

# **7CCS4PRJ Final Year Individual Project**

## **Real-Time Visualisation of Bus Delays in London**

### **iBus Disruption Monitor**

A project in collaboration with Transport for London

Final Project Report

Author: Konstantin Vladimirov Draganov

Supervisor: Dr Steffen Zschaler

Course: MSci Computer Science

Student ID: 1101314

September 2014 - April 2015

## **Abstract**

Automatic Vehicle Location (AVL) systems for bus fleets have been deployed successfully in many cities. They have enabled improved bus fleet management and operation as well as wide range of information for the travelling public. However there are still processes that can be improved or automated by utilising the data made available by the different systems including the AVL one. This report explores and analyses the tools and applications currently available at Transport for London (TFL) bus emergency and command unit. The report then proposes a prototypical tool for monitoring the bus delays in real time in the network. This tool offers objective source of processed information to bus operators and control room staff. However further work need to be done in order to place this tool in production environment. This is because arterial urban delay detection is very complex and unpredictable as will the report justify. The report concludes with some suggestions of how this project can be improved.

### **Originality Avowal**

I verify that I am the sole author of this report, except where explicitly stated to the contrary. I grant the right to King's College London to make paper and electronic copies of the submitted work for purposes of marking, plagiarism detection and archival, and to upload a copy of the work to Turnitin or another trusted plagiarism detection service. I confirm this report does not exceed 25,000 words.

Konstantin Vladimirov Draganov

September 2014 - April 2015

## **Acknowledgements**

First and foremost I offer my sincerest gratitude to my supervisor Dr Steffen Zschaler for the continuous support, guidance, encouragement, insightful comments and hard questions throughout the course of this exciting project.

I would also like to thank Andrew Highfield and Keith Elliot from TFL for providing me with all the needed data, information and feedback which has been immensely helpful.

Last but not the least, I would like to thank my parents Vladimir and Veselka also my sister Lilia and my partner Simona for their support and patience not only during the project, but throughout my life as well.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Scope . . . . .	4
1.2	Aims . . . . .	5
1.3	Objectives . . . . .	5
1.4	Report Structure . . . . .	6
<b>2</b>	<b>Background</b>	<b>7</b>
2.1	London Bus Network . . . . .	7
2.2	CentreComm . . . . .	8
2.3	iBus AVL . . . . .	9
2.4	Related Work . . . . .	11
2.5	Time Series . . . . .	14
2.6	Other . . . . .	25
2.7	Summary . . . . .	25
<b>3</b>	<b>Specification</b>	<b>26</b>
3.1	User Requirements . . . . .	26
3.2	Functional Requirements . . . . .	27
<b>4</b>	<b>Design</b>	<b>29</b>
4.1	Use cases . . . . .	31
4.2	System architecture . . . . .	32
4.3	State Machine . . . . .	33
4.4	Class organisation . . . . .	33
4.5	User Interface . . . . .	36
<b>5</b>	<b>Implementation</b>	<b>38</b>
5.1	iBus AVL Data . . . . .	38
5.2	Disruption Engine . . . . .	41
5.3	Graphical User Interface . . . . .	48
5.4	Problems . . . . .	52

<b>6</b>	<b>Testing</b>	<b>55</b>
6.1	Unit Testing . . . . .	55
6.2	White Box Testing . . . . .	55
6.3	Functional Testing . . . . .	56
6.4	Stress Testing . . . . .	57
<b>7</b>	<b>Results &amp; Evaluation</b>	<b>61</b>
7.1	Evaluation . . . . .	61
7.2	Results . . . . .	64
<b>8</b>	<b>Professional &amp; Ethical Issues</b>	<b>65</b>
<b>9</b>	<b>Conclusion</b>	<b>67</b>
9.1	Future Work . . . . .	68
	<b>References</b>	<b>75</b>

# Chapter 1

## Introduction

London bus network is one of the largest and most advanced bus networks in the worlds. It is responsible for more than 2.4 billion passenger journeys a year [15]. The constant population growth of England's capital has been also driving the expansion and improvement of the transport networks across the city. Transport for London (TFL) is in charge of its operation and its bus network is recognised as one of the top in the world in terms of reliability, affordability and cost-effectiveness [15]. The capital's bus network is continually expanding along with the city's population [30]. This leads to more pressure being put on the infrastructure which include not only the road network and the bus fleet, but on the technological systems that aid its operation.

Maintaining such a large scale network requires careful planning and monitoring. Being able to maintain such high reliability service 24/7 364 days in the year requires employing new technologies. This also helps keep costs down and thus keeping the service more affordable and accessible for the general travelling public. Each bus in the TFL network has been equipped with state of the art GPS enabled automatic vehicle location (AVL) system named iBus[35]. This AVL system has led to improved fleet management and has enabled the creation and improvement of multiple applications [50]. The system is generating large sets of data both in real-time as well as historical data. This information aids the bus operators and the emergency control room at TFL,

responsible for maintaining the bus network (CentreComm), to better manage and maintain the smooth operation of the bus network. This includes both planning for future demand and growth as well as emergencies and innervation during the daily operation of the network.

However there are still some situations and problems which require CentreComm staff to carry out manually analysis of the available data. This means that there is lack of readily available preprocessed information. Having to manually monitor thousands of buses continuously is very impractical. That is the reason why currently CentreComm operators rely heavily on individual bus operators and drivers to alert them of possible problems. Once alerted of a possible disruptions in the network they (CentreComm) can start their own investigation first verifying what they have been told by the bus drivers/operators and then into finding the cause and the actual severity of the problem. This often could lead to spending time and resources into investigating non existent problems. Worse it often leads to time and resources being spend on investigating and dealing with problems which actually are less important than others just because some drivers or bus operators have exaggerated the issue. Our project tries to address this inefficiencies and to propose a prototypical tool for real time monitoring of the bus delays in the network.

## 1.1 Scope

The scope of this project is to analyse the current work flow of CentreComm operators and their needs. The main goal is to design and implement a prototype which to automate and improve the work flows currently in place. This tool has to work and analyse the data that has been made available by TFL. This analysis is required to happen in real time as more data is being made available. The reason for this being that it would be used as an objective source of information for the delays in the bus network at each point in time. In addition to this the project needs to perform analysis of what visualisation of the output will be useful and suitable.



## 1.2 Aims

The main aim of this project is to design and implement a real-time visualisation tool which to highlight disrupted routes or parts of the TFL bus network which experience delays or are already experiencing. Potentially the system could alert of possible delays even before the bus drivers or operators have noticed and contacted CentreComm for assistance. This aim could be subdivided into two smaller aims:

- The first one which is independent of the other is to enable the processing of the data generated by the buses in the TFL's bus network. The tool need to be able to analyse the input data sets and calculate and output a list of the disruption that are observed in the network. It has to present information regarding the location (route section) in the transport network and their severity.
- The second part of the main aim above is to visualise the generated output in an easy to use and understand way. It is also important to note that the visualisation should be capable of updating itself whenever the list of delays have changed. This needs to happen in real time as well.

## 1.3 Objectives

The objectives that have been followed in order to successfully meet the above stated aims are:

- Obtain an in depth understanding of the problem and current work-flows that are in place at CentreComm.
- Researching similar work in the literature that has already been done and how it relates to our problem.
- Obtaining samples of the available data and understanding what it means.

- Gather and refine the user requirements during discussions and meetings with CentreComm staff and stakeholders.
- Design and develop initial prototype which to be further refined and improved after obtaining feedback from TFL.
- Test and evaluate that the tool works according to the user requirements and the design specifications.

## 1.4 Report Structure

In order to help the reader I have outlined the project structure here. The report will continue in the next chapter by providing the reader with a detailed background knowledge needed for the rest of the report as well as an in-depth review of the related work that is found in the literature. This will include brief of background on the current work-flow CentreComm operators follow and its inefficiencies. I will also give background on the iBus system and the data that the tool will need to operate with. I then explore related work that has already been done and how ours differs. This is followed by alternative approaches and models that could be utilised. Afterwards the report focusses on the specific requirements (Chapter 3) that have been identified and gathered from CentreComm. The report then goes on to outline the design (Chapter 4) and the implementation (Chapter 5) of the proposed system. This is followed by Chapter 6 and 7 which address testing and evaluation of the prototype respectively. I conclude the report with a summary of what has been achieved and guidance how the work presented in this report could be further developed and improved.

## Chapter 2

# Background

This chapter aims to introduce some concepts surrounding our problem domain, which to help the reader to understand and follow easier the following more technical chapters. It also presents a review of the related work that has been done in this area. The below sections look at some of the key aspects and problems that arise. I then conclude by providing a number of alternative methods for solving our problem.

### 2.1 London Bus Network

London bus network is one of the most advanced and renowned in the world. It runs 24 hours and it is extensive and frequent. Every route in the network is tendered to different bus company operator [12]. Each of this bus operator companies is then responsible for abiding to the contracts with TFL. This means that they (the different bus companies) are responsible to ensure that the services they are operating run according to the timetable as per the respective contract. There two main types of schedules that are being used:

- Fixed schedule - a bus stop need to be served at specific predefined times (e.g. 1:00pm, 1:20pm, 2:00pm etc.).
- Headway based - this means that buses should server bus stops at regular intervals (e.g. a stop need to be served by a bus every 5 minutes).

However under different circumstances some delays occurring on a given route are beyond the control of the different bus operator companies. A simple example could be a burst water/gas pipe on a street used by a bus route or any other incident (even terrorist attacks [14]) and even simply a severe congestion. In situations like this bus operators have no authority or power to respond or overcome such problems on their own. This is where CentreComm comes into place to respond and deal with such issues. In situations like this the bus drivers or orators would need to alert and ask CentreComm to intervene. The emergency command and control room at TFL can do so by for example implementing a short/long term diversions or curtailments (short turning) some of the buses on the affected routes. They also can seek assistance from London Traffic Control Centre<sup>1</sup> or even the Police under given circumstances.

Buses in the network can be classified by multiple factors, however for the purpose of this report the main distinction we need to consider apart from fixed schedule and headway based are **high** and **low** frequency bus routes [12]. High frequency routes are routes where there are 5 or more buses per hour attending a given bus stop. Low frequency routes are those that have 4 or less buses at a stop.

## 2.2 CentreComm

CentreComm is TFL's emergency command and control room responsible for all public buses in London. It has been in operation for more than 30 years [14] and it employs a dedicated team of professionals who work 24 hours 364 days in the year. They are dealing with more than a 1000 calls on a daily basis. The majority of these calls come from bus drivers or bus company operators regarding problems and incidents happening within the bus network. CentreComm staff implement planned long and short term changes in the bus network in response to different events taking place in the capital (including the 2012 Olympics). They are also responsible for reacting in real time to any unexpected and unpredicted changes and disruptions, maintaining the smooth,

---

<sup>1</sup><http://www.tfl.gov.uk/corporate/about-tfl/what-we-do/roads>

reliable and sage operation of London busy bus network.

London bus network consists of around 680 bus routes operated by more than 8000 buses [2]. Each of this buses is equipped with state of the art iBus system to help monitor and manage this enormous fleet. CentreComm's way of operation has been transformed beyond recognition since it has first opened and today. It started more than 30 years ago [14] and it consisted of a couple of operators equipped with two way radios and pen and papers. Today CentreComm operators make use of numerous screens each displaying interactive maps (displaying each bus location) and CCTV cameras in real-time. However there is still a lot of room for automation and improvement in their way of operation in order to effectively and efficiently deal with the growing bus network and its demands.

## 2.3 iBus AVL

Automatic vehicle location (AVL) systems provide real-time vehicle tracking most often this is achieved by the integration of Global Positioning System (GPS), wireless communications (e.g. SMS, GPRS) and geographic information system (GIS) [38]. AVL systems employ wireless communications for the transmission of the GPS coordinates and other data of the vehicle as it moves in the transport network. This once received by the a central server or computer allows the GIS software to map the location of the vehicle and it also enables further analysis to be performed based on this data.

All of London buses operating on the TFL bus network have been equipped with state of the art and award wining [11] AVL system named iBus [13]. This system has opened a range of new applications that could be built on top of it using the information that is made available. The iBus system consist of a number of computer and communication systems, sensors and transmitters as described in [23] and [50]. One of the key components of the system on-board unit (OBU) which mounted on each of buses in the TFL bus fleet and consists of a computational unit connected to sensors and GPS transmitters

(see figure 2.1 below taken from [23]). This OBU is responsible for a number of

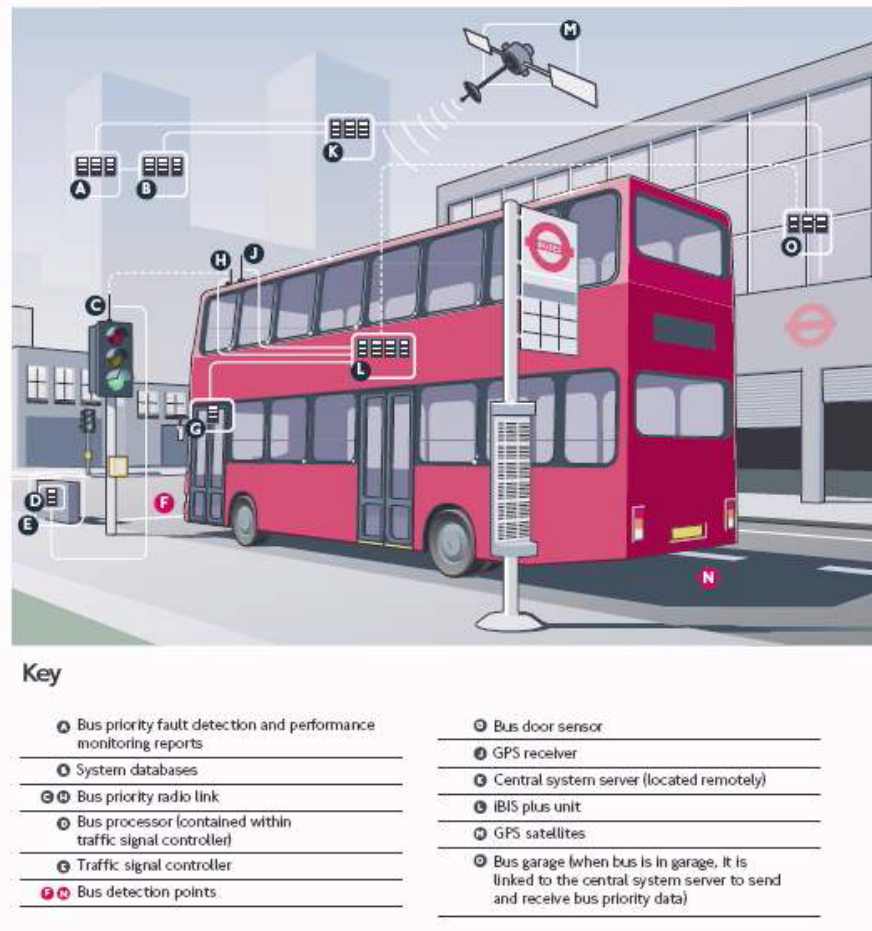


Figure 2.1: Overview of iBus System [23]

tasks including a regular (approximately every 30 seconds) transmission of the bus location. This information is currently used by the different bus operators for fleet management as well as by CentreComm for real-time monitoring of the buses and their locations. There have been a number of other applications and systems that have been implemented and put in to use as a result of the data that is generated by the iBus AVL. Some examples include Countdown (real-time passenger information), improved bus priority at traffic signals and more [23]. This has led to improved and more affordable transport service.

