to process Kafka messages in a predefined format by defining the method *processMessage()*. This method receives the identifier of the consumer that is going to process the current message, the topic that the message belongs to and the key-value pair of the message.

The package also contains our implementations of Kafka Consumers and Producers. Kafka is explained in greater detail in chapter 3, however a simple way of understanding producers and consumers in Kafka is to see them as the components responsible for creating messages that can be sent via Kafka topics and sending them, and for subscribing to topics and receiving messages from them, respectively. By implementing our own versions of such components, we gain greater control over their configurations and functionalities. The Producer class supports 3 message sending methods:

- *fireAndForget()* - where messages can be lost or arrive reordered on the consumer side.
- *sendAsync()* - where the producer asynchronously awaits for a delivery confirmation with the use of callbacks.
- *sendSync()* - where the producer waits for a delivery confirmation before proceeding to processing another message, ensuring the correct order on the consumer side.

Producer also has an inner class called ProducerCallback that, as the name suggests, deals with the callbacks related to the executions of the method *sendAsync()*. The Consumer class on the other hand works as a java thread by implementing the Runnable interface. Once instantiated, it listens to the topics where it is subscribed and processes arriving messages by calling *processMessage()*. Each consumer has an instance of RebalanceListener, our implementation of Kafka's ConsumerRebalanceListener for partition management. This listener is explained in greater detail in section 3.2.

