

# Keep Austin on time: A study of real-time bus tracking in MetroRapid

Ryan Young

October 26, 2014

## 1 Introduction

### 1.1 Life in the slow lane

Of all the reasons given for not using public transportation, “it takes too long” is probably the most widespread. After all, compared to the convenience of the private automobile, even the simplest grocery store errand demands a complex, multi-step planning process; transfers must be timed, stops must be located, and page after page of schedules must be carefully inspected. On top of that, the transit vehicles themselves run slow—and quite often, they run late, especially during periods of heavy traffic and road construction. Because of these factors, the choice between stewing in a city bus and blazing home in a muscle car in half the time is all too clear for the vast majority of Americans. Public transportation is inconvenient, unreliable, *and especially slow*.

The bus is the most widely employed mode of transit, but it is also the most disadvantaged in terms of speed. Condemned to use the same lanes used by general traffic, the bus faces the dual challenge of making stops and negotiating rush hour, while subways, metros, and even light rail vehicles (historically known as “streetcars”) glide from station to station on dedicated rights of way. But trains are outrageously expensive. American transit agencies have turned to a lower cost option: *make the buses go faster*. A new concept, known as **Bus Rapid Transit** (BRT), promises the benefits of fast and reliable rail systems combined with the cost effectiveness and flexibility of ordinary buses. Dedicated transit lanes segregate buses from other traffic. Off-vehicle fare payments, upgraded stops, and bigger doors allow passengers to board faster. Traffic signal tweaks give buses green lights. All of these perks shave minutes off the schedule and improve the rider experience.

However, every BRT system still faces reliability issues. Even the fanciest bus with the best BRT infrastructure in place can run late. *For the riders, a faster bus matters not if they find themselves waiting forever for it*. Twenty first century technology offers a potential solution: real-time bus tracking. GPS trackers fitted onto buses allow the transit agency to predict precisely when they will arrive at the stop. This information is made available to customers via hotlines, text messages, and smartphone apps. For many riders, real-time “next bus” estimates are game changers. One can imagine, for example, sitting in an air-conditioned coffee shop and deciding whether to stay inside or brave the hot summer

weather on the chance that the bus is actually keeping its schedule. By decreasing waiting times—and passenger frustration—real-time tracking makes public transit feel faster and more convenient.

Austin’s own Capital Metro transit authority recognized BRT’s enormous potential and implemented two new BRT routes this year. Route 801, established in January, and route 803, established in August, now carry passengers along Austin’s busiest north-south corridors. Branded as “MetroRapid,” the new routes feature:

- Bus lanes in downtown Austin.
- Fewer stops with wider spacing.
- Traffic signal prioritization outside of downtown Austin, which lengthens green lights so that buses can catch them.
- Premium, high-capacity buses with larger doors and the ability to board through any door, not just the front one, thereby streamlining the boarding process.
- Pre-boarding payments via smart cards and the Capital Metro smartphone app, also reducing boarding times.
- Upgraded bus stops, with raised boarding platforms at select ones.

MetroRapid (along with Capital Metro’s commuter rail service, MetroRail) also includes a real-time bus tracking system provided by the Trapeze Group. GPS positions for each bus are updated within 90 seconds; predicted arrival times are then displayed at digital signs at bus stops, within the Capital Metro smartphone app, and on the agency’s online trip planner.

## 1.2 A real-time problem

Despite studying computer science, I have always held an interest in public transportation. My friends and family roll their eyes (and call me outright crazy) whenever I insist on riding the bus around town. I have often played transportation computer games such as *SimCity*, *Transport Tycoon*, and *Simutrans*, and I enjoy learning about the latest transit technologies such as bus rapid transit.

One day, while exploring Austin, I was waiting for the 803 MetroRapid bus at the UT Pickle Research Campus. I noticed that the displayed arrival time was fluctuating between two to five minutes. I wondered how bus arrival times are estimated, and if the algorithms could be improved.

Sure enough, I uncovered a local news story from Columbus, Ohio complaining about deficiencies in Trapeze’s software. The transit agency had invested millions of dollars and years of development time into a system that simply didn’t work. The agency claimed that the the accuracy rates were “as low as 25 percent” and that the GPS tracking technology was unreliable (Rouan).

A search of the UT library system revealed a few of papers on the subject of bus prediction. The first, by Wei-Hua Lin and Jian Zeng, was an experimental study of real-time bus tracking in Blacksburg, Virginia, a small rural town. The researchers used GPS-based bus locations and on-time performance statistics from the Blacksburg transit agency to develop and compare four bus arrival prediction algorithms. The algorithms were measured on the metrics of deviation (how closely the algorithms predicted the actual time) and stability (how much the algorithm predictions changed over time). The study concluded that the best-performing algorithm took into account the location of the bus, its adherence to the schedule, and whether or not the stop is a time point (at time point stops, the bus must wait until the listed departure time to proceed if it is running early) (Lin and Zeng). However, limitations in the quality of the data were noted. GPS tracking data have a limited resolution, and tracking units can go for minutes at a time without updates. For cases like these, location data have to be interpolated, sacrificing accuracy (104).

Although interesting, the study has little direct application to MetroRapid due to the small population of Blacksburg. The study authors did not take traffic into account because Blacksburg has little traffic congestion (Lin and Zeng 106). Obviously, this would not be the case for Austin, as anybody driving down Guadalupe Street during rush hour can attest to.

The second paper, by Ali Abdelfattah and Ata Khan, was about models for predicting bus delays. Like the Blacksburg study, the researchers created a number of bus arrival time algorithms and measured their performance. However, unlike the Blacksburg study, they used a computerized traffic simulation based on an Ottawa, Canada bus route to develop the models instead of real-world data. After tuning the models within the simulation program, the authors then field-tested them by predicting bus arrival times at a number of intersections. The study concluded that the most effective algorithm used a combination of distance, number of bus stops, and travel time as inputs (Abdelfattah and Khan).

For my inquiry, I decided to apply Abdelfattah and Khan's study to MetroRapid. Using their models, could I predict bus arrival times better than Capital Metro could? While the question of making better predictions for bus arrival times has been studied by others, to my knowledge nobody has examined the problem with today's widely deployed bus arrival prediction systems—and certainly not with Capital Metro's.

## 2 The experiment

The general procedure was:

1. Develop a model to predict the estimated arrival time for a bus.
2. Collect location and estimated arrival time data for MetroRapid buses over a week.
3. Tune the model to fit this data and attempt to predict bus arrival times for the next few days.

To simplify the experiment, I decided to focus on four bus stops instead of the entire MetroRapid system. Half were located in suburban areas, and the other half were located in downtown:

- Crestview Station (suburban)
- UT West Mall Station (downtown)
- Republic Square Station (downtown)
- Pleasant Hills Station (suburban)

All four stations were on the southbound side for route 801. Route 801 was selected because it is the more mature of the two MetroRapid routes. The southbound direction was chosen by coin toss.

## 2.1 The scraper

First, I wrote a computer program (see A.1) to record bus location data from Capital Metro. This data is publicly available through a hidden service. For help, I turned to the authors of the MetroRappid app, a smartphone app created by local Austinites that provides an alternative to the official Capital Metro app. They had written documentation for requesting locations and arrival times from Capital Metro’s servers (Dawoodjee).

My program, written in the Python programming language, collects and saves bus locations every 30 seconds. Note that, as in the Blacksburg study, the intervals between GPS updates can vary. Capital Metro quotes within 90 seconds, while MetroRappid authors claim between 30 to 90 seconds. I chose 30 seconds to achieve high precision while minimizing the number of requests made to Capital Metro servers.

Due to the way Trapeze’s software is designed, estimated arrival times for each station must be retrieved separately. Data for each station was retrieved every 60 seconds. With four stations in the experiment, this worked out to one request every 15 seconds.

Once I got the program written and working, I let it run for a week while I turned my attention to developing a model.

## 2.2 My algorithm

For ease of implementation, I opted to implement Abdelfattah and Khan’s linear bus delay model:

$$DELAY = 0.4855 + 0.0287 \times DENS.LT + 0.0168 \times DENS.TH + 0.9654 \times LENGTH - 1.1969 \times M/T + 0.1130 \times STATION$$

- *DENS.LT* is the left-turn vehicle density in vehicles per lane per kilometer. For this value, I estimated that about  $\frac{1}{8}$  of all traffic would be making left turns. Thus, I used  $\frac{1}{8} \times DENS.TH$ .
- *DENS.TH* is the straight-through vehicle density in vehicles per lane per kilometer. For this value, I estimated about 15 vehicles/lane/km for suburban stations and double that for downtown stations. This was based on analysis of Google Earth imagery.

- *LENGTH* is distance to the station in kilometers. For this value, I used the straight-line difference between the station’s coordinates and the bus’s coordinates. While this is not as ideal as the roadway distance, the route of the 801 is mostly straight, so this approximation seemed “good enough.”
- *M/T* is the ratio of moving time to traveling time. (Traveling time is time spent between stations, while moving time is time spent actually moving, i.e. not at a stoplight or stuck in traffic.) Based on my own experience riding MetroRapid, I estimated this at 0.5.
- *STATION* is the number of bus stops that the bus must pass before arriving at the target stop. This was computed based on the stations per kilometer ratio; that is,  $STATION = LENGTH \times (stations\ per\ km)$ . For suburban areas, this is 2. For downtown areas, this is 3.
- *DELAY* is the delay, in minutes, that the bus is expected to face.

According to Adelfattah and Khan, from *DELAY* we can estimate the arrival time with the equation:

$$travel\ time = \frac{distance}{bus\ speed\ without\ delay} + DELAY$$

Where (*bus speed without delay*) is assumed to be 40 km/h. This is the same value that I used in my program.

Finally, by adding (*travel time*) to the current time, we arrive at the predicted arrival time.