CentreComm operators have access to an online GIS system showing each bus location in the network on a map in real-time. This system also allows them

to see whether a given bus is behind, ahead or on time according to its schedule. However this does not show or alert the control room staff if a bus or a route is disrupted. Control room staff also can see when was a given bus expected to arrive at a given bus stop and when it actually arrived, but this again is only per individual bus and there are no automatic alerts set in case buses on given route or a given stop is being attended late. Currently the work-flow is such that CentreComm need to go and analyse all this information manually (once bus drivers or bus company operator have contacted/alerted the control room) in order to figure out if there is a problem and how severe actually it is. This is very inefficient, tedious and error prone process. Here is where our project comes into place to address the lack of preprocessed information and automated alerts.

## 2.4 Related Work

The literature is full of research towards accurately predicting bus arrival times with various computational models being used [1]. Predicting bus arrival is complex as many factor need to be considered like for example bus dwell time at stops, general congestion and others [24]. This is also closely related to the issues of bus prioritisations for which we can find numerous examples in the literature. Some examples of work done towards bus prioritisation at traffic signals and junctions can be found in [22] and [34]. However these are not directly related to our problem domain and thus are not discussed further in this report. This is because we are not interested to know when is the next bus due to arrive at stop or should we give priority to a bus at traffic signal. We want to know if buses experience increased travel time and thus get delayed travelling through some parts of the network.

The main problem posed by this project of detecting disruptions in the bus network can be also translated to short term traffic congestion detection or/and travel time calculation. This is valid as we are not interested in individual bus delays. Some examples of such single instances of bus disruptions

are customer incident on board of a single bus or technical fault with this bus or other issue which is affecting a single vehicle rather than the route or network. We are interested into finding routes/sections in the network which are delayed/disrupted and this beyond the control of individual bus operators. Most often such problems are due some sort of congestion or road closures/repairs. However road problems are in most cases linked with increased traffic congestions as roads are used by other vehicles as well. In order to address this problem posed by this project we have focussed our attention towards work which relates to calculating the bus travel times or tries to give short-term traffic congestion predictions in a given transport network.

The literature contains plenty of work done toward detecting calculating travel times in non urban environment. This includes approaches based on AVL probe data and automatic vehicle identification (AVI) as well as induction loops [47]. There are also plenty research done towards traffic congestion detection based on AVL probe data [47]. However there are significant challenges due to the nature of the urban environment itself. Densely populate areas are influenced by many factors which can be affecting the general traffic flow. Another problem posed by urban environments is the irregularity of the AVL data transmission because of for example weak or lost signal at times (e.g. due to high buildings) or even there can be some noise which reduces the accuracy of the transmitted location.

There is however little research to my knowledge which focusses on the issues of detecting and short-term forecasting of traffic congestion/disruptions in arterial urban environment [47] [5]. From the tables shown in figures ?? to ?? in Appendix A taken from [47] we can easily see that most of that has been done has focussed on motorways and also has employed data from static detector points (e.g automatic vehicle identification). Only in recent years we can see that more attention has been given to the use of GPS and AVL data. This is probably due to increased popularity and usage of these technologies.

In the literature various approaches to measure and predict travel time can be found. These models are categorised in four type groups according to [51]



as follows:

- Statistical models - this type of models employ statistical tools and methods for analysis and forecasting. Some models of this type include:
  - Historical
  - Time Series
  - Nonparametric regression
  - Hybrid
- Computer Simulations - main models of this type are traffic simulations. They allow simulation of the traffic flow in a network resembling the characteristic of moving vehicles. Main advantage of these models is that they allow for the simulation of different scenarios. The main drawback however is that they require traffic flow prediction information in advance [42]. Due to the optimisation nature of this approaches usually high performance computers are employed. In [25] they make use of parallel computing.
- Mathematical Optimisation - this include dynamic traffic assignment (DTA) models. A good review of dynamic traffic assignment and simulation models can be found in [29]
- Artificial Intelligence (AI) - neural networks are an example of AI approach. They have received a lot of attention in terms of transportation systems applications. Some example include traffic signal control, traffic flow modelling and transportation planning[10, 18, 42].

The advantages and disadvantages of the listed types and models are summarised in table showed in figure 2.2 below.

From the above the most widely used and well defined are the statistical models. From them the historical approaches are relatively easy for implementation and have fast execution speed, but have difficulty with dealing with incidents. Time series models have many applications and are well formulated. Because of this reasons and also the nature of the data that has been made

Type	Model	Advantages	Disadvantages
Statistical models	Historical Profile Approaches	-Relatively easy to implement -Fast execution speed	-Difficult to respond to traffic incidents
	Time series models -ARMA/ARIMA -State Space/Kalman filter	-Many applications -Well-defined model formulation	-Difficult to handle missing data
	Nonparametric regression -Dynamic clustering/pattern recognition	-Pattern recognition -No assumption of underlying relationship	- Complexity of search for “neighbors”
	Hybrid models -Clustering+linear regression -ARIMA+SOM -Fuzzy logic+GA	-Smaller and more efficient network	-Not yet many implementations
Computer simulation	Traffic simulation	-Possible to simulate various situations	- Requires traffic flow prediction in priori
Mathematical optimization	Dynamic Traffic Assignment	-Various types of models available and well known	-Not suitable for micro-simulation
Artificial Intelligence	Neural Networks	-Suitable for complex, non-linear relationships	-Forecasting in black box -Training procedure

Figure 2.2: Traffic forecasting types and models [51]

available (detailed description of which can be found in Chapter 5) for this project we will focus our attention on time series analysis for the rest of this report.

## 2.5 Time Series

Time series is a sequence of data readings taken during successive time intervals [41]. This could be a continuous recording of readings or a set of discrete readings. In the context of our project we have the continuous process of bus readings (generated by the AVL system) being transmitted which generate a discrete set of observations. This results in a data set of measurement values which consists of the actual values with some noise. Time series data contains four main components (illustrated in figure 2.3)[7]:

- **Trend** - this is the long term pattern that the given time series data follow. The trend can have positive or negative value depending on the data exhibits an increase or decrease respectively in the long term patter. Time series data with no trend (it does not show nether increase or decrease) is said to be stationary.

- **Cyclical** - this when we can see that the data show up and down movement around a given trend is referred to as cyclical pattern. Main characteristic of the cycle is its duration which can depend on the type of measurement.
- **Seasonal** - seasonality occurs when the time series exhibits regular repetitive fluctuations. For instance, temperatures peak during summer months (in the northern hemisphere) and drop during winter months.
- **Irregular** - also known as the error component. These are random increases or decreases for a specific time period.

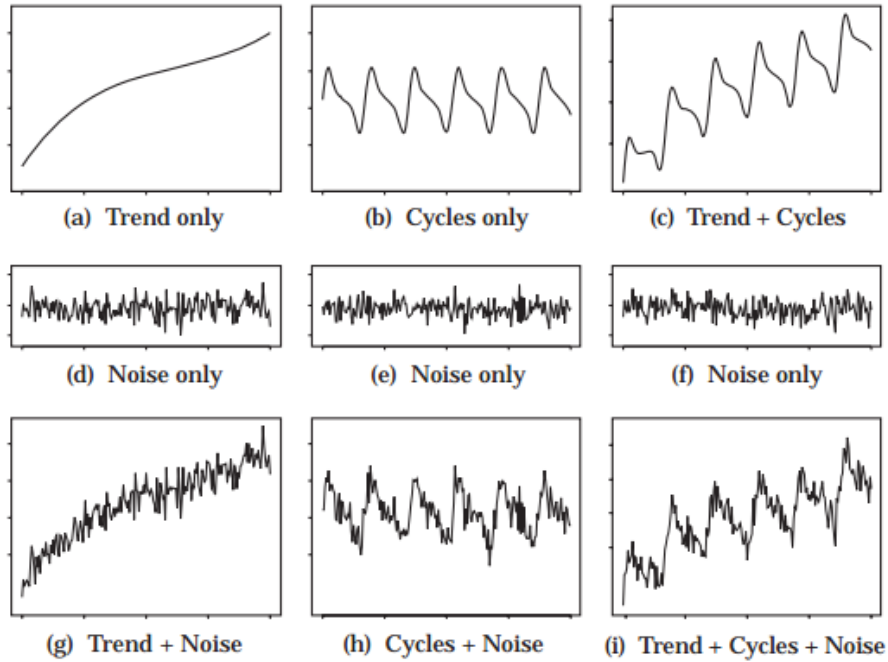


Figure 2.3: Time series data components

Analysis of time series data could be performed in order to extract and calculate some meaningful statistics from the data [41]. This could also result in producing a forecast of the data of interest for the next (future/unobserved) period of time based on the past observations. In order to highlight trends and make predictions we need to employ time series analysis techniques. Below I have presented some of the available techniques that could be employed when

analysing time series data. It is not an exhaustive review of all available methods and models as I have tried to keep the discussion relevant to this project.

### 2.5.1 Moving Averages

One technique commonly used in time series analysis is moving averages which is a form of smoothing. Smoothing means to dampen the effect of noise and irregularities in the original time series. Moving average also called rolling or running average is a statistical calculation method. It helps to analyse data series by calculating a series of averages of subsets of the data. Moving averages is commonly used in time series data analysis when the data is fairly stable and does not have significant trend, cyclical or seasonal effects. It can be used in order to smooth out a time series data with the aim of highlighting or estimating the underlying trend of the data. The other main usage is as a forecasting method again for time series. The main strength of these methods is that they are easy to understand. Moving averages are often used as the building block for more complex time series analysis. Below we present some of the main types of moving averages that are used in practice. [8, 41]

#### Simple Moving Average

Simple moving average (SMA) is calculated by adding all the observations for a given period of time and dividing this sum by the number of observations. This is popular statistical technique which is mainly used to calculate the trend direction. The simple average is only useful for estimating the next forecast when the data does not contain any trends. Each observation is weighted equally. If we consider shorter period window (meaning we only consider less observation e.g. only the last 5 or 10) for our averages they would react quicker to changes. While if we work with bigger period windows the averages would have greater lag. The equation for calculating SMA is given below in equation 2.1. In this  $n$  is the size of the window (e.g. the number of reading we are considering) and  $Value(i)$  is the actual value of observation  $i$ .

$$SMA = \frac{\sum_{0 \leq i \leq n} \text{Value}(i)}{n} \quad (2.1)$$

In table on figure 2.4 we can see sample time series data with window size  $n$  equal to 3. The graph in figure 2.5 shows the plotted actual observation values and the SMA calculated predictions.

Time period	Observation Value	SMA(n=3)
2007	10	N/A
2008	12	N/A
2009	16	N/A
2010	13	12.67
2011	17	13.67
2012	19	15.33
2013	15	16.33
2014	20	17.00
2015	22	18.00
2016	N/A	19.00

Figure 2.4:

Another form of SMA is centered moving average (CMA). Both are very similiar in terms that they use the same method for calculating the average value, but the differ in that the CMA calculates an average of  $n$  periods' data and associates it with the midpoint of the periods. An example can be seen in figures 2.6 and 2.7.