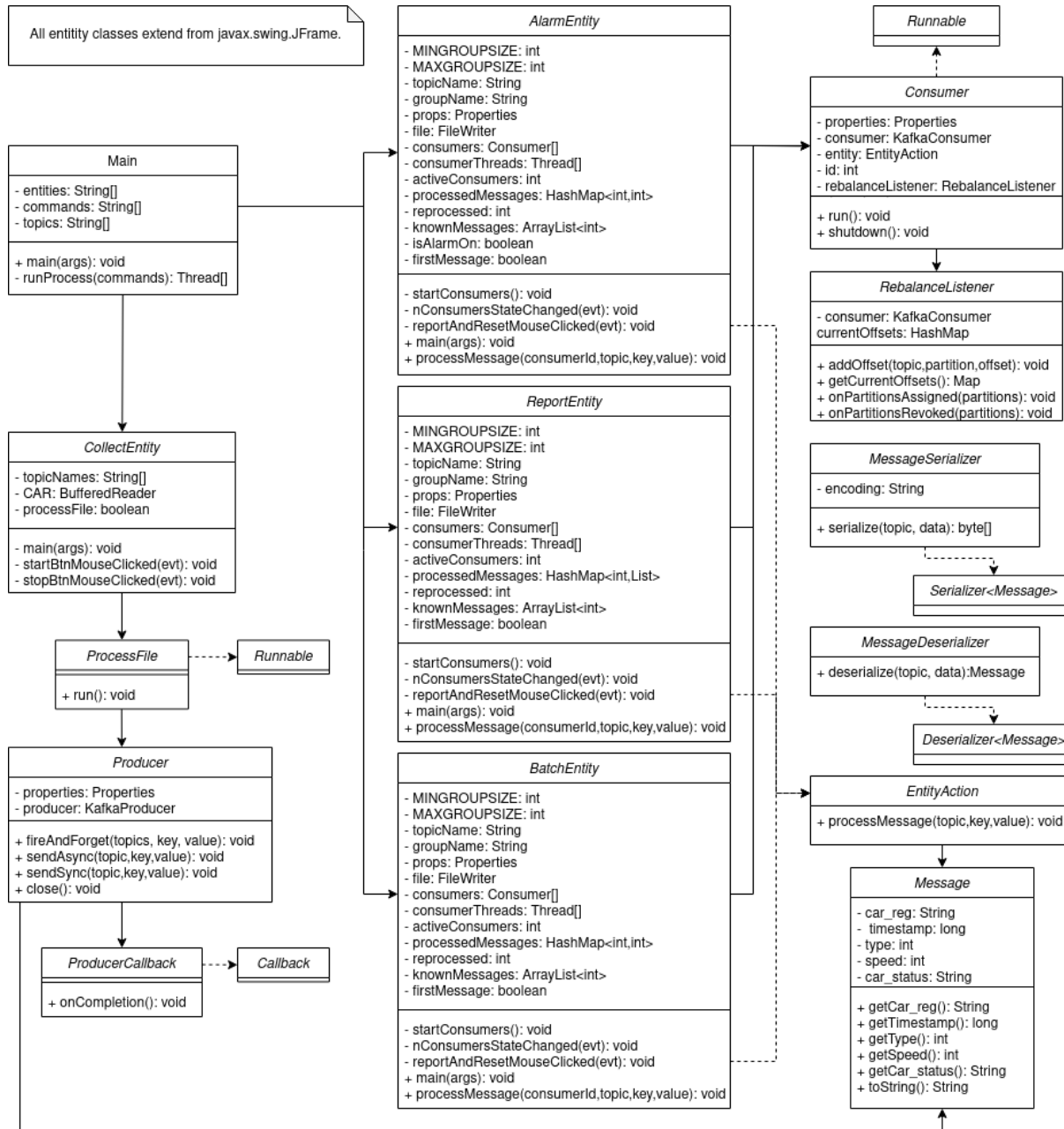


Figure 2: Java application's class diagram.

Figure 2 presents the project's class diagram with the relations between components. In order to more easily visualize interactions between entities, we designed an interaction diagram as well, present in Figure 3.

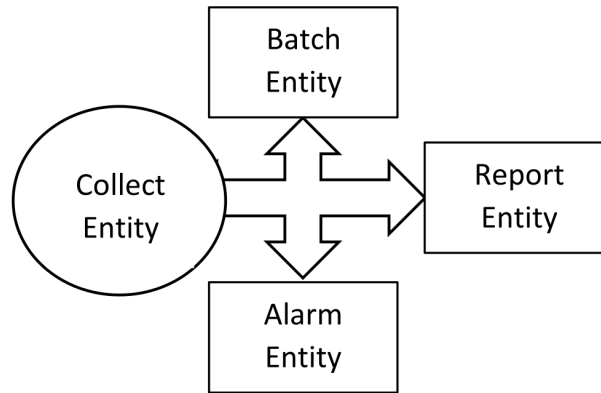


Figure 3: Entities interaction diagram.

2.3 User Interface

In any management platform there needs to be some sort of interface for users to interact with the system. Our choice was to develop a graphical user interface (GUI) for each entity using Java Swing (4).

All entities have a window dedicated to presenting relevant information, including processed messages status information and possible execution errors. Below we see a sequence of figures showing all GUIs in different program execution states.

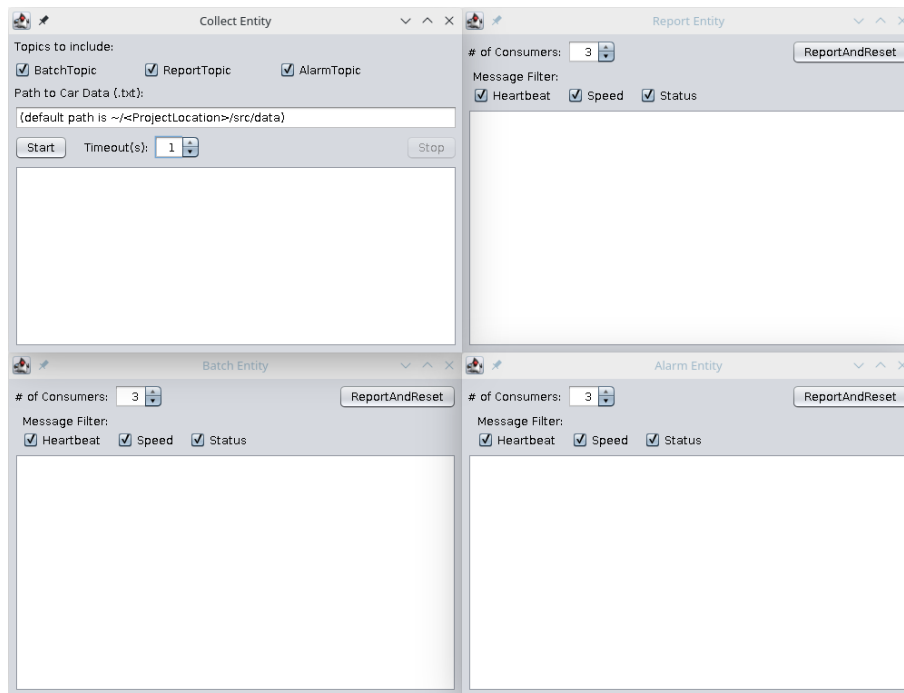


Figure 4: GUIs before file reading.

