

Max Bowman NEA

Optimised Safe Route Finder For Cyclists

Last updated 14:11, Tuesday 22nd March, 2022

Pages: 26

Written in L^AT_EX



Contents

1 Analysis	3
1.1 Problem Area	3
1.2 Client	3
1.3 Similar Systems	5
1.4 Features	5
1.4.1 Map data	5
1.4.2 Accident Data	5
1.4.3 Traffic Data	6
1.4.4 Finding the shortest route	6
1.5 Class Layout	7
1.6 Critical Path	7
1.7 Specification	8
1.8 Specification Justification	9
2 Design	10
2.1 Accident Download System	10
2.1.1 Pulling from the TFL API	10
2.1.2 Parsing the Data	10
2.2 Route Finding System	12
2.2.1 Overview and Design Choices	12
2.2.2 Open Street Map Data	12
2.2.3 Representing and Querying the Graph	17
2.2.4 Preprocessing and Subsequent Querying of the Graph	20
2.2.5 Saving the work done to disk	22
2.2.6 User Interface	22
2.2.7 Route Output	22
3 Technical Implementation	23
3.1 Overview	23
3.2 Accident Download System	23
3.3 Route Finding System	25
3.3.1 Parsing Open Street Map Data	25
4 Testing	26
4.1 Reliability Testing	26
4.2 Performance Testing	26
5 Evaluation	26



Figure 1: Route suggested by google maps

1 Analysis

1.1 Problem Area

There exist a plethora of route finding services for finding the shortest route between two points. These typically optimise for the shortest time taken to travel between two points. Provision for cyclists can be lacking, because it typically consists of an alteration to the existing code for cars, slightly modified to allow for cycle paths and different timings for cyclists.

Many of these route finding applications don't take into account safety considerations. This is shown in Figure 1 which shows the route Google Maps suggests that I cycle to school by. The route includes the A4, which is a very dangerous stretch of road for cyclists as it is a 3 lane road. My plan is to use accident statistics to work out which roads are dangerous and then avoid them. This would be done by imposing a cost for going somewhere that there had been an accident.

1.2 Client

I interviewed Yuvraj Dubey, one of my classmates, to talk about whether they would be interested in a product that attempted to calculate safe routes to cycle.

Do you own a bicycle?

Yes

Do you cycle regularly?

Yes

Do you ever feel in danger while cycling?

Yes

When do you feel in most danger while cycling?

At a Junction near Victoria, where there is a right turn and no cycle lane to do it in. I attempted to go this way once and was almost hit by a bus.

Are you aware of your surroundings while cycling?

Not really

Do you feel aware of the places which to cycle in safely, does this change when you are in places you commonly go?

No, but especially where I don't know where I am because then I do not know where it is safe to cycle. When I cycle home I have learnt a safe route, but if asked to cycle somewhere I did not know I would most likely end up in a dangerous place.

Would you use an application that tells you safe routes to get places?

Yes

Would you be happy using a command line application for this?

Yes

If said application took more than 60 seconds to calculate a route would you lose patience?

Yes, probably

Have you found that cycling directions generated by current products are safe and efficient?

They are definitely efficient, but they often take me onto busy roads and junctions which can be less safe.

From this interview and my own experiences cycling in London I determined that there was an application for my idea, but only if it ran in a reasonable amount of time and was otherwise easy to use, as people struggle to find safe routes themselves, especially when going somewhere they had not previously been, and commercial products such as Google Maps don't generate safe routes.

1.3 Similar Systems

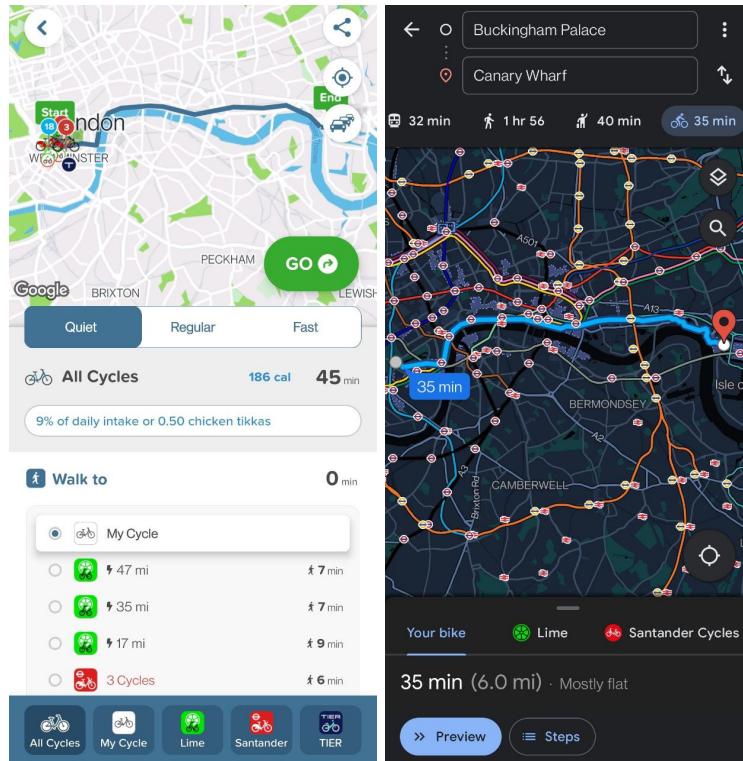


Figure 2: Citymapper and Google Maps respectively, generating the same directions

There exist similar systems for route finding such as Google Maps and Citymapper. These have advantage's over what I will be able to offer such as being able to use live traffic data, and having a good user interface. Citymapper has a option to choose between "Quiet", "Regular", and "Fast" routes, which seems to be based on the type of road you are taken down. However they don't really take accident density into account which will be my aims.

1.4 Features

1.4.1 Map data

I need a data source which can provide data on roads that are legal to cycle on as well as being freely available for me to use. I settled on Open Street Map, a project which combines data gathered by volunteers into one massive freely available map. The map is downloadable in the form of a large XML file or PBF file. A PBF file is just a binary version of the same thing. I will write code to parse one of these myself.