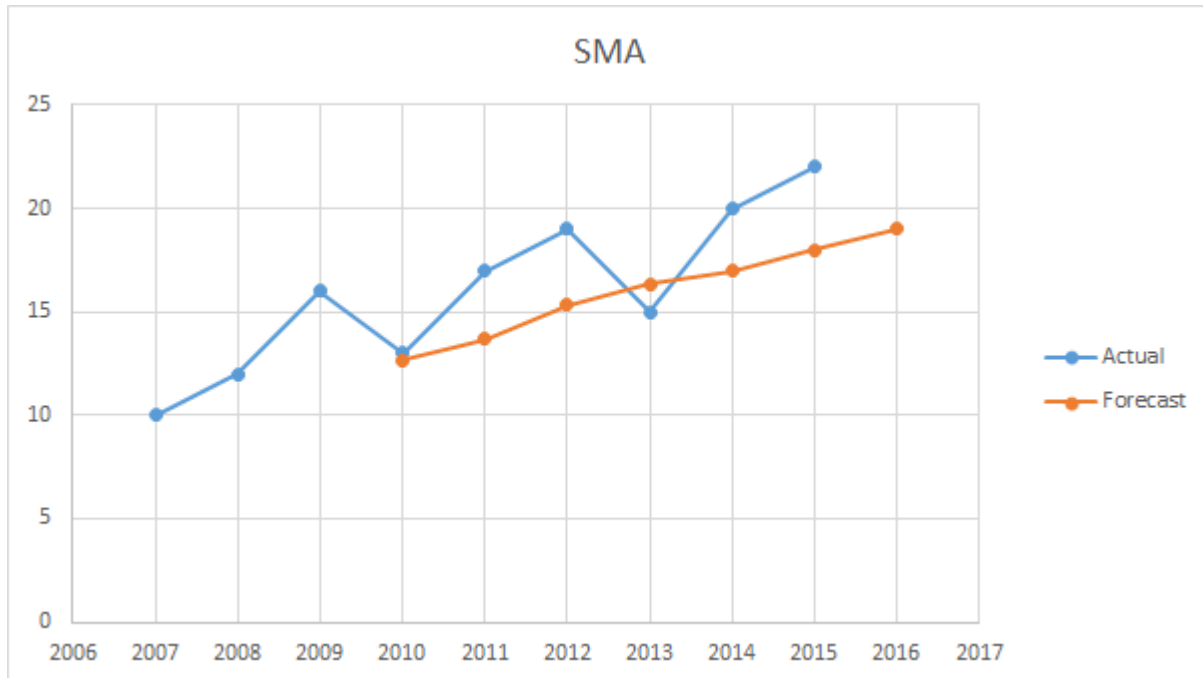


Figure 2.5:

Time period	Observation Value	CMA(n=3)
2007	10	N/A
2008	12	12.67
2009	16	13.67
2010	13	15.33
2011	17	16.33
2012	19	17.00
2013	15	18.00
2014	20	19.00
2015	22	N/A
2016	N/A	N/A

Figure 2.6:

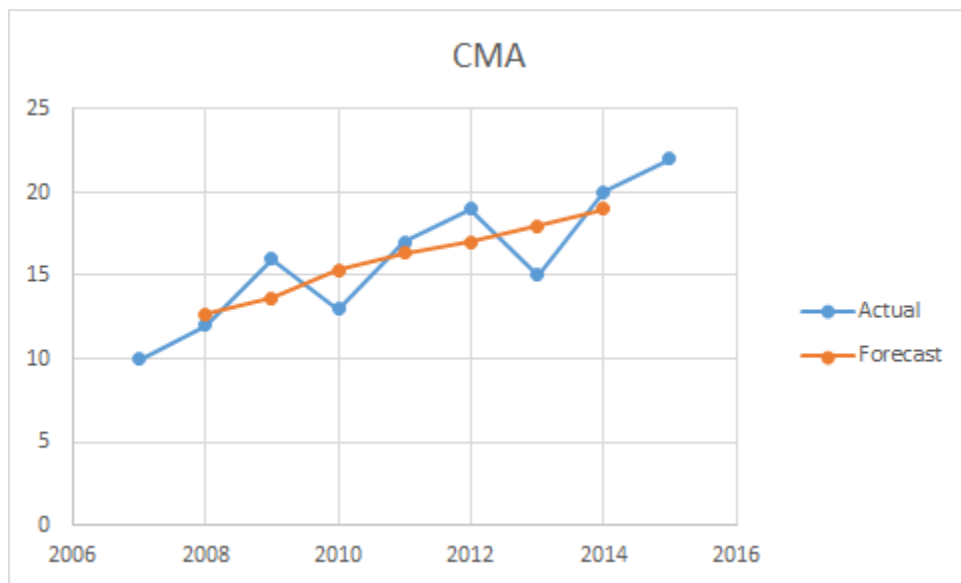


Figure 2.7:

### Weighted Moving Average

The problem with the simple moving average is that it weighted all data points equally. Meaning that both older and newer data would have the same effect on the average. This however is not the case when using weighted moving average (WMA). In WMA model each data point would be weighted differently according to period of when the observation was made. For example if we consider a  $n$  period moving average we can calculate the weight for the value taken in period  $i$  where  $0 \leq i \leq n$  by the following formula:

$$\text{Weight}(i) = \frac{2i}{n(n+1)}$$

This would mean that recent data have bigger impact on the result. However it should be noted that weighting formula given is only an example as it is the most natural and widely used weighing scheme for WMA. It is possible to use different weighting formula one example could be:

$$\text{Weight}(i) = \frac{2^i}{\sum_{0 \leq x \leq n} 2^x}$$

this would result in putting more weight on more recent data (e.g. older data is having less effect). The general equation for calculating WMA is given below as equation 2.2 where  $n$  is the number of observations (the size of the window).

$$WMA = \frac{\sum_{0 \leq i \leq n} (\text{Weight}(i) \text{Value}(i))}{\sum_{0 \leq i \leq n} \text{Weight}(i)} \quad (2.2)$$

An example of the application of WMA is shown in the table in figure 2.8. The forecast is plotted against the actual values and is shown in the graph in figure 2.9



Time period	Observation Value	Weight moving total (n=3)	WMA (n=3)
2007	10	N/A	N/A
2008	12	N/A	N/A
2009	16	N/A	N/A
2010	13	82	13.67
2011	17	83	13.83
2012	19	93	15.50
2013	15	104	17.33
2014	20	100	16.67
2015	22	109	18.17
2016	N/A	121	20.17

Figure 2.8:

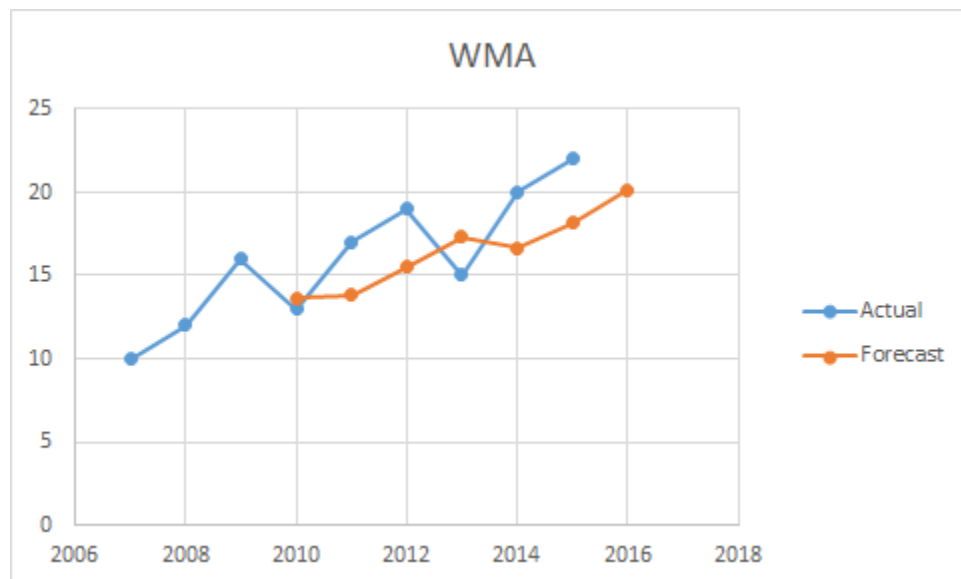


Figure 2.9:

## Exponentially Weighted Moving Average

Exponential smoothing was first suggested by Robert Goodell Brown [17]. Exponentially weighted moving average (EWMA) or also called exponential smoothing or simply exponential moving average (EMA) is very similar to WMA. The main difference is that in order to calculate it we do not need to keep all the data, but we could only store the latest value and the previous forecast only. Exponential moving average weights the data points exponential which means that the oldest data would have minimalistic effect on the result. There exist few exponential smoothing techniques including single, double and triple exponential moving average. Equations 2.3 and 2.4 give the simplest form for calculating single exponential smoothing. In this equation  $\alpha$  is called the smoothing factor and it is usually a value between 0 and 1. The closer  $\alpha$  is to 0 have greater smoothing factor, but are less responsive to recent changes thus have greater lag. Values of  $\alpha$  that are near to 1 have less smoothing effect, but are very reactive to recent changes in the data.

$$EMA_1 = Value_1 \quad (2.3)$$

$$\text{for } t > 1, EMA_t = \alpha Value_t + (1 - \alpha)EMA_{t-1} \quad (2.4)$$

Example of the application of EMA with different values of  $\alpha$  (0.2 and 0.8) is shown in figures 2.10 and 2.11. From this simple example it can clearly be seen the effect the value of  $\alpha$  has on the calculated value. From the graph it can be seen that the smaller  $\alpha$  value of 0.2 has greater smoothing effect, but greater lag. While the bigger value of this constant increases the reactivity to recent changes, but produces less smoothed line.

Time period	Observation Value	Previous EMA ( $\alpha = 0.2$ )	EMA( $\alpha = 0.2$ )	Previous EMA( $\alpha = 0.8$ )	EMA( $\alpha = 0.8$ )
2007	10	N/A	10.00	N/A	10.00
2008	12	10	10.40	10	11.60
2009	16	10.4	11.52	11.6	15.12
2010	13	11.52	11.82	15.12	13.42
2011	17	11.816	12.85	13.424	16.28
2012	19	12.8528	14.08	16.2848	18.46
2013	15	14.08224	14.27	18.45696	15.69
2014	20	14.265792	15.41	15.691392	19.14
2015	22	15.4126336	16.73	19.1382784	21.43
2016	N/A	16.73010688	N/A	21.42765568	N/A

Figure 2.10:

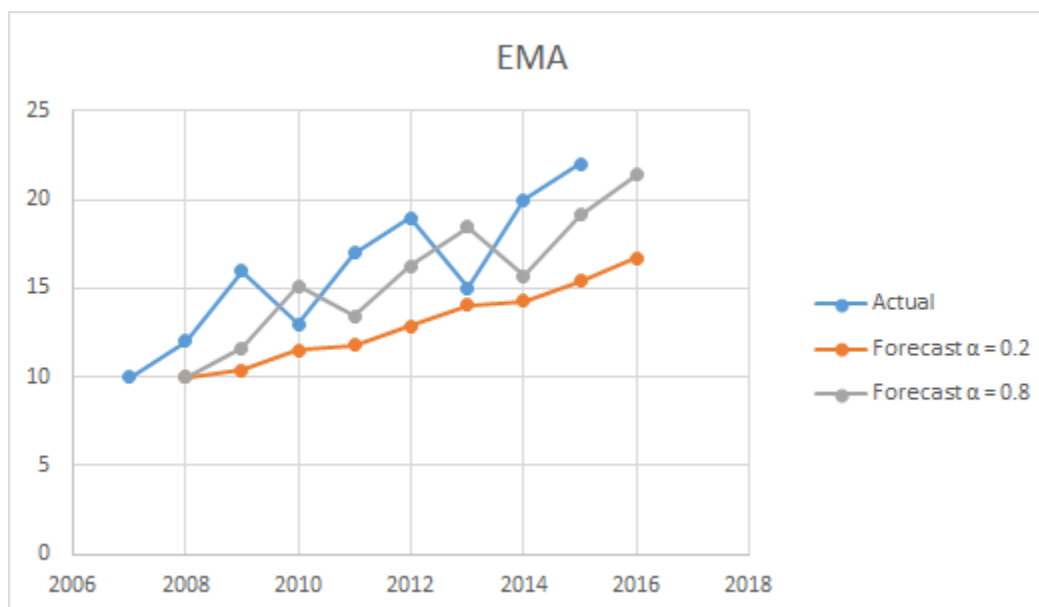


Figure 2.11:

## Summary

If we compare the presented moving average methods (SMA, CMA, WMA and EMA) we can clearly see that the SMA and CMA offer the most smoothing. However this comes with the trade-off of an increased lag (e.g. it takes longer to reflect recent changes).

The weighted moving average performance is influenced by the choice of window size as well as the choice of weights. There is no single rule what weights one should use and most often this is based on intuition and simulations in order to get optimal results.

As we have seen in the given examples The exponential moving average performance depends heavily on the chosen values for the  $\alpha$  constant. EMA offers the advantage of not having to keep all data point values in memory for the periods of our window. Whereas all the other presented techniques require us to specify a window size for our moving average and also to have the data for these periods available in memory. The choice of window size for the SMA, CMA and WMA has direct impact on the sensitivity (speed of reaction) of the method to changes. Increased size of the window results in less reactive moving average and increase in the opposite.

If a trend indication with better smoothing and little reaction for shorter movements is required then the simple average should produce the best results. However if smoothing is desired where you can still see shorter movements then it is better to use either WMA or EMA. Using either of the requires us to make some choices regarding what parameters (window size and weight for WMA and value of  $\alpha$  for EMA) to use in order to obtain best results.

### 2.5.2 Peak Detection

Detection and analysis of peaks (spikes) and valleys in time series is important in many applications (e.g signal processing, bioinformatics and many more). Peaks and valleys usually represent either significant events or errors in time series data. In our domain we are mostly interested in detecting high sudden changes in the traffic congestion conditions (e.g along route/sections in the bus

network). Peaks could be easily identified by visualising the data, however we are interested in automating this process. In the literature has many examples of peak detection application one such for example is [40] used for spike detection in microarray data (another example could be found in [3]). We do not go into details of any particular algorithms here because of the time constraints of this project (a study peak detection algorithms could be found in [46]). However it is worth to note that spike detection and analysis could be used to classify events that are detected. A simple example could be to distinguish if given detected congestion/disruption is incident related (e.g. sudden) or it is general traffic jam (e.g. rush hour).

## 2.6 Other

Other approaches that could be used include Kalman filtering [26] [20], Markov Chains [33] [37], Machine learning [21], Bayesian networks [49].

## 2.7 Summary

In this chapter I have aimed to provide the reader with broad view of the background of our problem and its domain. Detailed overview of the operations and the technologies and work flows used by TFL's bus command and control centre has been provided. This has led to detailed review and discussion of the related work that has been found in the literature. The chapter concluded with discussion on some of the available approaches. The exact approach taken and any implementation details are described in detail in Chapter 5 along with presentation and discussion of the data that has been provided by TFL for this project.