## 2.3 The first analysis

To measure the performance of arrival time estimates, the following metrics were used:

- Mean error, the mean of the differences between the predicted and actual arrival times. This gives a good measure of how close the algorithm was.
- Standard deviation error, the standard deviation of the differences between the predicted and actual arrival times. This measures the spread of the predicted times. Lower spread is more desirable.
- Stability. Taking a cue from the Blacksburg study, I also included this metric, which Lin and Zeng defined mathematically as the average of  $|e_n - e_{n-1}|$ ; that is, the average of the differences between successive predicted times. This metric is designed to detect algorithms that produce vastly different times in a short period of time. For example, predicting “5:45 PM” one minute and then “5:57 PM” the next is confusing and undesirable from a passenger’s perspective.

After a week had elapsed, I extended my program to determine the actual bus arrival times for each station and compare them with Capital Metro’s predicted times. Taking another hint from the Blacksburg study, I also compared the arrival times with the times listed on the schedule.

Capital Metro’s system stops predicting arrival times once the bus has actually arrived at the stop. Therefore, we can assume that the actual arrival time is the time at which Capital Metro no longer gives an arrival estimate for that bus.

### 2.3.1 Tables and figures

Current time	Actual arrival time	Scheduled arrival time	Capital Metro est. arrival time	Capital Metro est. error (min)
11:13 AM	11:45 AM	11:44 AM	11:44 AM	-1
11:14 AM	11:45 AM	11:44 AM	11:44 AM	-1
11:15 AM	11:45 AM	11:44 AM	11:45 AM	0
11:16 AM	11:45 AM	11:44 AM	11:45 AM	0
11:17 AM	11:45 AM	11:44 AM	11:45 AM	0
11:18 AM	11:45 AM	11:44 AM	11:45 AM	0
11:20 AM	11:45 AM	11:44 AM	11:44 AM	-1
11:21 AM	11:45 AM	11:44 AM	11:44 AM	-1
11:22 AM	11:45 AM	11:44 AM	11:43 AM	-2
11:23 AM	11:45 AM	11:44 AM	11:43 AM	-2
11:24 AM	11:45 AM	11:44 AM	11:43 AM	-2
11:25 AM	11:45 AM	11:44 AM	11:44 AM	-1
11:26 AM	11:45 AM	11:44 AM	11:44 AM	-1
11:28 AM	11:45 AM	11:44 AM	11:44 AM	-1
11:29 AM	11:45 AM	11:44 AM	11:44 AM	-1
11:30 AM	11:45 AM	11:44 AM	11:44 AM	-1
11:31 AM	11:45 AM	11:44 AM	11:44 AM	-1
11:32 AM	11:45 AM	11:44 AM	11:45 AM	0
11:33 AM	11:45 AM	11:44 AM	11:45 AM	0
11:34 AM	11:45 AM	11:44 AM	11:44 AM	-1
11:36 AM	11:45 AM	11:44 AM	11:44 AM	-1
11:37 AM	11:45 AM	11:44 AM	11:45 AM	0
11:38 AM	11:45 AM	11:44 AM	11:45 AM	0
11:39 AM	11:45 AM	11:44 AM	11:45 AM	0
11:40 AM	11:45 AM	11:44 AM	11:46 AM	1
11:41 AM	11:45 AM	11:44 AM	11:46 AM	1
11:42 AM	11:45 AM	11:44 AM	11:46 AM	1
11:44 AM	11:45 AM	11:44 AM	11:45 AM	0
11:45 AM	11:45 AM	11:44 AM	11:46 AM	1

Table 1: Sample data for the Wednesday 11:44 AM bus at UT West Mall Station.

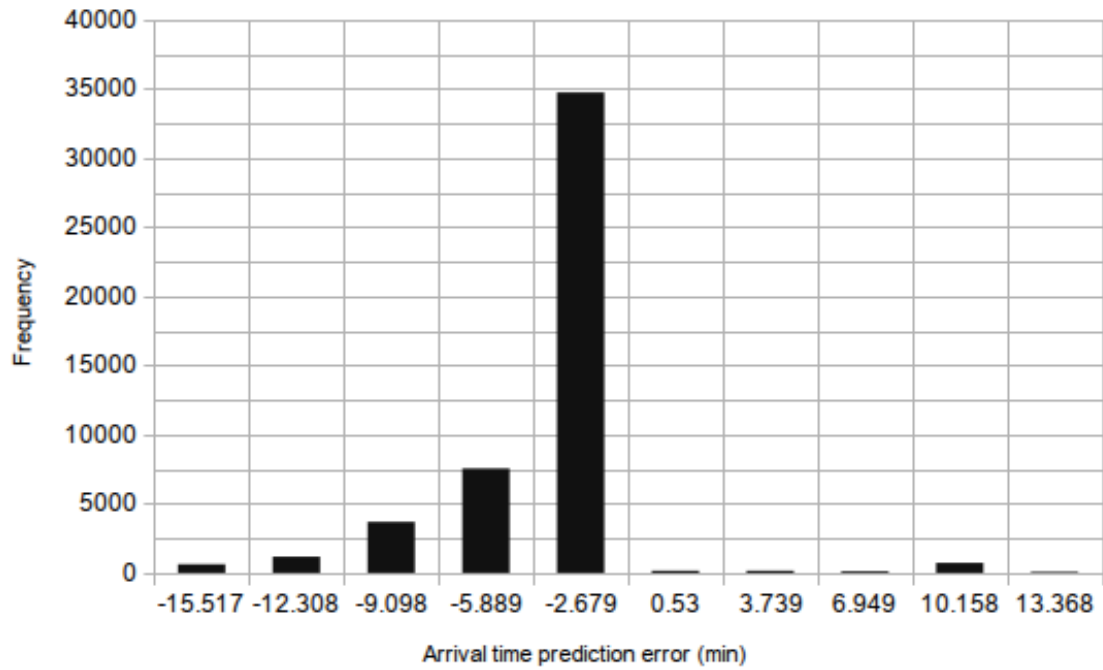


Figure 1: Histogram for schedule-based arrival time predictions.

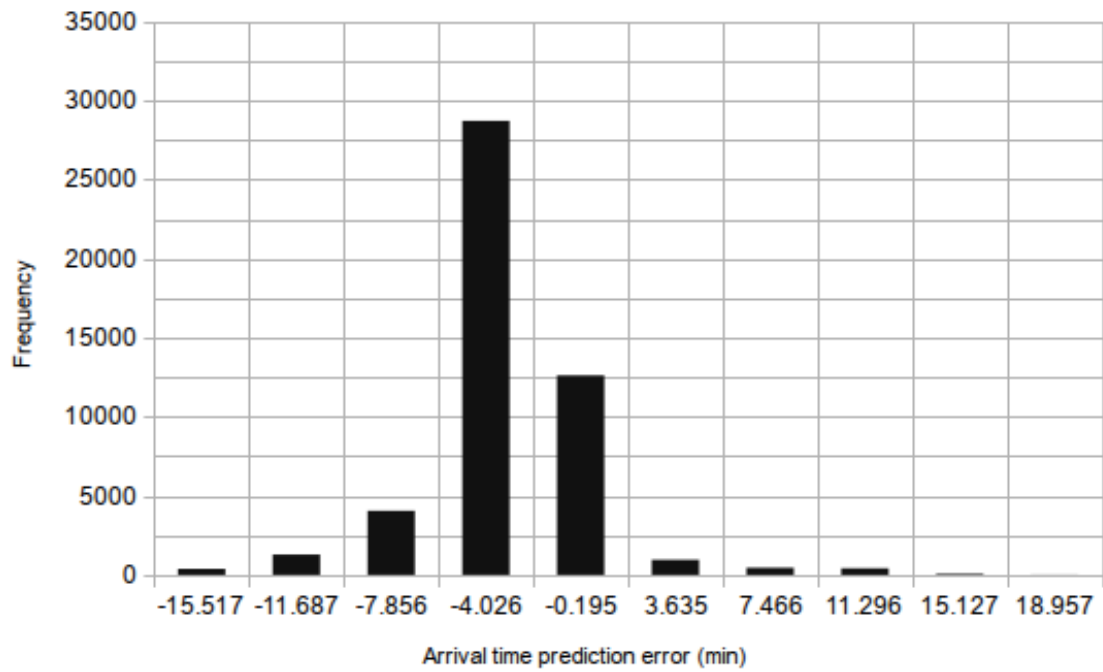


Figure 2: Histogram for Capital Metro arrival time predictions.

	Schedule	Capital Metro predicted
Error mean	-1.785	-1.137
Error std. deviation	2.342	2.614
Stability	0.0	0.390

Table 2: Summary of bus arrival prediction data over a six-day period. All units in minutes.

### 2.3.2 Comments

This analysis was very surprising.

First, note that the “stability” value for the schedule is zero. This is to be expected, because the arrival time listed on the schedule—and its difference from the actual arrival time—is constant.

More intriguing, however, are the mean and standard deviation values. On average, buses arrived within a minute of the time predicted by Capital Metro and within two minutes of the scheduled time. From the standard deviation values and the histogram plots, we can conclude that in the majority of cases the error was within five minutes. First, this indicates that Capital Metro ought to be congratulated, because most of the MetroRapid buses were running on time. Second, this calls into question the usefulness of real-time bus tracking. The expense of onboard GPS trackers and the computer systems required to utilize them amount to what is effectively a *0.65-minute difference* in the average prediction error, compared to simply handing each passenger a copy of the bus schedule.

## 2.4 Tuning my algorithm

From the initial data I had collected, I was now ready to adjust and fine-tune my own arrival time model. However, I soon realized that I had made a major mistake.

When requesting bus locations, Capital Metro’s server returns not one but a list of locations for each bus. The locations go back in time, creating a “track” for the bus as it moves along its route. Based on the design of the MetroRappid app, I had assumed that the *last* location in the list is the latest one.

