

1 TRB Paper Manuscript #15-0309  
2 **Develop New York City Surface Transit Boarding and Alighting Ridership Daily Production Application**  
3 **Using Big Data**

4  
5 Qifeng Zeng, Alla Reddy\*, Alex Lu, and Brian Levine

6 \* *Corresponding author*

7  
8  
9 Qifeng Zeng

10 Computer Specialist

11 New York City Transit Authority

12 2 Broadway, Cubicle A17.092, New York, N.Y. 10004-2208

13 Tel: (646) 252-2400

14 Email: [Qifeng.Zeng2@nyct.com](mailto:Qifeng.Zeng2@nyct.com)

15  
16 Alla Reddy

17 Senior Director, System Data & Research (SDR), Operations Planning

18 New York City Transit Authority

19 2 Broadway, Office A17.92, New York, N.Y. 10004-2208

20 Tel: (646) 252-5662

21 Email: [Alla.Reddy@nyct.com](mailto:Alla.Reddy@nyct.com)

22  
23 Alex Lu

24 Transit Analyst

25 P.O. Box 684, Ossining, N.Y. 10562-0684

26 Tel: (212) 340-2684

27 Email: [me@lexciestuff.net](mailto:me@lexciestuff.net)

28  
29 Brian Levine

30 Principal Transportation Planner, System Data & Research (SDR), Operations Planning

31 New York City Transit Authority

32 2 Broadway, Office A17.100, New York, N.Y. 10004-2208

33 Tel: (646) 252-5541

34 Email: [Brian.Levine@nyct.com](mailto:Brian.Levine@nyct.com)

35  
36  
37 Submitted for Consideration for Publication in

38 *Transportation Research Records: Journal of the Transportation Research Board*

39  
40  
41 **Word Count:** 274 (Abstract) + 6,273 (Text) + 4 \* 250 (Figures) = 7,547 Words.

42 **Submittal Date:** November 15, 2014.

43

44

45

46

**1 ABSTRACT**

2 New York City Transit Authority (NYCTA) has put into production a ridership application that  
3 determines surface transit boarding and alighting locations for each one of approximately 2.8 million  
4 daily passenger trips on board 218 bus routes, to support bus service scheduling and planning. The  
5 application combines data from an Automated Vehicle Location (AVL) system, a multimodal entry-only  
6 non-geographic Automated Fare Collection (AFC) system, and General Transit Feed Specification (GTFS)  
7 schedule and shape file streams. To accomplish this, NYCTA developed a highly optimized network  
8 generation tool to estimate bus boarding/alighting locations and link loads, by creating a scaled-down  
9 custom network based on first-bus trajectory (from AVL boarding location data) and a few possible  
10 AFC/AVL-inferred second-leg pickup stops. Solving for the shortest walking path on this sub-network  
11 yields connection points and disembarkation points much more efficiently than solving for all 128  
12 million potential O/D pairs citywide. The optimized program executes in less than three hours in an  
13 automated environment and is responsive to daily detours, special events, and weather-driven  
14 ridership, which can support operations management’s needs for next-day reporting as well as the  
15 monitoring of fast-changing patterns and trends in ridership behavior. It also allows multiple days of  
16 route-level program output to be aggregated for schedule making purposes, providing a significantly  
17 more representative understanding of typical passenger loads that were historically obtained from a  
18 few labor-intensive on-board observations (ridechecks) collected over a multi-year period. Results were  
19 validated and found to be consistent with manual ridechecks, limited O/D surveys, and earlier  
20 estimation methodologies. Obtained accuracy was sufficient to achieve acceptance of AVL-AFC data in  
21 lieu of traditional on-board observations for NYCTA schedule making and vehicle assignment purposes.

22

23

## 1 INTRODUCTION

2 This paper presents New York City Transit Authority’s (NYCTA’s) production algorithm to determine  
3 surface transit boarding and alighting locations for each one of approximately 2.8 million daily passenger  
4 trips on 218 bus routes by combining data from:

- 5
- 6 • an Automated Vehicle Location (AVL) system with explicitly limited scope in providing analytics.
- 7 • a legacy regional multimodal entry-only Automated Fare Collection (AFC) system, designed over  
8 twenty years ago, that covers six distinct bus operators, seven bus rapid transit lines, two  
9 subway systems, an inter-state rail transit line, a tram, and an airport shuttle.
- 10 • a General Transit Feed Specification (GTFS) schedule-and-shape file stream generated by  
11 commercial-off-the-shelf (COTS) crew and vehicle scheduling system.
- 12

13 The goal of this project is to supply detailed ridership reports to operational management on daily basis.  
14 The reports, based on AVL and AFC data, will be compiled at the trip and stop level, with roll-up  
15 summaries. This will provide invaluable operational and analytical tool that can be used to assess the  
16 effects of planned service changes (General Orders, special events) and unplanned disruptions (weather,  
17 traffic, etc.) and will permit the monitoring of fast changing patterns and trends in ridership behavior.  
18 AFC data (over 9 million total daily transactions by over 3 million unique cards from an estimated 40  
19 million outstanding farecards in active circulation) and AVL data (over 4.8 million time-and-location  
20 reports from 4,800 vehicles) are available as early morning bulk download for the previous 24-hour  
21 period. In order to process AFC/AVL data by the next day, the program must execute totally unattended  
22 within a maximum running time of five hours.

23

24 To account for frequent short-term route changes to surface transit network—publicized in Authority’s  
25 GTFS feeds—the algorithm uses these changes as input and assigns passenger accordingly. It must also  
26 accept and make “best guess” decisions without human intervention for common error conditions, such  
27 as GTFS data not matching actual operations; wrong vehicle head signs; corrupt or missing AFC data; and  
28 busted or non-reporting Global Positioning System (GPS) devices.

29

30 To meet these goals and requirements, NYCTA took the following approach:

- 31 • For the most complex case, two consecutive bus trips, AFC data provides boarding timestamps  
32 to the nearest six minutes, allowing AVL to geocode to the nearest 0.5 miles (about 3-5 possible  
33 stops). A highly optimized network generation tool then produces a limited custom network  
34 based on first bus trajectory (from AVL boarding location) and the few possible second-bus pick-  
35 up stops; solving for shortest walking path on this small network thence yields transfer point  
36 and consequently disembarkation point from first bus.
- 37 • For the trips from bus to subway, determining the bus disembarkation stop requires solving for  
38 the shortest walk between route-specific possible alighting stops and actual subway entrance  
39 logged by AFC, which provides definitive geographic locations. Similarly, a small network is  
40 generated then solved for each transfer.
- 41 • For the trips from subway to bus, a daily lookup table is generated by finding the closest subway  
42 station to each bus stop within given distance. Among 3-5 possible bus boarding stops implied  
43 by AVL data, the one closest to subway is selected. An earlier methodology that solves  
44 schedule-based shortest path to determine exit stations was scrapped due to excessive  
45 computing time.
- 46
- 47

## 1 **Relationship to Prior Work**

2 Development of AFC origin-destination (O/D) modelling capability in New York City dates back to late  
3 20<sup>th</sup> Century when Barry et al. (2) published an inference method assigning the destination of one trip  
4 based on next-trip origin. This approach had been known at NYCTA since 1997 when AFC system first  
5 began reporting card serial number data. Variation of this approach plus various extensions utilizing  
6 schedule data (5) was subsequently popularized and applied in Chicago (3), Boston (4), Beijing (5),  
7 London (6), Hong Kong (7) and no doubt other rail transit systems.

8  
9 A decade later, Freimer et al. (8) reported that NYCTA had extended this approach into multimodal  
10 region-wide O/D models,—with bus boarding and disembarkation localization provided by schedule-  
11 based dead-reckoning creating a linked transit trip database from over 7 million daily AFC swipes.  
12 Combining New York City’s famous traffic gridlock and lack of functioning AVL, bus O/D estimations and  
13 consequent vehicle-level passenger load determinations were insufficiently accurate for scheduling  
14 frequency determination and vehicle assignment. Few bus routes were validated against empirical data,  
15 and datasets required manual processing and “massaging” before becoming usable for O/D estimations.  
16 Only a few days’ data from 2004 was processed, and model quickly became outdated. Stop localization  
17 was so problematic that NYCTA elected to exploit travel pattern symmetries in subsequent work (20)  
18 instead of guessing where bus boarding occurred using AFC data alone.