## Chapter 3

# Specification

In this chapter I have introduced and formalised the user and functional requirements. This is an important step of any project, especially computer science project, as it formalises the problems that the project is trying to address as outlined by the project aims and objectives in chapter 1. It also allows to be used as a measure for evaluating the success of the project once it has been completed. The requirements presented below have evolved and have been refined throughout the project lifetime in response to feedback and discussions carried out with the key stakeholders.

### 3.1 User Requirements

The user requirements provide a list of the functionalities that the user(s) expects to be able to perform and see the according results, in the end product. These are what is expected from the system, but are not concerned how they are designed or implemented. The main user requirements are listed as follows:

1. The tool must be able to produce a prioritised list of the disruption in the bus network that it has knowledge of.
2. The tool must be prioritising the disruptions according to the user defined rules (these are still discussed and gathered from the user) The tool must

be updating this list of disruptions whenever there is more data. This should happen as real-time as possible.

3. The tool must be able to provide detailed information for every detected disruption. This has to include the specific section and route that are affected and its severity.
4. The user must be able to interact with the system in order to lower or increase the priority of a given disruption (even ignore one).

## 3.2 Functional Requirements

These requirements specify in more details what the expected behaviour and functionality of the system/tool is. They are built on top of the user requirements as an input and are detailed list of what the system should be able to accomplish technically. The tool/system must:

1. Have an appropriate and useful representation of the bus network in order to be able to monitor and detect problems in it.
2. Be able to read and process CSV files as this is the primary input of the AVL feed files (more detailed discussion on the exact input and its format is presented in Chapter 5).
3. Listen/monitor for new incoming data and process it in real-time.
4. Be able to update itself whenever new data is detected and processed.
5. Be able to run without intervention 24/7.
6. Keep track of the disruptions detected and track how they evolve and develop.
7. Be able to keep information for a given window of time (e.g. data feeds from the last 2 hours).
8. Be able to output a prioritised list of disruptions.

9. Visualise the generated output appropriately. It should be compatible to run under Firefox or Internet Explorer as this are the main browsers used by CentreComm/TFL staff.
10. Display on request detailed information for the requested disruption. This should include a graph representing the route/section average disruption time.
11. Be easily configurable and maintainable.



## Chapter 4

# Design

For the purpose of the successful completion of this project I have decided to employ agile software development methodologies with evolutionary prototyping. The reason for taking this approach are the strengths of the agile software developmental methodology which is that it is incremental, cooperative, flexible and adaptive [28]. Our project is addressing an issue which does not have well specified set of requirements and the clients do not have a clear view of what they actually expect. This led us to use evolutionary prototyping [9] as this helps minimise the impact of misunderstanding or miscommunication of the requirements. The risk of which are relatively high as the goals of this project are relatively new and there is not much similar work done. This technique would also give better idea of what the end product would be capable of and would look like to the client. With evolutionary prototyping the system is continuously refined and improved. Each iteration builds on top of the previous thus meaning that with each increment we add more functionality and features or/and refine/improve what has already been implemented. Simply stated this means that with each iteration we are one step closer to the end product. This allows us to add features which were not previously considered or remove ones that are no longer viable or needed. In addition this approach also allows us to engage with the key stakeholders very early in the project life-cycle. This would provide us with valuable feedback which again brings a

lot of advantages.

- The delivery of the tool is speeded up and also minimises the risks of failing to deliver a working product before the project deadline [9].
- Users would engage with the product early in the project lifetime. This however poses some risk like the users requesting more features which were not previously mentioned or discussed. This means that we need to maintain some balance as this project has a fixed deadline and limited resources.
- Increased chances of fully understanding and meeting the user requirements and expectations from the tool.

Each increment (iteration) consists of the following stages:

1. Requirements specification & refinement
2. Design
3. Prototype implementation & Testing
4. User testing and feedback provision
5. Evaluate

The requirements and specification step would document what the tool should look and work like at the end of each iteration. These would be following the aims and objectives we have defined in Chapter 1 of this report. After a prototype has been implemented and tested by the user we will evaluate the progress and make any changes in our requirements, design and project plan accordingly. Thus after a number of iterations we should have a prototype which to resemble the desired product as required by the user.

In the below sections I have presented the system at an abstract level. This follows from careful analysis and consideration of the requirements stated in the previous chapter. I have tried to highlight all major key components and classes which are described formally. This allows us to better organise and

structure our problem. The diagrams presented follow the unified modelling language (UML) paradigm [32] and are platform-independent models (PIM) of the system. This allows us to focus on the design of the system itself without distracting our attention with platform specific decisions. Once these models are created we can then easily transform them into platform specific models (PSM) using the desired technologies.

## 4.1 Use cases

Based on the user requirements stated in the previous chapter a use case diagram has been derived and presented in figure A.1 in Appendix A. The use case diagram depicts the way users(actors) interact with the tool (system). There are three types of users of the system (every next type is extension of the previous as it can be seen from the diagram):

- **Normal** users are people with general access to the system. They have read-only right and can interact with the system by requesting a list of disruptions and more details for a selected from them disruption. They have also access to the disruption history of the network. This use case was not in the initial requirements and was identified as a useful feature during demonstrations of the prototype to TFL.
- **Operator** is assumed to be a CentreComm staff member who is required to authenticate into the system. This would allow them the extra functionalities of adding comments to disruptions for others to see. Such users also have the functionality, of hiding and showing disruptions that are currently detected in the bus network, available to them.
- **Administrator** users have the most privileges of all type of users of the system. In addition to the above actions they are also allowed to view and change the configuration settings of the system. This is to allow easier configuration of the application.

## 4.2 System architecture

In figure A.2 in Appendix A the architecture diagram of the system could be seen. The overall architecture is following a four-tier architecture with an model view controller (MVC) pattern for the user interface. This architecture allows us to decompose to system into separate subsystems where lower tiers do not depend on higher tiers. This system allows for the implementation details of subsystems to be changed without affecting other components if the interfaces do not change.

As it can be seen from the diagram the system is divided in four main layers:

- Representation Layer - responsible for visualisation of the user interface. It consists of a number of user views.
- Representation Control Layer - responsible for the control/transition between user interface windows. In this layer I have made use of the front controller pattern [16] as it allows us to combine the common logic in one controller.
- Application Logic Layer - this layer is the functional core of the system. This is where all the business logic is encoded and corresponding calculations are done. The most important part of the system is the Disruption Engine which is responsible for:
  1. Monitoring for new feeds and processing them.
  2. Updating the bus network status (calculating delays and detecting disruptions).
  3. Writing the changes to the system database.

The Disruption Model is the other major component of the system. It is responsible for retrieving the disruptions and their details from the database and providing them to the representation control tier.

- Data Layer - this the the data repository layer. It consists of comma separated values (CSV) files representing the AVL feeds that are being

pushed to the system. Here we also have the system database which contains all the configuration settings parameters and the output of the engine (disruptions and their details that are detected by the tool).

This architecture diagram represents coarse grained view of the system to be implemented. Each of the components presented above could be implemented as a number of smaller components and modules depending on the specific technologies.

### 4.3 State Machine

State machine diagrams are useful in explaining what in what states the system could be and how it transitions from one state to the other. The state machine diagram of the disruption engine component of our system is presented in figure A.5 in Appendix A. This diagram represent the states in which the disruption engine could reside and the available transitions. As it can be seen from the figure once the engine is initialised correctly it enters a continuous loop. This loop represents the engine waiting for new feeds to be detected by the system. Once detected they are processed and the bus network state is updated. If the change of the bus network has not changes from what was previously observed the tool does not make any changes to the database, else it would write all the changes that have been detected and calculated. In case the system fails to connect or write the changes to the database the system terminates producing the appropriate error alerting the maintenance personnel. This behaviour is appropriate for this project as its scope is to produce a proof concept working prototype rather than a fully functional ready to deploy production tool.

### 4.4 Class organisation

Class decomposition diagrams are the main building block of object oriented programming paradigm. They are widely used tool for the organisation and design of a software system. The diagram consist of classes, which encapsulate

some attributes (state) and methods (functionalities), and the associations between the individual class. Each class is depicted by a box with the name of the class on top and its member attributes and methods below. Associations are represented by lines connecting those classes where a relation exists. In figure A.3 in Appendix A we have depicted the class diagram of the disruption engine component from the architecture diagram. We have only presented the class diagram of this component as this is the component which captures the main business logic with regards to detecting the disruptions in the bus network. In the below subsections I have provided some explanation of the most important classes in the class diagram.

#### **4.4.1 iBusMonitor**

This could be viewed as the entry point of the tool engine. Its most notable attribute is the link to the configuration file specifying the database connection properties. It can have a number of bus network and is responsible for initialising the tool and monitoring. This means it encapsulates the logic for listening for new feed files being published.

#### **4.4.2 BusNetwork**

This class represents a given bus network. In the context of this project this can be viewed as the TFL bus network. This object encapsulates all attributes that relate to the network state and its behaviour. Each bus network consists of at least one bus stop and at least one route otherwise it does not make any sense to have a network without any stops or routes.

#### **4.4.3 BusStop**

The bus stop class is representation of a bus stop in the bus network. It consists of a number of expected attributes that a stop would have. We also make the assumption that a single stop can belong to only one bus network.

#### **4.4.4 Route**

This class is one of the most important in the context of this project as it encapsulates the state of a single route in the network. Each route is associated with at least one run (in most cases each route would have two runs In/Out-bound) and a number of observations.

#### **4.4.5 Observation**

The observation class captures the state and functionality of a single observation. By observation we mean a single reading extracted from the AVL data input. This reading is expected to be coming from a single bus logged on a given bus route, thus it would belong to this route.

#### **4.4.6 Run**

This object represent a route's run state and methods. Its main properties consist of list of consecutive readings made on this run for each logged bus. It also provides interface for detecting and updating disruptions on this run, thus it needs to keep track of the disruptions that were previously seen along this run.

#### **4.4.7 Section**

This is the most basic part of a bus route apart from the bus stop. Each section represents the part of the route between two consecutive bus stops along this route. This means it is characterised by a start and end stop and the sequence of this section along the route. In this class we calculate the delay per individual section (more on how this is done in the following chapter).

#### **4.4.8 Disruption**

This class simply captures all the attributes of a disruption. Each disruption would have an identification number, sections between the disruption is observed and the corresponding delay and trend. It also provides methods for

updating and saving the details to the database.

## 4.5 User Interface

The user interface of our system is addressing the second main aim of this project the one of visualising the calculated list of delays. The design and rationale behind the user interface has been developed throughout discussions and meetings with CentreComm staff. The main requirement for the design is to be easy to identify the most important issues in the network. One of the uses of the user interface would be on a large screen walls which to be used not only by CentreComm, but by Traffic Management and even the Metropolitan Police. Another usage would be by the individual operators to access it through their personal computers. This has led to the decision of using a web based application for the purpose of satisfying those requirements. Using web rather than a standalone application we allow for our system to be accessible from any device capable of running a browser (e.g. computer, laptop, smart TV, smart-phone, tablet etc.). Another advantage is that we only need to deploy the web application once and it can be universally accessed through the local network or even through the internet.

In order to improve separations of concerns we have also decided to have separate application for the disruption engine and the user interface. This means that we can change each one without affecting the other (considering we maintain the correct interfaces). This also allows us to implement and add more user interfaces apart from the web application if needed. For example we may later want to create a dedicated mobile (tablet or smart-phone) application using the output from the disruption engine. This separation also allows to have different dedicated specialised people for maintaining each of the applications.

We have decided to use tabular approach for visualising the prioritised list of disruption. Each entry in the list would give detail of the route and section which are delayed. It also provides the time when the disruption was first



detected. Any additional information is only provided on request from the user.

The overall architecture of the structure of the user interface can be seen in the architecture diagram in figure A.2 in Appendix A. An early mock-up of the graphic user interface can be seen in figure 4.1. This has however evolved a lot through the project. In chapter 5 below we will give detailed description of the implemented visualisation.

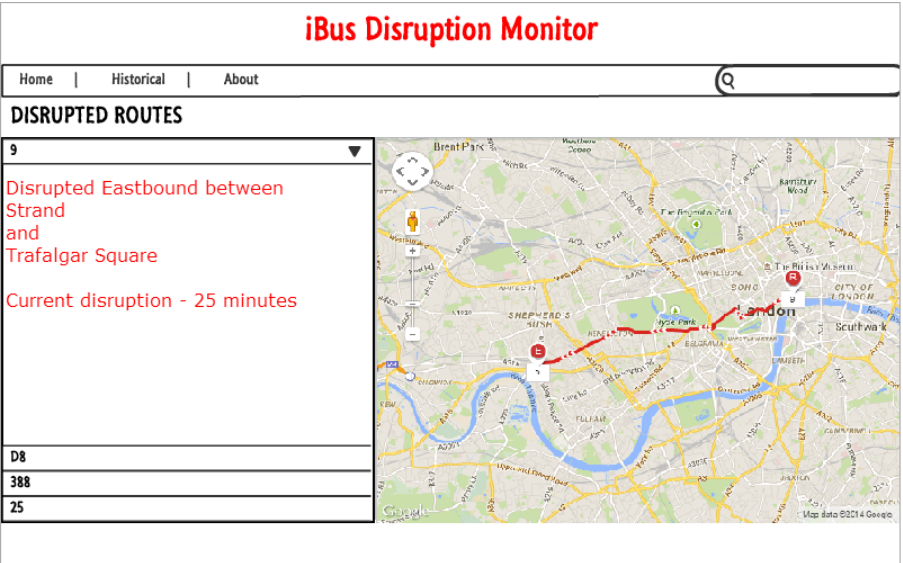


Figure 4.1: Initial GUI mock-up

## Chapter 5

# Implementation

This chapter aims to present the reader with explanation of the key implementation aspects. These include major challenges, decisions and problems that have been encountered and taken during the course of this project. I have tried to avoid going into too much technical details except where this is essential and provides the reader with better insight and understanding of the material.