1.4.2 Accident Data

Transport for London has an excellent API which you can download accident data from. The data comes in JSON files and contains information about what type of vehicle was involved in each accident, the severity, and the coordinates of the accident. I need to find a way of mapping the coordinates onto the road network so that the danger of roads can be calculated as accurately as possible. Originally I was looking at adding the accident to the nearest edge in the graph but I decided against this after considering how intersections are typically represented in OSM. Road intersections typically look something like what

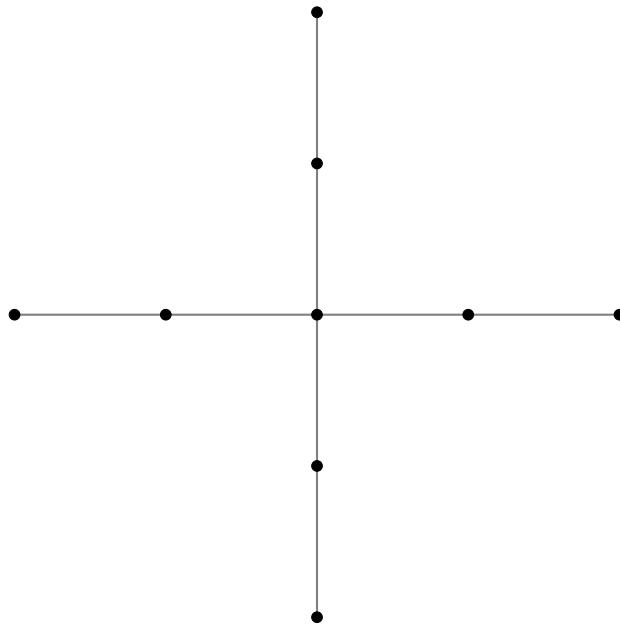


Figure 3: Typical OSM representation of an intersection

is represented in Figure 3. If an accident were to occur at the intersection it would end up being added to one of the segments leading into the intersection, so the danger would not be properly calculated if not passing through that segment. Instead I thought about adding the danger to the nearest node, so that the accident would always be counted when a route passes through that intersection. The other problem that I have to deal with, if adding data to the nearest node is that the density of nodes is not constant. Straighter roads will not need to use as many OSM nodes as curved roads, so i might incorrectly add cost to the wrong node. My proposed solution to this is to interpolate in nodes to a very high density to deal with this problem.

1.4.3 Traffic Data

The main problem with this is that the accident data is absolute and can thus not be used to calculate probabilities. For example, more accidents happen on King's Street than the dangerous road I showed earlier, but this doesn't mean that King's Street is more dangerous merely that more cyclists travel on it. This means that I need to get accurate cyclist traffic data for the whole of London in order to turn my accident statistics into accident probabilities. The Department for Transport and the Office for National Statistics both keep data on traffic, but it isn't applicable because cycle data is only given as a total ¹ and at specific count points. This means that I will either need to work out traffic data or get it from some dataset, such as Strava's Global Heatmap.

Luckily the rest of the application can work without including traffic data so my plan is to deal with this problem at a later time or hope that the avoidance of accidents alone will be enough.

1.4.4 Finding the shortest route

Algorithms for finding the shortest path in a Graph are abundant. The most well known is Dijkstra's algorithm, and it's variant A*. In large Graphs, both Dijkstra's algorithm and A* can be very slow. I will most likely use preprocessing based algorithms such as ALT*/Contraction Hierachies to make

¹0.6 billion miles per year in London

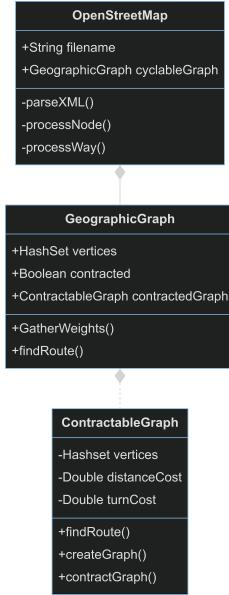


Figure 4: Rough idea of the class diagram I will use

my program run a lot faster. This preprocessing will take a long time, but when completed it will significantly reduce the time taken to complete searches.

1.5 Class Layout

As seen in Figure 4, I think I will use 3 major classes which will deal with parsing, storing, and preprocessing Graph data respectively. These will be linked by object composition rather than inheritance. Where possible I will use private fields and methods.

1.6 Critical Path

- Download accidents from the TFL API, parse for accidents that are of interest, and store it in a format that can be easily parsed later
- Import OSM map of London as a Graph which can be easily queried.
- Write code that imports the accident data and connects it to given Graph nodes.
- Write a simple Dijkstra's Algorithm method of finding the shortest path.
- Write a preprocessing based method of finding the shortest path.
- Write a frontend for the input of user options.

1.7 Specification

1. Accident Download System
 - 1.1 The System must be capable of interfacing with the TFL API to download accidents
 - 1.2 The System must be capable of parsing accident data to determine which accidents are relevant
 - 1.3 The System must be capable of storing accident data in a easily accessible file for the route finding algorithm to parse
2. Route Finding System
 - 2.1 The System should be capable of processing an Open Street Map map file into a suitable data structure
 - 2.2 The System should be capable of processing the accident data from the Accident Download System
 - 2.3 The System should be capable of attaching the data from the Accident Download System to the Suitable Data Structure from specification point 2.1
 - 2.4 The System should be capable of using Accident Data combined with other suitable cost estimation functions to find a directed route between two points that are places on roads in London
 - 2.5 The System should always suggest a route that can be followed while observing all currently known laws of physics
 - 2.6 The System should always suggest a route that can be followed while observing all international, national, local laws and all relevant bylaws
 - 2.7 The System should suggest a safe route wherever possible
 - 2.8 The System should have guards in place for certain types of roads which are deemed too dangerous to consider
 - 2.9 The System should be able to generate this route quickly
 - 2.10 The System should have a mode for preprocessing to make the route generation even quicker
 - 2.11 When using the same parameters, the route generated using a preprocessed graph and associated algorithms should be the same as that generated by the non-preprocessing based approach
 - 2.12 The System should be easy to use
 - 2.13 The System should not crash, and should give appropriate non crashing errors if user data is determined to be bad
 - 2.14 The parameters for the cost estimation functions should be user definable, but sensible defaults should also be defined
 - 2.15 The Route should be output in a format that is easy for the user to follow on a mobile device
 - 2.16 The System should have a system for saving all data from intensive calculation to disk

1.8 Specification Justification

1. Accident Download System
 - 1.1 This point is required in order to obtain the necessary accident data
 - 1.2 This point is required to reduce the file sizes needed to be packaged with the project, and make it easier to use later
 - 1.3 This point is required as otherwise the system would not do much good in making the Route Finding System simpler
2. Route Finding System
 - 2.1 The suitable data structure is required for all further queries
 - 2.2 The accident data is required for route finding applications
 - 2.3 This is required to tell where the costs for the accidents should be applied in shortest path queries
 - 2.4 This is the main function of the project; if it can't do that it will be worthless
 - 2.5 The Route should be realistic and connected; it can't ask people to teleport to their destination or phase through walls
 - 2.6 The Route should only take people onto roads which are publicly accessible and legal to cycle on, as obviously I do not want to encourage people to break laws
 - 2.7 The Route should attempt to be as safe as possible. Of course in some situations there is only one route which could be unsafe. In that situation that route would be suggested
 - 2.8 Some roads can be very dangerous to cyclists, such as canal paths, and these won't necessarily be represented through the data
 - 2.9 If the route generation takes too long the user will get bored. Other products offering the same features can generate routes very quickly
 - 2.10 London is very large, and route generation may be slow without preprocessing, so the preprocessing mode will help with that
 - 2.11 This makes sure that the preprocessing based approach is not losing information about the best route in any way
 - 2.12 If the system was not easy to use, people would not use it
 - 2.13 This feeds into the previous specification point
 - 2.14 The customisability would allow users to make their own decisions about tradeoffs between time and safety, and sensible defaults should generate sensible routes
 - 2.15 Especially with longer routes, the probability of the user being able to remember the route is slim, so they will need a mobile device to be able to guide them
 - 2.16 The parsing of XML files may take a while, and the preprocessing for query speedup definitely will. If the results of this can be saved to disk repeating this every time the program is run will not be necessary.

```
1 [{
2     "id": 0,
3     "lat": 0.0,
4     "lon": 0.0,
5     "location": "string",
6     "date": "string",
7     "severity": "string",
8     "borough": "string",
9     "casualties": [
10         {
11             "age": 0,
12             "class": "string",
13             "severity": "string",
14             "mode": "string",
15             "ageBand": "string"
16         },
17         {
18             "type": "string"
19         }
20     }
}]
```

Figure 5: The Default output from the API [1]

2 Design

2.1 Accident Download System

As outlined in the Specification, this system should be capable of interfacing with the TFL API, downloading the accidents, deciding which ones are relevant, and storing this in a useful format.