But as I examined the data by hand, things weren’t adding up. A bus would be considered “arrived” when it seemed from the location data that it was actually still several blocks away. Thinking it was an anomaly, I discarded the set. But then bus after bus showed the same peculiarity. Finally, it dawned on me; the *first* location is the latest, not the last one.

I could still analyze the performance of Capital Metro’s arrival predictions, because I did not need good bus location data to do that. However, my week’s worth of data were now useless for developing my own algorithm. My only choice was to quickly correct the program’s flaw and collect a weekend’s worth of good location data to use.

After collecting the weekend data set, I extended my program (see A.2) to analyze it in retrospect, applying my new algorithm to determine its performance. After some observations, I made the following tweaks to my model:

- Do not predict arrival times if the scheduled arrival time is at least 30 minutes away. This is because predicted times varied substantially when the bus was still 30 minutes



out, and they would not be of much use to the passenger anyway. From the raw data, it appears that the Capital Metro system does this as well.

Current time	Actual arrival time	Scheduled arrival time	Capital Metro est. arrvl. time
<i>09:52 AM</i>	<i>10:51 AM</i>	<i>10:52 AM</i>	<i>10:52 AM</i>
<i>09:53 AM</i>	<i>10:51 AM</i>	<i>10:52 AM</i>	<i>10:52 AM</i>
<i>09:54 AM</i>	<i>10:51 AM</i>	<i>10:52 AM</i>	<i>10:52 AM</i>
<i>09:56 AM</i>	<i>10:51 AM</i>	<i>10:52 AM</i>	<i>10:52 AM</i>
...	...	...	...
<i>10:21 AM</i>	<i>10:51 AM</i>	<i>10:52 AM</i>	<i>10:52 AM</i>
<i>10:22 AM</i>	<i>10:51 AM</i>	<i>10:52 AM</i>	<i>10:52 AM</i>
10:23 AM	10:51 AM	10:52 AM	10:53 AM
10:24 AM	10:51 AM	10:52 AM	10:54 AM
10:25 AM	10:51 AM	10:52 AM	10:53 AM
...	...	...	...

Table 3: Estimated arrival times remain constant until 30 minutes before the scheduled arrival time.

- Reduce the prediction frequency from every one minute to every two minutes. This smooths out abrupt changes in the predicted arrival time.

After these modifications, the results seemed promising:

	Schedule	Capital Metro predicted	My algorithm predicted
Error mean	-1.954	-1.425	-0.794
Error std. deviation	3.956	3.729	4.384
Stability	0.0	0.196	0.202

Table 4: Summary of bus arrival prediction data over a two-day period. My algorithm tested retrospectively. All units in minutes.

Although standard deviation and stability were slightly higher than the Capital Metro data, indicating that my algorithm’s predicted times varied slightly more, the mean error was lower. With my model tested, it was time to see how it performed in the real world.

## 2.5 The second analysis

I modified the scraper program (see A.3) to predict bus arrival times using my algorithm in addition to collecting Capital Metro’s predicted arrival times. Thus, the algorithm was tested by interpreting data in real time instead of reexamining data that had been previously collected.

After running for two days (on Sunday and Monday), the results were in.

### 2.5.1 Tables and figures

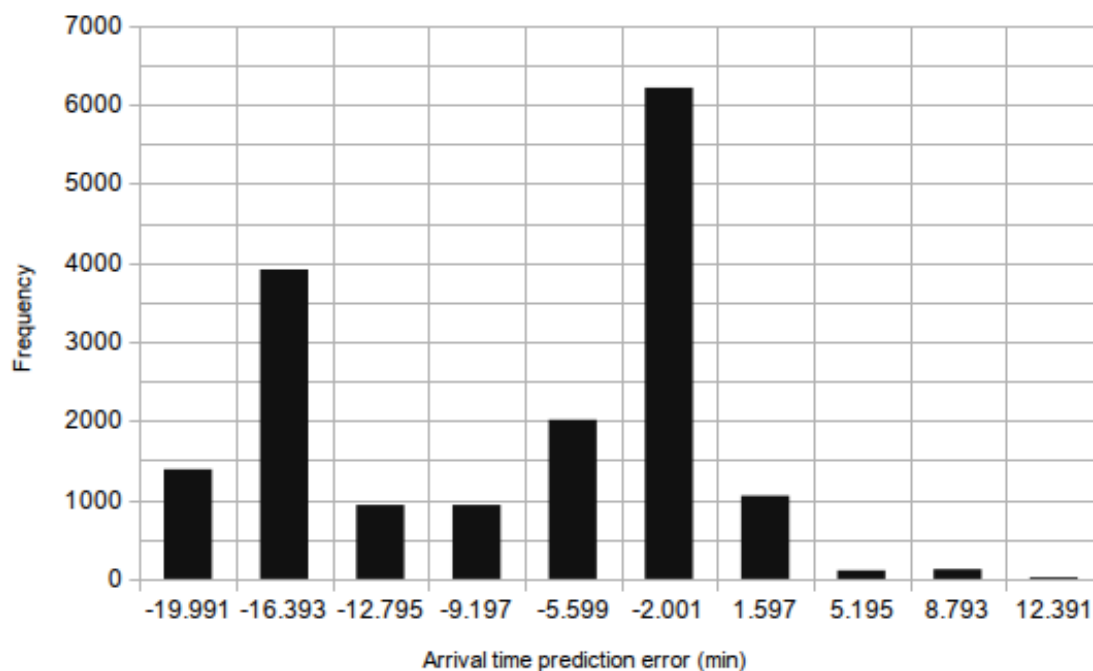


Figure 3: Histogram for my algorithm’s arrival time predictions.

	Schedule	Capital Metro predicted	My algorithm predicted
Error mean	-1.632	-0.741	-6.372
Error std. deviation	3.230	2.795	7.048
Stability	0.0	0.296	0.491

Table 5: Summary of bus arrival prediction data over a two-day period.

### 2.5.2 Comments

Needless to say, the algorithm’s performance was poor. On average, it produced arrival times that were approximately six minutes too early. Considering the high standard deviation and a histogram that indicates widely dispersed values, the predicted times were not very precise either. By every measure, the algorithm was a worse predictor of bus arrival times than both Capital Metro and the bus schedule.

## 3 Conclusions

The biggest problem with the experiment was a simple lack of resources. Attempting to come up with a good model for predicting bus arrival times proved far too difficult and complex

for me to execute well. Consider the equation that I was using; I had to guess almost all of the variables based on my own intuition, and despite my enthusiasm for the subject, I am no transportation engineer. Also, without time and money, I could not determine these values precisely using surveys or by consulting with the city of Austin.

Another major issue was the quality of the GPS-based bus location data. As in the Blacksburg study, oftentimes buses would not send updated location information for minutes at a time, as demonstrated by a typical recording of a bus's track:

Time	Location (latitude, longitude)
05:50 PM	30.222734, -97.766434
05:51 PM	30.222488, -97.766525
05:51 PM	30.218885, -97.766525
05:52 PM	30.218885, -97.766525
05:53 PM	30.20541, -97.774696
05:53 PM	30.20541, -97.774696
05:54 PM	30.20541, -97.774696
05:54 PM	30.20541, -97.774696
05:55 PM	30.20541, -97.774696
05:55 PM	30.20541, -97.774696
05:56 PM	30.200048, -97.776283
05:56 PM	30.200048, -97.776283
05:57 PM	30.192375, -97.779228
05:57 PM	30.192375, -97.779228
05:58 PM	30.192375, -97.779228
05:58 PM	30.171007, -97.786095
05:59 PM	30.171007, -97.786095
05:59 PM	30.171007, -97.786095
06:00 PM	30.171007, -97.786095

Hence, there was always a high degree of uncertainty as to the exact location of the bus.

Other flaws in my methodology were more subtle. For instance, I had assumed that a bus had arrived at the stop when Capital Metro no longer offered an arrival estimate for it. Close inspection of the location data, however, revealed that this would sometimes result in a computed arrival time that was a minute or two later than the actual arrival time. As the GPS data are not precise enough to determine exactly when a bus has arrived, the ideal solution would be to have observers recording arrival times at the stop. Unfortunately, this was not practical concept for my inquiry.

A secondary issue was that I measured the distances between buses and stops in terms of straight-line distances, not roadway distances (as was done in the Blacksburg and Ottawa studies), which would be more accurate. The latter requires a detailed computer representation of the bus route, something that was out of the question given my limited skill set.

While my bus arrival time prediction model was a complete failure, there was one neat thing that I learned from my inquiry: real-time arrival estimates, at least for route 801, do

not actually make much of a difference. The difference compared to the published schedule was a minute or two at best. Since the buses are almost always running early or on time, Capital Metro might consider tightening the schedule to reduce the travel time along the route. Also, I suggest a follow up study focusing on the effectiveness of current real-time tracking systems. Can their costs be justified? Would it be more effective, for example, to reserve their use for detours and extreme traffic jams?

As for me, the next time I'm waiting to catch a ride on the MetroRapid, I won't be staring nervously at the estimated arrival time.

## Works Cited

- Abdelfattah, Ali M and Ata M Khan. "Models for predicting bus delays". *Transportation Research Record: Journal of the Transportation Research Board* 1623.1 (1998): 8–15. Print.
- Dawoodjee, Luqmaan. "The CapMetro API". 2014. Web. 17 Oct. 2014.
- Lin, Wei-Hua and Jian Zeng. "Experimental study of real-time bus arrival time prediction with GPS data". *Transportation Research Record: Journal of the Transportation Research Board* 1666.1 (1999): 101–109. Print.
- Rouan, Rick. "COTA says its real-time bus-tracking system doesn't work". The Columbus Dispatch, 24 July 2014. Web. 17 Oct. 2014.