19  
20 Meanwhile, as smartcard-based AFC, AVL, and even Automated Passenger Counters (APC) became  
21 popular elsewhere, algorithms for processing and combining AFC-AVL-APC data streams and generating  
22 O/D matrices thrived in both proprietary (30) and public space (40), with research being carried out in  
23 locales including: Minneapolis (9,1); Columbus, Ohio (16); Chicago (39); Montreal (12); London (11,12a);  
24 Taipei, Taiwan (10); Jinan, China (41); Shenzhen, China (17); Santiago (14); and São Paulo (15). A recent  
25 TCRP study reviewed and catalogued many common analytical models utilized within the industry to  
26 process archived data (19).

27  
28 Despite the popularity of AFC-AVL-APC methodological research, few references exist in the literature  
29 for automatically generating surface transfer points. Relevant research focuses on grouping bus stops  
30 together using geographic proximity, text similarity, and land-use pattern to automatically provide  
31 transfer information (1). Like many studies, databases were used for implementation, and the need for  
32 grouping highlights the challenges of creating and maintaining enormous lookup tables.

33  
34 While O/D analysis of passively collected transit data is a crowded field, our work extends an operations-  
35 support, production-oriented philosophy of zero manual intervention (ZMI) generation of next-day  
36 “flash” reports—first accomplished in New York with AFC-based passenger-miles computation (20),  
37 followed by track-circuit based performance monitoring (21)—to much more complex tasks of  
38 calibrating and generating full O/D matrices. Results were convincing and comparable to field  
39 observations at sufficiently high granularity and fidelity at individual vehicle load, passenger origins and  
40 destinations, peak-load points, and time-period levels—to be accepted by NYCTA scheduling dept. for  
41 use over labor-intensive onboard on/off counts (colloquially, ‘ridecheck’). While newer or smaller transit  
42 properties could run SQL queries from composite smartcard-APC-GPS data, NYCTA accomplished this  
43 with legacy entrance-only AFC, simple low-cost AVL, and varying bus routes.

44  
45

## 1 **Summary of Innovations**

2 This algorithm takes a unique approach and makes a major contribution as it amounts to an automatic,  
3 error-tolerant, and analyst-free daily re-calibration of a classic shortest-path passenger routing model  
4 utilizing daily GTFS file as base network and combined AFC-AVL data as constraint to model each  
5 individual farecard itinerary for every trip every day. It also demonstrates problem-specific optimization.

6 Key innovations are:

- 7 • Large scale production implementation
- 8 • Zero manual intervention
- 9 • Each actual bus-bus and bus-subway consecutive passenger trips observed in AFC data (about  
10 2.5 million daily transfers, 1.0 million unique customers) is solved discretely, rather than solving  
11 transfers between all 218 routes at over 16,000 stops system-wide.
- 12 • Dynamic transfer links between two bus routes are generated only from feasible disembarkation  
13 points for that passenger on that first bus to the few known possible second bus boarding  
14 locations (and not all bus routes in the vicinity).
- 15 • Tight on-the-fly integration of previously unconnected AVL, AFC, and GTFS data sources.

16

## 17 **BIG DATA**

18 Input data is taken from three major sources: MetroCard (NYCTA farecard) transaction data, MTA  
19 BusTime AVL data, and Authority GTFS feed. MetroCard's AFC data format was well-documented in  
20 prior Authority papers (2,8,20,22,23). This discussion will focus on two new sources of data used by this  
21 algorithm.

22

### 23 **Daily Schedule GTFS Feed**

24 GTFS format, introduced in 2007, contains a wealth of information describing transit schedules, and  
25 routes in Geographic Information Systems (GIS) format. The Authority started publishing bus and  
26 subway GTFS data generated by industry-standard scheduling software on its website for third party  
27 developers in January 2010 (31). Bus scheduling processes are described in detail by Holmes (24). MTA  
28 GTFS feeds contain sufficient detail for systematically and accurately constructing traditional computer  
29 transit networks consisting of nodes and links elements and their properties such as running time and  
30 wait time.

31

32 Most transit systems publish system-wide schedule changes several times a year (colloquially, a 'Pick').  
33 However, NYCTA routes sometimes undergo major route/schedule changes between Picks, which are  
34 very difficult for developers to keep track on daily basis. MTA GTFS files are constantly updated to  
35 include between Pick reroutes—providing impetus for building a robust application capturing everyday  
36 operational activities.

37

### 38 **MTA BusTime AVL Data**

39 Under MTA's BusTime project, NYCTA installed experimental AVL devices on B63 route in Brooklyn  
40 beginning April 2011. Almost all of over 4,800 buses running >300 NYCTA and MTA Bus local and express  
41 routes were equipped with AVL devices by March 2014. BusTime uses GPS hardware and wireless  
42 communications to track bus locations in real-time.

43

44 Due to legacy hardware and proprietary software issues onboard NYCTA buses, AVL computers actually  
45 "sniff" serial communication links between fareboxes and headsign units to determine basic trip

1 information like operator pass number (employee ID), route code, direction, and destination. GPS  
2 location is reported to central servers approximately every 30 seconds. Server software integrates this  
3 information with map, route, schedule data, previous real-time updates, headsign information, and real-  
4 time operational data like operators' timekeeping system—linking vehicle fleet numbers with a  
5 particular trip within scheduled runs. Sophisticated inferential algorithms then provide a “best guess” in  
6 real-time of (a) whether bus is in-or-out of service, (b) what route variant—‘path’ in TA parlance—is  
7 being served, (c) route direction, and (d) if off-route—i.e. detour in effect. Based on inferred  
8 information it predicts what stops bus will make and distance from each stop. BusTime server software  
9 utilizes the open source OneBusAway package (27); output data is publicly accessible to developers as a  
10 feed (28).

11  
12 However, AVL archive data requires significant post-processing and scrubbing to make usable. Because  
13 BusTime ‘path’ is a real-time “best guess”, it can be inconsistent from stop to stop. Only after trip is  
14 complete can path be conclusively determined— path is not traced out until bus reaches its terminus.  
15 Another issue is dealing with buses with head signs set to “Not In Service” or “Next Bus Please” but yet  
16 remain in service. Also, about 10% of GPS devices are not functioning on an average basis. This process  
17 and associated AVL issues are described in detail in Levine et al. (25) and Glikin et al. (26).

## 18 19 **METHODOLOGY**

### 20 ***Dynamic Transfer Link Network***

21  
22 Most bus origin-destination estimation studies (12a, 38,39) use pre-generated bus-bus transfer lookup-  
23 table to infer alighting bus stop(s) or area. Transfer lookup-tables consist of permissible transfers  
24 systemwide. Each study delicately defines policies and crafts methodologies to create lookup tables.  
25 It's the most important factor influencing quality of results.

26  
27 Bus-bus lookup approaches are not appropriate for New York City surface transit daily production  
28 application. The method generally involves great deal of human judgment and tedious time-consuming  
29 manual operation to generate transfer tables. Even if systematic approaches is found for table  
30 generation, plenty of resulting exceptions still requires manual editing. Therefore, it is most suitable for  
31 research studies where dimensions and complexities of transport networks can be contained. New York  
32 City transit network has over 16,000 bus stops and 426 subway stations with various transfer exceptions  
33 that can change frequently, making it impossible, due to size and complexities, to create and maintain a  
34 production oriented lookup table.

35  
36 Fundamental network computing method are explored for this application. Shortest path methods are  
37 used for deriving bus alighting locations. Network computing methods usually prepare static transit  
38 networks for shortest-path (35a). Surface transit networks include bus routes links, and walking links  
39 connecting bus routes via bus stops. Walking links are conceptually same as transfer tables except for  
40 each link having different impedances. Underlying structure of walking links is street network; the street  
41 network, in theory, can connect any two bus stops. Even with moderate surface networks and strict  
42 rules in governing walking link creation, size of walking links permuted by bus stops could easily grow  
43 out of control.

44  
45 In addition, classic Dijkstra's shortest-path algorithm has order of growth proportional to vertices  
46 squared, i.e.  $O(|V|^2)$  (29). The entire New York City surface network couldn't be solved in reasonable

1 time (Figure 1(a)) with 128 million unique potential O/D pairs on 24 hours bus network with classical  
2 method of transit network generation and shortest path.

