

# **7CCS4PRJ Final Year Individual Project**

## **Real-Time Visualisation of Bus Delays in London using AVL Data**

### **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 is still a lot of area of utilisation of the data that these AVL systems generate and produce. This report explores and analyses the tools and applications currently available at Transport for London (TFL) bus operation 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 as urban delay detection is very complex and unpredictable as will the report justify.

### **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	Motivation . . . . .	3
1.2	Scope . . . . .	4
1.3	Aims . . . . .	5
1.4	Objectives . . . . .	5
1.5	Report Structure . . . . .	6
<b>2</b>	<b>Background research</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 . . . . .	12
2.6	Other . . . . .	16
2.7	Summary . . . . .	16
<b>3</b>	<b>Specification</b>	<b>18</b>
3.1	User Requirements . . . . .	18
3.2	Functional . . . . .	19
<b>4</b>	<b>Design</b>	<b>21</b>
4.1	Use cases . . . . .	23
4.2	System architecture . . . . .	24
4.3	State Machine . . . . .	25
4.4	Class organisation . . . . .	25
<b>5</b>	<b>Implementation</b>	<b>28</b>
5.1	iBus AVL Data . . . . .	28
5.2	Database . . . . .	29
5.3	Disruption Engine . . . . .	29
5.4	Graphical User Interface . . . . .	31

<b>6</b>	<b>Testing</b>	<b>32</b>
6.1	Unit Testing . . . . .	33
6.2	Functional Testing . . . . .	33
6.3	Integration Testing . . . . .	33
6.4	Stress Testing . . . . .	33
<b>7</b>	<b>Results &amp; Evaluation</b>	<b>34</b>
7.1	Evaluation . . . . .	34
7.2	Results . . . . .	34
<b>8</b>	<b>Professional &amp; Ethical Issues</b>	<b>35</b>
<b>9</b>	<b>Conclusion &amp; Future Work</b>	<b>37</b>
9.1	Future Work . . . . .	37
9.2	Conclusion . . . . .	37
	<b>References</b>	<b>42</b>
<b>A</b>	<b>User Guide</b>	<b>43</b>
A.1	Installation . . . . .	43
A.2	Execution . . . . .	43
A.3	Dependencies . . . . .	44
<b>B</b>	<b>Source Code</b>	<b>45</b>
B.1	Engine . . . . .	46
B.2	Web Front End . . . . .	46
B.3	Database . . . . .	46

# Chapter 1

## Introduction

### 1.1 Motivation

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 [12]. 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 [12]. The capital's bus network is continually increasing along with the city's population [21]. 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[25]. This AVL system has led to improved fleet management and has enabled the creation and improvement of multiple applications [32]. The system is generating large sets of data both in real-time as well as historical data. This helps, 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. However there are currently problems for which CentreComm staff are required to manually analyse these data sets in order to discover in which parts of the networks problems are occurring. This is because there is lack of readily available preprocessed information. This is very impractical and time consuming and currently they ultimately rely on individual bus operators and drivers to notify them of possible problems. Once alerted of a possible disruption in the network they can start their own investigation into 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. This is where this project comes in place to address this inefficiency and to propose, implement and evaluate a prototypical tool for real time monitoring of the bus network.

## 1.2 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 aid the control room staff. This tool has to work and analyse the data that has been made available in real time. The main two problems that are in the scope of this project and need to be solved are:

- Analysing and identifying the disruptions in the bus network in real-time using the provided data.
- Producing a prioritised list of the disruptions classified using rules and heuristics gathered during meetings and discussions with the key stakeholders from TFL.
- Visualising the calculated list in an appropriate format.
- Evaluating the performance of the tool. This is essential as it need to be proven that the results could be trusted.



### 1.3 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 where disruption might happen or have already happened even before the bus drivers or operators have noticed and alerted CentreComm. This 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 (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.

### 1.4 Objectives

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

- Getting in depth understanding of the problem and current work-flows that are in place.
- Researching similar work that has already been done and how it can relate 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.