2.1.1 Pulling from the TFL API

The Transport For London API is an excellent API which can provide information on many different aspects of the Transport For London network. The specific API which I used is the AccidentStats API. This API is very simple, you simply request a given year and a JSON object is returned which contains all of the accidents that happened in london that year. These consist of all the accidents that were reported to the police as happening in that year. Of course, there will be many more accidents than are on the API, but these will mostly be more minor accidents. The Data is returned as a list of accidents, formatted as shown in Figure 5.

As the data about accidents that happened in the past is not going to change any time soon, and TFL only updates this API every year, it is simplest just to download the files once, parse it once, and then use that result in the route finder. I used two scripts to do this.

TFL started gathering this data in 2005, and the most recent update was in 2019, so the first script downloads all the data from 2005 to 2019 and save it in a subfolder called `accidents`. The script loops through all the years between 2005 and 2019 and creates a JSON file for that year.

2.1.2 Parsing the Data

The next step is to go through the data from all the years², and save all the accidents that are pertinent to my project. As seen in Figure 5, each accident listing has information on casualties, of which there may be many. As this is a cycling application, I only want data on accidents where at

²I later decided to remove some years see Section 3.2

```
1 [  
2     ["Latitude", "Longitude", "Severity"]  
3 ]
```

Figure 6: Intended format for parsed accident file

least one of the casualties was a cyclist. The API provides a lot of data, but all that is relevant is the latitude and longitude, as well as the severity of the incident. All of this data was ultimately saved to a file called `output.json`. The intended format of the file data will be saved in is shown in Figure 6.

2.2 Route Finding System

2.2.1 Overview and Design Choices

I decided to use Kotlin for this project. Kotlin is related to Java in that it allows access to all the Standard Java Libraries, as well as external ones, and it runs in the Java Virtual Machine. Kotlin does not need to maintain backwards compatibility with old Java code, so it has a much cleaner API and has nice syntax. As this was a new project not dependent on using Java code Kotlin was an obvious choice. Furthermore, Java code can be easily translated into Kotlin code, so I could code things in Java then translate over if necessary. Like Java, kotlin allows Object-Oriented-Programming, so that is the approach I will be taking in this project.

As seen in Figure 7, I decided to define 3 main classes for my NEA. `OpenStreetMap` contains methods for parsing OSM files and creating a graph. `OpenStreetMap` contains a single instance of the `GeographicGraph` class by object composition. This is because there will only be one `GeographicGraph` class associated with a given OSM file, and this system could be easily extended to hold information on two separate graph areas at the same time, which could be used if I wanted to extend the system to more cities. The `GeographicGraph` can contain an instance of the `ContractableGraph` class, if it has been contracted. If it has not been contracted, this class will not exist.

I also created two other classes, `IntTuple` and `DoubleTuple`, which are used in priority queues in both the `GeographicGraph` and `ContractableGraph` classes.

2.2.2 Open Street Map Data

Open Street Map is a world map generated by user mapping. It contains multitudes of data on all sorts of mappable things, from the height of stories to the roads that encompass them. Files can be in one of two formats; PBF³ which is a highly efficient binary format, and XML⁴ which is a markdown language. I tried to use a library to parse the PBF file into a Graph, but the main library for that did not work. So instead I decided to parse it manually as an XML file.

Getting the relevant file Openstreetmap has an API for requesting parts of the map, but it does not allow requesting large areas, such as the whole of London. There are many mirrors from which you can download large parts of the world. I decided to use geofabrik [2] because it allows you to download single countries. Next I extracted London from this dataset. This is not strictly necessary, as my program can still deal with large areas such as the whole of the UK, but it does not have accident data for such an area. Furthermore it would increase the file size needed on disk and increase the running time of any preprocessing. In order to extract London, I used a GeoJSON[3] file, which is essentially a list of coordinates. I found a GeoJSON file on the internet [4] of the M25 boundary, and used the following command⁵ to cut out London:

```
$osmium extract -p course_m25_boundary.json united_kingdom.osm.pbf -o london.osm
```

This file could be used by my program, as it intelligently avoids parsing non-routable ways, but I decided to further reduce the file size by removing all nodes and ways that were not part of highways. This was just for quality of life, as it makes parsing the file much faster.

```
$osmium tags-filter london.osm nw/highway -o ways.osm
```