3  
4 A highly optimized network generation tool is developed, restricting shortest-path problem sizes when  
5 solving for a specific farecard's exact trajectory given fuzzy information about boarding locations. A  
6 novel approach, called "Dynamic Transfer Link" method, is developed for this study. Instead of creating  
7 walking links statically, the method takes advantage of known next bus trip's fuzzy boarding locations  
8 (due to six-minute farebox timing error margin) or subway stations, and creates walking links  
9 dynamically during shortest path computing process. Only possible walking links between two bus  
10 routes or bus route and subway system meeting predefined criteria (usually "distance") are generated.  
11 Number of walking links created dynamically during shortest path process is dramatically reduced to a  
12 fraction of static network size. Onerous transfer exceptions constraints are automatically and implicitly  
13 considered by exploring only paths traced by an actual farecard. It takes merely 1½ hours to compute  
14 over 2 million daily bus passengers using "Dynamic Transfer Link" method with 1.9 miles as maximum  
15 walking distance allowed.

16  
17 The illustrated example of "Dynamic Transfer Link" is shown in Figure 1(b). Four consecutive MetroCard  
18 transactions for a subway-bus-bus-subway itinerary (Station 0, Route 1, Route 2, and Station 4):

- 19
- 20 • From Station 0 to Boarding Stop 2, a transfer link is created by lookups based on the
  - 21 geographically closest subway station to bus stops on Route 1.
  - 22 • From Routes 1 to 2, dynamic transfer links are created by linking all bus stops downstream of
  - 23 Boarding Stop 2 to all possible Boarding Stops 3, so long as transfer links are less than 1.9
  - 24 miles—all possible walking transfers.
  - 25 • From Route 2 to Station 4, stops downstream of Stop 3 are linked with subway stations to create
  - 26 the bus-subway transfers.

27  
28 With dynamic transfer links, a partial transit network specific to this farecard itinerary is completed and  
29 fed into Dijkstra's algorithm for path finding using shortest path assumption. Output results include  
30 boarding and alighting stops on each individual bus trip.

### 31 32 **Implementation Flow Chart**

33  
34 This application takes five general major steps (Figure 1(c)):

- 35
- 36 1. Generate subway-bus lookup table and bus route network from GTFS.
  - 37 2. Determine the arrival time at all stops for all trips from BusTime data.
  - 38 3. Attach each AFC transaction to bus stops groups.
  - 39 4. Build and solve custom shortest-path network for each bus-bus AFC transaction sequence.  
40 (Individual passenger travel path is the basic output)
  - 41 5. Results are factored to account for O/D indeterminate AFC transactions, cash passengers, and  
42 fare evaders.
  - 43 6. For longer-period summaries, data from multiple days are blended.

44  
45 Numerous proof-of-concept methods to generate O/D matrices have been tested both by NYCTA and  
46 outside research groups (16) using familiar synthetic methodologies such as iterative proportional fitting  
47 (IPF) that can generate a full O/D matrix from marginal on/off's obtained from APC or AFC, but to date no  
48 other groups have reported an automated data merge of entry-only AFC and AVL data by treating each

1 farecard discretely, thereby preserving as much of the actual O/D information as possible without  
2 analyst manual intervention and model calibration. Complete O/D matrix is generated from all 2.8  
3 million daily swipes.

4

#### 5 **Step 1: Construct Network**

6 Computer transit networks have links (bus, subway, walk for transfer, O/D access and egress) and nodes  
7 (subway stations, bus stops, transfer points, etc.). Transit routes are usually static and can be  
8 automatically generated from GTFS. However, coding of transfer, access, and egress links are more  
9 challenging. Transfers between subway stations are usually straightforward, as physical connections  
10 exist between stations without incurring additional fare. Bus-bus transfers occur on the street, which  
11 presents numerous ways of connecting. Developers usually have to decide which transfers are practical  
12 feasible links. Generally, certain criteria are set for feasible connections, normally based on threshold  
13 distance, and links are generated systematically. With this approach, no transfer links are created at this  
14 stage.

15

16 Subway to bus transfer matrices are built from GTFS for that specific day based on geographic proximity.  
17 GIS analysis was conducted using the shortest distance between geocodes of subway stations and bus  
18 stops. Separate pre-processing generates a system wide bus network that describes nodes and links  
19 from GTFS schedule data and the distance between nodes using GTFS bus shape files and GIS linear  
20 referencing feature.

21

#### 22 **Step 2: Deriving Bus Stop Arrival Time**

23 GIS-based bus route tracking was developed to identify and validate route variant for each individual  
24 bus trip. Each trip is compared to system wide bus network and GTFS route shape. After correct path is  
25 assigned, entire bus trip is re-processed to establish arrival times at each stop. If no valid path was  
26 found, the trip is removed.

27

#### 28 **Step 3: Merging AFC and AVL Data**

29 Prior modelling efforts (8) were done without AVL data and had tremendous challenges in determining  
30 bus location. AFC data provides the farebox number when passengers board the bus, but not exact  
31 geographical location due to the constant movement of the bus. The AFC transaction time is linked to  
32 bus stop arrival time, developed in Step 2 to determine boarding stops. Since MetroCard is entry-only,  
33 only boarding points can be determined at this stage. The limitation of bus farebox memory when  
34 MetroCard was introduced constrained timestamp resolution to one-tenth of an hour (six minutes),  
35 during which bus might have covered 3-5 stops. Thus, a number of possible boarding stops are selected.

36

#### 37 **Step 4: Generating Dynamic Bus Transfer Links and Running Shortest Path**

38 With numerous known boarding points for individual journeys, this step applies the “Dynamic Transfer  
39 Links” and Dijkstra’s shortest-path algorithm method developed for the application to determine bus  
40 boarding and alighting points as well as walking links, and consequently obtaining the paths of the  
41 journey. As heart of the application, the “Dynamic Transfer Links” method is detailed in “Dynamic  
42 Transfer Link” section appearing earlier.

43

44



**1 Step 5: Non-Farecard Passenger Trip Estimation**

2 Almost all the transit system has non-farecard passengers. In New York City transit, about 9% (May  
3 2014) of bus passengers still use cash fare. Cash to AFC passenger ratio by bus route and direction are  
4 calculated and corresponding O/D patterns are used for allocation. Fare evasion ratio (20) at bus route  
5 level is used to further adjust passenger trips. Using the same approach, individual fare evader trips are  
6 synthesized based on AFC O/D patterns. Figure 2(a) shows NYCT borough-level adjustment factor  
7 summary.  
8

**9 Step 6: Blending/Trending Data**

10 For schedule production purposes, medium-term ridership trends are more helpful than daily results.  
11 When combining data from multiple days, as bus assignments differ each day, schedule run number and  
12 scheduled trip start time are used as keys to merge multiple days' data in a database. Since AVL  
13 malfunctions tend to be randomly distributed throughout the fleet, multiple-day averages provide  
14 reliable data on all trips specifically even if some trips ran with busted GPS on some dates. Medium-  
15 term averages also smooth out effects of non-recurring noise like street congestion, special events,  
16 weather conditions, or simply non-routine customer travel.  
17