## 1.5 Report Structure

In order to help the reader I have outlined the project structure here. The report would continue in the next chapter by providing the reader with all the background knowledge needed for the rest of the report. This would 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 would 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 that have been identified and gathered from CentreComm. The report then goes on to discuss the design and the implementation of the proposed system. This is followed then by Chapter 5 and 6 which address testing and evaluation of the prototype. 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 research

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. I also present 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 [9]. These operators agree with TFL that the routes they are operating would be served either according to a predefined fixed schedule (e.g. a bus stop need to be served at 1pm, 3pm etc.) or on a headway (e.g. a bus stop need to be served 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 [11]) and even simply a severe congestion. In situations like this bus operators have no authority or power to overcome such

problems on their own. They can only ask CentreComm to intervene. CentreComm can do so by for example implementing a short/long term diversions or curtailments (short turning) some of the buses on the affected routes.

Buses in the network can be classified by multiple factors, however for the purpose of this report the main distinction we need to consider are high and low frequency bus routes. 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 alighting 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 [11] 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, reliable and safe operation of London busy bus network.

London bus network consists of around 680 bus routes operated by more than 8000 buses [3]. Each of these 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 [11] 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 maintain the growing bus network.

## 2.3 iBus AVL

Automatic vehicle location systems make use of the Global Positioning System (GPS) to enable the remote (using the internet) tracking of the locations of the vehicles in a given fleet. This system combines a number of technologies including GPS, cellular communications and more with the goal of improving and cutting the cost of fleet management.

All of London buses operating on the TFL bus network have been equipped with state of the art and award winning [8] AVL system named iBus [10]. 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 [18] and [32]. 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 1 below taken from [18]). This OBU is responsible for a number of 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 are 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 [18]. This has led to improved and more affordable transport service.

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

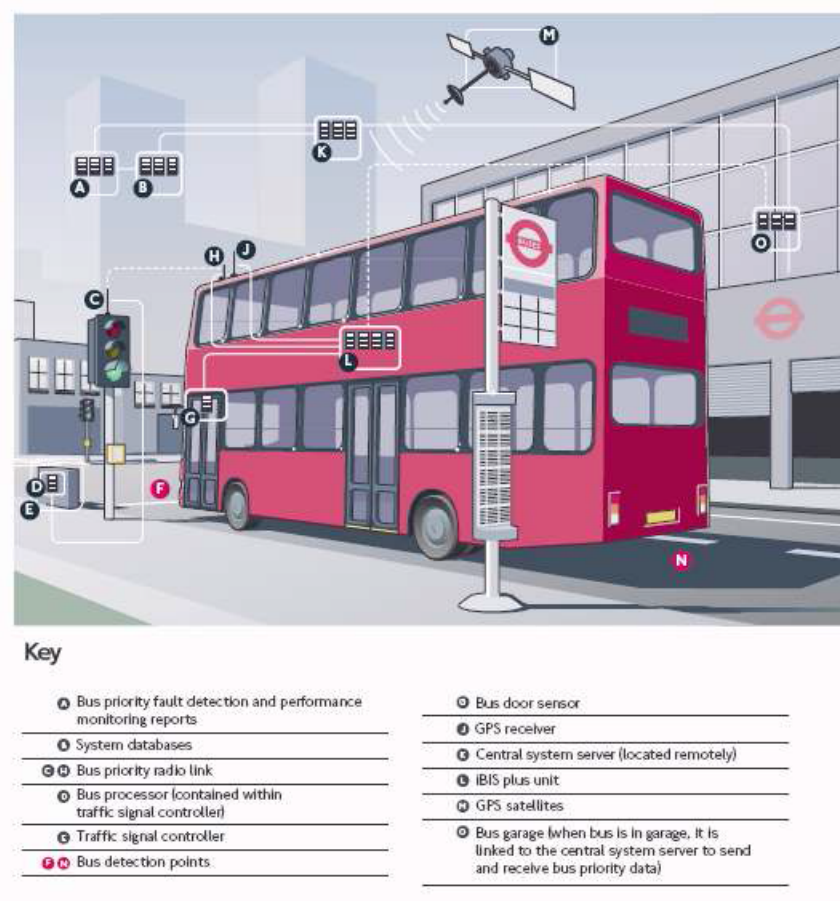


Figure 2.1: Overview of iBus System (Source: Transport for London, 2006a)

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. 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 a there is a problem and how severe 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.