## A Source code listings

### A.1 phaseI.py

```
1 import sys
2 import time
3 import datetime
4 import pickle
5 import sched
6 import urllib.request
7 import urllib.parse
8 import xml.dom.minidom
9 import math
10 import statistics
11
12 DATA_FILE = "phaseI.pickle"
13
14 CAPMETRO_NEXTBUS = "http://www.capmetro.org/planner/s_nextbus2.asp"
15 CAPMETRO_BUSLOCS = "http://www.capmetro.org/planner/s_buslocation.asp"
16 POLL_INTERVAL = 30
17 SAVE_INTERVAL = 15 * 60
18
19 # Split trips at points that are this amount of time apart.
20 TRIP_TIMEDELTA_SPLIT = datetime.timedelta(hours = 12)
21 ACCEPTABLE_ERROR_MIN = 20
22 DEFAULT_SPEED_KMH = 40
23 ARRIVAL_ESTIMATE_INTERVAL = datetime.timedelta(minutes = 2)
24
25 class Station:
26     def __init__(self, stop_id, name, latitude, longitude):
27         self.stop_id = stop_id
28         self.name = name
```

```

29     self.trips = {}
30     self.latitude = latitude
31     self.longitude = longitude
32 def record_trip(self, timestamp, trip_id, vehicle_id, sched_arrival, est_arrival):
33     def fix_time(time):
34         time = time.strip()
35         # Error parsing estimated arrival time '12:02 XM': time data '12:02 XM'
36         # does not match format '%I:%M %p'
37         time = time.replace("X", "A")
38         # Error parsing estimated arrival time '00:29 AM': time data '00:29 AM'
39         # does not match format '%I:%M %p'
40         time = time.replace("00", "12")
41         return time
42     if not trip_id in self.trips:
43         self.trips[trip_id] = []
44     sched_arrival_time = None
45     sched_arrival = fix_time(sched_arrival)
46     try:
47         sched_arrival_time = datetime.datetime.strptime(sched_arrival,
48                 "%I:%M%p").time()
49     except Exception as e:
50         print("Error_parsing_scheduled_arrival_time_" + sched_arrival + "':" +
51                 str(e))
52     est_arrival_time = None
53     est_arrival = fix_time(est_arrival)
54     try:
55         est_arrival_time = datetime.datetime.strptime(est_arrival,
56                 "%I:%M%p").time()
57     except Exception as e:
58         print("Error_parsing_estimated_arrival_time_" + est_arrival + "':" +
59                 + str(e))
60     self.trips[trip_id].append({
61         "timestamp": timestamp,
62         "sched_arrival_time": sched_arrival_time,
63         "est_arrival_time": est_arrival_time,
64         "vehicle_id": vehicle_id
65     })
66
67     stations = [
68         Station(497, "UT_West_Mall", 30.286064, -97.741815),
69         Station(5867, "Republic_Square", 30.267751, -97.746857),
70         Station(5606, "Crestview", 30.337852, -97.719035),
71         Station(5872, "Pleasant_Hill", 30.192391, -97.779203)
72     ]
73     vehicles = {}
74
75     # Credit John D. Cook, http://www.johndcook.com/python\_longitude\_latitude.html
76     # (public domain)
77     def distance_on_unit_sphere(lat1, long1, lat2, long2):
78
79         # Convert latitude and longitude to
80         # spherical coordinates in radians.
81         degrees_to_radians = math.pi/180.0
82
83         # phi = 90 - latitude
84         phi1 = (90.0 - lat1)*degrees_to_radians
85         phi2 = (90.0 - lat2)*degrees_to_radians
86
87         # theta = longitude
88         theta1 = long1*degrees_to_radians
89         theta2 = long2*degrees_to_radians
90
91         # Compute spherical distance from spherical coordinates.
92
93         # For two locations in spherical coordinates
94         # (1, theta, phi) and (1, theta, phi)
95         # cosine( arc length ) =
96         # sin phi sin phi' cos(theta-theta') + cos phi cos phi'

```

```

97     # distance = rho * arc length
98
99     cos = (math.sin(phi1)*math.sin(phi2)*math.cos(theta1 - theta2) +
100            math.cos(phi1)*math.cos(phi2))
101     arc = math.acos( cos )
102
103     # Remember to multiply arc by the radius of the earth
104     # in your favorite set of units to get length.
105     return arc
106
107 def dom_value(node):
108     if node.firstChild == None:
109         return ""
110     else:
111         return node.firstChild.data
112
113 def combine_time_date(time, dt):
114     # Combines a time with date information from a datetime. We need some extra
115     # logic to handle cases straddling midnight.
116     midnight = datetime.time(0, 0)
117     if time > midnight and dt.time() < midnight:
118         dt_correct = datetime.timedelta(days = 1)
119     elif time < midnight and dt.time() > midnight:
120         dt_correct = datetime.timedelta(days = -1)
121     else:
122         dt_correct = datetime.timedelta()
123     return datetime.datetime.combine(dt.date(), time) + dt_correct
124
125 def load_data():
126     global stations
127     global vehicles
128     try:
129         prev_data = pickle.load(open(DATA_FILE, "rb"))
130         stations = prev_data["stations"]
131         vehicles = prev_data["vehicles"]
132     except Exception as e:
133         print("Error_loading_data_file:_ " + str(e))
134
135 def save_data():
136     global stations
137     global vehicles
138     print("Saving_data...")
139     pickle.dump({ "stations": stations, "vehicles": vehicles },
140                open(DATA_FILE, "wb"))
141
142 def retrieve_station_data(station):
143     url = CAPMETRO_NEXTBUS + "?" + urllib.parse.urlencode({ "route": "801",
144                                                            "stopid": str(station.stop_id) })
145     http_handle = None
146     try:
147         print("Downloading_" + url)
148         http_handle = urllib.request.urlopen(url)
149     except Exception as e:
150         print("Error_retrieving_data_for_station_" + str(station.stop_id) + ":_ " +
151              str(e))
152     return False
153 xml_tree = None
154 now = datetime.datetime.now()
155 try:
156     xml_tree = xml.dom.minidom.parse(http_handle)
157     runs = xml_tree.getElementsByTagName("Run")
158     for trip in runs:
159         vehicle_id = dom_value(trip.getElementsByTagName("Vehicleid")[0]).strip()
160         if dom_value(trip.getElementsByTagName("Valid")[0]) == "Y":
161             station.record_trip(
162                 timestamp = now,
163                 trip_id = dom_value(trip.getElementsByTagName("Tripid")[0]).strip(),
164                 vehicle_id = vehicle_id,

```

```

165         sched_arrival = dom_value(trip.getElementsByTagName("Triptime")[0]),
166         est_arrival = dom_value(trip.getElementsByTagName("Estimatedtime")[0])
167     )
168 except Exception as e:
169     print("Error_parsing_tracking_data_for_station_" + str(station.stop_id) +
170           ":", e)
171     return False
172 return True
173
174 def retrieve_vehicle_data():
175     global vehicles
176     url = CAPMETRO.BUSLOCS + "?" + urllib.parse.urlencode({ "route": "801" })
177     try:
178         print("Downloading_" + url)
179         http_handle = urllib.request.urlopen(url)
180     except Exception as e:
181         print("Error_retrieving_bus_tracking_data:", e)
182         return False
183     xml_tree = None
184     now = datetime.datetime.now()
185     try:
186         xml_tree = xml.dom.minidom.parse(http_handle)
187         xml_vehicles = xml_tree.getElementsByTagName("Vehicle")
188         for vehicle in xml_vehicles:
189             vehicle_id = dom_value(vehicle.getElementsByTagName("Vehicleid")[0]).strip()
190             latest_pos = [float(x) for x in dom_value(
191                 vehicle.getElementsByTagName("Position")[0]).split(",")]
192             if not vehicle_id in vehicles:
193                 vehicles[vehicle_id] = []
194             vehicle_data = vehicles[vehicle_id]
195             # Do not record duplicate data (data that hasn't been updated within our
196             # update interval).
197             if (len(vehicle_data) == 0 or
198                 vehicle_data[-1]["latitude"] != latest_pos[0] or
199                 vehicle_data[-1]["longitude"] != latest_pos[1]):
200                 vehicle_data.append({
201                     "timestamp": now,
202                     "block": dom_value(
203                         vehicle.getElementsByTagName("Block")[0]).strip(),
204                     "reliable": dom_value(
205                         vehicle.getElementsByTagName("Reliable")[0]).strip() == "Y",
206                     "off_route": dom_value(
207                         vehicle.getElementsByTagName("Offroute")[0]).strip() == "Y",
208                     "stopped": dom_value(
209                         vehicle.getElementsByTagName("Stopped")[0]).strip() == "Y",
210                     "in_service": dom_value(
211                         vehicle.getElementsByTagName("Inservice")[0]).strip() == "Y",
212                     "speed": float(dom_value(vehicle.getElementsByTagName("Speed")[0])),
213                     "heading": int(dom_value(vehicle.getElementsByTagName("Heading")[0])),
214                     "latitude": latest_pos[0],
215                     "longitude": latest_pos[1]
216                 })
217     except Exception as e:
218         print("Error_parsing_bus_tracking_data:", e)
219         return False
220     return True
221
222 def collect():
223     load_data()
224     scheduler = sched.scheduler(time.time, time.sleep)
225     def collect_stations(i):
226         retrieve_station_data(stations[i])
227         i = (i + 1) % len(stations)
228         now = datetime.datetime.now()
229         now_time = now.time()
230         # Cease collecting between 12:45 AM and 4:45 AM.
231         if now_time >= datetime.time(0, 45) and now_time <= datetime.time(4, 45):
232             scheduler.enterabs(time.mktime(datetime.datetime.combine(now,