**18 NEW YORK CITY SURFACE TRANSIT APPLICATION**

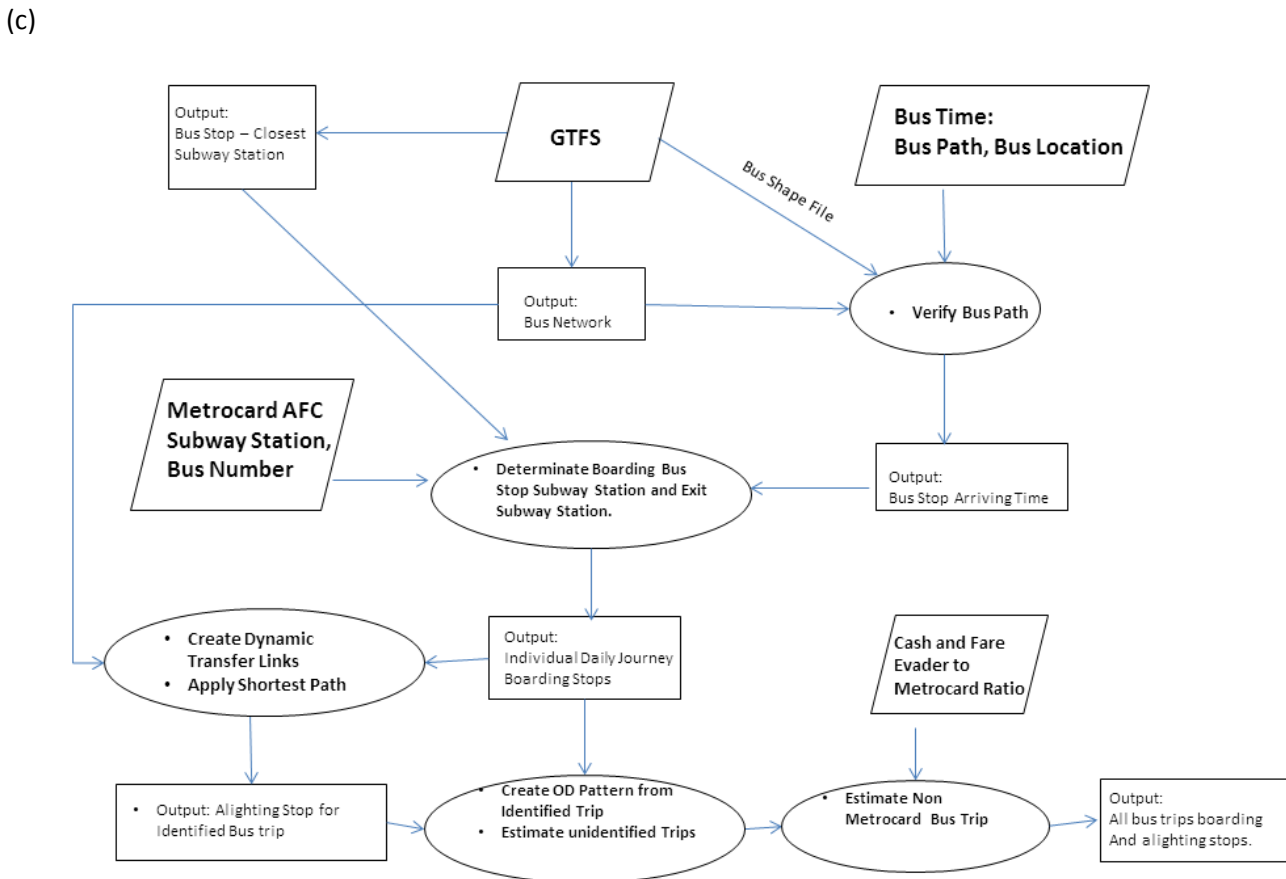
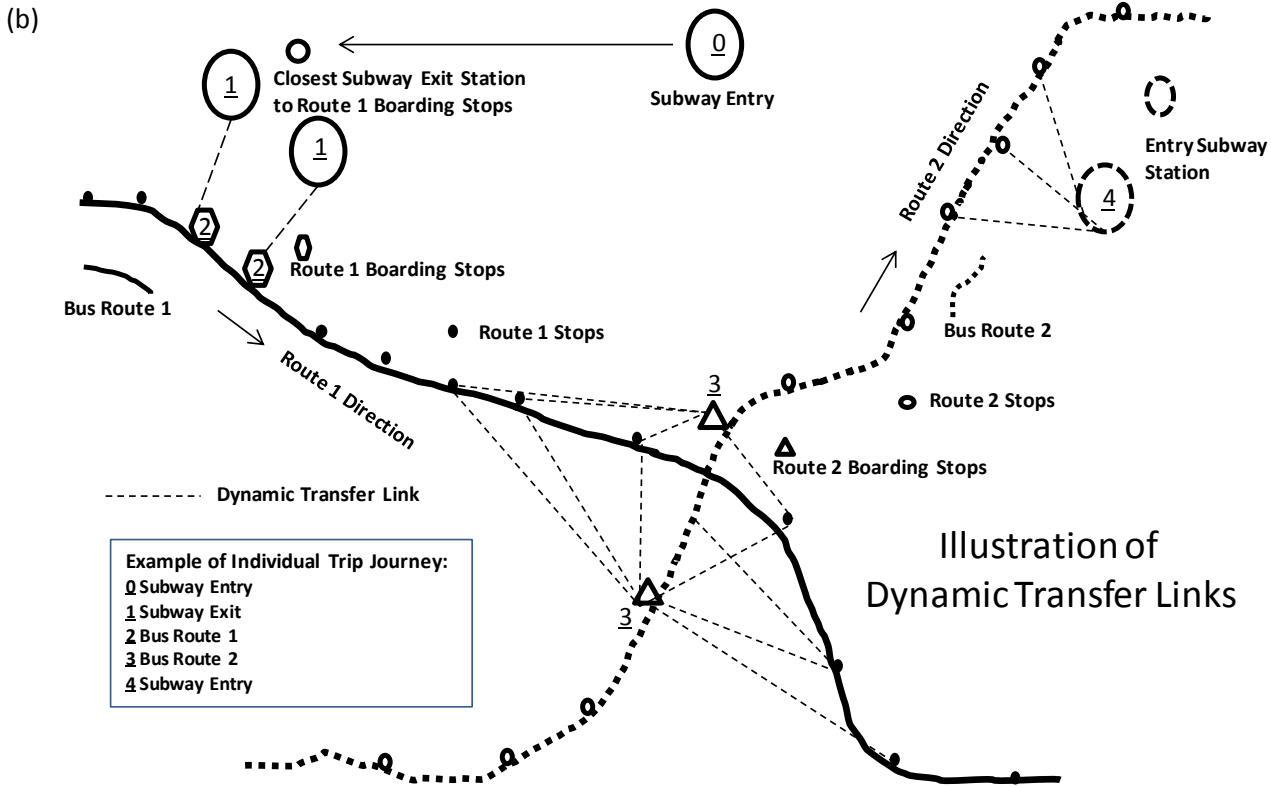
19 Several key assumptions were needed to customize general methodology described to New York City  
20 surface transit daily production application. Most are typical reassurances of observations and policies  
21 widespread in AFC analysis practice. Length of maximum transfer is a key parameter affecting results  
22 and outcome. Sensitivity analysis was conducted to determine appropriate transfer length in New York  
23 City's environment.  
24

**25 One-to-One Passenger-Farecard Relationship**

26 Individual passengers are identified by their MetroCard, assuming that a passenger uses one card  
27 throughout the day. However, this assumption is incorrect in some scenarios. Most common and  
28 troublesome scenario is when passengers purchase and begin to use new MetroCards midday.  
29 Connection between old and new MetroCards then cannot be established. In March 2013, Authority  
30 Tariff was modified to levy \$1.00 surcharge for each new farecard issued—reducing of plastic waste,  
31 decreasing required MetroCard Vending Machine (MVM) servicing frequencies (to re-fill farecards), and  
32 reducing trash in stations. This fare policy was effective in encouraging customers to add value or time  
33 to existing MetroCards and thus significantly improving passenger-farecard assumption's validity. By  
34 May 2014, 11% of total in-system sales transactions were new cards, compared to 39% before March  
35 2013.  
36

1

O/Ds solved	All O/Ds systemwide	AFC observed O/Ds only	AFC observed O/Ds only	AFC observed O/Ds only	AFC observed O/Ds only
<b>Network Generation Approach</b>	<b>Classical lookup</b>	<b>Route-based stop masking</b>	<b>Distance-based stop masking</b>	<b>AVL-based origin masking</b>	<b>Directional dest. masking</b>
<i>Potential Origins</i>	~16,000	~90	~32	3-5	3-5
<i>Potential Destinations</i>	~16,000	~90	~32	~90	~23
<i>O/D permutations (vertices) solved</i>	$\sim 16,000^2 / 2 = \sim 128$ million	$\sim 90^2 / 2 = \sim 4,050$	$\sim 32^2 / 2 = \sim 512$	$4 * 90 / 2 = \sim 180$	$4 * 23 / 2 = \sim 46$
<i>Number of sub-problems solved</i>	1 (to populate lookup table)	<2.8 million (worst case)	<2.8 million (worst case)	<2.8 million (worst case)	<2.8 million (worst case)
<i>Total operations</i>	$(128m)^2 = 16,384 * 10^{12}$	$4,050^2 * 2.8m = 45.9 * 10^{12}$	$512^2 * 2.8m = 734$ bn	$180^2 * 2.8m = 90.7$ bn	$46^2 * 2.8m = 5.9$ bn
<i>Est. execution time</i>	317 years	324 days	5.2 days	15 hours	1 hour



2 **Figure 1** Algorithm Methodology Detail: (a) Theoretical Permutation Estimation for Various Algorithmic  
 3 Approaches; (b) Illustrative Example of Dynamic Transfer Link Generation; (c) Final Algorithm Processing  
 4 Flowchart.