## 2.4 Related Work

The main problem posed by this project of detecting disruptions in the bus network can be also translated to short term traffic congestion detection. 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 etc. We are interested into finding routes/sections in the network which are delay/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.

The literature is full of research towards accurately predicting bus arrival times [REFERENCE]. Also there is a lot of work done toward detecting calculating travel times in non urban environment. This include approached based on AVL probe data and automatic vehicle identification (AVI) as well as induction loops [REFERENCE]. There are also plenty research done towards traffic congestion detection based on AVL probe data [30]. However there are significant challenges due to the nature of the urban environment it self. Densely populate areas are influenced by many factors which could 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 could be some noise which reduces the accuracy of the stated 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 [30] [6]. As it can be seen from the tables in figures A.1 to A.4 taken from [30]. From these tables we can easily see that most of the 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 mass usage

of such technologies. The literature is abundant of research on addressing the issues of bus prioritisation at traffic signals and junctions [17] and [2], but this is not directly related to our problem domain and thus is not discussed further in this report.

The approaches and research that we have examined show us that the the main model that have been used for detecting and forecasting short-term traffic congestion/disruptions are: According to [33] traffic forecasting models could be categorised as follows:

- Statistical models
- Computer Simulations
- Mathematical Optimisation
- Artificial Intelligence

The statistical models could be further subdivided to:

- Historical
- Time Series
- Nonparametric regression
- Hybrid

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 their wide usage and advantages and the available AVL data for this project we 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. 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 generated and transmitted. Analysis of time series data could be performed in order to extract and calculate some meaningful statistics from the data. 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**

Moving average also called rolling or running average is a statistical calculation method. It is universal analysis technique. It is type of mathematical convolution [REFERENCE]. It helps to analyse data series by calculating a series of averages of subsets of the data. It can also be referred to as rolling or moving mean and it acts like a filter. Moving averages is commonly used in time series data analysis. 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 are the three main type of moving averages.

#### **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 term averages they would react quicker to changes while long term averages would have greater lag. The equation for calculating SMA is given below as equation 2.1.

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

### 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 when the observation was made. For example if we consider a n period moving average we could multiply the oldest observation by

$$\frac{2 * 1}{n(n + 1)}$$

and the newest by

$$\frac{2 * n}{n(n + 1)}$$

. This would mean that recent data have bigger impact on the result. It is important that the weights add up to 1. The general equation for calculating WMA is given below as equation 2.2.

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

### Exponential Moving Average

Exponential smoothing was first suggested by Robert Goodell Brown [14].

Exponential moving average (EMA) or also called exponential smoothing 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.

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

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

## Summary

If we compare the three types (simple, exponential and weighted) we can clearly see that the simple moving average offers the most smoothing however there is the trade-off of the biggest lag. The exponential moving average will react faster to changes and sits closer to the actual readings, however it might overreach at times. The weighted moving follows the movement even more closely than the exponential moving average. Determining which one to use depends on the objective. If a trend indication with better smoothing and little reaction for shorter movements 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.

### **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 [1] used for spike detection in microarray data (another example could be found in [4]). 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 [29]). 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 have could be used include Kalman filtering [19] [15], Markov Chains [24] [26], Machine learning [16], Bayesian networks [31].

## **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 below:

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

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 it self whenever new data is detected and processed.
5. Be able to run without intervention 24/7.
6. Keep track of the disruption 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 [20]. The research nature (as seen in chapter 2 of this report) of this project itself makes it sensible to use evolutionary prototyping [7] 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 [7].
- 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 resembles the desired product as required by the user.