The system I have developed as a prototype which to satisfy our aims consists of two sub systems. These are the disruption engine which address the first aim of detecting disruptions in the bus network and a web front end application which to visualise the calculated disruptions. Below I have presented the major implementation decisions and challenges that were faced during the development of this system. We begin with first describing the data that has been provided by TFL for this project. Afterwards we go into discussion of the implementation of the disruption engine which captures the core functionality and business logic of the system. Then we cover the visualisation part which can be viewed as an extension of the disruption engine.

### 5.1 iBus AVL Data

The data that is used for this project is provided by the Technical Service Group (TSG) at TFL. This data consists of comma separated value (CSV) files.

There is an individual file for every different bus operator which contains the data for all buses currently operated by this company. Initially every bus in the network would transmit its unique identification number and GPS coordinates approximately every 30 seconds [23]. This information is then preprocessed by a central server. This results in more information being derived as the central server has knowledge of the whole network and the bus schedules and headways. This results in the CSV feed files that have been provided to us for this project. An example of the content of the raw feed file and a formatted version is presented in figure 5.1 and 5.2 respectively. Below I have provided a detailed explanation of each field in these files [39].

- **Vehicle Id** - this is a unique id of the vehicle.
- **Bonnet Code** - this is the bonnet code of the bus.
- **Registration Number** - this is the number of the registration plate of the bus.
- **Time of Data** - this refers to the date and time of when this data is received from the respective vehicle.
- **Base Version** - this the version of the system that is run by the respective bus.
- **Trip Id** - stores the internal trips id <sup>1,2</sup>. This would increment every time a bus starts new run either at end of its current run or if it is curtailed.
- **LBSL<sup>3</sup> Trip Number** - this is LBSL trip number<sup>1</sup>. This is similar to the Trip Id however this is a global trip id thus it is incremented whenever a bus in the network start a new trip.
- **Trip Type** - the type of the trip as follows:

1. - From depot to start stop of the block.

---

<sup>1</sup>This is only valid when bus is properly logged.

<sup>2</sup>Not available in route variant.

<sup>3</sup>London bus services limited [45]

2. - To new starting point.
3. - Normal trip with passengers.
4. - From the last stop of the block to the depot.
5. - Without passengers.
6. - Route variant.
7. - Vehicle not logged in either block or route.

- **Contract Route** - the route name<sup>1,2</sup>.
- **Last Stop ShortDesc** - this is the lbl code of the last stop visited by the respective bus<sup>1,2</sup>.
- **Schedule Deviation** - this is the standard deviation from the schedule, calculated using the bus position telegram, for the respective bus<sup>1,2</sup>.
- **Longitude** - this is longitude of the place from where the vehicle is sending the telegram<sup>4</sup>.
- **Latitude** - this is latitude of the place from where the vehicle is sending the telegram<sup>4</sup>.
- **Event Id** - the last event Id.
- **Duration** - currently not being populated.

The information we are interested is the deviation from the schedule. This value is calculated the same way for both low<sup>5</sup> and high<sup>6</sup> frequency buses by knowing the bus schedule. It must be noted that it is possible vehicles would start the route run already with some deviation from the schedule.

---

<sup>4</sup>GPS raw data divided by 3,600,000.

<sup>5</sup>Less than 5 buses per hour.

<sup>6</sup>5 or more buses per hour.

```

1 VEHICLE_ID;BONNET_CODE;REGISTRATION_NUMBER;TIME_OF_DATA;BASE_VERSION;BLOCK_NUMBER;TRIP_ID;LBSL_TRIP_NUMBER;TRIP_TYPE;CONTRACT_ROUTE;LAST_STOP_SHORT_DESC;SCHEDULE_DEVIATION;LONGITUDE;LATITUDE;EVENT_ID;DURATION;
2 18351;AE11;CJ61SUL;2015/02/19 14:01:52;20150213;85004;457551;49;3;298;17775;-60;-0.12806;51.63210;0;;
3 18356;AE16;SB61SUL;2015/02/19 14:01:47;20150213;85002;451147;45;3;298;OC04;-1080;-0.19187;51.69676;0;;
4 18366;TPL927;EY03FNL;2015/02/19 14:01:00;20150213;-2147483645;-2147483645;-2147483645;7;UL8;SN;-2147483645;-0.20822;51.49595;0;;
5 18357;ELV2;PN02XCR;2015/02/19 14:01:39;20150213;-2147483645;-2147483645;-2147483645;7;UL8;RR15;-2147483645;-0.20710;51.49124;0;;
6 18354;AE14;KS61SUL;2015/02/19 14:01:51;20150213;85001;451047;43;3;298;OC04;-2280;-0.19187;51.69676;0;;
7 18355;AE15;TW61SUL;2015/02/19 14:01:07;20150213;85005;455779;41;3;298;OC06;-3300;-0.19859;51.69806;0;;
8 18352;AE12;DS61SUL;2015/02/19 14:01:40;20150213;85005;455780;48;3;298;16915;-1050;-0.15202;51.65436;0;;
9 18353;AE13;KR61SUL;2015/02/19 14:01:12;20150213;85003;455689;47;3;298;994;-1200;-0.12851;51.63257;0;;

```

Figure 5.1: Sample Raw iBus Data

	VEHICLE_ID	BONNET_CODE	REGISTRATION_NUMBER	TIME_OF_DATA	BASE_VERSION	BLOCK_NUMBER	TRIP_ID	LBSL_TRIP_NUMBER	TRIP_TYPE	CONTRACT_ROUTE	LAST_STOP_SHORT_DESC	SCHEDULE_DEVIATION	LONGITUDE	LATITUDE	EVENT_ID	DURATION
1	18351	AE11	CJ61SUL	19/02/2015 14:01	20150213	85004	457551	49	3	298	17775	-60	-0.12806	51.6321	0	
2	18356	AE16	SB61SUL	19/02/2015 14:01	20150213	85002	451147	45	3	298	OC04	-1080	-0.19187	51.69676	0	
3	18366	TPL927	EY03FNL	19/02/2015 14:01	20150213	-2147483645	-2.15E+09	-2147483645	7	UL8	SN	-2147483645	-0.20822	51.49595	0	
4	18357	ELV2	PN02XCR	19/02/2015 14:01	20150213	-2147483645	-2.15E+09	-2147483645	7	UL8	RR15	-2147483645	-0.2071	51.49124	0	
5	18354	AE14	KS61SUL	19/02/2015 14:01	20150213	85001	451047	43	3	298	OC04	-2280	-0.19187	51.69676	0	
6	18355	AE15	TW61SUL	19/02/2015 14:01	20150213	85005	455779	41	3	298	OC06	-3300	-0.19859	51.69806	0	
7	18352	AE12	DS61SUL	19/02/2015 14:01	20150213	85005	455780	48	3	298	16915	-1050	-0.15202	51.65436	0	
8	18353	AE13	KR61SUL	19/02/2015 14:01	20150213	85003	455689	47	3	298	994	-1200	-0.12851	51.63257	0	

Figure 5.2: Formatted Sample iBus Data

## 5.2 Disruption Engine

The disruption engine is implemented based on the design given in the previous chapter. The main implementation language used to the implementation is Scala. Scala is both functional and object oriented language [31]. It is a type-safe Java Virtual Machine (JVM) language [31]. This means that it is compatible with existing Java<sup>7</sup> code which allows for reuse of existing Java libraries. Scala was first introduced back in 2003, however it has been only in the past few years that it had gained more popularity. In addition Scala enables the programmer to write more concise and clear code than Java. The decision of using Scala has also been influenced by the fact that I have good knowledge of the object oriented programming paradigm as well as experience in Java. This allowed me to quickly learn Scala and put it into use.

The disruption engine need to have an accurate internal representation of the bus network. TFL's bus network and any other bus network usually consists of bus routes. Each bus route often has multiple runs (directions - e.g. inbound and outbound). In turn each run consists of a sequence of bus stops that the bus passes through. In addition to these typical bus network components for our implementation we also have the notion of a section. By this we mean a pair of consecutive bus stops along a given run of a given route. In figure 5.1 below we can see an example for route 15 outbound where we have

<sup>7</sup>[https://www.java.com/en/download/faq/whatis\\_java.xml](https://www.java.com/en/download/faq/whatis_java.xml)

depicted two section X (between Leman Street and Tower Of London Stops) and Y (from Tower Of London to Great Tower Street). This means that if we have  $n$  stops on a given run then we have  $n - 1$  sections on the same run.

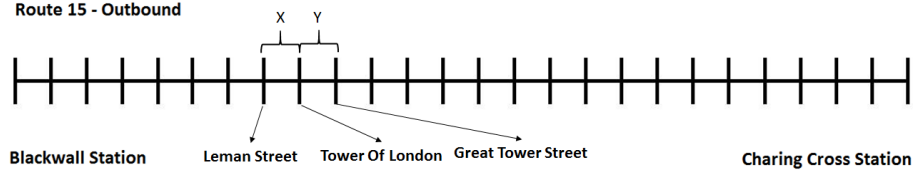


Figure 5.3: Example of a section

Our tool makes use of a PostgreSQL<sup>8</sup> database for storing all configuration parameters and also for storing the required information for the bus network representation. In figure A.8 in Appendix A I have presented the relational model on which the database is based. The initial versions of the prototype used a CSV files as a means for storing the output. The decision to switch the flat file storage for a database is based on a number of things. The most important being is that using a database rather than CSV files we could keep historical data of the detected disruption for future analysis which is accomplished much easier using a proper database. Another advantage for the database approach is that it offers better concurrency support out of the box. Unlike the flat files which if manipulated concurrently could result in inconsistent state. Also having a database means that our disruption engine would require only details for establishing a connection to the respective database where it can read all other information that it requires. Otherwise it would require a number of other files and parameters to be defined which is not as maintainable as having one single database containing all of the configuration parameters as well as data for the bus network representation. For the above reasons I have decided to use PostgreSQL as it is advanced open source relation database management system. PostgreSQL has very good document and community support which big advantage for any technology. Another advantage for using PostgreSQL for our implementation is that it ensures reliability and data in-

<sup>8</sup><http://www.postgresql.org/about/>

tegrity [27]. Also it supports Listen<sup>9</sup> and Notify<sup>10</sup> functionality which could be used for real-time updating of the web application along with Server Sent Events (SSE) [48].

### 5.2.1 Bus Network representation

In order for the engine to have full bus network representation it requires information of the bus routes and the bus stops in the network. This information is freely available to anyone on TFL website <sup>11</sup>. It consists of two CSV files, one containing information on all bus routes and one for all bus stops in the network. The bus file contains a list of bus routes respectively has a one or more runs which consists of sequence of bus stops. In our implementation we assume that this information is preloaded into the database. This pre load consists of simply extracting the information from the CSV files obtained from TFL website into the respective tables (BusStops and BusRouteSequences). In addition to this we need to pre load also the sections which will allows us to store individual section information. Currently sections are being generated manually however this could be easily automated, but this is not in the scope of this project. Then once the engine is started it would read and load from the database all bus routes and respective sections. This results in the BusNetwork class storing a HashMap which maps a bus route name to the corresponding Route object. In turn the Route class maintains a list of all runs for the respective route. BusStop information (apart from the bus stop LBSL code which we use as an id) is only loaded from the database on request. This whole process is part of the initialisation of the monitoring tool along with pre loading some other environment configuration parameters from the database. It happens only once throughout the execution of the tool and takes place just after the engine is started.

---

<sup>9</sup><http://www.postgresql.org/docs/9.1/static/sql-listen.html>

<sup>10</sup><http://www.postgresql.org/docs/9.1/static/sql-notify.html>

<sup>11</sup><http://www.tfl.gov.uk/info-for/open-data-users/our-feeds>

### 5.2.2 Monitoring and processing new feeds

Once the system is initialised the tool will continuously monitor a specified directory (configurable from the database) for new feeds being written (pushed). Once new feed files are detected they are picked up and processed by the engine. The processing consists of extracting the data of interest and calculating the time lost by buses on average for each section. Extracting the data means reading the CSV feed file line by line. Each line would be a data (observation) for a given bus thus we associate each observation with the route that is currently logged on. Once the observations are extracted they are sorted by the time of the data field and any data that is older than a given predefined threshold (e.g. 120 minutes - this is configurable) is discarded. Once a given feed file is processed it is moved to a predefined processed feed directory. This completes the feed file processing step. Feed processing is performed on batches of feeds (see the state machine diagram on figure A.5 in Appendix A). This means that once the engine detects new feed(s) in the directory it is monitoring it will process all new feed files.

### 5.2.3 Bus network state update

Once feed processing has finished the system needs to update the bus network state. This consists of a number of steps. First we need to calculate the lost time per section. This process is done iteratively for every route in the network that has active buses (readings have been transmitted in some predefined interval of time e.g. 90 minutes). Each bus observation is then taken on a given route (see Figure 5.4). At least two or more observations are required in order to calculate the time loss for a section. In the example below (figure 5.4) we can take the first reading  $x_1$  and the second reading  $x_2$ . The difference in the

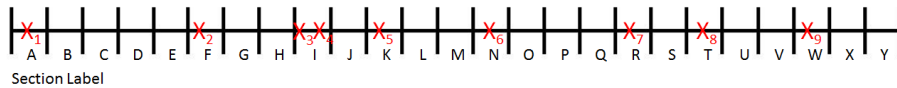


Figure 5.4: Example of observations

schedule deviation between  $x_2$  and  $x_1$  is then calculated (see table figure 5.5).