³Protocol Buffer Format

⁴Extensible Markup Language

⁵I used the OSM manipulation tool Osmium

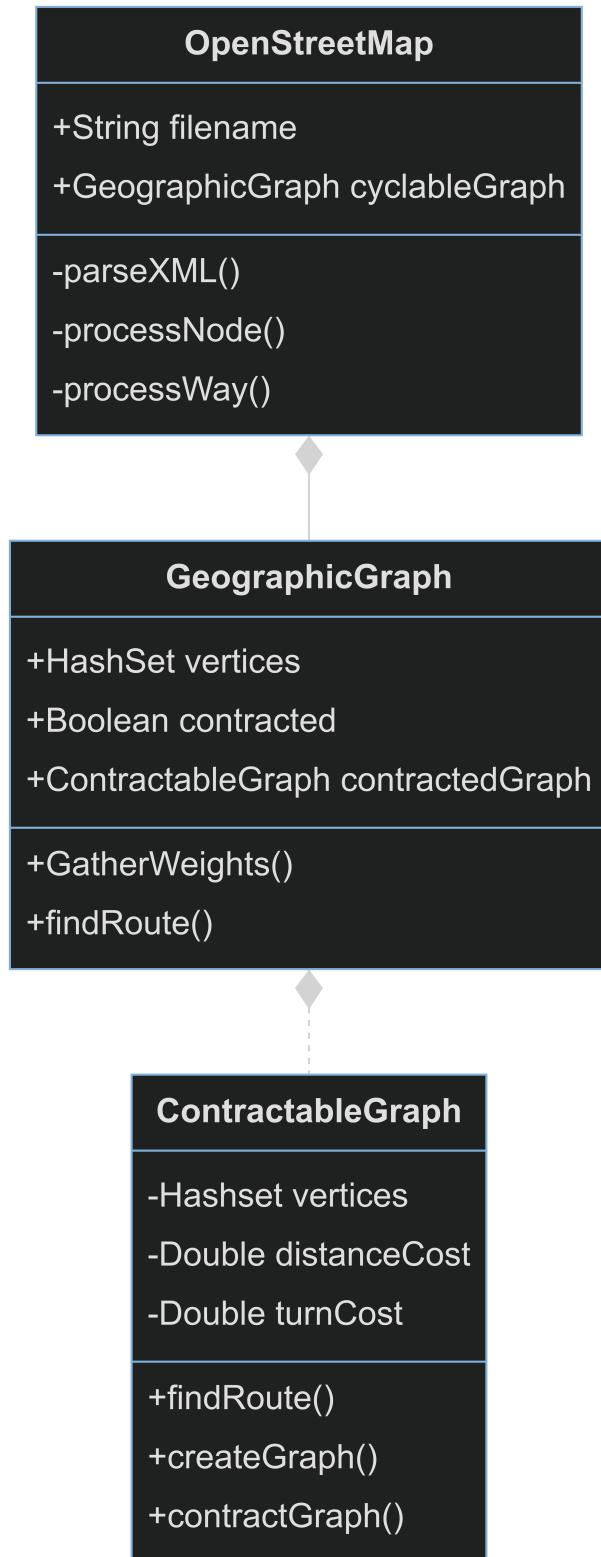


Figure 7: Rough UML diagram for the project

```

1 <way id="1202" version="32" timestamp="2019-05-15T08:24:07Z" uid="7105697" user="
2   _Garrison_ changeset="70264933">
3     <nd ref="5335693253"/>
4     <nd ref="5335693250"/>
5     <nd ref="104429"/>
6     <nd ref="3330244696"/>
7     <nd ref="5335691516"/>
8     <tag k="bicycle" v="no"/>
9     <tag k="foot" v="no"/>
10    <tag k="highway" v="trunk"/>
11    <tag k="horse" v="no"/>
12    <tag k="lanes" v="3"/>
13    <tag k="lit" v="yes"/>
14    <tag k="maxspeed" v="40 mph"/>
15    <tag k="maxspeed:enforcement" v="average"/>
16    <tag k="name" v="North Circular Road"/>
17    <tag k="oneway" v="yes"/>
18    <tag k="operator" v="Transport for London"/>
19    <tag k="ref" v="A406"/>
20    <tag k="sidewalk" v="none"/>
21    <tag k="surface" v="asphalt"/>
21 </way>
```

Figure 8: One of the Ways in `ways.osm` extracted

This created a 370MB file called `ways.osm`. I could have reduced this further by removing the ways that are irrelevant to me, but I decided to do that in the parser so that it would be more easily configurable.

Parsing the file Open Street Map files are built out of 3 main elements [5]:

- Nodes, which contain longitude and latitude.
- Ways, which contain between 2 and 2000 ordered nodes. These are used to form roads, but also all kinds of other polygons, such as buildings, rivers, and walls.
- Relations, which contain many ways, nodes and relations. These are used to connect things together, like all roads on a certain bus route. These are not relevant to my project.

All of these elements can also contain additional information on what type of node or way it is, as well as a unique id. As I only care about nodes that are part of roads I can ignore the relation element completely. In a `.osm` file, the nodes all come first. This means that I can parse the nodes, and then figure out how all the nodes are connected using the ways.

After some research, I determined that the best way to parse the XML was to use the library `dom4j` [6]. I decided that the `OpenStreetMap` class needs 3 functions: `parseXML`, `processNode`, and `processWay`. The function that is called with the direct input file is `.parseXML`. Essentially what it does is iterate through the items in the XML file, calling either `processNode` or `processWay`. `processNode` simply adds a new node to the `GeographicGraph`. The `processWay` function is more complex as it has to work out whether ways should be added to the graph, and if additional data needs to be stored about them. Within each way stored in the `.osm` file, there will be additional tags that hold more data about the way [7]. An example way, which is part of the North Circular road, is shown in Figure 8. First there are a ordered list of nodes which are sequentially connected, then there is a collection of optional Key Value pairs which form tags. The `processWay` function analyzes these tags in order to determine if the road is acceptable for cyclists. The example in Figure 8 would probably not be a good way to include, because the tag `bicycle` is set to `no`, and it is a 3 lane road with a 40mph speed limit.

The program uses accident data to work out which roads are dangerous, but some roads will not have many accidents on them simply because they are so dangerous nobody would ever cycle on them.

Parsing the ways In order to convert from XML to a data representation in memory, we simply iterate through all the subelements in the way. If it is a reference to a node, the id is appended to a ordered list, and if it is a tag, it is added to a HashMap so it can be more easily queried.

Deciding what is allowed

- Highway tag
 - This tag is present in all roads mapped in OSM, so its presence must be checked for.
 - This tag represents the “the importance of the highway within the road network as a whole”[7].
 - This does not normally represent anything about the road quality, safety, usage, layout, or maximum speed. The exceptions to this are
 - * `motorway` which is obviously not wanted.
 - * `living_street` which represents places where pedestrians and cyclists have legal priority, such as Low Traffic Neighbourhoods.
 - * `cycleway` which encompasses segregated areas for cyclists.
 - * `bridleway` which is a segregated area for horses, where cyclists may be allowed.
 - * `footway` which is a footpath which cannot be cycled on.
- Access tag. This tag represents what the access arrangements for the area in question are. If it is `no` or `private`, the way would not be wanted.
- Bicycle and Motor tags. These tags are the same as the access tag, but for specific types of vehicle.
- Surface tag. This represents the surface the route is made out of. Ways that are made out of dirt or similar materials will not be parsed.
- Note tag. This can contain lots of different notes, but the one `processWay` looks for is the `towpath`, as these can be dangerous in a way that is not represented in the TFL accident data
- Oneway tag. If true, the route will only be connected oneway.
- Maxspeed tag. This can be used to exclude dangerous roads.

After analyzing these tags `processWay` has multiple outcomes.

- Way is included in the graph in both directions.
- Way is included in the graph only in the same direction as the ordered list of nodes, because `oneway` is set to `yes`.
- Way is considered too dangerous to include or access to cyclists is not guaranteed to be legal so specification point 2.5 would not be followed.
- Way is included in the graph and the nodes are added to a set of safe nodes because it is in a cycle route or footpath.
- Way is included in the graph and the nodes are added to a set of nodes that will be slower to travel on because of a footpath.

Some of these outcomes can overlap, such as being added to the group of safe nodes and having edges only added in one direction.

1. The way is parsed in some form if:

The `highway` value is not null, and is within a set of allowed highway types

The `access` and `bicycle`⁶ value are not within a set of disallowed access types.