In the below sections I have presented the system at an abstract level as a. 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 [23] 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 create we can than easily transform them into platform specific models (PSM) using the desired technologies. This technique is known as Model-Driven Architecture (MDA) software design approach [22].

## 4.1 Use cases

Base 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 a 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 supposed to be a CentreComm staff members who need 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 to them of hiding and showing disruptions that are currently detected in the bus network.
- **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) patter 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 [13] 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 for 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 would terminate.

### 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. The class diagram for our system is presented in figure A.3 in Appendix A. 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. It 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 need 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.

## 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 not to go into too much technical details except where this would provide better insight and understanding.

### 5.1 iBus AVL Data

INCLUDE SAMPLE OF THE DATA IN THE APPENDIX iBus system generate very short telegram messages with the bus location [18]. This information is then processed on a server and more information is calculated and derived.

The data that has been provided by Technical Service Group (TSG) at TFL for this project consists of preprocessed iBus feeds. This feeds currently are being generated every 5 minutes. Each file consists of the following fields:

1. BusId

The information we are interested is the deviation from the schedule. This is calculated by knowing the expected arrival time of the bus at a stop from the schedule and is compared to the observed time. This value is calculated the same way for both low and high frequency buses(Low frequency buses are supposed to run according to a fixed schedule (e.g. a bus should arrive at stop



at predefined time) and usually routes which have less than 5 buses an hour. High frequency bus routes should maintain headways - this meaning a bus should be arriving at stop at predefined intervals (e.g. each 2 minutes)). More about the supplied data for this project would be covered in the requirements section.

Key property of time series data is stationarity. This means that the behaviour of the time series data does not change over time. In our case the data generally speaking the time series data is not stationary. However if we consider short window size it might be possible to treat the data as close to stationary.

## 5.2 Database

Why have I have decided to use database vs flat file data store. Pros and cons Why postgresql. Reference to the database model Explanation and discussion of the database model.

## 5.3 Disruption Engine

However as other research has pointed out [28] calculating travel time as measure for congestion is difficult task and it is very dependable on the environment and its conditions (e.g. weather, time of day, public demand etc.). For this reason and because of the data available this project would not try to measure disruptions by calculating travel time or bus speeds. TFL has provided us with example of the AVL data which among other things contains a the GPS coordinates of the bus at a given point in time and preprocessed deviation from the schedule value. For the rest of this chapter we assume this value is accurately calculated and that we would receive this value for each bus in the network at some regular interval. The provided data is discussed in further details in Chapter 5.

For this project we monitor and measure the schedule deviation value as calculating the congestion is very challenging and still not very well understood

in the case of arterial urban traffic.

However from the literature [28] we can see that it could be difficult to precisely define what we mean by congestion in a transport network. It seems that congestion could mean various things to different studies and people.

What is the general approach I have taken

### **5.3.1 Technologies**

SCALA, XML database connection file, Database is Postgresql Why Scala and DB Postgresql.

### **5.3.2 Bus Network representation**

What other data is needed?

### **5.3.3 Monitoring for new feeds**

### **5.3.4 Feed File Processing and observation extraction**

How is the input processed etc. ...

### **5.3.5 Updating network state**

Graphics to depict how the algorithm works

### **Applying WMA**

The important factors to consider are the weights and the period/window size to use. We could also apply exponential smoothing on top of the WMA. Problems: The bus could have started the journey ahead of schedule and thus to intentionally be losing time. The buses could curtail anywhere on a route without notification The buses could be diverted (this could be short term or long term) it can also be for few stops or it could be a longer diversion.

**Applying EMA**

**Detecting disruptions**

### **5.3.6 Problems & Optimisations**

## **5.4 Graphical User Interface**

Include a figure showing the layout. It is responsive etc. Why Ruby on Rails as front end. Foundation as CSS layout framework as it is lightweight and it enables easy responsive design. Ajax short polling - as this is only prototype and is supported by all major browsers.

### **5.4.1 Problems**

Data available, data frequency, no knowledge when buses do curtail, not taking into account bus dwell time, bus drivers who are running ahead of schedule could be driving slower on purpose and thus. Tool gives an upper bound of the of the WMA lost time per section.