```

```

233         datetime.time(4, 45)).timetuple()), 2, collect_stations, (0,))
234     else:
235         scheduler.enter(POLLINTERVAL / len(stations), 2, collect_stations,
236                         (i,))
237 def collect_vehicles():
238     retrieve_vehicle_data()
239     now = datetime.datetime.now()
240     now_time = now.time()
241     # Cease collecting between 12:45 AM and 4:45 AM.
242     if now_time >= datetime.time(0, 45) and now_time <= datetime.time(4, 45):
243         scheduler.enterabs(time.mktime(datetime.datetime.combine(now,
244                                     datetime.time(4, 45)).timetuple()), 1, collect_vehicles, ())
245     else:
246         scheduler.enter(POLLINTERVAL, 1, collect_vehicles, ())
247 def collect_save():
248     save_data()
249     scheduler.enter(SAVEINTERVAL, 3, collect_save, ())
250
251 collect_vehicles()
252 collect_stations(0)
253 collect_save()
254 try:
255     scheduler.run()
256 except KeyboardInterrupt:
257     print("Stopped.")
258     save_data()
259     sys.exit(0)
260
261 def split_unique_trips(trip):
262     # Trips are distinguished by trip ID, but this is only unique for one day.
263     # This function splits one "trip" into separate trips for each date.
264     trips = []
265     i = 0
266     this_ts = None
267     last_ts = None
268     last_split = 0
269     while i < len(trip):
270         this_ts = trip[i]["timestamp"]
271         if last_ts != None and this_ts - last_ts > TRIP_TIMEDELTA_SPLIT:
272             trips.append(trip[last_split:i])
273             last_split = i
274         last_ts = this_ts
275         i += 1
276     if trip[last_split:] != []:
277         trips.append(trip[last_split:])
278     return trips
279
280 def get_vehicle_position(vehicle_id, dt):
281     global vehicles
282     if vehicle_id in vehicles:
283         for point in vehicles[vehicle_id]:
284             if (abs(point["timestamp"] - dt) < datetime.timedelta(seconds = 30) or
285                 point["timestamp"] > dt):
286                 return point
287         return None
288     else:
289         return None
290
291 def trip_error_data(trip, station):
292     error_data = {
293         "schedule": [],
294         "capmetro": [],
295         "me": []
296     }
297     arrival_time = trip[-1]["timestamp"]
298     this_ts = None
299     this_dt_correct = None
300     this_sched = None

```



```

301     this_est = None
302     for point in trip:
303         this_ts = point["timestamp"]
304         vehicle_id = point["vehicle_id"]
305         this_est = combine_time_date(point["est_arrival_time"], this_ts)
306         this_sched = combine_time_date(point["sched_arrival_time"], this_ts)
307         # Calculate error and append to the lists.
308         schedule_error_min = (this_sched - arrival_time).total_seconds() / 60
309         capmetro_error_min = (this_est - arrival_time).total_seconds() / 60
310         if abs(schedule_error_min) < ACCEPTABLE_ERROR_MIN:
311             error_data["schedule"].append(schedule_error_min)
312             error_data["capmetro"].append(capmetro_error_min)
313     return error_data
314
315 def compute_stability(data):
316     diffs = []
317     i = 1
318     while i < len(data):
319         diffs.append(abs(data[i] - data[i - 1]))
320         i += 1
321     return statistics.mean(diffs)
322
323 def analyze():
324     def append_keys(dict1, dict2):
325         for key, value in dict1.items():
326             dict2[key].append(value)
327     global stations
328     load_data()
329     metrics = {
330         "schedule": {
331             "_data": [],
332             "stability": []
333         },
334         "capmetro": {
335             "_data": [],
336             "stability": []
337         }
338     }
339     for station in stations:
340         for trip_id, trip_data in station.trips.items():
341             trips = split_unique_trips(trip_data)
342             for trip in trips:
343                 metrics["schedule"]["_data"] += this_error_data["schedule"]
344                 metrics["capmetro"]["_data"] += this_error_data["capmetro"]
345                 if len(this_error_data["schedule"]) > 1:
346                     metrics["schedule"]["stability"].append(
347                         compute_stability(this_error_data["schedule"]))
348                     metrics["capmetro"]["stability"].append(
349                         compute_stability(this_error_data["capmetro"]))
350
351     print("\nBUS_SCHEDULE_ONLY")
352     print("Error_mean:_ " + str(statistics.mean(metrics["schedule"]["_data"])))
353     print("Error_std_dev:_ " + str(statistics.stdev(metrics["schedule"]["_data"])))
354     print("Stability:_ " + str(statistics.mean(metrics["schedule"]["stability"])))
355     print("\nCAPITAL_METRO_ESTIMATES")
356     print("Error_mean:_ " + str(statistics.mean(metrics["capmetro"]["_data"])))
357     print("Error_std_dev:_ " + str(statistics.stdev(metrics["capmetro"]["_data"])))
358     print("Stability:_ " + str(statistics.mean(metrics["capmetro"]["stability"])))
359
360 if __name__ == "__main__":
361     if sys.argv[1] == "collect":
362         collect()
363     elif sys.argv[1] == "analyze":
364         analyze()

```

## A.2 phaseII.py

```

1 import sys
2 import time
3 import datetime
4 import pickle
5 import sched
6 import urllib.request
7 import urllib.parse
8 import xml.dom.minidom
9 import math
10 import statistics
11
12 DATA_FILE = "phaseII.pickle"
13
14 CAPMETRO_NEXTBUS = "http://www.capmetro.org/planner/s_nextbus2.asp"
15 CAPMETRO_BUSLOCS = "http://www.capmetro.org/planner/s_buslocation.asp"
16 POLLINTERVAL = 30
17 SAVEINTERVAL = 15 * 60
18
19 # Split trips at points that are this amount of time apart.
20 TRIP_TIMEDELTA_SPLIT = datetime.timedelta(hours = 12)
21 ACCEPTABLE_ERROR_MIN = 20
22 DEFAULT_SPEED_KMH = 40
23 ARRIVAL_ESTIMATE_INTERVAL = datetime.timedelta(minutes = 2)
24
25 class Station:
26     def __init__(self, stop_id, name, latitude, longitude):
27         self.stop_id = stop_id
28         self.name = name
29         self.trips = {}
30         self.latitude = latitude
31         self.longitude = longitude
32     def record_trip(self, timestamp, trip_id, vehicle_id, sched_arrival, est_arrival):
33         def fix_time(time):
34             time = time.strip()
35             # Error parsing estimated arrival time '12:02 XM': time data '12:02 XM'
36             # does not match format '%I:%M %p'
37             time = time.replace("X", "A")
38             # Error parsing estimated arrival time '00:29 AM': time data '00:29 AM'
39             # does not match format '%I:%M %p'
40             time = time.replace("00", "12")
41             return time
42         if not trip_id in self.trips:
43             self.trips[trip_id] = []
44         sched_arrival_time = None
45         sched_arrival = fix_time(sched_arrival)
46         try:
47             sched_arrival_time = datetime.datetime.strptime(sched_arrival,
48                 "%I:%M%p").time()
49         except Exception as e:
50             print("Error_parsing_scheduled_arrival_time_" + sched_arrival + "':" +
51                 str(e))
52         est_arrival_time = None
53         est_arrival = fix_time(est_arrival)
54         try:
55             est_arrival_time = datetime.datetime.strptime(est_arrival,
56                 "%I:%M%p").time()
57         except Exception as e:
58             print("Error_parsing_estimated_arrival_time_" + est_arrival + "':" +
59                 str(e))
60         self.trips[trip_id].append({
61             "timestamp": timestamp,
62             "sched_arrival_time": sched_arrival_time,
63             "est_arrival_time": est_arrival_time,
64             "vehicle_id": vehicle_id
65         })
66
67 stations = [
68     Station(497, "UT_West_Mall", 30.286064, -97.741815),

```

```

69     Station(5867, "Republic_Square", 30.267751, -97.746857),
70     Station(5606, "Crestview", 30.337852, -97.719035),
71     Station(5872, "Pleasant_Hill", 30.192391, -97.779203)
72 ]
73 vehicles = {}
74
75 # Credit John D. Cook, http://www.johndcook.com/python\_longitude\_latitude.html
76 # (public domain)
77 def distance_on_unit_sphere(lat1, long1, lat2, long2):
78
79     # Convert latitude and longitude to
80     # spherical coordinates in radians.
81     degrees_to_radians = math.pi/180.0
82
83     # phi = 90 - latitude
84     phi1 = (90.0 - lat1)*degrees_to_radians
85     phi2 = (90.0 - lat2)*degrees_to_radians
86
87     # theta = longitude
88     theta1 = long1*degrees_to_radians
89     theta2 = long2*degrees_to_radians
90
91     # Compute spherical distance from spherical coordinates.
92
93     # For two locations in spherical coordinates
94     # (1, theta, phi) and (1, theta, phi)
95     # cosine( arc length ) =
96     #   sin phi sin phi' cos(theta-theta') + cos phi cos phi'
97     # distance = rho * arc length
98
99     cos = (math.sin(phi1)*math.sin(phi2)*math.cos(theta1 - theta2) +
100            math.cos(phi1)*math.cos(phi2))
101     arc = math.acos( cos )
102
103     # Remember to multiply arc by the radius of the earth
104     # in your favorite set of units to get length.
105     return arc
106
107 def dom_value(node):
108     if node.firstChild == None:
109         return ""
110     else:
111         return node.firstChild.data
112
113 def combine_time_date(time, dt):
114     # Combines a time with date information from a datetime. We need some extra
115     # logic to handle cases straddling midnight.
116     midnight = datetime.time(0, 0)
117     if time > midnight and dt.time() < midnight:
118         dt_correct = datetime.timedelta(days = 1)
119     elif time < midnight and dt.time() > midnight:
120         dt_correct = datetime.timedelta(days = -1)
121     else:
122         dt_correct = datetime.timedelta()
123     return datetime.datetime.combine(dt.date(), time) + dt_correct
124
125 def load_data():
126     global stations
127     global vehicles
128     try:
129         prev_data = pickle.load(open(DATA_FILE, "rb"))
130         stations = prev_data["stations"]
131         vehicles = prev_data["vehicles"]
132     except Exception as e:
133         print("Error_loading_data_file:_" + str(e))
134
135 def save_data():
136     global stations