### 1 **Passenger May Travel Maximum of 1.9 Miles While Connecting**

2 To determine connection validity between consecutive waypoints, maximum thresholds for  
 3 time/distance are required. In New York City, people are more likely to use longer “walk link” for mid-  
 4 distance travel than other cities. Besides walking, customers have rich selections of transportation  
 5 alternatives including bicycle, taxi, skateboarding, black car (gypsy cab), third-party bus shuttle, ethnic  
 6 jitney, horse-drawn carriage, pedicab, commuter rail, and ferry. With high density commercial and  
 7 residential development throughout the City, higher probabilities of having one or more activities in  
 8 between consecutive trips (i.e. trip-chaining) occurs. With timespans between consecutive system  
 9 entries varying from 6 minutes to 24 hours, connections between consecutive trips (i.e. between two  
 10 waypoints), may not necessarily be short walks made solely for transfer purposes, as is typical for other  
 11 cities.

12  
 13 Rather than setting thresholds from theoretical walking distance tolerance, sensitivity analysis of  
 14 farecard trajectory observations was carried out to determine threshold from percentage of passenger  
 15 trips identifiable. If distance between inferred alighting and known boarding stops is greater than  
 16 maximum transfer length, it is assumed that an error condition has occurred and passenger’s alighting  
 17 stop must be determined without reference to next trip.

18  
 19 Figure 2(b) shows percent of total passenger trips identified as function of maximum transfer threshold  
 20 from ¼ to 3 miles for selected May 2014 weekdays. At one-quarter mile, 70% of trips are identified. The  
 21 percentage increases rapidly to 79% as transfer threshold increases to 10,000ft (1.9 miles). However,  
 22 from 1.9 to 3 miles, additional trips identified rapidly approach point of diminishing return. Thus, 1.9  
 23 miles, in the upper ranges of reasonable walking distances in New York, is selected as transfer threshold.

24

### 25 **Every MetroCard Returns to the Same Location**

26 This assumption is necessary when using 24-hour blocks of MetroCard data. It’s assumed last trip of the  
 27 day terminates where first trip originated (2). Given nighttime activity dynamics, Barry argued 3am is  
 28 when days start in New York City. In this study, we used midnight as start-of-day for data processing  
 29 convenience with most daily input files. The 3am schedule was originally proposed for subway-only  
 30 trips, where 1.3% of daily ridership occurs between midnight and 3am—much higher than equivalent  
 31 0.6% figure for buses. Furthermore, a recent study (37) shows that subway O-D estimating improve only  
 32 ~1% if starting time moves from midnight to 3am in New York City.

33

### 34 **Dealing with Non-Contiguous Metrocard Itineraries**

35 Consecutive recorded locations of system entries (‘swipes’) do not make logical sense in some  
 36 situations. Quite common in AFC data are consecutive swipes on two bus routes miles apart—perhaps  
 37 in different boroughs. Such transfers are either rare or geographically impossible. Possible explanations  
 38 for this occurrence are:

39

- 40 1. Most frequently, passenger used another transport mode between two consecutive swipes.
- 41 2. Passenger may have used different MetroCards mid-journey. When passengers change cards,  
 42 inferred destinations of both old and new cards are incorrect—resulting in irrational  
 43 connections between swipes.
- 44 3. More than one Pay-per-Ride passengers can share one card, resulting in non-sequential system  
 45 entry transactions and causing non-contiguous itineraries.

- 1 4. Loss of transaction data occurs due to farebox failures; sequence of swipes for MetroCards  
2 having lost transactions is affected.
- 3 5. AVL/GPS devices may not function as expected, resulting in bus trips failing to provide passenger  
4 boarding locations. When geographic fixes are not available, travel paths for customers having  
5 unknown boarding locations cannot be built.

6  
7 In general, when connection between two consecutive swipes is irrational, passenger is treated as  
8 having paid cash or evaded the fare. If boarding stop can be determined by AVL, boarding stop is  
9 assigned; otherwise both boarding and alighting stops are factored from O-D patterns generated from  
10 other identifiable trips.

### 11 ***MetroCard and Non-AFC Passengers Have Similar Trip Patterns***

12  
13 Since inception in 1997, MetroCard market share has steadily increased. MetroCard accounted for over  
14 95% of NYCTA subway/bus usage in May 2014. Average weekday bus ridership in May (excluding  
15 evaders) is 2.3 million, yet only 8.5% (190,000 passengers) paid cash or used paper transfers. Non-  
16 MetroCard boardings like cash, single ride tickets (essentially tokens), senior return tickets and others  
17 (32) do not generate transaction information. While non-MetroCard passengers might have different  
18 O/D patterns, since MetroCard accounts for majority of travel (and generated O/D matches land-use  
19 patterns well), it's reasonable to assume differences are *deminimis*. Non-MetroCard passengers are  
20 thus assumed to have similar trip pattern as AFC passengers.

### 21 ***Special AFC Processing***

22  
23 NYCTA is a complicated system; one general modelling methodology could not provide the entire  
24 solution. Also, due to differing implementations of MetroCard AFC on transport carriers or modes post-  
25 dating early 1990s AFC design, some exception processing is required:

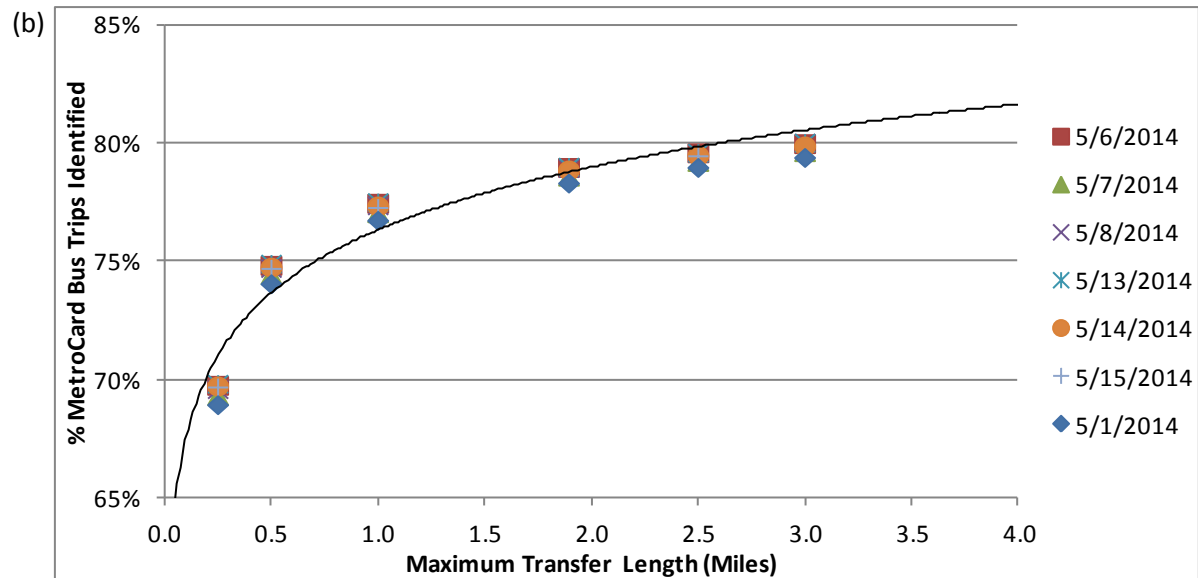
26  
27 **Bus Rapid Transit**—Since 2008, seven new rapid bus routes ([+selectbusservice](#), SBS) were  
28 implemented in New York. Except for S79 on Hylan Boulevard, Staten Island (33), all remaining six  
29 routes including Bx12 (36), M15 (35), M34, B44, Bx41 (34), and M60 utilize roadside proof-of-payment  
30 fare collection. Passengers use MetroCard to register at wayside MetroCard Fare Collectors (MFCs) and  
31 receive receipts retained for random on-board inspections. SBS physical buses no longer have  
32 passenger information, and thus special AFC data processing is required. Assuming SBS passenger  
33 always boards next arriving bus, onboard passenger loads can be determined. Since roadside farebox  
34 locations are static, specific boarding stops are determined. Dynamic transfer links are created similar  
35 to the regular bus transfers.