## Chapter 6

# Testing

The disruption engine consists of a number of individual components and algorithms. These are tested using a number of unit tests. Due to the nature of the engine, stress tests were carried out to test it for any memory leaks and performance issues. Evaluation of the output of the disruption engine will be performed by first carrying out further literature review to find some guidance on what window for the data to consider to use and what weights to use for optimal results. This will be followed by running the system on a given set of data (e.g. a week worth of AVL data) and comparing the output with the actual state of the network during this period.

The user interface consists of a simple web application capable of displaying list of disruptions. Testing and evaluation of this web application will be performed as user testing. As I mentioned above, some feedback and problems were identified which will be addressed and fixed. Follow up user tests will be carried out by giving access to friends and family to the web site to use and give feedback on. This should give reasonable confidence in the correctness of the user interface as there is no complex logic incorporated in the web front end application.

## **6.1 Unit Testing**

What is it? Why was needed? How was it performed?

## **6.2 Functional Testing**

Focus on the user interface. What is it? Why was needed? How was it performed?

## **6.3 Integration Testing**

What is it? Why was needed? How was it performed?

## **6.4 Stress Testing**

Why was this needed? How was it carried out? - Results/proof

## Chapter 7

# Results & Evaluation

### 7.1 Evaluation

Averaging methods: These techniques could be evaluated by calculating the error (the difference between the prediction and the actual value), the squared error and also the the sum of the squared errors (SSE) and respectively the mean of the squared errors (MSE). The model which minimizes the MSE is the best. It can be shown mathematically that the one that minimizes the MSE for a set of random data is the mean.

Consider metrics for Mean Absolute Difference and Mean Absolute Error.

Peak analysis. Determining the optima Weights and the data window size.

root mean square error (RMSE) and root mean square percentage error (RMSPE)

### 7.2 Results

## 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 [5]. These rules and professional standards that govern the individual decisions and behaviour. The main rule 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."
- "have due regard for the legitimate rights of Third Parties"

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 ides 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 main with our tool as

discussed in the previous chapter is the fact that this system is expected to run 24/7 with a large number of files being processed every day. I have also used appropriate method to safeguard the database from any SQL injection [27] 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.



## Chapter 9

# Conclusion & Future Work

### 9.1 Future Work

Distinction between incidents and congestion (use peak detection algorithm). Data from other source could be used (taxis AVL, couriers services AVL) etc. Historical data could be employed in order to make further analysis and correlations with weather data, time or the day/week/year etc. Increasing the frequency of the data means that we could make use of the actual geo location (alternative approach to one taken) information in order to calculate and monitor the bus speeds rather than the preprocessed schedule deviation value.

Using data from taxis [Cite - Requirements and Potential of GPS-based Floating Car Data for Traffic Management Stockholm Case study]

### 9.2 Conclusion

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 could be applied not only in the domain of this project, but in every problem 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. One thing that did not work so well is the evaluation of the software system. This is mainly owned 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.

# References

- [1]
- [2] *Mass deployment of bus priority using real-time passenger information systems in London*, Leeuwenhorst, Netherlands, 2007.
- [3] 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.
- [4] Hamed Azami and Saeid Sanei. Spike detection approaches for noisy neuronal data: Assessment and comparison. *Neurocomputing*, 133(0):491 – 506, 2014.
- [5] BCS. Bcs code of conduct. <http://www.bcs.org/category/6030>, June 2011. Online; accessed 6-April-2015.
- [6] 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.
- [7] J. L. Connell and L. Shafer. *Structured Rapid Prototyping: An Evolutionary Approach to Software Development*. Yourdon Press, Upper Saddle River, NJ, USA, 1989.
- [8] 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.
- [9] 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.
- [10] 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.
- [11] 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.
- [12] 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.
- [13] Martin Fowler. *Patterns of enterprise application architecture*. Addison-Wesley, Boston, 2003.
- [14] Everette S. Gardner. Exponential smoothing: The state of the art. *Journal of Forecasting*, 4(1):1–28, 1985.
- [15] 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.
- [16] Ryan Jay Herring. Real-time traffic modeling and estimation with streaming probe data using machine learning. 2010.