The `surface` value is not within a list of disallowed surfaces

The `maxspeed` value is not more than 30 miles per hour

The `note` value does not say that the path is a towpath

2. The way is only parsed in one direction if:

The `oneway` tag is set to yes

3. The way is marked as a safe node if it is

Marked as a `cycleway`, `footway`, or `pedestrian` or `motor_vehicle` is set to private

4. The way is marked as a slow way if it is

Marked as a `footway` or `pedestrian` and `bicycle` is not set to `designated`

Of course, it will be very easy to iteratively develop this when the whole system is finished. If a certain way that looks like it should be routed down is not routed down or vice versa the function can be adjusted. It is important to note that this function is designed to be quite open in deciding which ways are safe. It only excludes ways that are plainly very dangerous. The actual decision of where to go based on danger will be made by the route finding algorithms.

⁶In places where cycling is not permitted but walking with a bike is this will be dismount instead of no

2.2.3 Representing and Querying the Graph

As shown in Figure 7, my design uses a class called `GeographicGraph` to represent London⁷. There are many graph libraries available to java programs, but creating my own implementation will allow me more flexibility in implementation of methods that relate to real world applications. As shown in Figure 7, the `GeographicGraph` class has some major methods that need to be implemented. `gatherWeights` will read in a accident file in the format specified in Figure 6, and attach the accident to the nearest node in the graph. `findRoute` will find the best route between two nodes in the graph. There will also need to be other functions for maintaining and querying the data structure.

Representing the Graph There are many possible methods of representing Graphs, such as using an adjacency matrix, or an adjacency list. An adjacency matrix would not work in this case, as the graph is very sparse. An adjacency list could work, but maintaining information about the nodes separately from the adjacency list would be difficult. Instead, my design uses an OOP approach, where every relevant node in the OSM representation is converted into an equivalent node represented by a subclass called `GeographicNode`. This is advantageous because it allows associated data, such as the number of accidents attached and the coordinates, to be stored in the same place. As shown in Figure 9, the connections would be represented as a `HashSet`. Ideally, this would contain other `GeographicNode` objects that the node in question was connected to. However this would lead to problems with serialisation as explained in Section 2.2.5 so instead the unique ids from the OSM representation are used, along with a lookup table for converting these back to `GeographicNode` objects.

Another feature that the `GeographicGraph` class needs to provide is facility for the conversion of coordinates to the id of the closest node. This will be useful for allowing the user to ask for directions in terms of coordinates, as well as for attaching accidents to the nearest node. In order to do this quickly and efficiently my design calls for the use of an R*Tree[8]. This is a special type of data structure for holding coordinates that groups coordinates into a tree of rectangles. This allows much faster queries when looking for nodes that satisfy certain geometric properties, such as asking for the closest node. The * in R*Tree means that a heuristic is used to try and optimise query times, by minimising overlap between rectangles. The implementation of an R*Tree is beyond the bounds of this project, and any I could implement would be less powerful than one from a library. To that end, I decided to use the library `rtree2`[9]. I have used it in other projects before and found it to be very powerful. The R*-Tree will be set up when the graph is created, and a function called `getNearestNode` will convert between coordinates and OSM node ids.

Creating the Graph The `OpenStreetMap` class deals with parsing the XML, but it has to interface with the `GeographicGraph` class. When adding new nodes to the Graph, it can simply create a new instance of `GeographicNode` and add it to the relevant lookup table. For adding an edge, the functionality can be provided by a `addEdge` function.

Pruning the Graph Due to the sequential order of parsing the OSM representation of the Graph, some nodes may be left disconnected because the way that connects them together was not parsed. If the system is asked to find a route from a node that is not connected to any other node, it will of course fail. This could lead to a crash, which would be in contravention of specification point 2.13, or at least a bad user experience, if given an error for seemingly sensible inputs. Another problem could be that the Graph might have multiple components, perhaps because access to an area was cut off by safety protocols. Whatever the cause, both of these problems can be averted by removing all nodes that cannot be routed to by one location in the main component. This can be done by picking any node from which all other nodes can be connected to, finding the whole component, then removing any `GeographicNode` objects that are not in that component from the lookup table they are stored in.

⁷Other areas could be represented as well

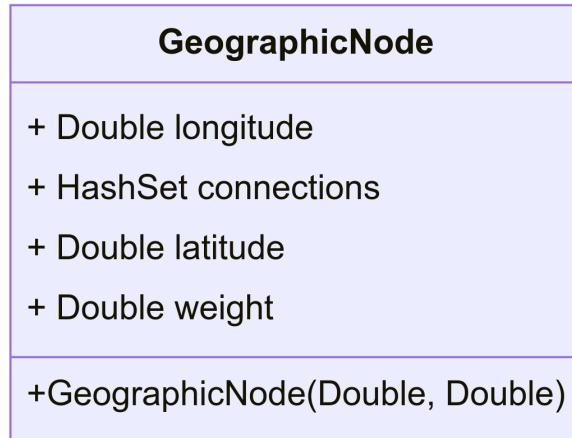


Figure 9: Rough UML class diagram for the **GeographicNode** subclass

Querying the Graph The main query that has to be provided for is finding the shortest route between two points. This functionality can be provided either by a simple Dijkstra's Algorithm implementation, or through a different query for the preprocessed Graph. If the graph has been preprocessed, that mode should be used. Otherwise Dijkstra's shortest path algorithm should be used.

Dijkstra's Algorithm Dijkstra's Algorithm[11] is a immensely popular algorithm for finding routes in graphs. It builds up a tree of nodes that it knows the shortest route to until said tree contains the target node. It does this by iterative "settling" of nodes. The known unsettled node that has the shortest route cost to the start node is settled, and all of it's neighbours are added to a list of nodes to be settled. Every node has a tentative `dist` value and a tentative `prev` value, defined by the currently known shortest path to that route. When the node is settled this value is known to be true. This process is repeated until the algorithm is completed. Pseudocode for Dijkstra's algorithm is shown in Figure 10. A Priority Queue can be used to find the next node to process. This is detailed in the Final Design section.

A variant of Dijkstra's algorithm called A*[12] uses heuristics to reduce the required search space. The heuristic estimates the distance to the end node, meaning that nodes are vaguely settled in the right direction. However this heuristic has to be "admissible" if the shortest possible route is to be found. This means that it never makes overestimates. As the density of accidents cannot be predicted between any point and the end node, the heuristic would be extremely loosely fitting, and would fail to significantly reduce the search space.

```

1  function Dijkstra(Graph, source, target):
2
3      for each vertex v in Graph.Vertices:
4          dist[v] ← INFINITY
5          prev[v] ← UNDEFINED
6          add v to Q
7          dist[source] ← 0
8
9      while Q is not empty:
10         u ← vertex in Q with min dist[u]
11         remove u from Q
12         if u == target: break
13
14         for each neighbor v of u still in Q:
15             alt ← dist[u] + Graph.Edges(u, v)
16             if alt < dist[v]:
17                 dist[v] ← alt
18                 prev[v] ← u
19
20     S ← empty sequence
21     u ← target
22     if prev[u] is defined or u = source:
23         while u is defined:
24             insert u at the beginning of S
25             u ← prev[u]

```

Figure 10: Dijkstra's algorithm modified from Wikipedia[10]