Reading	Schedule Deviation Value (Minutes)	Change
$x_1$	1	
$x_2$	7	6
$x_3$	10	3
$x_4$	15	5
$x_5$	21	6
$x_6$	20	-1
$x_7$	20	0
$x_8$	24	4
$x_9$	22	-2

Figure 5.5: Example schedule deviation calculations for the example in figure 5.4

Once we have calculated the schedule deviation change between two stops we need to assign it to the respective sections. From our example we can see that reading  $x_1$  has been sent somewhere from section  $A$  and that  $x_2$  from section  $F$ . This means that the bus have travelled through 6 sections ( $A$  to  $F$ ) and during that time it has lost 6 minutes. The problem is that we do not know where exactly this delay has happened. It is possible that there was a delay just in one of the section or it could have been distributed along all or few sections. For this reason what we do is to distribute this lost time evenly along all sections in-between the two readings. In our example this would mean that we would assign 1 minute to each of sections  $A$ ,  $B$ ,  $C$ ,  $D$ ,  $E$  and  $F$ . Then we take the next pair of observation in this case  $x_2$  and  $x_3$  and do the same procedure. Every time we have new time loss value for a section we add it to the existing lost time for the section and associate it with the newest observation time of the data. In this case this means that section  $F$  will have time loss value of  $1 + \frac{3}{4}$  and time-stamp equal to the time of data of observation  $x_3$ . If we however take the case of reading  $x_3$  and  $x_4$  we know for sure that the delay has occurred in section  $I$  thus we assign the 5 minutes lost time only to section  $I$ . Repeating the above steps for each bus on a route we end up with a

list of values representing the delay (time lost) and the time (time of the data of the latter observation in the examples given above this means we take the time of the data of  $x_2$  and  $x_4$  respectively) when this has occurred.

Once we have a list of these values for each section of each bus route we need to calculate the weighted moving average of this data. To do this we firstly need to sort this list of values for each section of a route run by the time of the observation in ascending order. Then we calculate the weighed moving average by assigning the oldest data weight of 1 and the newest data weight of  $n$  ( $n$  is the number of data entries of a section). Calculating the weighted moving average instead of a simple average we put more weight and the newer data and also help dampen the effect of a single irregularity (e.g. a bus has experienced a technical fault).

Doing the above steps we obtain a value representing the WMA time loss (delay) for each section (see figure 5.6). The next step then is to examine each route run and its sections in particular and check if any of them are disrupted. To accomplish this we firstly check if the total cumulative lost time, of the sections of the respective run, is greater than or equal to some predefine minimal threshold (this is configurable). In case this does not hold (e.g it is below the minimal threshold) the algorithm will move onto the next run and will consider this one as clear (without any problematic sections). Any negative values (e.g where buses have gained time) are treated as 0 as we are only interested if there are any problematic sections of the route.

However if for example we assume that the minimal threshold is 20 minutes and take the example from figure 5.6 we can clearly see that the cumulative lost time is greater than or equal 20. In such cases we need to look more closely at this route run and try to identify the sections which are causing the most delay. We do this by searching for a number of consecutive sections which have their sum of the delay time greater or equal to some predefined threshold (e.g. 20 minutes). However small delays (e.g of 1-2 minutes) are treated as 0 as there always some slight variations and we are only interested in major disruption which are beyond the control of individual bus operator companies.

If we take the example presented in figure 5.6 we can see that there seems to be some significant problem between sections *L* and *O*. The sum of the delay is  $7 + 14 + 12 + 5 = 36$  minutes which is greater than our example threshold of 20 minutes. This results in the engine detecting and outputting a disruption between those sections.

Lost time in minutes																			
0	0	1	1	2	1	1	2	0	1	3	7	14	12	5	2	0	0	1	3
A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
Section Label																			

Figure 5.6: Example of a disruption

However we have to limit the number of consecutive section we look at a single time as for example we consider the below case (see Figure 5.7). We can have a long route run with very small loss of time per section (as in the below example of the order of 1 to 3 minutes) but overall they might add up to 20 or more minutes in particular if the route run is long (e.g has 30 or more stops). This however does not represent a single problematic hotspot and is responsibility of the bus operators to manage such cases and not CentreComm. For this reason our algorithm would mark the start of a disruption as the first section it encounters with a greater WMA delay value and it would continue expanding this disruption until it encounter a section with value less than the minimal threshold for section time loss value. Then it checks if the total delay for the detected sections is greater than or equal to the minimal disruption threshold. If true it outputs it as a disruption affecting those sections. This is done for the rest of the route run even if some disruptions are already detected at the beginning of the run of a particular route.

Lost time in minutes																			
1	0	1	1	2	1	1	2	0	1	1	3	2	1	0	2	0	0	1	3
A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
Section Label																			

Figure 5.7: Example of no disruption

Every detected disruption or one that has changes its state is updated and written to the database. This allows for the user interface to update itself

with the latest information. The engine also writes a snapshot of the sections state of every route run in the network which has active disruptions. We only write data for the sections which belong to runs with disruption due to performance considerations. For example if we take the London bus network which has 680 routes most of which have 2 runs (outbound and inbound) and each run on average has 30 stops it means we have more than 40000 sections. Updating them each time would take significant computational time (e.g. couple of seconds depending on the machine and the location of the database) and also would be waste of space. For this reason our implementation only updates sections which belong to disrupted runs as this information is utilised for displaying detailed graphs in the user interface.

### 5.3 Graphical User Interface

[1. Technologies used and why those in particular; 2. What has been achieved and how?]

The engine implementation described above address only half of our requirements. In order to fully meet the project goals we need to be able to visualise the detected disruption in the network. For this reason we need to produce some kind of user interface which to be able to meet the requirements for visualisation of the detected bus delays. As described in section 4.5 of this report we have decided to have individual application for the back-end disruption engine and the user interface. For the implementation of the user interface we have chosen to create a web application using Ruby<sup>12</sup> as the programming language running on Rails<sup>13</sup> framework.

The Rails framework is full stack open source web framework implemented in Ruby which allows for quick development and deployment [19]. It is relatively easy to use and maintain and it supports a wide range of software engineering paradigms and patterns. The main ones being Model View Controller (MVC), Don't Repeat Yourself (DRY), Convention Over Configuration

---

<sup>12</sup><https://www.ruby-lang.org/en/>

<sup>13</sup><http://rubyonrails.org/>

(CoC) and the Active Record pattern [16]. Ruby is object-oriented scripting language which is famous for its conciseness. It enables the programmer to express his intentions quickly in very few lines code and it is also easy to read this code latter. Another advantage that Rails framework provide is the automatic code generators which provide code skeletons which save effort and time. Rails also allows for agile software development as this is in its roots [19]. It is also important to note that Rails has a strong support community and documentations as well as expansive list of add-ons and third-party libraries.

Other languages and framework have been also considered as well to be used instead of Ruby and Rails for this project. However the chosen ones provide us with technologies which are to be used for quick prototyping and agile development which our project is following. In addition to the above this was also an opportunity for me to learn a new language and framework which I have not used before.

I have also made use of Foundation<sup>14</sup> Cascading Style Sheets<sup>15</sup> (CSS) front-end framework. This framework has enabled quick prototyping as it has a rich library of predefined components. However its main advantage over some other similar frameworks is its notion of grid which allows for quick and easy implementation of responsive websites. Making the web application with responsive design allows us to reach more platforms and client systems with a single application. This allows us with little additional implementation effort and time to achieve results which look good on both large screens (desktop computers, laptops etc.) as well as small mobile devices (e.g. smart phones, tablets etc.). This framework also has the advantages of good support in terms of documentation and add-ons and it is lightweight.

The ruby on rails web application is implemented following the MVC pattern. MVC allows the programmer to structure their application code better with clear separation of concerns. In our case it consists of few simple models representing the corresponding database tables, controllers and a number of

---

<sup>14</sup><http://foundation.zurb.com/>

<sup>15</sup><http://www.w3.org/Style/CSS/Overview.en.html>

views. The models are implemented as Active Records<sup>16</sup> which provide interfaces to the respective database tables. They provide create, read, update and delete (CRUD) functionalities. The controllers are responsible for processing the client request by parsing and checking for any parameters, credentials where necessary and responding with the requested information. The responses are in most cases rendered views containing the requested information.

The user interface can be seen in figures ?? to ??. It consists of three main views which are the disruption view (figure ??), the history view (figure ??) and the settings view (figure ??). The main view is the disruption one which is responsible for visualising the disruptions in the network in any point in time. The list is displayed in tabular form (figures ??). This provides instant awareness of what the network state is to the CentreComm operators. Disruptions are prioritised by their severity and are color coded.

The application also has a basic authentication in place which enables to distinguish between the three type of users as described in section 4.1. Guest users do not have to authenticate, but they have limited access. They have read only view of the system. The system also allows users to login in (figure ??) and depending on their status they could either be operators or administrator users. The administrator users have full access to the functionalities which includes view of the settings and the ability to change them, this is not possible if you are guest or operator user. Logged in users are allowed to hide/show (the result is that guest users would not see hidden disruptions) disruptions (figure ??) and also to add comments (figure ??). All users including guest are enabled to request further information for a given disruption which includes a graph of the lost time on the disrupted route see figure ??. This graph is a combo chart<sup>17</sup> (combination of line and bar chart) and is created using Google Charts<sup>18</sup>. The bars of the graph denote the WMA delay per section while the line depict the cumulative lost time along the route (this treats all negative values as 0). The Google Charts API<sup>19</sup> allows us to make the graph

---

<sup>16</sup>[http://guides.rubyonrails.org/active\\_record\\_basics.html](http://guides.rubyonrails.org/active_record_basics.html)

<sup>17</sup><https://developers.google.com/chart/interactive/docs/gallery/combochart>

<sup>18</sup><https://developers.google.com/chart/>

<sup>19</sup>API - Application Program Interface

interactive such that on hover we can display more detailed information for the selected section which allows us to create a clear and organised layout with all the information available. This graph is very useful as it provides the CentreComm operator or other users a detailed view of what is the state of the given route. The real-time disruption list and the history table are by default sorted by severity of the section delay and the total delay. However users are enable to sort by any other column. This is achieved using Ajax <sup>20</sup> in order to minimise network load by only reloading the changed part of the web page (we do not need to reload the navigation menus, footers etc.). We also have added a data filter for the history view. The history view has not been part of the initial requirements or aims of the project, but it has been added latter in the project just as a proof of concept.

One important aspect of the user interface is how to update the disruption list in real time whenever there are changes. This is achieved by implementing Ajax short polling [6]. This means that we use asynchronous JavaScript on the client side to make request to the server at predefined intervals. Alternative methods for implementing the updating of the list include Web Sockets<sup>21</sup>, Ajax Long Polling [6], Server Sent Events (SSE)[48]. The main advantages of using Ajax short polling over Web Sockets and SSE are that Ajax Short Polling is supported by all major web browsers natively unlike SSE which lack Internet Explorer support and it is easy and quick to implement it. The drawback of using Ajax polling is that in order to achieve near real time update the clients will need to make frequent request to the server which wastes network bandwidth and server resources. SSE is probably the best approach in this scenario as it establishes a persistent long-term connection on which the server is able to push new data once it becomes available to all connected clients. However SSE is relatively new standard and it is has been standardized only as part of HTML5<sup>22</sup>. Currently Internet Explorer does not support SSE and this is a problem as this is one of the main web browsers CentreComm staff

---

<sup>20</sup>Short for asynchronous JavaScript and XML

<sup>21</sup><https://tools.ietf.org/html/rfc6455>

<sup>22</sup><http://www.w3.org/TR/html5/>

use. In addition to this SSE require the use of a concurrency enabled server. Web Sockets provide a persistent two way connection between the client and the server however in our case we are mainly interested in pushing new data from the server to the client. There is very little information the clients need to send to the server and thus we have excluded this approach as viable one.

## 5.4 Problems

During the implementation of the project a number of obstacles and problems have arisen. In this section we discuss the major ones. All that are not mentioned below are assumed to be solved by our implementation and not of enough significance for the reader.

The main challenge during the implementation has been how to assign the lost time between two consecutive observations to the respective sections that the bus has travelled through. This is because the data we work with is sparse (currently every 5 minutes) which means that during the time between those two observations the bus could have passed a number of sections and bus stops and we will only know the last bus stop it has attended. To solve this problem our the disruption engine keeps an ordered list of observations for each unique vehicle which is active on a given route. To assign the delays accordingly the engine takes two observations at a time and does a number of steps. The first step is to check if the observation have been made along the same run (e.g. both reading come from the bus when it was travelling outbound along the respective route). To do this we need check if the last stop from the earlier observation and the last stop from the latter reading are both from the same run and the earlier is preceding the latter last stop.

One problem that was encountered is that there is no explicit information in the data that has been provided if a bus is on diversion. Diversions vary greatly in terms of length (few or many stops) and duration (e.g. it can last half an hour or few weeks/months) of the implemented diversion . What happens in such cases is that the bus would still transmit its position however the



central iBus server would not be able to calculate the schedule deviation. This results in readings in the feed file which have abnormal values for the schedule deviation and other fields. These abnormal values however are always the same and are equal to the negative integer  $-2147483645$ . However this does not happen only when the bus is on diversion, it can happen if there is problem with the GPS (e.g. weak signal due to high buildings) of the respective bus or if the bus is not logged on properly in the iBus AVL system. This means that we are unable to know for sure what such readings mean. For these reasons when such readings are observed by the engine they would simply be ignored. The implications of this are that our implementation loses some accuracy and it may produce delayed alerts. If we consider the example shown on figure 5.8 where buses on this route are on diversion from section  $F$  to section  $P$ . In such scenarios what will happen is that we will get correct data before and after the diversion which would be processed. However during the diversion the data we will get from the buses currently travelling on this diverted path will have meaningless values and thus the disruption engine will ignore it. This means that if buses experience delay during the diversion this would only be picked up by the system once they return on the normal route. However the delay that the system will observe would be distributed along the whole diversion (as described above in section 5.1). This means that in case we have long diversion some significant delay might not be picked up by the system. This can also happen if there is problem during part of the diversion, but during the rest of the diverted route the buses actually manage to get back on schedule.