- [17] 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.
- [18] 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.
- [19] 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.
- [20] J.A. Livermore. Factors that impact implementing an agile software development methodology. In *SoutheastCon, 2007. Proceedings. IEEE*, pages 82–86, March 2007.
- [21] 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.
- [22] Object Management Group (OMG). Mda. <http://www.omg.org/mda/>, August 2014. Online; accessed 2-April-2015.
- [23] Object Management Group (OMG). Uml. <http://www.uml.org/>, April 2015. Online; accessed 2-April-2015.
- [24] 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.
- [25] Clarke R., Bowen T., and Head J. Mass deployment of bus priority using real-time passenger information systems in london, 2007.
- [26] 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.

- [27] Zhendong Su and Gary Wassermann. The essence of command injection attacks in web applications. *SIGPLAN Not.*, 41(1):372–382, January 2006.
- [28] S. Tantiyanugulchai and R.L. Bertini. Arterial performance measurement using transit buses as probe vehicles. In *Intelligent Transportation Systems, 2003. Proceedings. 2003 IEEE*, volume 1, pages 102–107 vol.1, Oct 2003.
- [29] 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.
- [30] Eleni I. Vlahogianni, Matthew G. Karlaftis, and John C. Golias. Short-term traffic forecasting: Where we are and where we’re going. *Transportation Research Part C: Emerging Technologies*, 43, Part 1(0):3 – 19, 2014. Special Issue on Short-term Traffic Flow Forecasting.
- [31] 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.
- [32] Alan Wong and Nick Hounsell. Using the ibus system to provide improved public transport information and applications for london. Paper 01753, July 2010.
- [33] Jinsoo You and Tschangho John Kim. Empirical analysis of a travel-time forecasting model. *Geographical Analysis*, 39(4):397–417, 2007.

# Appendix A

## User Guide

You must provide an adequate user guide for your software. The guide should provide easily understood instructions on how to use your software. A particularly useful approach is to treat the user guide as a walk-through of a typical session, or set of sessions, which collectively display all of the features of your package. Technical details of how the package works are rarely required. Keep the guide concise and simple. The extensive use of diagrams, illustrating the package in action, can often be particularly helpful. The user guide is sometimes included as a chapter in the main body of the report, but is often better included in an appendix to the main report.

### A.1 Installation

Link to gitHub repository

### A.2 Execution

Link to gitHub repository

## A.3 Dependencies

Link to gitHub repository Scala and Java versions, scalatest library, PostgreSQL, scala xml library, scala geo coordinate library, scala logger Ruby on Rails, Ruby version, Gems that have been used Foundation, jQuery - Include the specific versions



# Appendix B

## Source Code

Complete source code listings must be submitted as an appendix to the report. The project source codes are usually spread out over several files/units. You should try to help the reader to navigate through your source code by providing a “table of contents” (titles of these files/units and one line descriptions). The first page of the program listings folder must contain the following statement certifying the work as your own: “I verify that I am the sole author of the programs contained in this folder, except where explicitly stated to the contrary”. Your (typed) signature and the date should follow this statement.

All work on programs must stop once the code is submitted. You are required to keep safely several copies of this version of the program - one copy must be kept on the departmental disk space - and you must use one of these copies in the project examination. Your examiners may ask to see the last-modified dates of your program files, and may ask you to demonstrate that the program files you use in the project examination are identical to the program files you had stored on the departmental disk space before you submitted the project. Any attempt to demonstrate code that is not included in your submitted source listings is an attempt to cheat; any such attempt will be reported to the KCL Misconduct Committee.

**You may find it easier to firstly generate a PDF of your source code using a text editor and then merge it to the end of your report.**

There are many free tools available that allow you to merge PDF files.

**B.1 Engine**

**B.2 Web Front End**

**B.3 Database**