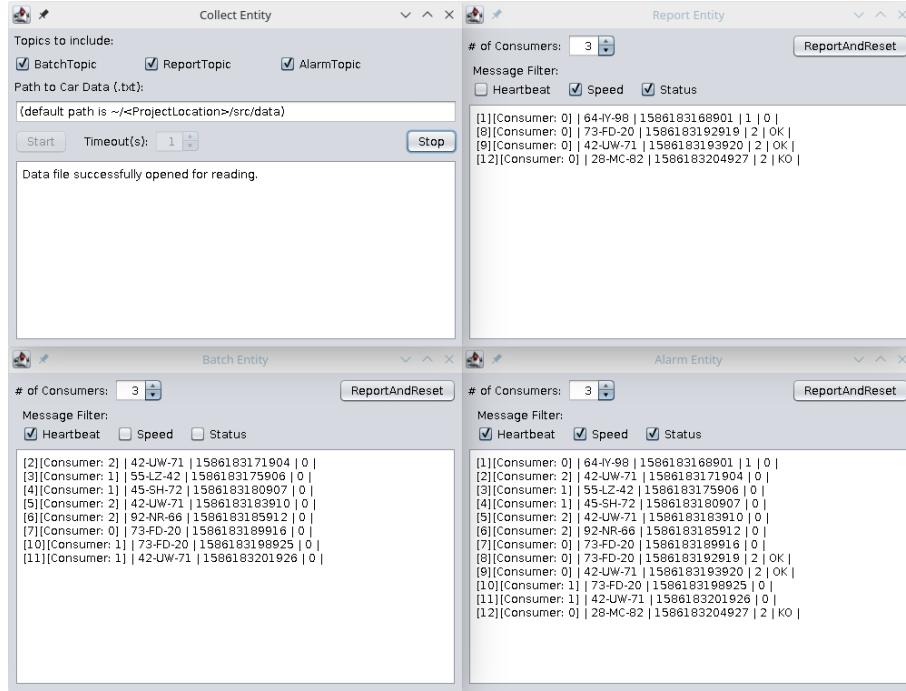


Figure 5: GUIs during message transmission.

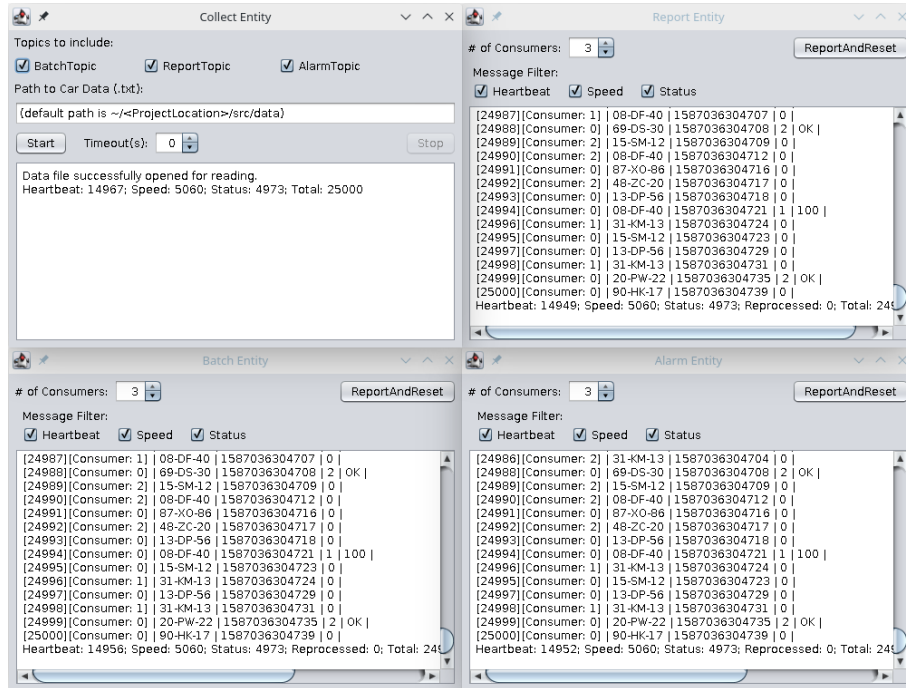


Figure 6: GUIs with counter reports.