36  
37 **Connections with Regional Transit Operators**—Numerous non-NYCTA transit operators in the region  
38 also accept MetroCard as form of payment. AFC data gathered using their terminals also require special  
39 processing, detailed in Figure 2(c).

40

(a)

Category	Manhattan	Bronx	Brooklyn	Queens	Staten Is.	Systemwide
<b>AFC Passengers</b>	85.85%	77.38%	82.35%	87.42%	75.37%	82.26%
<b>Non-AFC Passengers</b>						
Cash	3.82%	5.86%	5.78%	5.37%	7.61%	5.40%
Partial Fare	1.90%	3.32%	2.30%	1.92%	4.01%	2.50%
Front Door	2.12%	5.22%	3.17%	1.76%	6.67%	3.42%
Rear Door	0.32%	1.46%	0.44%	0.14%	0.29%	0.63%
Child Over 44"	1.61%	2.27%	1.81%	0.87%	1.87%	1.75%
Child Under 44"	3.06%	3.23%	2.71%	1.32%	2.38%	2.70%
Flash Pass/Wheelchair	1.01%	0.98%	1.15%	0.89%	1.50%	1.05%
Broken Farebox	0.30%	0.30%	0.30%	0.30%	0.30%	0.30%
<b>Total</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>



(c)

Operator	Required Special Processing
Staten Island Railway (SIR)	Staten Island Railway (SIR) is a separate rapid transit line that doesn't connect to NYC subways. It has 22 stations, two of which (St. George and Tompkinsville, near ferry terminal) have fare control facilities. Fares are charged for both entries and exits, similar to downtown terminal ticket barriers in typical commuter rail practice. SIR entry AFC transaction is no different from subway entries. Since most boardings from suburban stations have destinations of either St. George or Tompkinsville, exit transactions are considered as boardings at any other non-fare control stations. For exit transactions, dynamic transfer links from buses include all possible transfers from any downstream bus stops within 1.9 miles of suburban SIR stations.
Staten Island Ferry (SIF)	Staten Island Ferry is a free shuttle operated by NYC Dept. of Transportation between South Ferry, Manhattan and St. George, Staten Island (SI). Named after such famous SI politicians as Marchi and Molinari, SIF serves mainly Staten Island residents' commuting needs to Downtown Manhattan. SIF route is five miles long. Many SI residents make daily intermodal trips of subway-ferry-bus or bus-ferry-SIR daily. Although physical distance between South Ferry and St George is >1.9 miles, transfers via ferry are definitely legitimate. Therefore, a pseudo St. George bus stop is created with South Ferry geocodes. When dynamic transfer links are created from St. George to Downtown transit services, pseudo locations are used, representing SIF.
Port Authority Trans-Hudson (PATH) Train	PATH is a rail transit service between New York and New Jersey, with six stations in Manhattan. MetroCard can be used to pay PATH fares. Some individuals' itineraries include travel on PATH. To model their activities, the six PATH stations in Manhattan are replaced with closest TA subway stations. All N.J. PATH stations are replaced with either World Trade Center or Herald Square subway stations, which are likely Manhattan destinations of PATH passengers.
Nassau Inter County Express (NICE) Bus	NICE (Nassau Inter County Express) bus mainly serves Nassau County, to the east of Queens. NICE adopted MetroCard as primary fare media and have service closely coordinated with TA facilities and routes. Over 8,000 transfers were made between NICE and TA/MTA buses in October 2013. NICE bus makes GTFS and GIS data available, but does not provide similar AVL data. Based on NICE route provided in AFC data, we consider all stops on that route possible boarding stops. When transfers are made between TA/MTA and NICE bus, dynamic transfer links are generated between TA/MTA route's all possible alighting stops within 1.9 miles of NICE route's every possible entry stop.
Bee-Line (Westchester County Bus System)	Bee-Line is Westchester County bus, to the north of Bronx, with over 7,000 monthly transfers to/from TA/MTA buses. Bee-Line bus does not provide GTFS data; we simply use Bx16, which operates along the Bronx-Westchester county line, as approximate geographic representation of all Bee-line routes.

4  
5  
6  
7  
8

**Figure 2** Algorithm Application Detail based on New York City empirical data: (a) Non-AFC Passenger Adjustment Factors; (b) Maximum Transfer Length Sensitivity Analysis; (c) Special AFC Processing due to regional integration of farecard systems.

## 1 RESULTS VALIDATION

2 Three different approaches were used to rigorously validate model outputs:

3

- 4 1. Route-level average trip length
- 5 2. Bus passenger load at each stop:
  - 6 a. versus classic manual ridecheck, time-period level average
  - 7 b. versus experimental APC data, trip level
- 8 3. Aggregate origin-destination matrix

9

### 10 ***Unlinked Trip Length Validation***

11 Validation process starts with unlinked passenger trip length. This is a route characteristic that depends  
12 on the land-use and activity patterns at the stops—provided routing does not change, average  
13 passenger lengths-of-haul by route should remain relatively constant. Conventionally, this is considered  
14 a key model quality measurement. Since model output is individual bus passenger boarding and  
15 alighting stops, lengths-of-haul is easily determined.

16

17 NYCTA had developed an AFC-based prior model (20) that directly measures unlinked trip length but  
18 makes a constant-speed bus operations assumption. That model was accepted by Federal Transit  
19 Administration (FTA) for official passenger-miles computation. Figure 3(a) and 3(b) show comparisons  
20 between AFC-only and improved AFC-AVL models for local and express bus routes derived from March  
21 6, 2013 data.

22

23 Average trip length difference for local routes is <2.3%, with errors appearing randomly distributed and  
24 approximately equal number of overestimated routes compared to underestimated ones, largely  
25 validating new model for local bus. However, validation results show that AFC-AVL model estimates  
26 express bus trip lengths on average 34% higher than AFC-only model. To understand large differences in  
27 express bus results, recall that:

28

29 Constant average speed assumptions implicit in distance-time conversions are a little more  
30 problematic, as NYCTA has routes where average speeds vary considerably over the whole route.  
31 Express Bus average speeds might be 15 mph at the residential end, 30~45 mph in non-stop  
32 express portions, and <10 mph in downtown—all within one trip. [...] Using average speeds leads  
33 to biases in these cases. (20)