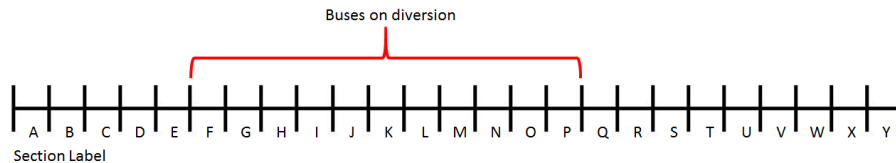


Figure 5.8: Example of diversion

Another problem with the available data is that we do lose some reading when the bus turns at the end of a run. We cannot simply assume that every bus would travel along a route from start to end and turn as there are occasions

when buses are curtailed<sup>23</sup> along the route. This can happen for a number of reasons including service regulation, heavy congestion/disruption along some sections of a route, bus driver exceeding the allowed by law work hours and more.

Our system relies on monitoring for changes in the schedule deviation value as indicator for the delays in the network. However there are some scenarios where the engine can detect delays where in fact there are no real life disruptions taking place. To understand why this can happen consider the case when a bus is initially ahead of schedule. This can result in the bus driver or operator company to deliberately force this service to lose time in order to be right on time. This is especially a problem for low frequency buses as they run according to fixed schedule and not based on headways. Thus we can claim that our implementation provides an upper bound of the disruption delays in the bus network.

Another problem of the approach taken is that because the system relies on monitoring the schedule deviation value and it does not treat differently calculations where during the first observation the bus is ahead or schedule and when it is already behind. This means that it is possible that the bus driver is losing time on purpose for service regulations. Thus we can assume that our tool provides an upper bound of the disruption.

Some other problem with the implementation include inconsistencies in the bus route sequence data. This however is affecting small number of routes and does not prevent our system from delivering a proof of concept. This issues has been discussed with members of the TSG and have been considered irrelevant for the purpose of this project. However it has been agreed that they would need to be addressed in case this is tool is put into production.

---

<sup>23</sup>To cut short.

# Chapter 6

## Testing

Testing is important and integral part of any software development project. It is necessary to carry out testing throughout the implementation as well as once the system is completed in order to make sure it functions correctly according to the requirements and the design.

### 6.1 Unit Testing

As described in the previous chapters the system is composed of a number of classes each of them representing different part of the bus network with its own functionalities and responsibilities. Throughout the development of the product a number of unit test have been carried out in order to ensure that each of the classes function correctly as per requirements and carries out the required task. Each module was tested separately before being connected to the rest of the system.

### 6.2 White Box Testing

Apart from the unit testing through the development of the tool we have made use of the debugging tool provided by IntelliJ IDEA<sup>1</sup> integrated development environment (IDE). This was the main IDE used for the development of this

---

<sup>1</sup><https://www.jetbrains.com/idea/>

project (both the disruption engine and the user interface web application) and has provided many useful tools and support. The built-in debugger allowed us to carry out inspection of the code during its execution while still developing the application in order to make sure that desired outcome is produced. The debugging tool allows us to halt the execution (in real-time) at specific point during the calculations and examine the state of the program (the values of all variables). This has proved very useful as it allowed us to quickly fix and rectify any issues early in the implementation which otherwise could have remained hidden.

## 6.3 Functional Testing

Functional testing is type of black box testing. This is because we test the software product by feeding it some input and examining the generated output. Functional testing does not mean we are carrying out test on individual methods/functions or classes/components. It means that we are testing parts of the functionalities stated in the requirements. As our system consists of two sub systems we are able to perform separate functional testing of each of the systems.

The disruption engine was tested by having an automated functional tests which feed the system with a feed file and then examine the produced results. Some of the example test cases include:

- Correct, but empty feed file is pushed to the system.
- Single correct feed file is pushed, containing a single bus reading.
- Two correct feed files are pushed containing a single bus reading from the same bus each.
- A malformed feed file is pushed to the system.
- A batch of feed files is pushed to the system.

The web application part of the system does not contain much complex logic. The most complicated part is that of retrieving the right data. This

required some more complex SQL queries which were tested manually using and SQL editor program. However the rest of the functionalities of the user interface have been tested by me throughout the implementation. As the user interface is basically visualising the generated output from the disruption engine a two main scenarios have been employed during its testing. The first is testing the functionalities while the disruption engine is either not running or it does not produce any output. The other is while new output is being generated by the back end and the user interface is continuously updating itself. In addition to these scenarios due to the fact the system distinguishes between three type of user roles I have assumed each of the roles and carried out manual tests to verify the correct behaviour from the system.

Some feedback has been received regarding the user interface during demonstration and discussions with CentreComm and my supervisor. All remarks and suggestions from the received feedback has been addressed. In addition to that I have asked family and friends to spend time and test the functionality of the user interface for this I provided them with a list of use cases and the expected behaviour. This has proved useful in identifying some issues with the functionality and the compatibility of the product with different web browsers and systems. All of the identified problems have been addressed and rectified after that.

## 6.4 Stress Testing

One of the main requirements for detecting the disruptions in the bus network is for this to be calculated in real-time. Because of this we need to make sure that the implementation is able to cope with large amounts of data quickly. Each bus reading we receive in the AVL feed file is on average 100 bytes. There could be roughly up to 7000 buses active in the bus network at each point in time. This means that the system should be capable of processing 1MB of data every 30 seconds as close to real-time as possible. Also as the system will run continuously with more data being pushed to the system throughout

its execution it need to be ensured that any memory that is occupied by data that is irrelevant (old) be discarded appropriately without causing any memory leaks.

The approach that I have taken in order to test the system for such performance issues consisted of firstly creating an automated script for simulating the feed pushing to the tool. This program takes two main input parameters. One is the directory in which the feeds for the simulation are being stored and the second main input is the rate at which these files are copied to the directory which is used by the disruption engine for monitoring. Currently as discussed previously in this report the AVL data files which were provided for this project are generated every 5 minutes. However this could be changed in future if the tool goes into production. For this reason we need to make sure that the proposed prototype is capable of dealing with all of the data that is being pushed to the engine every 30 seconds (as this is the current rate iBus equipped buses transmit data). Apart from the rate at which new data is pushed to the tool the performance of our system depends mainly on how much data is pushed (e.g. during the night ours there are less active buses thus less data generated compared to daytime) and also on how many disruptions are detected in the system. The earlier is clearly obviously that if there is more data to be processed the system would naturally take longer. However the second statement is because we only update the database information for a route, run and section only if there is changes. This means that the update time could vary greatly depending on the number of changes at each network update. In order to minimise the impact of this the system is implemented in such a way that it makes all database updates after a network update in one transaction. This means that the main performance bottleneck of the system becomes the updating of the database. The initial versions of the prototype did not make use of a database, but rather wrote the output into flat CSV files. The reason for adding a database in place of the flat files is that the implementation change, so that it could keep historical state of the network which could later be used for further analysis and reports which are features

that have received a lot of positive feedback from TFL.

Simulations have been run to make sure the system is capable of processing and updating itself in real-time. This tests have been carried out on a laptop running Intel Core i7-3610QM with 16GB DDR3 1600Mhz of ram and Windows 10 x64 bit operating system. The simulations consisted of feeding the tool with a week worth of data that was provided by TFL. The data was made of 21629 separate (each one for a single bus operator) AVL feeds which have been group into 2814 batches according to their timestamps. The total size of the sample is 1.07 Gigabyte (1076394754 bytes) and the average size of a single feed file is 49.77 Kilobytes (49766.27 bytes). We ran a number of simulations with this data keeping all parameters the same and re-initialising the system before each run. The only parameter we altered, to make sure the system is capable of handling data at higher rates, is the rate at which the feed simulation program pushed the AVL files. The results could be seen on figure A.10 in Appendix A. From the results we can conclude that even on a normal computer in a development environment the system is capable of producing output in almost real-time. The memory usage also could be seen in the results shown on figure A.10 in Appendix A. The memory was measured through the execution of the tool and readings were taken every minute. The memory usage is highly dependent on the amount of the data being processed and also depends on when the Java Virtual Machine(JVM<sup>2</sup>) garbage collector executes. The latter is because the system might have already removed all references to some object, however this memory would not be unallocated until the JVM garbage collector is called.

The system current implementation is updating each route in the network consecutively. However the system has been implemented with concurrent execution of this calculation in mind. This means that it could easily be updated such that each route is concurrently processed. This is also the case for the parsing of the CSV feed files and observation extraction. As the aim of this project is to built a prototype I have decided not to spend the limited

---

<sup>2</sup><http://www.oracle.com/webfolder/technetwork/tutorials/obe/java/gc01/index.html>

project time on such optimisations as it is more important to prove if this is viable solution. As it can be seen from the stress tests carried out even without such parallelism in place the system is capable of generating output in less than a second. This means that the bottleneck for achieving real time updates to the users is the user interface implementation (e.g. how often the client web browser will poll the server for updates).



## Chapter 7

# Results & Evaluation

Evaluation is important step of every software project yet it is sometimes neglected. This chapter aims to discuss the how we have evaluated our implementation. It concludes by discussing the results that were obtained.

### 7.1 Evaluation

#### 7.1.1 Disruption Engine

In order to evaluate and verify that our tool is able to detect disruptions we need to run the system with some data for which we know what the output should be. This means that we can analyse the output once the system has processed the sample input data. Analysing the data need to focus on whether the system is able to:

- Detect all disruptions present in the data it is fed with.
- Are they being detected in real-time (there no or little lag).
- Is the information calculated representing the an accurate estimate of what the real situation is.

detect all disruptions

However achieving this is problematic in our case as the data provided by TFL is only the real AVL feed files. There is no formal information or list of

all delays, their severity and duration for a given period of time. We have been provided with some links to web blogs which contain some of the diversions that are being implemented on daily basis in response to disruptions in the network. However we cannot formally use this data to evaluate our product as this is neither an exhaustive list nor precise and accurate source of information (the delays stated there vary greatly and are somewhat subjective). We can analyse the provided bus feeds manually for a give period to extract what the actual state of the network is during that time. However if we want to cover a number of scenarios this approach becomes infeasible as it will require great effort and time which are limited. And if we want more cases for evaluation it will requires us to go over the same process again and again.

In order to overcome this problem we have decided to generate some test data on our own. We have achieved this by implementing a simple program for generating iBus AVL like feed files. This allows us to compile test data with different scenarios and test if our system achieves the expected results. We have limited our test data to four scenarios as follows:

1. The schedule deviation value remain fairly constant with very minor changes across the route for every bus on that route. This scenarios should yield no disruption detections.
2. Single bus incident (e.g. bus breakdown, customer incident, error reading etc.). This again should not be picked up by the system.
3. General traffic scenario. This means that there is gradual schedule deviation increase at some point of the route run. This is the most typical real life scenario (e.g. rush hour traffic build up) and it should be detected by the system.
4. Incident (e.g. traffic collision) along a route. The data for this would represent a sudden increase in the schedule deviation. This needs to be detected by the tool.

Once we have generated the data for the above scenarios we need to run the system and feed it with the artificially generated input. We run each of

the above scenarios separately and the system is re-initialised before each run. This allows us to have the system in the exactly same state for each test.

This tests allowed us to run multiple simulation in order to adjust the weighted moving average parameters. These are the weights and the moving average window. This parameter values are critical for the accuracy of our tool. They have direct impact on what is being detected by the tool and with what lag. We also altered the data validity parameter which represents when we discard any observations.

In order to calibrate our prototype and its output we ran a numerous simulation with the describer scenarios and different values for our parameters. Using our test generated data we have obtained best results with having the data validity set to 90 minutes, but then only have moving average window of 5. This means that we only consider the last 5 observations for a given section which have occurred in the last hour and a half. We also use weight of 1 to 5 (1 for the oldest and 5 for the most recent). This allowed for a balance output, as the system tended to overreact if we used an exponentially growing weights or it lagged behind if we used a greater moving average window. Also keeping data for the last 90 minutes in memory should work well for both low and high frequency buses, as from the data provided we do not have any indication of the type of the service. This however needs to be further tested and evaluated by either generating further test data or using real AVL data for which the bus network state at that point is know.

### **7.1.2 Visualisation**

The web application part of our system was evaluated by simply using the list of requirements. Using the define requirements from Chapter 3 of this report we were able to verify that all expected functionality is place and produces the correct results. This also included testing that the application is behaving and displaying the same way on the most widely used browsers (Internet Explorer, Firefox and Chrome) as well one some mobile devices (Android tablet, Ipad and Iphone). All of the test verified that the system is consistent across the

various devices and no functional issues have been identified.

## 7.2 Results

During the testing and evaluation of the proposed prototype we have ensured that all user and functional requirements have been met. Our simulations using artificially generated data proved that the proposed system is capable of effectively monitor changes in the schedule deviation value. However further evaluations and analysis is required into whether using this value is accurate and reliable measure of the actual delays in the bus network. As it has been discussed in the background chapter of this report there is very little or no studies which address the problems this particular project is trying to solve. This means that we are unable to compare our results with those of others simply because we could not find any.

## Chapter 8

# Professional & Ethical Issues

Throughout every stage of this project I have made every effort to follow the rules and guidelines that are set out by the British Computer Society (BCS) Code of Conduct & Code of Practice [4]. These are rules and professional standards that govern the individual decisions and behaviour. The main rules that apply almost to every software development project states the individual should:

- "have due regard for public health, privacy, security and wellbeing of others and the environment."[4]
- "have due regard for the legitimate rights of Third Parties"[4].

The whole project has been planned, designed and developed with both of these rules, as well as other rules and standards, in mind. The system makes use of a number of third party libraries and framework. However I have made explicitly the use of any such libraries and provided the according reference to the source of the original idea/product. I have also given references to any work or ideas that I have made use of throughout the project.