```

```

137     global vehicles
138     print("Saving_data...")
139     pickle.dump({ "stations": stations, "vehicles": vehicles },
140                 open(DATA_FILE, "wb"))
141
142 def retrieve_station_data(station):
143     url = CAPMETRO_NEXTBUS + "?" + urllib.parse.urlencode({ "route": "801",
144                                                             "stopid": str(station.stop_id) })
145     http_handle = None
146     try:
147         print("Downloading_" + url)
148         http_handle = urllib.request.urlopen(url)
149     except Exception as e:
150         print("Error_retrieving_data_for_station_" + str(station.stop_id) + ":" +
151               str(e))
152         return False
153     xml_tree = None
154     now = datetime.datetime.now()
155     try:
156         xml_tree = xml.dom.minidom.parse(http_handle)
157         runs = xml_tree.getElementsByTagName("Run")
158         for trip in runs:
159             vehicle_id = dom_value(trip.getElementsByTagName("Vehicleid")[0]).strip()
160             if dom_value(trip.getElementsByTagName("Valid")[0]) == "Y":
161                 station.record_trip(
162                     timestamp = now,
163                     trip_id = dom_value(trip.getElementsByTagName("Tripid")[0]).strip(),
164                     vehicle_id = vehicle_id,
165                     sched_arrival = dom_value(trip.getElementsByTagName("Triptime")[0]),
166                     est_arrival = dom_value(trip.getElementsByTagName("Estimatedtime")[0])
167                 )
168     except Exception as e:
169         print("Error_parsing_tracking_data_for_station_" + str(station.stop_id) +
170               ":" + str(e))
171         return False
172     return True
173
174 def retrieve_vehicle_data():
175     global vehicles
176     url = CAPMETRO_BUSLOCS + "?" + urllib.parse.urlencode({ "route": "801" })
177     try:
178         print("Downloading_" + url)
179         http_handle = urllib.request.urlopen(url)
180     except Exception as e:
181         print("Error_retrieving_bus_tracking_data:" + str(e))
182         return False
183     xml_tree = None
184     now = datetime.datetime.now()
185     try:
186         xml_tree = xml.dom.minidom.parse(http_handle)
187         xml_vehicles = xml_tree.getElementsByTagName("Vehicle")
188         for vehicle in xml_vehicles:
189             vehicle_id = dom_value(vehicle.getElementsByTagName("Vehicleid")[0]).strip()
190             latest_pos = [float(x) for x in dom_value(
191                 vehicle.getElementsByTagName("Position")[0]).split(",")]
192             if not vehicle_id in vehicles:
193                 vehicles[vehicle_id] = []
194             vehicle_data = vehicles[vehicle_id]
195             # Do not record duplicate data (data that hasn't been updated within our
196             # update interval).
197             if (len(vehicle_data) == 0 or
198                 vehicle_data[-1]["latitude"] != latest_pos[0] or
199                 vehicle_data[-1]["longitude"] != latest_pos[1]):
200                 vehicle_data.append({
201                     "timestamp": now,
202                     "block": dom_value(
203                         vehicle.getElementsByTagName("Block")[0]).strip(),
204                     "reliable": dom_value(

```

```

205         vehicle.getElementsByTagName("Reliable")[0]).strip() == "Y" ,
206     "off_route": dom_value(
207         vehicle.getElementsByTagName("Offroute")[0]).strip() == "Y" ,
208     "stopped": dom_value(
209         vehicle.getElementsByTagName("Stopped")[0]).strip() == "Y" ,
210     "in_service": dom_value(
211         vehicle.getElementsByTagName("Inservice")[0]).strip() == "Y" ,
212     "speed": float(dom_value(vehicle.getElementsByTagName("Speed")[0])),
213     "heading": int(dom_value(vehicle.getElementsByTagName("Heading")[0])),
214     "latitude": latest_pos[0],
215     "longitude": latest_pos[1]
216     })
217 except Exception as e:
218     print("Error parsing bus tracking data:_" + str(e))
219     return False
220 return True
221
222 def collect():
223     load_data()
224     scheduler = sched.scheduler(time.time, time.sleep)
225     def collect_stations(i):
226         retrieve_station_data(stations[i])
227         i = (i + 1) % len(stations)
228         now = datetime.datetime.now()
229         now_time = now.time()
230         # Cease collecting between 12:45 AM and 4:45 AM.
231         if now_time >= datetime.time(0, 45) and now_time <= datetime.time(4, 45):
232             scheduler.enterabs(time.mktime(datetime.datetime.combine(now,
233                 datetime.time(4, 45)).timetuple()), 2, collect_stations, (0,))
234         else:
235             scheduler.enter(POLLINTERVAL / len(stations), 2, collect_stations,
236                 (i,))
237     def collect_vehicles():
238         retrieve_vehicle_data()
239         now = datetime.datetime.now()
240         now_time = now.time()
241         # Cease collecting between 12:45 AM and 4:45 AM.
242         if now_time >= datetime.time(0, 45) and now_time <= datetime.time(4, 45):
243             scheduler.enterabs(time.mktime(datetime.datetime.combine(now,
244                 datetime.time(4, 45)).timetuple()), 1, collect_vehicles, ())
245         else:
246             scheduler.enter(POLLINTERVAL, 1, collect_vehicles, ())
247     def collect_save():
248         save_data()
249         scheduler.enter(SAVEINTERVAL, 3, collect_save, ())
250
251     collect_vehicles()
252     collect_stations(0)
253     collect_save()
254     try:
255         scheduler.run()
256     except KeyboardInterrupt:
257         print("Stopped.")
258         save_data()
259         sys.exit(0)
260
261 def split_unique_trips(trip):
262     # Trips are distinguished by trip ID, but this is only unique for one day.
263     # This function splits one "trip" into separate trips for each date.
264     trips = []
265     i = 0
266     this_ts = None
267     last_ts = None
268     last_split = 0
269     while i < len(trip):
270         this_ts = trip[i]["timestamp"]
271         if last_ts != None and this_ts - last_ts > TRIP_TIMEDELTA_SPLIT:
272             trips.append(trip[last_split:i])

```

```

273         last_split = i
274         last_ts = this_ts
275         i += 1
276     if trip[last_split:] != []:
277         trips.append(trip[last_split:])
278     return trips
279
280 def get_vehicle_position(vehicle_id, dt):
281     global vehicles
282     if vehicle_id in vehicles:
283         for point in vehicles[vehicle_id]:
284             if (abs(point["timestamp"] - dt) < datetime.timedelta(seconds = 30) or
285                 point["timestamp"] > dt):
286                 return point
287     return None
288 else:
289     return None
290
291 def estimate_arrival_time(vehicle_id, dt, scheduled_dt, station):
292     if scheduled_dt - dt <= datetime.timedelta(minutes = 30):
293         pos = get_vehicle_position(vehicle_id, dt)
294         if (pos != None and not pos["off_route"] and
295             pos["latitude"] != 0 and pos["longitude"] != 0):
296             if station.name == "Pleasant_Hill" or station.name == "Crestview":
297                 traffic_density = 15
298                 station_density = 2
299             elif station.name == "UT_West_Mall" or station.name == "Republic_Square":
300                 traffic_density = 30
301                 station_density = 3
302             dist_km = distance_on_unit_sphere(pos["latitude"], pos["longitude"],
303                 station.latitude, station.longitude) * 6373
304             delay_min = (0.4855 + 0.0287 * traffic_density / 8 + 0.0168 *
305                 traffic_density + 0.9654 * dist_km - 1.1969 * 0.5 + 0.1130 *
306                 dist_km * station_density)
307             return dt + datetime.timedelta(hours = dist_km / DEFAULT_SPEED_KMH +
308                 delay_min / 60)
309     else:
310         return scheduled_dt
311 else:
312     return scheduled_dt
313
314 def trip_error_data(trip, station):
315     error_data = {
316         "schedule": [],
317         "capmetro": [],
318         "me": []
319     }
320     arrival_time = trip[-1]["timestamp"]
321     this_ts = None
322     this_dt_correct = None
323     this_sched = None
324     this_est = None
325     this_my_est = None
326     this_my_est_last_dt = None
327     for point in trip:
328         this_ts = point["timestamp"]
329         vehicle_id = point["vehicle_id"]
330         this_est = combine_time_date(point["est_arrival_time"], this_ts)
331         this_sched = combine_time_date(point["sched_arrival_time"], this_ts)
332         if (this_my_est_last_dt == None or
333             this_ts - this_my_est_last_dt > ARRIVAL_ESTIMATE_INTERVAL):
334             this_my_est = estimate_arrival_time(vehicle_id, this_ts, this_sched,
335                 station)
336             this_my_est_last_dt = this_ts
337         # Calculate error and append to the lists.
338         schedule_error_min = (this_sched - arrival_time).total_seconds() / 60
339         capmetro_error_min = (this_est - arrival_time).total_seconds() / 60
340         me_error_min = (this_my_est - arrival_time).total_seconds() / 60