The central entity, Collect, is the one with more control over the platform. In its GUI, the user can: define the path to the data file CAR.TXT (a standard path is defined by default, directed towards the "data" package); select which topics shall be used, allowing the program to run even if one of the consumer entities isn't working; define a timeout variable that sets

the time period between message dispatches; initiate the collecting and transmission of data by clicking the `Start` button. During execution, at any time, the user can also terminate data transmission by clicking the `Stop` button (enabled only during execution).

Batch, Report and Alarm entities all have similar interfaces with the same operations available to the user. These operations are: defining the number of consumers ; filtering which messages to present on the information window according to their type; presenting the number of messages of each type, of reprocessed messages and in total processed so far by clicking the `ReportAndReset` button (as the name suggests, these counters are set to zero every time the button is clicked).

3 Kafka Infrastructure

As a distributed streaming platform, Apache Kafka has 3 main capabilities: publish and subscribe to record streams, similarly to a message queue or messaging system; store those record streams in a fault-tolerant way; process streams of records as they occur. Its general use is for building real-time streaming data pipelines to establish communication between applications, and building real-time streaming applications to transform and react to streams of data.

For our specific case, the messaging system scenario seemed the most appropriate. As stated in their online documentation (2), normally these type of systems follow 2 possible models: message queueing or publish-subscribe. The queue model is characterized by a pool of consumers that read from the messaging server and each record goes to one of the consumers; this has several positive aspects such as allowing the division of the records in a balanced fashion by the registered consumers and enabling the processing rate and power to scale; on the other hand, being a queue, when a record is processed it may not be reprocessed. Comparatively, the publish-subscribe model is characterized by broadcasting each record to all registered consumers; this broadcasting feature can be seen as a positive aspect, but is counterweighted by the fact that there is no way to scale the process, since every message is sent to every consumer.

With this in mind, we established Kafka's infrastructure as the message system and the backbone for the entire platform, making use of both models mentioned above. The usage of this platform was not only a constraint in the proposed assignment, but was also justified by our own personal analysis that showed that Kafka is a solution complete enough for our purposes and guarantees a number of quality metrics that were found important. It was these metrics we explored during development in order to understand their influence on the quality attributes of a system with a standard IoT architecture.

3.1 Initialization and Destruction Scripts

Obviously, to use this platform we first needed to instantiate the services and infrastructure composing Apache Kafka. To do so, we had 2 approaches: use the instantiation scripts provided with the source code available online (2) or use docker as a middleware framework that would enable us to abstract the instantiation process and use Kafka in a container. We chose the former option, which requested more work from our side but enabled us to instantiate the platform with more precision and control.

Manually running each necessary command was obviously out of question. The solution was to develop 2 scripts to automate both the creation and deletion of the infrastructure necessary for this project. These scripts are available together with the source code, in the package folder `scripts`.

The initialization script is the first action to be executed when the project is ran. Also, as you might guess, the deletion script is the last thing to be executed before the system is terminated. In relation to *initKafka.sh*, the initialization script, it is responsible for: starting Zookeeper, the entity in charge of managing the Kafka Brokers; initializing the Kafka Brokers themselves, entities with the logic of the platform, for service providing (3 brokers are deployed); creating the necessary Kafka Topics. The final order of actions performed is:

1. Instantiate Zookeeper
2. Instantiate Kafka Brokers
3. Store process IDs of the Kafka Brokers in a auxiliary file
4. Create necessary topics for this project

In relations to *deleteKafka.sh*, the deletion script, its responsibility is to kill the processes previously stored in the auxiliary files and delete the logs generated by the same processes. The final order of actions performed is:

1. Fetch process IDs from the auxiliary file
2. Kill the fetched processes
3. Delete the generated logs

3.2 Topics and Constraints

To establish constraints to the infrastructure consequently to the project, we had 3 methods, either in the definition of properties for the consumers, in the definition of the producers or in the definition of properties for the topics themselves.

For this project it was established that there should only exist 3 topics, the *BatchTopic*, *ReportTopic* and the *AlarmTopic*, each one with the objective of establishing communication between the *CollectEntity* and the *BatchEntity*, *ReportEntity* and *AlarmEntity*, respectively.