34

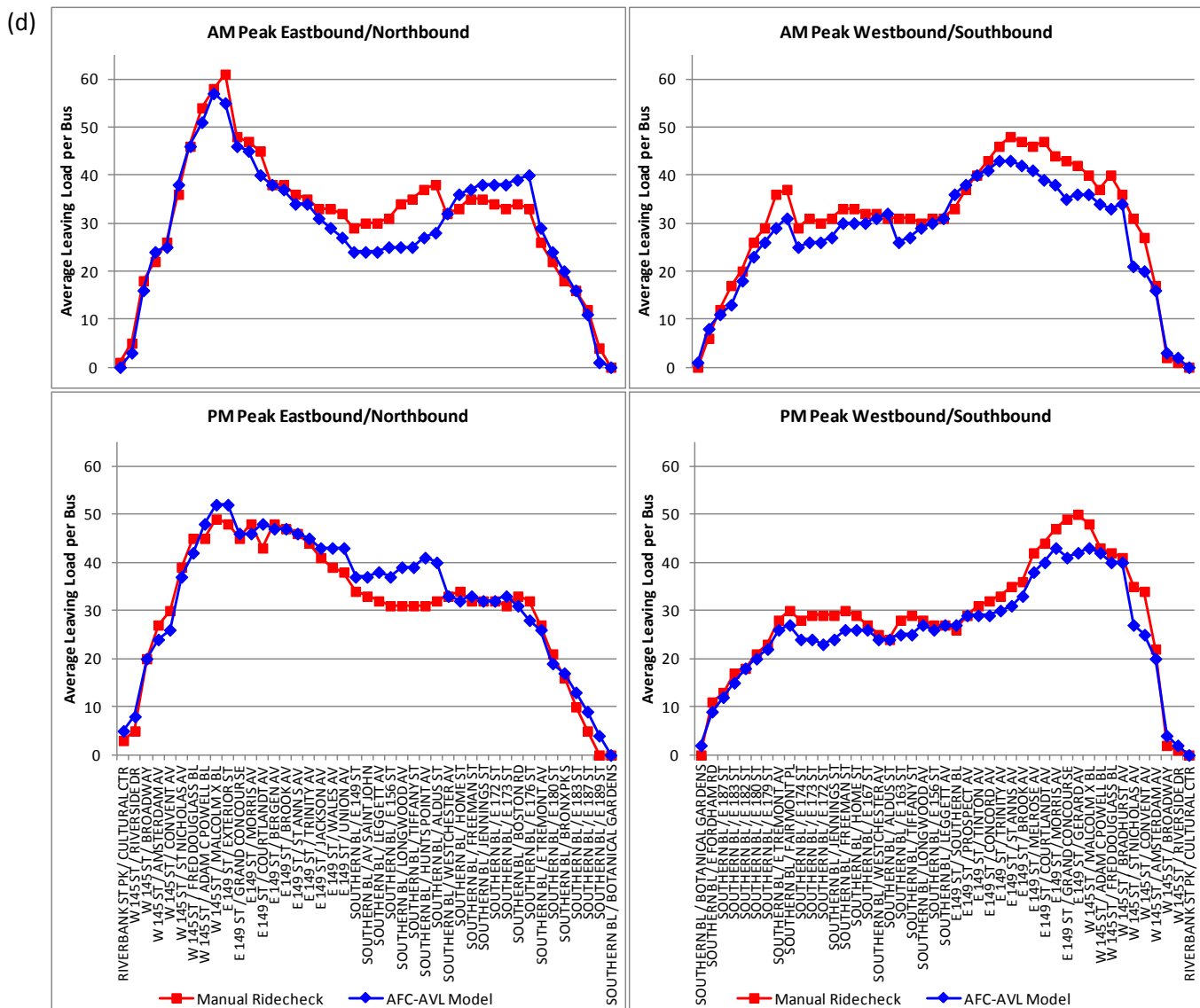
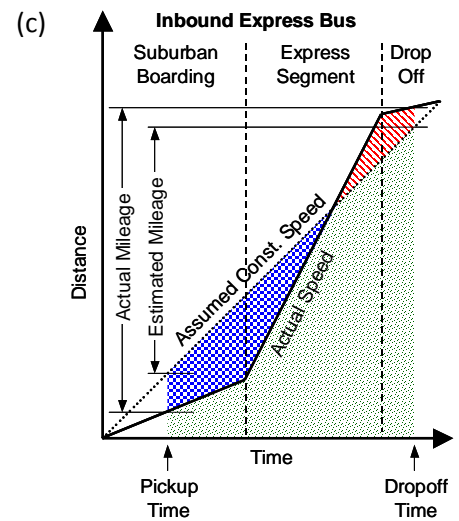
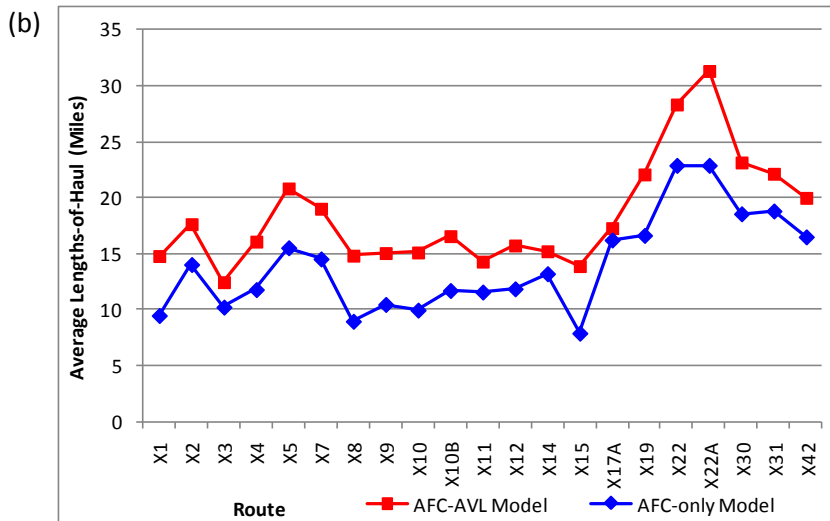
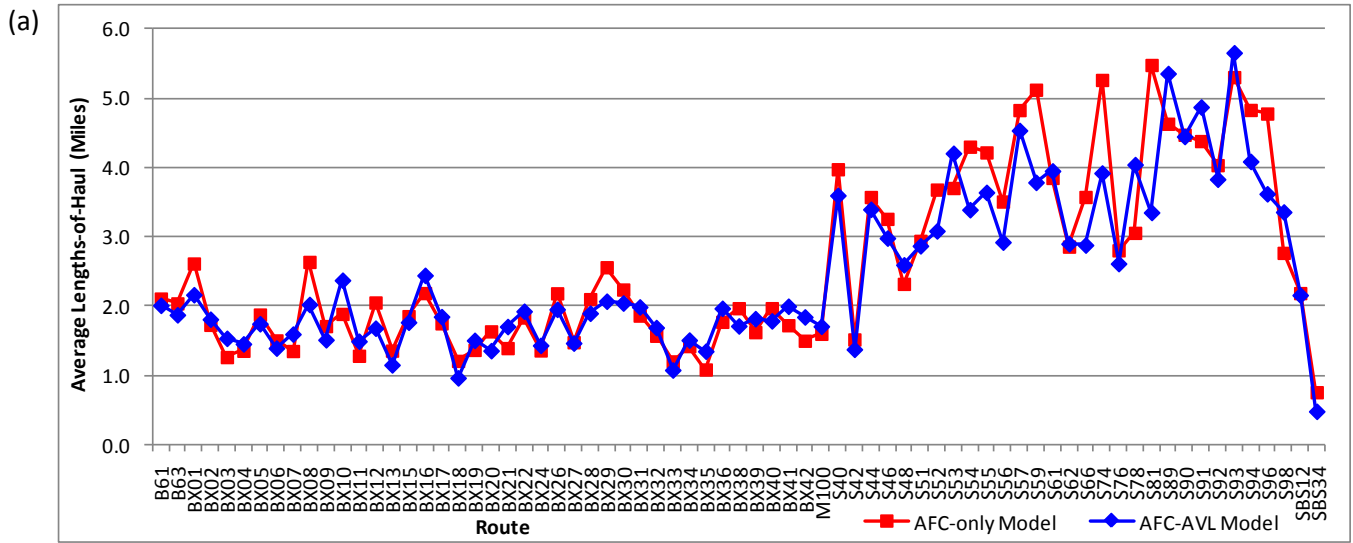
35 Figure 3(c) (from (20)) clearly shows the authors anticipated that for inbound express buses, passenger-  
36 mileage underestimation would occur for all boardings. As no ridechecks are performed on express  
37 buses, that assumption could not be validated then. The new AFC-AVL model supports that assertion.