Cost Function Any route finding algorithm needs a function or set of functions for the cost that would be incurred if going a certain way. If going from node u to node v , the function c would be:

$$c(u, v) = w(v) + a \times s(u, v) \times d(u, v) + b \times t(p(u), u, v)$$

$$w(x) = \text{slight}(x) + 2 \times \text{serious}(x) + 3 \times \text{fatal}(x)$$

$$t(a, b, c) = \begin{cases} 1 & \text{if } \angle abc < c \\ 0 & \text{otherwise} \end{cases}$$

$$s(u, v) = \begin{cases} d & \text{if both nodes are safe nodes} \\ e & \text{if both nodes are slow nodes} \\ f & \text{otherwise} \end{cases}$$

where a, b, c, d, e , and f are variables that can be tweaked. The turn function $t(a, b, c)$ can be used to impose a cost on turning too often, as this makes routes harder to follow and increases danger. $\angle abc$ can be calculated using the cosine rule. For the triangle shown in Figure 11

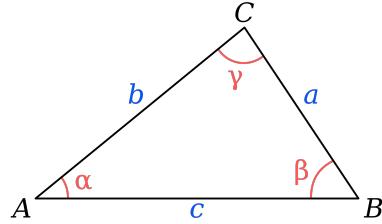


Figure 11: Triangle used to illustrate the cosine rule [13]

$$\begin{aligned} a^2 &= b^2 + c^2 - 2bc \cos \alpha \\ \therefore \cos \alpha &= \frac{b^2 + c^2 - a^2}{2bc} \\ \therefore \alpha &= \arccos\left(\frac{b^2 + c^2 - a^2}{2bc}\right) \end{aligned}$$

Calculating Distances To calculate real world distances⁸ from coordinates, Great Circle distance is required. This is the shortest possible distance across the surface of a sphere. The haversine formula[14] is one formula that calculates this distance. For two points (ϕ_1, λ_1) and (ϕ_2, λ_2) where ϕ and λ are latitude and longitude respectively, the great circle distance can be calculated as

$$d = 2r \arcsin \sqrt{\sin^2\left(\frac{\phi_2 - \phi_1}{2}\right) + \cos \phi_1 \cos \phi_2 \sin^2\left(\frac{\lambda_2 - \lambda_1}{2}\right)}$$

where r is the radius of the earth, so approximately 6371 km.

2.2.4 Preprocessing and Subsequent Querying of the Graph

There are many different methods of route finding[15]. In the Analysis section I considered both Contraction Hierachies[16] and ALT*. Upon further research, I determined that contraction hierachies would be more efficient than ALT* in this situation.

Contraction Hierachies In most road networks there is an inherently hierachical structure. When a human plots a long distance route, they don't consider small dirt tracks, and instead choose the best motorway. This is because in most cases motorways are faster. I do not expect the same level of hierachy in figuring out where is cyclable, because cyclists do not have the same limiting factor of the speed limit. However there will still be some routes that are obviously the best, at least to my program.

Contraction hierachies work by the process of iterative contraction of nodes. Some order of contraction is defined, in order of least to most important, and the nodes are contracted in this order. When a node is contracted the shortest paths are calculated between all pairs of neighbour nodes. If the shortest path goes through the node to be contracted a new shortcut is added. This maintains all shortest paths in the Graph. An example in an undirected graph is shown in Figure 12.

Bidirectional Dijkstra Queries The bidirectional variant of dijkstra is a useful speedup technique. A forward search from the source node is interlaced with a backward search from the target node. When a node is settled by both searches all nodes on the shortest path are known to have been settled. This variant can reduce the search space, and ideas from this search are used when querying the contracted graph.

⁸Assuming the world is a sphere

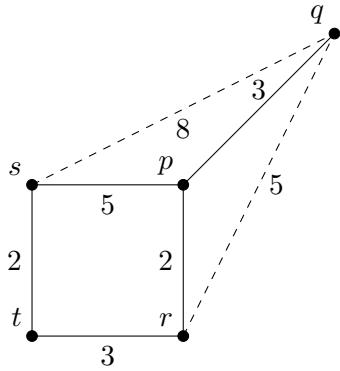


Figure 12: The shortcuts that would be added while contracting node p

Querying the Contracted Graph When contraction is complete, a new graph is created called G^* . This Graph contains all of the nodes and edges in the original graph, as well as all of the shortcuts created during the contraction process. All nodes have a new item of information attached to them, known as the hierarchy. The first node to be contracted will have a hierarchy of 1, and the last node to be contracted will have a hierarchy of n . Two new graphs can be developed from the hierarchy: $G^* \uparrow$ and $G^* \downarrow$. $G^* \uparrow$ contains only edges where the hierarchy increases, and $G^* \downarrow$ contains the opposite. In order to find the shortest path between two nodes u and v , a bidirectional search consisting of a forward search in $G^* \uparrow$ and a backward search in $G^* \downarrow$. In an undirected graph this is the same as two forward searches in $G^* \uparrow$, but in a directed graph the backwards search must go backwards on directed edges. Implementing this will require the `GeographicNode` equivalent class in `ContractableGraph` to contain incoming connections and shortcuts. It has been proven[16] that this method of query will leave at least one node on the shortest path settled by both the forward and backward searches.

Finding the shortest route Once the query is completed, a node on the shortest path will be in the intersection of both searches. The total distance of a route containing a given node is the sum of the distance from the forward search and the distance from the backward search. So to find a node on the shortest path, the minimum total distance cost in the intersection of both searches must be found. From there the `prev` values can be used to backtrack to the start and end of the route. At this point the route will almost certainly contain shortcuts, which must be somehow expanded into their original form.

Limitations When using contraction hierachies, turn costs cannot be as easily integrated because the cost function must be per edge not for collections of edges. It is possible to use a second graph which has the original graph's edges as nodes, but this approach would slow down contraction, and add more complexity to the project. Therefore the ability to consider turn costs will only be added to the simple dijkstra searches and any turn costs will be ignored in the preprocessing mode.

Node Order The optimal node ordering is the one that minimises the expected query time. This can be estimated to be the same as the ordering that leads to the lowest total amount of shortcuts being added, which will in turn minimise the expected search space size. Finding the best ordering is an NP-hard problem, as it requires simulating every possible node ordering to find which one gives the best results. Instead the contraction algorithm will use a heuristic based approach. There are many heuristics of varying complexity, but the two that my design uses are the difference between the amount of shortcuts added and edges removed, and the number of deleted neighbours. Keeping the number of deleted neighbours low makes sure the contraction is spread out across the graph. When contraction is not uniform, very poor performance is attained. The edge difference is essentially a metric of how simple

the graph would be after the contraction, so it makes sense to greedily make the graph as simple as possible. The edge difference heuristic can be calculated by simulating the contraction of a node. The deleted neighbours heuristic can only go up as the graph is contracted, and the same tends to happen for the difference of edges. Therefore it makes sense to load all of the nodes into a priority queue by their heuristic values, and pull out the minimum node. If the heuristic value has changed since the heuristic was placed into the priority queue the node is reinserted with the new value, otherwise it is contracted.