As already mentioned, there are three different types of messages, and all types should be sent to every one of the topics. This condition is quite relevant since this means that the properties and conditions assured in shared resources across message types, such as the topics, will be dictated by the most restricting message type. The higher restrictions on a given constraint can be defined by different message types, since each message type as a different set of constraints. Those constrains are:

- | | | |
|--|--|---|
| <ul style="list-style-type: none"> ● HEARTBEAT - can be lost - can be reordered - can be reprocessed | <ul style="list-style-type: none"> ● SPEED - can be lost - cannot be reordered - cannot be reprocessed | <ul style="list-style-type: none"> ● STATUS - cannot be lost - cannot be reordered - can be reprocessed |
|--|--|---|

4 Additional Remarks

4.1 Documentation

Our attitude towards the developed code was to ensure it could be applied to other similar scenarios and reused in systems intended to be deployed in real scenarios. With this in mind, we took great care with regards to code readability. By maintaining a code style equal throughout the project and defining intuitive and self-explaining variable and method names, we made the code easy to understand by someone already contextualized with Kafka.

Nevertheless, we wanted to make sure this was also true to someone looking at our project for the first time, so we resorted to the well-known Javadoc (5) tool to manage all code documentation. Comments were also added in key points throughout the code, including the scripts.

4.2 Assignment Contributions

As the entire development phase took place in a time where on-site cooperation was not possible, we resorted to online communication platforms to debate decisions and discuss difficulties. Team scheduling allowed us to work on the project simultaneously, so no member suffered from unbalanced workloads. The dimension of the project did not appeal to the usage of repository pull requests and other synchronization tools. However, each small solution was verified and agreed by both team members.

Having said this, it is difficult to isolate what each member actually implemented, as the influence of both is present in all components. Nevertheless, one might say that each had stronger responsibilities on a set of project aspects: Filipe took care of the execution of the individual Java processes and of the Shell scripts, while João developed the Kafka-related classes such as Consumer, Producer and EntityAction; Filipe developed the Python script for generation of `CAR.TXT`, while João developed the Shell scripts for Kafka initialization and deletion; each implemented 2 entities and each wrote a portion of this report; Filipe made sure everything was coherent throughout the report and the code documentation, while João solved the most critical issues regarding the configuration of the topics. In terms of work percentage, we believe it was about 50% for each student.

Conclusions

It is only natural for an informatics engineer to eventually come in contact with Apache Kafka. The wide range of applications of such technology, its high performance in several dimensions and the ease of configuration makes it a valuable tool for several scenarios and with a variety of purposes. This is true for previous software development projects we were responsible for, both in the academical and the professional environments. Nevertheless, prior to this assignment we were never given the chance to truly explore Kafka's capabilities more deeply.

Given such opportunity, there was room for learning valuable details about Kafka. The first thing we understood was the use of different topic configurations according to their purpose and how this could have been useful in past scenarios. Also, the idea of load balancing through the use of multiple partitions within a topic was found very interesting. Having to manage different topics and different message types forced us to find flexible solutions, which was perhaps the most rewarding aspect of the assignment once our tests proved it to be working as intended.

Regarding the overall perspective over the work delivered, we believe it fulfills all of the defined requirements while ensuring some level of redundancy and providing enough resources for a correct reuse of the code in future projects. For future work, it would be interesting to more realistically simulate the smart vehicles and apply some sort of more complex computations in, for example, the Batch Entity.

References

1. Smarter With Gartner, *Staying on Track with Connected Car Security*, <https://www.gartner.com/smarterwithgartner/staying-on-track-with-connected-car-security/>, accessed in April 2020.
2. Apache Kafka, *Apache Kafka: A Distributed Streaming Platform*, <https://kafka.apache.org/>, accessed in April 2020.
3. Óscar Pereira, SA: *Practical Assignment no.2*, University of Aveiro, 2019/20.
4. Oracle, *Swing*, <https://docs.oracle.com/javase/8/docs/technotes/guides/swing/index.html>, accessed in April 2020.
5. Oracle, *Javadoc Technology*, <https://docs.oracle.com/javase/8/docs/technotes/guides/javadoc/index.html>, accessed in April 2020.