I have also tried to make sure that the applications that were developed as part of this project do not pose any harm neither to the computers they

are running on or interacting with nor to their users. The tool is expected to run 24/7 with a large number of files being processed every day. This means that we need to make sure does not contain any memory leaks as discussed in previous chapters. I have also used appropriate method to safeguard the database from any SQL injection [43] which could potentially alter the data unintentionally or without the appropriate permission. However it should be noted that this is a prototypical system and not a fully working and security proof production version.

The web application displays the last time and date when the disruption engine has updated its state. This value represent the latest time of data value from the iBud AVL feeds. This allows the users to be aware how old the data they are seeing actually is.

In case of failure during the execution of the disruption engine an email alert system can easily be set. This will allow for the responsible maintenance personnel to be alerted in time. However if the system malfunctions in manner such that it continues its execution and continue to generate output it can lead to confusion among the users of the application. We cannot be responsible if this is caused by the input data. However if the the input is correct and the calculations are wrong it will possibly lead to lost of trust in the application. In order to minimise such risks thorough testing should be carried out once the system is deployed. It is recommended that the system undergoes trial runs with real data being pushed in real-time. This would be another test building on the rest of the testing that has been done to prove the correct system behaviour in real-life scenarios.

## Chapter 9

# Conclusion

During the course of this exciting project we have examined carefully the needs of CentreComm for automation of their current work flow which could lead to better operation, management and cut in costs of controlling London's bus network. This has naturally led to a in depth literature review of the related work and approaches that could be undertaken in order to solve the problem posed by our project. This has resulted in the design and implementation of a prototypical tool for detecting bus delays in real time using iBus AVL data.

This report proposed a prototypical tool for detecting disruptions in London bus network. This is achieved by employing moving average smoothing technique for analysis of the time series data presented. The main strength of this approach is its simplicity. This is growing area of interest for intelligent transportation systems (ITS) and I expect to receive much more attention in near future with the rise in AVL data availability from different sources. I have also given some directions and proposal for improving and driving this work further.

This project has taught me a lot about transportation systems and networks. It has given me great insight into the complex operations that go into operating and managing large bus networks as is London's one. In addition to that I have also learnt a lot about different statistical techniques available for analysis of time series data. Such techniques are very useful and can be applied

not only in the domain of this project, but in a broad range of problems where there is time series data which needs to be analysed.

The software development approach taken for developing the prototype for this project seemed to work well. It has allowed me to gather useful feedback early on and to refine the project requirements as initially the user requirements were very broad and vague. This is probably due to the fact that there is not much closely related work previously done. Significant amount of this project was devoted to gathering and refining the user requirements. This included numerous meetings, discussions and even shadowing at CentreComm. Another major part of this project was carrying out proper background research which even needed to be revisited latter in the project.

One thing that did not work so well is the evaluation of the software system. This is owned partly due to the lack of readily available information which could be used as well as due to delay in planning and carrying out this evaluation. However this is something I have now learnt and would help me better plan my future projects I undertake.

## **9.1 Future Work**

The proposed system has provided some initial results as seen in the previous chapter. There is however a lot that could be improved and built upon on it. In this section I provide some suggestions and direction for taking this work further.

One simple extension that could take this project one step ahead is to try to implement a peak/valley detection algorithm which could distinguish between incidents and just increased traffic congestion. This was briefly discussed in chapter 2 of this report however due to time limits of this project has not been looked in more depth or considered for design and implementation.

Another feature that has received a lot of positive feedback from different levels in TFL is the ability to generate complex historical reports of what has been the network state. According to TFL this would provide more objective



view of what has happened in the transport network and who should take the responsibility.

More historical data could also be employed and correlated with the real-time data received. Examples of such data may include weather data, time of the day/week, workdays compared to weekends and public holidays etc. This could improve the accuracy of the system as well as bring to light some persisting problems under given environments.

The increased popularity and usage of various AVL systems being used by different fleets could be used along with the data present in this report. This could include taxi fleets [36], delivery service fleet, emergency services and many more. Having more information and especially from different sources could result in one central congestion monitoring system which could more accurately and in shorter-terms calculate and predict traffic conditions in the arterial road network of the city. There is the potential for employing data from GPS enabled devices that most people of the general travelling public own these days. This something is already being investigated [44] and I can see this will get more attention in the near future.

# References

- [1] Mehmet Altinkaya and Metin Zontul. Urban bus arrival time prediction: A review of computational models. *International Journal of Recent Technology and Engineering (IJRTE)*, 2(4):164–169, 2013.
- [2] Greater London Authority. Number of buses by type of bus in london. <http://data.london.gov.uk/dataset/number-buses-type-bus-london>, 2014. Online; accessed 22-October-2014.
- [3] Hamed Azami and Saeid Sanei. Spike detection approaches for noisy neuronal data: Assessment and comparison. *Neurocomputing*, 133(0):491 – 506, 2014.
- [4] BCS. Bcs code of conduct. <http://www.bcs.org/category/6030>, June 2011. Online; accessed 6-April-2015.
- [5] A.I. Bejan, R.J. Gibbens, D. Evans, A.R. Beresford, J. Bacon, and A. Friday. Statistical modelling and analysis of sparse bus probe data in urban areas. In *Intelligent Transportation Systems (ITSC), 2010 13th International IEEE Conference on*, pages 1256–1263, Sept 2010.
- [6] Engin Bozdog, Ali Mesbah, and Arie Van Deursen. A comparison of push and pull techniques for ajax. In *Web Site Evolution, 2007. WSE 2007. 9th IEEE International Workshop on*, pages 15–22. IEEE, 2007.
- [7] Peter J Brockwell. *Introduction to time series and forecasting*, volume 1. Taylor & Francis, 2002.

- [8] Peter J Brockwell and Richard A Davis. *Time series: theory and methods*. Springer Science & Business Media, 2009.
- [9] J. L. Connell and L. Shafer. *Structured Rapid Prototyping: An Evolutionary Approach to Software Development*. Yourdon Press, Upper Saddle River, NJ, USA, 1989.
- [10] Mark S. Dougherty and Mark R. Cobbett. Short-term inter-urban traffic forecasts using neural networks. *International Journal of Forecasting*, 13(1):21 – 31, 1997.
- [11] Transport for London. Transport for london’s ibus wins innovation award. <https://www.tfl.gov.uk/info-for/media/press-releases/2008/march/transport-for-londons-ibus-wins-innovation-award>, March 2008. Online; accessed 2-April-2015.
- [12] Transport for London. <https://www.tfl.gov.uk/cdn/static/cms/documents/uploads/forms/lbsl-tendering-and-contracting.pdf>, 2008. Online; accessed 2-April-2015.
- [13] Transport for London. All london’s buses now fitted with ibus. <https://www.tfl.gov.uk/info-for/media/press-releases/2009/april/all-londons-buses-now-fitted-with-ibus>, April 2009. Online; accessed 2-April-2015.
- [14] Transport for London. <https://www.tfl.gov.uk/info-for/media/press-releases/2009/may/centrecomm-celebrates-30-years-keeping-londons-buses-moving>. Online, May 2009. Online; accessed 2-December-2014.
- [15] Transport for London Media. <https://www.tfl.gov.uk/info-for/media/press-releases/2014/may/annual-passenger-journeys-on-london-s-buses-top-2-4-billion>. Online, May 2014. Online; accessed 2-December-2014.
- [16] Martin Fowler. *Patterns of enterprise application architecture*. Addison-Wesley, Boston, 2003.

- [17] Everette S. Gardner. Exponential smoothing: The state of the art. *Journal of Forecasting*, 4(1):1–28, 1985.
- [18] John F Gilmore and Naohiko Abe. Neural network models for traffic control and congestion prediction. *Journal of Intelligent Transportation Systems*, 2(3):231–252, 1995.
- [19] A Pragmatic Guide. Agile web development with rails. 2006.
- [20] Jianhua Guo, Wei Huang, and Billy M. Williams. Adaptive kalman filter approach for stochastic short-term traffic flow rate prediction and uncertainty quantification. *Transportation Research Part C: Emerging Technologies*, 43, Part 1(0):50 – 64, 2014. Special Issue on Short-term Traffic Flow Forecasting.
- [21] Ryan Jay Herring. Real-time traffic modeling and estimation with streaming probe data using machine learning. 2010.
- [22] N.B. Hounsell and B.P. Shrestha. Avl based bus priority at traffic signals: a review of architectures and case study. *European Journal of Transport and Infrastructure Research*, 5(1):13–29, June 2005.
- [23] N.B. Hounsell, B.P. Shrestha, and A. Wong. Data management and applications in a world-leading bus fleet. *Transportation Research Part C: Emerging Technologies*, 22(0):76 – 87, 2012.
- [24] Ran Hee Jeong. *The prediction of bus arrival time using automatic vehicle location systems data*. PhD thesis, Texas A&M University, 2005.
- [25] Thanavat Junchaya, Gang-Len Chang, and Alberto Santiago. Advanced traffic management system: real-time network traffic simulation methodology with a massively parallel computing architecture. *Transportation Research Record*, (1358), 1992.
- [26] Siem Jan Koopman. Exact initial kalman filtering and smoothing for nonstationary time series models. *Journal of the American Statistical Association*, 92(440):1630–1638, 1997.

- [27] Reuven M Lerner. Open-source databases, part iii: choosing a database. *Linux Journal*, 2007.
- [28] J.A. Livermore. Factors that impact implementing an agile software development methodology. In *SoutheastCon, 2007. Proceedings. IEEE*, pages 82–86, March 2007.
- [29] Hani S Mahmassani, Srinivas Peeta, Gang-Len Chang, and Thanavat Junchaya. A review of dynamic assignment and traffic simulation models for adis/atms applications, 1991.
- [30] BBC News. Extra 500 buses planned for growing capital before 2021. <http://www.bbc.co.uk/news/uk-england-london-30285777>, 2014. Online; accessed 2-December-2014.
- [31] Martin Odersky, Lex Spoon, and Bill Venner. *Programming in scala*. Artima Inc, 2008.
- [32] Object Management Group (OMG). Uml. <http://www.uml.org/>, April 2015. Online; accessed 2-April-2015.
- [33] Yan Qi and Sherif Ishak. A hidden markov model for short term prediction of traffic conditions on freeways. *Transportation Research Part C: Emerging Technologies*, 43, Part 1(0):95 – 111, 2014. Special Issue on Short-term Traffic Flow Forecasting.
- [34] "Clarke R., Bowen T., and J" Head. Mass deployment of bus priority using real-time passenger information systems in london. In *Proc. European Transport conference, 2007.*, Leeuwenhorst, Netherlands, 2007.
- [35] Clarke R., Bowen T., and Head J. Mass deployment of bus priority using real-time passenger information systems in london, 2007.
- [36] Mahmood Rahmani, Haris N Koutsopoulos, and Anand Ranganathan. Requirements and potential of gps-based floating car data for traffic management: Stockholm case study. In *Intelligent Transportation Systems*

- (ITSC), 2010 13th International IEEE Conference on, pages 730–735. IEEE, 2010.
- [37] Mohsen Ramezani and Nikolas Geroliminis. On the estimation of arterial route travel time distribution with markov chains. *Transportation Research Part B: Methodological*, 46(10):1576 – 1590, 2012.
  - [38] Stephen Riter and Jan McCoy. Automatic vehicle location—an overview. *IEEE Transactions on Vehicular Technology*, 26(1), 1977.
  - [39] Steve Robinson. Ticket info on buses in service, 2010. Proposed Change Paper Provided by TFL.
  - [40] A. Sboner, A. Romanel, A. Malossini, F. Ciocchetta, F. Demichelis, I. Azzini, E. Blanzieri, and R. Dell’Anna. Simple methods for peak and valley detection in time series microarray data. In Patrick McConnell, SimonM. Lin, and Patrick Hurban, editors, *Methods of Microarray Data Analysis V*, pages 27–44. Springer US, 2007.
  - [41] Robert H Shumway and David S Stoffer. *Time series analysis and its applications: with R examples*. Springer Science & Business Media, 2010.
  - [42] Brian L Smith and Michael J Demetsky. Traffic flow forecasting: comparison of modeling approaches. *Journal of transportation engineering*, 123(4):261–266, 1997.
  - [43] Zhendong Su and Gary Wassermann. The essence of command injection attacks in web applications. *SIGPLAN Not.*, 41(1):372–382, January 2006.
  - [44] Arvind Thiagarajan, James Biagioni, Tomas Gerlich, and Jakob Eriksson. Cooperative transit tracking using smart-phones. In *Proceedings of the 8th ACM Conference on Embedded Networked Sensor Systems*, pages 85–98. ACM, 2010.
  - [45] Trapeze. London bus services limited (lbsl). <http://www.trapezegroup.eu/london-bus-services-limited-lbsl>. Online; accessed 2-April-2015.

- [46] D Ventzas and N Petrellis. Peak searching algorithms and applications. In *The IASTED International Conference on Signal and Image Processing and Applications ~ SIPA 2011*, 2011.
- [47] Eleni I. Vlahogianni, Matthew G. Karlaftis, and John C. Golias. Short-term traffic forecasting: Where we are and where weâ€™re going. *Transportation Research Part C: Emerging Technologies*, 43, Part 1(0):3 – 19, 2014. Special Issue on Short-term Traffic Flow Forecasting.
- [48] W3C. Server-sent events. <http://www.w3.org/TR/eventsource/>. Online; accessed 2-April-2015.
- [49] Jian Wang, Wei Deng, and Yuntao Guo. New bayesian combination method for short-term traffic flow forecasting. *Transportation Research Part C: Emerging Technologies*, 43, Part 1(0):79 – 94, 2014. Special Issue on Short-term Traffic Flow Forecasting.
- [50] Alan Wong and Nick Hounsell. Using the ibus system to provide improved public transport information and applications for london. Paper 01753, July 2010.
- [51] Jinsoo You and Tschangho John Kim. Towards developing a travel time forecasting model for location-based services: A review. In *Methods and Models in Transport and Telecommunications*, pages 45–61. Springer, 2005.