```

```

341         if abs(schedule_error_min) < ACCEPTABLE_ERROR_MIN:
342             error_data["schedule"].append(schedule_error_min)
343             error_data["capmetro"].append(capmetro_error_min)
344             error_data["me"].append(me_error_min)
345     return error_data
346
347 def compute_stability(data):
348     diffs = []
349     i = 1
350     while i < len(data):
351         diffs.append(abs(data[i] - data[i - 1]))
352         i += 1
353     return statistics.mean(diffs)
354
355 def analyze():
356     def append_keys(dict1, dict2):
357         for key, value in dict1.items():
358             dict2[key].append(value)
359     global stations
360     load_data()
361     metrics = {
362         "schedule": {
363             "_data": [],
364             "stability": []
365         },
366         "capmetro": {
367             "_data": [],
368             "stability": []
369         },
370         "me": {
371             "_data": [],
372             "stability": []
373         }
374     }
375     for station in stations:
376         for trip_id, trip_data in station.trips.items():
377             trips = split_unique_trips(trip_data)
378             for trip in trips:
379                 metrics["schedule"]["_data"] += this_error_data["schedule"]
380                 metrics["capmetro"]["_data"] += this_error_data["capmetro"]
381                 metrics["me"]["_data"] += this_error_data["me"]
382                 if len(this_error_data["schedule"]) > 1:
383                     metrics["schedule"]["stability"].append(
384                         compute_stability(this_error_data["schedule"]))
385                 metrics["capmetro"]["stability"].append(
386                     compute_stability(this_error_data["capmetro"]))
387                 metrics["me"]["stability"].append(
388                     compute_stability(this_error_data["me"]))
389
390     print("\nBUS_SCHEDULE_ONLY")
391     print("Error_mean:_ " + str(statistics.mean(metrics["schedule"]["_data"])))
392     print("Error_std_dev:_ " + str(statistics.stdev(metrics["schedule"]["_data"])))
393     print("Stability:_ " + str(statistics.mean(metrics["schedule"]["stability"])))
394     print("\nCAPITAL_METRO_ESTIMATES")
395     print("Error_mean:_ " + str(statistics.mean(metrics["capmetro"]["_data"])))
396     print("Error_std_dev:_ " + str(statistics.stdev(metrics["capmetro"]["_data"])))
397     print("Stability:_ " + str(statistics.mean(metrics["capmetro"]["stability"])))
398     print("\nMY_ESTIMATES")
399     print("Error_mean:_ " + str(statistics.mean(metrics["me"]["_data"])))
400     print("Error_std_dev:_ " + str(statistics.stdev(metrics["me"]["_data"])))
401     print("Stability:_ " + str(statistics.mean(metrics["me"]["stability"])))
402
403 if __name__ == "__main__":
404     if sys.argv[1] == "collect":
405         collect()
406     elif sys.argv[1] == "analyze":
407         analyze()

```

## A.3 phaseIII.py

```
1 import sys
2 import time
3 import datetime
4 import pickle
5 import sched
6 import urllib.request
7 import urllib.parse
8 import xml.dom.minidom
9 import math
10 import statistics
11
12 DATA_FILE = "phaseIII.pickle"
13
14 CAPMETRO_NEXTBUS = "http://www.capmetro.org/planner/s_nextbus2.asp"
15 CAPMETRO_BUSLOCS = "http://www.capmetro.org/planner/s_buslocation.asp"
16 POLLINTERVAL = 30
17 SAVEINTERVAL = 15 * 60
18
19 # Split trips at points that are this amount of time apart.
20 TRIP_TIMEDELTA_SPLIT = datetime.timedelta(hours = 12)
21 ACCEPTABLE_ERROR_MIN = 20
22 DEFAULT_SPEED_KMH = 40
23 ARRIVAL_ESTIMATE_INTERVAL = datetime.timedelta(minutes = 2)
24
25 class Station:
26     def __init__(self, stop_id, name, latitude, longitude):
27         self.stop_id = stop_id
28         self.name = name
29         self.trips = {}
30         self.latitude = latitude
31         self.longitude = longitude
32         self._my_est_last = None
33         self._my_est_last_dt = None
34     def record_trip(self, timestamp, trip_id, vehicle_id, sched_arrival,
35                   est_arrival, my_est_arrival):
36         if not trip_id in self.trips:
37             self.trips[trip_id] = []
38         data = {
39             "timestamp": timestamp,
40             "sched_arrival": sched_arrival,
41             "est_arrival": est_arrival,
42             "vehicle_id": vehicle_id
43         }
44         if (self._my_est_last == None or
45             timestamp - self._my_est_last_dt > ARRIVAL_ESTIMATE_INTERVAL):
46             data["my_est_arrival"] = my_est_arrival
47             self._my_est_last = my_est_arrival
48             self._my_est_last_dt = timestamp
49         else:
50             data["my_est_arrival"] = self._my_est_last
51         self.trips[trip_id].append(data)
52
53 stations = [
54     Station(497, "UT_West_Mall", 30.286064, -97.741815),
55     Station(5867, "Republic_Square", 30.267751, -97.746857),
56     Station(5606, "Crestview", 30.337852, -97.719035),
57     Station(5872, "Pleasant_Hill", 30.192391, -97.779203)
58 ]
59 vehicles = {}
60
61 # Credit John D. Cook, http://www.johndcook.com/python_longitude_latitude.html
62 # (public domain)
63 def distance_on_unit_sphere(lat1, long1, lat2, long2):
64
65     # Convert latitude and longitude to
```



```

66     # spherical coordinates in radians.
67     degrees_to_radians = math.pi/180.0
68
69     # phi = 90 - latitude
70     phi1 = (90.0 - lat1)*degrees_to_radians
71     phi2 = (90.0 - lat2)*degrees_to_radians
72
73     # theta = longitude
74     theta1 = long1*degrees_to_radians
75     theta2 = long2*degrees_to_radians
76
77     # Compute spherical distance from spherical coordinates.
78
79     # For two locations in spherical coordinates
80     # (1, theta, phi) and (1, theta, phi)
81     # cosine( arc length ) =
82     #   sin phi sin phi' cos(theta-theta') + cos phi cos phi'
83     # distance = rho * arc length
84
85     cos = (math.sin(phi1)*math.sin(phi2)*math.cos(theta1 - theta2) +
86           math.cos(phi1)*math.cos(phi2))
87     arc = math.acos( cos )
88
89     # Remember to multiply arc by the radius of the earth
90     # in your favorite set of units to get length.
91     return arc
92
93 def dom_value(node):
94     if node.firstChild == None:
95         return ""
96     else:
97         return node.firstChild.data
98
99 def load_data():
100     global stations
101     global vehicles
102     try:
103         prev_data = pickle.load(open(DATA_FILE, "rb"))
104         stations = prev_data["stations"]
105         vehicles = prev_data["vehicles"]
106     except Exception as e:
107         print("Error_loading_data_file:_" + str(e))
108
109 def save_data():
110     global stations
111     global vehicles
112     print("Saving_data...")
113     pickle.dump({ "stations": stations, "vehicles": vehicles },
114               open(DATA_FILE, "wb"))
115
116 def parse_time(time_s, now_dt):
117     time_s = time_s.strip()
118     # Error parsing estimated arrival time '12:02 XM': time data '12:02 XM'
119     # does not match format '%I:%M %p'
120     time_s = time_s.replace("X", "A")
121     # Error parsing estimated arrival time '00:29 AM': time data '00:29 AM'
122     # does not match format '%I:%M %p'
123     time_s = time_s.replace("00", "12")
124     try:
125         time = datetime.datetime.strptime(time_s, "%I:%M%p").time()
126     except Exception as e:
127         print("Error_parsing_time_" + time_s + "':" + str(e))
128         return None
129     # We need some extra logic to handle cases straddling midnight.
130     if (time > datetime.time(0, 0) and
131         now_dt.time() < datetime.time(0, 0)):
132         dt_correct = datetime.timedelta(days = 1)
133     elif (time < datetime.time(0, 0) and

```

```

134         now_dt.time() > datetime.time(0, 0)):
135         dt_correct = datetime.timedelta(days = -1)
136     else:
137         dt_correct = datetime.timedelta()
138     return datetime.datetime.combine(now_dt.date(), time) + dt_correct
139
140 def get_vehicle_position(vehicle_id, dt):
141     global vehicles
142     if vehicle_id in vehicles:
143         for point in vehicles[vehicle_id]:
144             if point["timestamp"] > dt - datetime.timedelta(seconds = 30):
145                 return point
146         return None
147     else:
148         return None
149
150 def estimate_arrival_time(vehicle_id, dt, scheduled_dt, station):
151     if scheduled_dt - dt <= datetime.timedelta(minutes = 30):
152         pos = get_vehicle_position(vehicle_id, dt)
153         if (pos != None and not pos["off_route"] and
154             pos["latitude"] != 0 and pos["longitude"] != 0):
155             if station.name == "Pleasant_Hill" or station.name == "Crestview":
156                 traffic_density = 15
157                 station_density = 2
158             elif station.name == "UT_West_Mall" or station.name == "Republic_Square":
159                 traffic_density = 30
160                 station_density = 3
161             dist_km = distance_on_unit_sphere(pos["latitude"], pos["longitude"],
162                 station.latitude, station.longitude) * 6373
163             delay_min = (0.4855 + 0.0287 * traffic_density / 8 + 0.0168 *
164                 traffic_density + 0.9654 * dist_km - 1.1969 * 0.5 + 0.1130 *
165                 dist_km * station_density)
166             return dt + datetime.timedelta(hours = dist_km / DEFAULT.SPEED_KMH +
167                 delay_min / 60)
168         else:
169             return scheduled_dt
170     else:
171         return scheduled_dt
172
173 def retrieve_station_data(station):
174     url = CAPMETRO.NEXTIBUS + "?" + urllib.parse.urlencode({ "route": "801",
175         "stopid": str(station.stop_id) })
176     http_handle = None
177     try:
178         print("Downloading_" + url)
179         http_handle = urllib.request.urlopen(url)
180     except Exception as e:
181         print("Error_retrieving_data_for_station_" + str(station.stop_id) + ":" +
182             str(e))
183     return False
184     xml_tree = None
185     now = datetime.datetime.now()
186     try:
187         xml_tree = xml.dom.minidom.parse(http_handle)
188         runs = xml_tree.getElementsByTagName("Run")
189         for trip in runs:
190             vehicle_id = dom_value(trip.getElementsByTagName("Vehicleid")[0]).strip()
191             if dom_value(trip.getElementsByTagName("Valid")[0]) == "Y":
192                 station.record_trip(
193                     timestamp = now,
194                     trip_id = dom_value(trip.getElementsByTagName("Tripid")[0]).strip(),
195                     vehicle_id = vehicle_id,
196                     sched_arrival = parse_time(dom_value(
197                         trip.getElementsByTagName("Triptime")[0]), now),
198                     est_arrival = parse_time(dom_value(
199                         trip.getElementsByTagName("Estimatedtime")[0]), now),
200                     my_est_arrival = estimate_arrival_time(vehicle_id, now, parse_time(
201                         dom_value(trip.getElementsByTagName("Triptime")[0]), now), station)