Optimising Contraction In order to contract a node, and compute the edge difference, it needs to be ascertained which pairs of neighbour nodes have their shortest path going through the node. If checking individually for each pair, there would be $O(n^2)$ queries. Instead a variant of dijkstra's algorithm can be used, with a stopping condition of settling all the other neighbour nodes instead of just one. These queries might take slightly longer, but they would find all n neighbour nodes. Another useful optimisation technique is to refill the priority queue every so often, for instance when lots of re-insertions have to happen.

2.2.5 Saving the work done to disk

So that intensive workloads such as contracting the graph do not need to be repeated a facility for saving and loading this information must exist. It would be possible to write some manual serialisation code, but it is simpler to just save the whole class objects. Kotlin has a fairly easy to use serialisation library that can serialise whole objects to disk[17]. The standard mode for this is serialising to JSON, but as my files will likely be very large, my design uses the alternate CBOR encoding instead. As mentioned previously, `GeographicNode` objects cannot contain other `GeographicNode` objects. This is because the kotlin serialiser will get stuck recursively serialising the same nodes. CBOR⁹ can be used for more concise binary representation of JSON objects.

2.2.6 User Interface

My design calls for a simple terminal based user interface. Users should be able to load maps and preprocessed data from pre existing files, and trigger additional preprocessing at different levels. This can be simply implemented by means of a loop where users can choose options or end the program.

2.2.7 Route Output

In order to fulfil specification point 2.15, the route must be output in a format that is not just a list of coordinates. Writing software that can use a route to give turn by turn directions, along with route recalculation when straying off the route is beyond the bounds of this project. However such software does already exist. One example of this is Osmand, an application that uses Open Street Map data to provide a routing application. Osmand can generate its own routes, but it can also be used to follow routes stored in gpx files. GPX files are a subset of XML that can be used to store geographic data such as routes. Many other applications accept these, so a user would not be limited to Osmand. Writing to a GPX file manually would be boring, so my design calls for the use of a library.

⁹Concise Binary Object Representation

3 Technical Implementation

3.1 Overview

Figure 13 shows the complete file system. The `accidents` subfolder is used to generate `output.json`, and the `maps` folder contains different `osm` files for different areas. All of the `.py` files are used in the creation of `output.json` and all of the files in `src/main/kotlin` make up the program as a whole, with `Main.kt` as the entry point.

3.2 Accident Download System

As this system only really needs to be run once, I just wrote some hacky python scripts to do this. As mentioned in Section 2.1.1, I used a simple script called `download.py` to pull all the relevant accident data from the TFL API and store it in the `accidents` folder.

```

1 import requests
2 import json
3 for year in range(2005,2020):
4     print(year)
5     response = requests.get(f"https://api.tfl.gov.uk/AccidentStats/{year}")
6     with open(f"accidents/{year}.json","w") as yearfile:
7         json.dump(response.json(),yearfile)
```

In order to extract useful information from these files, I used `parse.py` to extract all those accidents where one of the listed casualties was a cyclist. While data exists between 2005 and 2019, when looking at some of the areas with high accident density, I found that the accident density had greatly decreased in recent years. This led to the realisation that the TFL accident data was not solely collected for the purpose of providing interesting data for A-Level NEA coursework, and that TFL was probably looking at the data as well and fixing those areas which were lethal. Thus I decided to only use accidents from after 2010.

```

1 import json
2 #we want to parse for accidents that happened to cyclists and store the severity and
3 #the location
4 accidents = []
5 for year in range(2010,2020):
6     print(len(accidents), year)
7     with open(f"accidents/{year}.json","r") as outfile:
8         data = json.load(outfile)
9     for x in data:
10         try:
11             if x["casualties"]["mode"] == "PedalCycle":
12                 accidents.append([x["lat"],x["lon"],x["severity"]])
13         except:
14             if len([person for person in x["casualties"] if person["mode"] == "PedalCycle"]) >= 1:
15                 accidents.append([x["lat"],x["lon"],x["severity"]])
16 with open("output.json","w") as outfile:
17     json.dump(accidents,outfile)
```

I also wrote a small visualisation script using `matplotlib` to verify that the data was dense enough to be useful and imported correctly. I called this `display.py`.

```

1 import json
2 import matplotlib.pyplot as plt
3 x = []
4 y = []
5 with open("output.json","r") as inputfile:
6     data = json.load(inputfile)
7 for accident in data:
```

```
$ tree
.
├── accidents
│   ├── 2005.json
│   ├── 2006.json
│   ├── 2007.json
│   ├── 2008.json
│   ├── 2009.json
│   ├── 2010.json
│   ├── 2011.json
│   ├── 2012.json
│   ├── 2013.json
│   ├── 2014.json
│   ├── 2015.json
│   ├── 2016.json
│   ├── 2017.json
│   ├── 2018.json
│   └── 2019.json
├── build.gradle.kts
├── download.py
├── gradlew
├── gradlew.bat
└── maps
    ├── course_m25_boundary.json
    ├── cyclable.osm
    ├── cyclable.osm.pbf
    ├── goldhawk.osm
    ├── testarea.json
    ├── testarea.osm
    └── ways.osm
├── output.json
├── parse.py
├── display.py
├── route.gpx
├── savedGraph.bin
├── settings.gradle.kts
└── src
    └── main
        └── kotlin
            ├── ContractableGraph.kt
            ├── ContractableGraphTest.kt
            ├── DoubleTuple.kt
            ├── GeographicGraph.kt
            ├── GeographicGraphTest.kt
            ├── IntTuple.kt
            ├── Main.kt
            ├── OpenStreetMap.kt
            └── OpenStreetMapTest.kt

```

6 directories, 42 files

Figure 13: Project file structure. Automatically generated files have been omitted

```

8     x.append(accident[1])
9     y.append(accident[0])
10    print(accident[0],accident[1])
11 plt.scatter(x,y,0.05)
12 plt.show()