```

```

202         )
203     except Exception as e:
204         print("Error_parsing_tracking_data_for_station_" + str(station.stop_id) +
205               ":",_ + str(e))
206         return False
207     return True
208
209 def retrieve_vehicle_data():
210     global vehicles
211     url = CAPMETRO.BUSLOCS + "?" + urllib.parse.urlencode({ "route": "801" })
212     try:
213         print("Downloading_" + url)
214         http_handle = urllib.request.urlopen(url)
215     except Exception as e:
216         print("Error_retrieving_bus_tracking_data:_ + str(e))
217         return False
218     xml_tree = None
219     now = datetime.datetime.now()
220     try:
221         xml_tree = xml.dom.minidom.parse(http_handle)
222         xml_vehicles = xml_tree.getElementsByTagName("Vehicle")
223         for vehicle in xml_vehicles:
224             vehicle_id = dom_value(vehicle.getElementsByTagName("Vehicleid")[0]).strip()
225             latest_pos = [float(x) for x in dom_value(
226                 vehicle.getElementsByTagName("Position")[0]).split(",")]
227             if not vehicle_id in vehicles:
228                 vehicles[vehicle_id] = []
229             vehicle_data = vehicles[vehicle_id]
230             # Do not record duplicate data (data that hasn't been updated within our
231             # update interval).
232             if (len(vehicle_data) == 0 or
233                 vehicle_data[-1]["latitude"] != latest_pos[0] or
234                 vehicle_data[-1]["longitude"] != latest_pos[1]):
235                 vehicle_data.append({
236                     "timestamp": now,
237                     "block": dom_value(
238                         vehicle.getElementsByTagName("Block")[0]).strip(),
239                     "reliable": dom_value(
240                         vehicle.getElementsByTagName("Reliable")[0]).strip() == "Y",
241                     "off_route": dom_value(
242                         vehicle.getElementsByTagName("Offroute")[0]).strip() == "Y",
243                     "stopped": dom_value(
244                         vehicle.getElementsByTagName("Stopped")[0]).strip() == "Y",
245                     "in_service": dom_value(
246                         vehicle.getElementsByTagName("Inservice")[0]).strip() == "Y",
247                     "speed": float(dom_value(vehicle.getElementsByTagName("Speed")[0])),
248                     "heading": int(dom_value(vehicle.getElementsByTagName("Heading")[0])),
249                     "latitude": latest_pos[0],
250                     "longitude": latest_pos[1]
251                 })
252     except Exception as e:
253         print("Error_parsing_bus_tracking_data:_ + str(e))
254         return False
255     return True
256
257 def collect():
258     load_data()
259     scheduler = sched.scheduler(time.time, time.sleep)
260     def collect_stations(i):
261         retrieve_station_data(stations[i])
262         i = (i + 1) % len(stations)
263         now = datetime.datetime.now()
264         now_time = now.time()
265         # Cease collecting between 12:45 AM and 4:45 AM.
266         if now_time >= datetime.time(0, 45) and now_time <= datetime.time(4, 45):
267             scheduler.enterabs(time.mktime(datetime.datetime.combine(now,
268                 datetime.time(4, 45)).timetuple()), 2, collect_stations, (0,))
269     else:

```

```

270         scheduler.enter(POLLINTERVAL / len(stations), 2, collect_stations ,
271                          (i,))
272     def collect_vehicles():
273         retrieve_vehicle_data()
274         now = datetime.datetime.now()
275         now_time = now.time()
276         # Cease collecting between 12:45 AM and 4:45 AM.
277         if now_time >= datetime.time(0, 45) and now_time <= datetime.time(4, 45):
278             scheduler.enterabs(time.mktime(datetime.datetime.combine(now,
279                                     datetime.time(4, 45)).timetuple()), 1, collect_vehicles, ())
280         else:
281             scheduler.enter(POLLINTERVAL, 1, collect_vehicles, ())
282     def collect_save():
283         save_data()
284         scheduler.enter(SAVEINTERVAL, 3, collect_save, ())
285
286     collect_vehicles()
287     collect_stations(0)
288     collect_save()
289     try:
290         scheduler.run()
291     except KeyboardInterrupt:
292         print("Stopped.")
293         save_data()
294         sys.exit(0)
295
296 def split_unique_trips(trip):
297     # Trips are distinguished by trip ID, but this is only unique for one day.
298     # This function splits one "trip" into separate trips for each date.
299     trips = []
300     i = 0
301     this_ts = None
302     last_ts = None
303     last_split = 0
304     while i < len(trip):
305         this_ts = trip[i]["timestamp"]
306         if last_ts != None and this_ts - last_ts > TRIP_TIMEDELTA_SPLIT:
307             trips.append(trip[last_split:i])
308             last_split = i
309             last_ts = this_ts
310             i += 1
311     if trip[last_split:] != []:
312         trips.append(trip[last_split:])
313     return trips
314
315 def get_station_arrival_time(search_start_dt, vehicle_id, station):
316     global stations
317     for this_station in stations:
318         for trip_id, trip in this_station.trips.items():
319             if len(trip) > 0 and trip[0]["vehicle_id"] == vehicle_id:
320                 for point in trip:
321                     ts = point["timestamp"]
322                     vehicle_pos = get_vehicle_position(vehicle_id, ts)
323                     if (ts > search_start_dt and
324                         ts - search_start_dt < TRIP_TIMEDELTA_SPLIT and
325                         vehicle_pos != None):
326                         dist_m = distance_on_unit_sphere(station.latitude, station.longitude,
327                                                         vehicle_pos["latitude"], vehicle_pos["longitude"]) * 6373000
328                         if dist_m < ARRIVAL_DISTANCE_M:
329                             return ts
330     return None
331
332 def trip_error_data(trip, station):
333     error_data = {
334         "schedule": [],
335         "capmetro": [],
336         "me": []
337     }

```

```

338 arrival_time = trip[-1]["timestamp"]
339 for point in trip:
340     # Calculate error and append to the lists.
341     schedule_error_min = (point["sched_arrival"] - arrival_time).total_seconds() / 60
342     capmetro_error_min = (point["est_arrival"] - arrival_time).total_seconds() / 60
343     my_error_min = (point["my_est_arrival"] - arrival_time).total_seconds() / 60
344     if (abs(schedule_error_min) < ACCEPTABLE_ERROR_MIN):
345         error_data["schedule"].append(schedule_error_min)
346         error_data["capmetro"].append(capmetro_error_min)
347         error_data["me"].append(my_error_min)
348 return error_data
349
350 def compute_stability(data):
351     diffs = []
352     i = 1
353     while i < len(data):
354         diffs.append(abs(data[i] - data[i - 1]))
355         i += 1
356     return statistics.mean(diffs)
357
358 def analyze():
359     def append_keys(dict1, dict2):
360         for key, value in dict1.items():
361             dict2[key].append(value)
362     global stations
363     load_data()
364     metrics = {
365         "schedule": {
366             "_data": [],
367             "stability": []
368         },
369         "capmetro": {
370             "_data": [],
371             "stability": []
372         },
373         "me": {
374             "_data": [],
375             "stability": []
376         }
377     }
378     for station in stations:
379         for trip_id, trip_data in station.trips.items():
380             trips = split_unique_trips(trip_data)
381             for trip in trips:
382                 this_error_data = trip_error_data(trip, station)
383                 metrics["schedule"]["_data"] += this_error_data["schedule"]
384                 metrics["capmetro"]["_data"] += this_error_data["capmetro"]
385                 metrics["me"]["_data"] += this_error_data["me"]
386                 if len(this_error_data["schedule"]) > 1:
387                     metrics["schedule"]["stability"].append(
388                         compute_stability(this_error_data["schedule"]))
389                     metrics["capmetro"]["stability"].append(
390                         compute_stability(this_error_data["capmetro"]))
391                     metrics["me"]["stability"].append(
392                         compute_stability(this_error_data["me"]))
393
394     print("\nBUS_SCHEDULE_ONLY")
395     print("Error_mean:_" + str(statistics.mean(metrics["schedule"]["_data"])))
396     print("Error_std_dev:_" + str(statistics.stdev(metrics["schedule"]["_data"])))
397     print("Stability:_" + str(statistics.mean(metrics["schedule"]["stability"])))
398     print("\nCAPITAL_METRO_ESTIMATES")
399     print("Error_mean:_" + str(statistics.mean(metrics["capmetro"]["_data"])))
400     print("Error_std_dev:_" + str(statistics.stdev(metrics["capmetro"]["_data"])))
401     print("Stability:_" + str(statistics.mean(metrics["capmetro"]["stability"])))
402     print("\nMY_ESTIMATES")
403     print("Error_mean:_" + str(statistics.mean(metrics["me"]["_data"])))
404     print("Error_std_dev:_" + str(statistics.stdev(metrics["me"]["_data"])))
405     print("Stability:_" + str(statistics.mean(metrics["me"]["stability"])))

```

```
406
407 if __name__ == "__main__":
408     if sys.argv[1] == "collect":
409         collect()
410     elif sys.argv[1] == "analyze":
411         analyze()
```