```

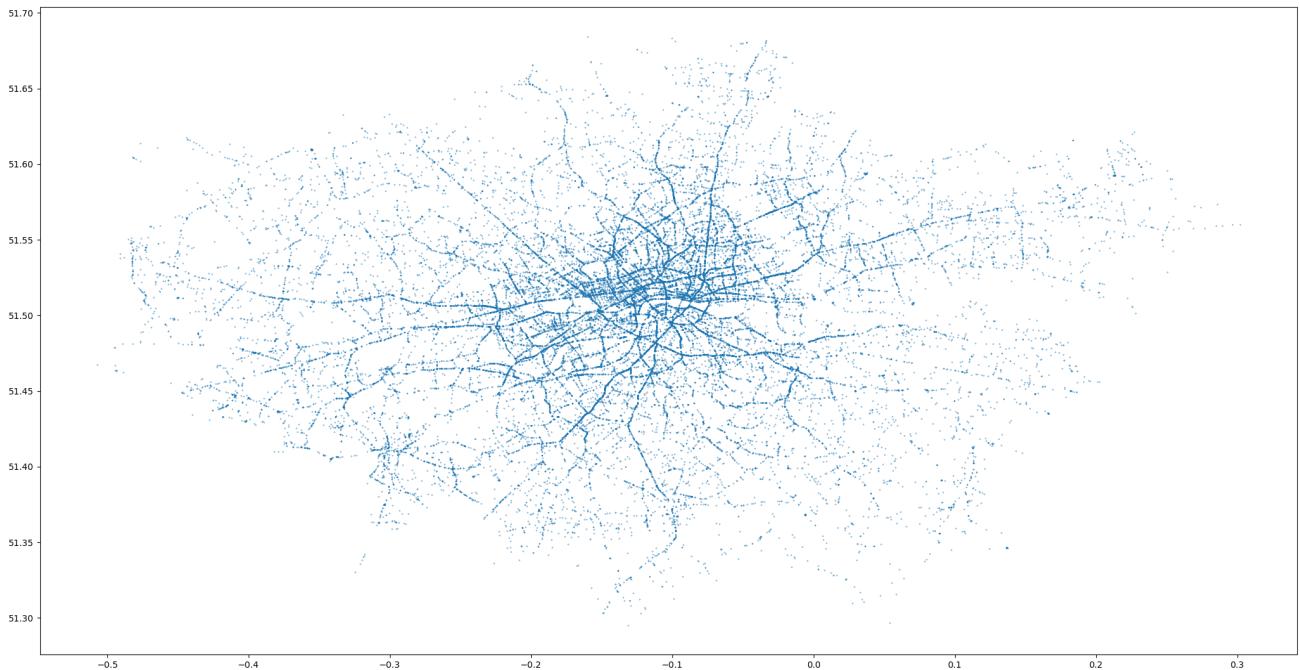


Figure 14: Plot of all accidents in which there were injured cyclists in London since 2010

This generated the plot shown in Figure 14. Major roads are clearly visible, so a route finding algorithm will be able to avoid these. One problem that was apparent from the data was that there are significantly more accidents in central London than in the suburbs. Central London may be more dangerous, but there are also far more cyclists, and as I cannot get road usage data easily routes will probably be biased against central London, possibly causing them to go around. Hopefully, the distance cost will pull routes back towards central London, and they will instead pick between roads for safety on a more local level.

3.3 Route Finding System

The accident download system of the previous section only really serves to provide the input for this system. This system *is the project*. **Add big UML class diagram thing here.**

3.3.1 Parsing Open Street Map Data

Add a UML class diagram of OpenStreetMap The `OpenStreetMap` class contains the methods for parsing `.osm` files and forming a new `GeographicGraph` class. Some of these methods cannot be called in the `init` method of `GeographicGraph` because they operate on the graph after it has been populated by `parseXML`, so they are placed in the constructor for `OpenStreetMap`. This constructor calls `parseXML`, which is shown in Figure 15.

```
1 private fun parseXML(filename: String) {
2     val stream = File(filename).inputStream()
3     val saxReader = SAXReader()
4     val cyclableDocument = saxReader.read(stream)
5     val root: Element = cyclableDocument.rootElement
6     val it: Iterator<Element> = root.elementIterator()
7     while (it.hasNext()) {
8         val element: Element = it.next()
9         when (element.qName.name) {
10             "node" -> processNode(element)
11             "way" -> processWay(element)
12         }
13     }
14 }
```

Figure 15: The `parseXML` function from OSM

4 Testing

4.1 Reliability Testing

4.2 Performance Testing

5 Evaluation

Therefore, the project was a success.

References

- [1] Transport For London. *AccidentStats API Documentation Page*. URL: https://api-portal.tfl.gov.uk/api-details#api=AccidentStats&operation=AccidentStats_Get. (accessed: 17.03.2022).
- [2] geofabrik. *Open Street Map mirror*. URL: <https://download.geofabrik.de>.
- [3] GeoJSON developers. *GeoJSON website*. URL: <https://geojson.org/>.
- [4] Stuart Grange. *GeoJSON download site*. URL: <https://skgrange.github.io/data.html>.
- [5] Open Street Map Wiki. *Open Street Map Wiki Elements Page*. URL: <https://wiki.openstreetmap.org/wiki/Elements>. (accessed: 17.03.2022).
- [6] dom4j. *dom4j documentation site*. URL: <https://dom4j.github.io>.
- [7] Open Street Map Wiki. *Open Street Map Wiki Key:highway Page*. URL: <https://wiki.openstreetmap.org/wiki/Key:highway>. (accessed: 18.03.2022).
- [8] Norbert Beckmann et al. “The R*-tree: an efficient and robust access method for points and rectangles”. In: *Proceedings of the 1990 ACM SIGMOD international conference on Management of data - SIGMOD '90*. ACM Press, 1990. DOI: [10.1145/93597.98741](https://doi.org/10.1145/93597.98741). URL: <https://doi.org/10.1145/93597.98741>.
- [9] David Moten. *Library Github Page*. URL: <https://github.com/davidmoten/rtree2>.
- [10] Wikipedia Editors. URL: https://en.wikipedia.org/wiki/Dijkstra's_algorithm#Pseudocode. Accessed (19.03.2022).
- [11] E. W. Dijkstra. “A note on two problems in connexion with graphs”. In: *Numerische Mathematik* 1.1 (Dec. 1959), pp. 269–271. DOI: [10.1007/bf01386390](https://doi.org/10.1007/bf01386390). URL: <https://doi.org/10.1007/bf01386390>.
- [12] Peter Hart, Nils Nilsson, and Bertram Raphael. “A Formal Basis for the Heuristic Determination of Minimum Cost Paths”. In: *IEEE Transactions on Systems Science and Cybernetics* 4.2 (1968), pp. 100–107. DOI: [10.1109/tssc.1968.300136](https://doi.org/10.1109/tssc.1968.300136). URL: <https://doi.org/10.1109/tssc.1968.300136>.
- [13] Dweisman. URL: https://commons.wikimedia.org/wiki/File:Triangle_with_notations.svg.
- [14] J. de Mendoza y R'ios. *Memoria sobre algunos metodos nuevos de calcular la longitud por las distancias lunares y explicaciones pr'acticas de una teor'ia para la soluci'on de otros problemas de navegaci'on*. Imp. Real, 1795. URL: <https://books.google.co.uk/books?id=030t00qlX2AC>.
- [15] Timothy 'zeevox' Langer. *How do travel planners work?* URL: <https://www.youtube.com/watch?v=wLNai11U2tI>.
- [16] Robert Geisberger et al. “Exact Routing in Large Road Networks Using Contraction Hierarchies”. In: *Transportation Science* 46.3 (Aug. 2012), pp. 388–404. DOI: [10.1287/trsc.1110.0401](https://doi.org/10.1287/trsc.1110.0401). URL: <https://doi.org/10.1287/trsc.1110.0401>.
- [17] *Serialization: Kotlin*. URL: <https://kotlinlang.org/docs/serialization.html>.