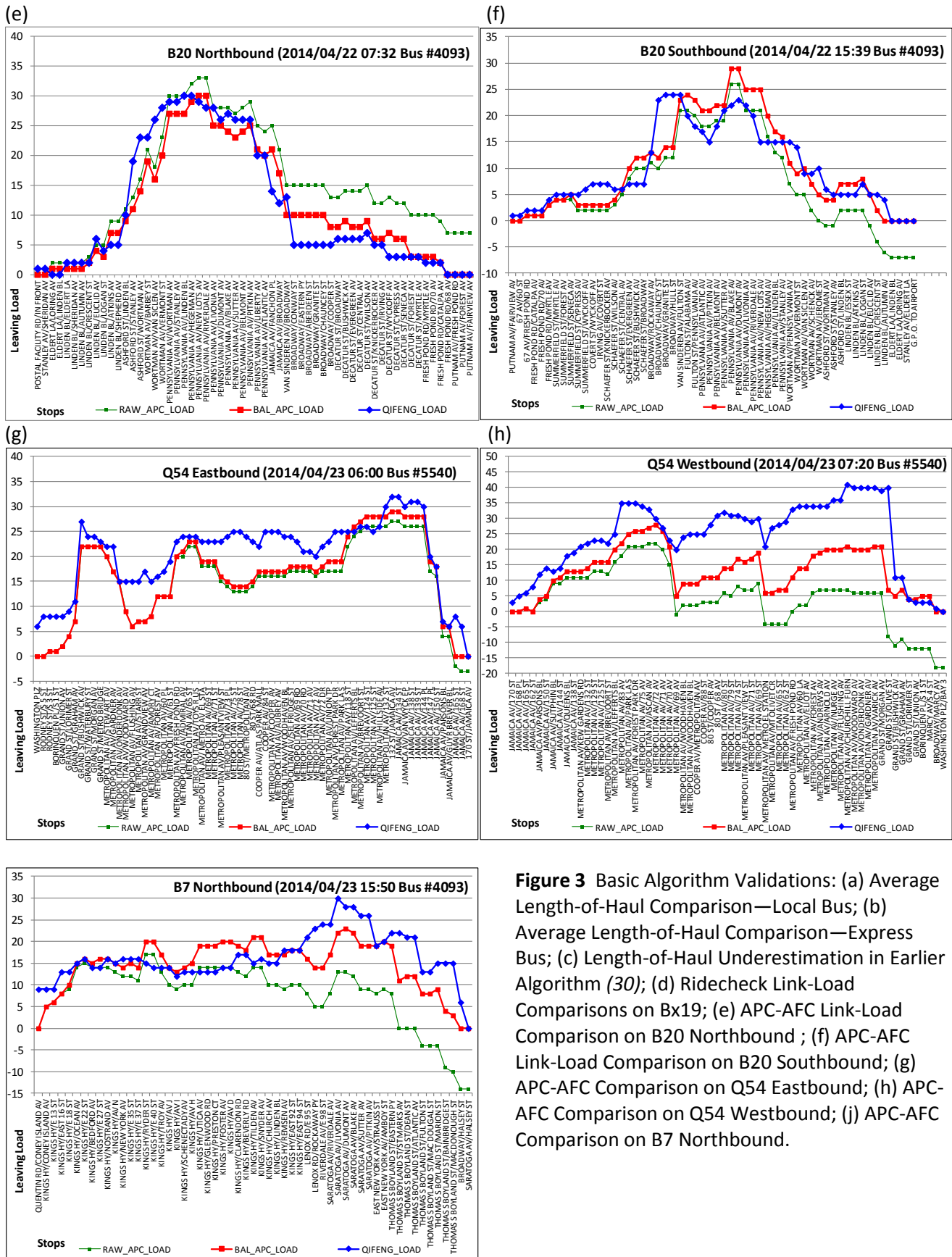
38

### 39 ***Ridecheck Validation—Ridership Profile***

40 Main validation effort is focused on comparison against ridecheck data. One goal of this study is to  
41 provide an all-electronic alternative for bus scheduling ridecheck. Comparing it with ridecheck,  
42 hithertofore considered standard, would essentially determine this model's practical worthiness.

43





**Figure 3** Basic Algorithm Validations: (a) Average Length-of-Haul Comparison—Local Bus; (b) Average Length-of-Haul Comparison—Express Bus; (c) Length-of-Haul Underestimation in Earlier Algorithm (30); (d) Ridecheck Link-Load Comparisons on Bx19; (e) APC-AFC Link-Load Comparison on B20 Northbound ; (f) APC-AFC Link-Load Comparison on B20 Southbound; (g) APC-AFC Comparison on Q54 Eastbound; (h) APC-AFC Comparison on Q54 Westbound; (j) APC-AFC Comparison on B7 Northbound.



1 Leaving load from each bus stop, averaged over time periods (i.e. AM Peak, Midday, etc.) is used for  
2 quality control comparison. Very few true parallel tests were possible because the Authority had  
3 already cut traffic checker forces in anticipation of system wide AVL deployment, even before AVL  
4 devices were reporting any data. However, Bx19 (Riverbank Park, Manhattan to Botanic Gardens, Bronx  
5 via 145 St. and Southern Blvd.) had ridechecks conducted on November 19, 2012 when most buses at  
6 West Farms Depot were reporting AVL data. A same-day comparison at time period- and stop-level  
7 (Figure 3(d)) allows for easy visual examination showing excellent load profile and absolute value  
8 correspondence between manually-collected ridecheck results and automatically-generated AFC-AVL  
9 model. More importantly, peak load point moves only by one or two stops no matter which data source  
10 is used, meaning model deployment would not produce any sudden shifts in frequency determination  
11 and scheduling considerations.

12

### 13 ***Ridecheck Validation—Correlation Studies***

14 To understand model usefulness in finding spot overcrowding or bunching (uneven loading) situations,  
15 ridecheck is matched to model data by unique bus number, run, and stop time (using manual fuzzy  
16 match allowing for inevitable timing calibration issues and/or typos in vehicle/run numbers.) Figure  
17 4(a,b,c) shows Bx19 comparison at individual trip-by-trip level using matched trips. AFC-AVL model  
18 tended to provide slightly lower boarding and disembarkation estimates than ridecheck figures (by up to  
19 5% on average, although there were plenty of cases where AFC estimated a higher number), but slightly  
20 higher total passenger-stops (by about 8%; a passenger-stop is one passenger travelling one stop.)  
21 Overall, R-square for on/off is 95% and for passenger-stops is 90%, demonstrating direct correlation  
22 between two data sources.

23

### 24 ***Comparison to Experimental APC Data***

25 In parallel with AFC model development, NYCTA installed experimental APC equipment in selected buses  
26 to gauge fleet wide deployment cost and feasibility. A comparison with APC data provided additional  
27 model validation opportunities but also highlighted needs for post-processing APC data.

28

29 Figure 3(g,j) reveals situations where APC software didn't correctly track beginning-of-trip, resulting in  
30 boardings at origin being missed and consequent negative loading near the terminus. Balancing  
31 algorithms fixed negative loading issues but had no basis for synthesizing missing boardings. Figure 3(e)  
32 shows classic passenger-remaining-onboard issues when disembarkations are not correctly counted.  
33 Figure 3(h) shows a very crowded trip where blocked front door sensors or other issues made it difficult  
34 track boardings correctly, whereas disembarkation was accurately counted leading to negative loads.

35

36 Figure 4(d) shows stop-level correlation analyses showing differences between AFC and APC data of  
37 about one-third. Likely reasons include situations shown in Figure 3(g,j) which might be corrected by  
38 keeping APC unit powered while bus engine is switched off and/or headsign not set in period prior to  
39 departure from origin.

40

41 While APC data streams, once stabilized, can be helpful, time-tested method of measuring farebox  
42 takings (i.e. AFC counts) continues to be an important data source for service adjustments. APC data  
43 requires substantial post-processing prior to use (30). APC hardware is not reliable compared to AFC  
44 equipment. APC dedicated maintenance costs are relatively high for data collection, whereas AFC costs  
45 are shared with fare collection business needs. Additionally, O/D matrices or individual travel itinerary  
46 are not available directly from APC. Future work could focus on limited fleet deployment of APC,  
47 integrating sample APC data streams as calibration factors in an AFC-AVL model.

1

**2 Aggregate Origin-Destination Matrix Validation**

3 In late 2013, NYCTA studied transit service in Co-op City, a community of high-rise apartments in  
4 northern Bronx with over 40,000 residents. In conjunction with limited ridecheck and origin-destination  
5 market surveys, AFC-AVL model was used and generated satisfactory results. One unanticipated benefit  
6 was incidental availability of limited O/D survey data, used for additional model validation. Conducted  
7 throughout Co-Op city during weekdays, 1,363 customers were interviewed. Distribution of  
8 destinations are shown in Figure 4(e). Intermodal linked trips derived from AFC-AVL model for over  
9 10,000 daily Co-op City originating trips yielded very similar results.

10

**11 General Observations**

12 Some bus routes have undergone intensive stop-by-stop comparison with ridecheck data. Few  
13 ridechecks were available for same-date comparison, thus most validation was completed by comparing  
14 monthly average data to ridecheck from earlier dates. Generally speaking:

15

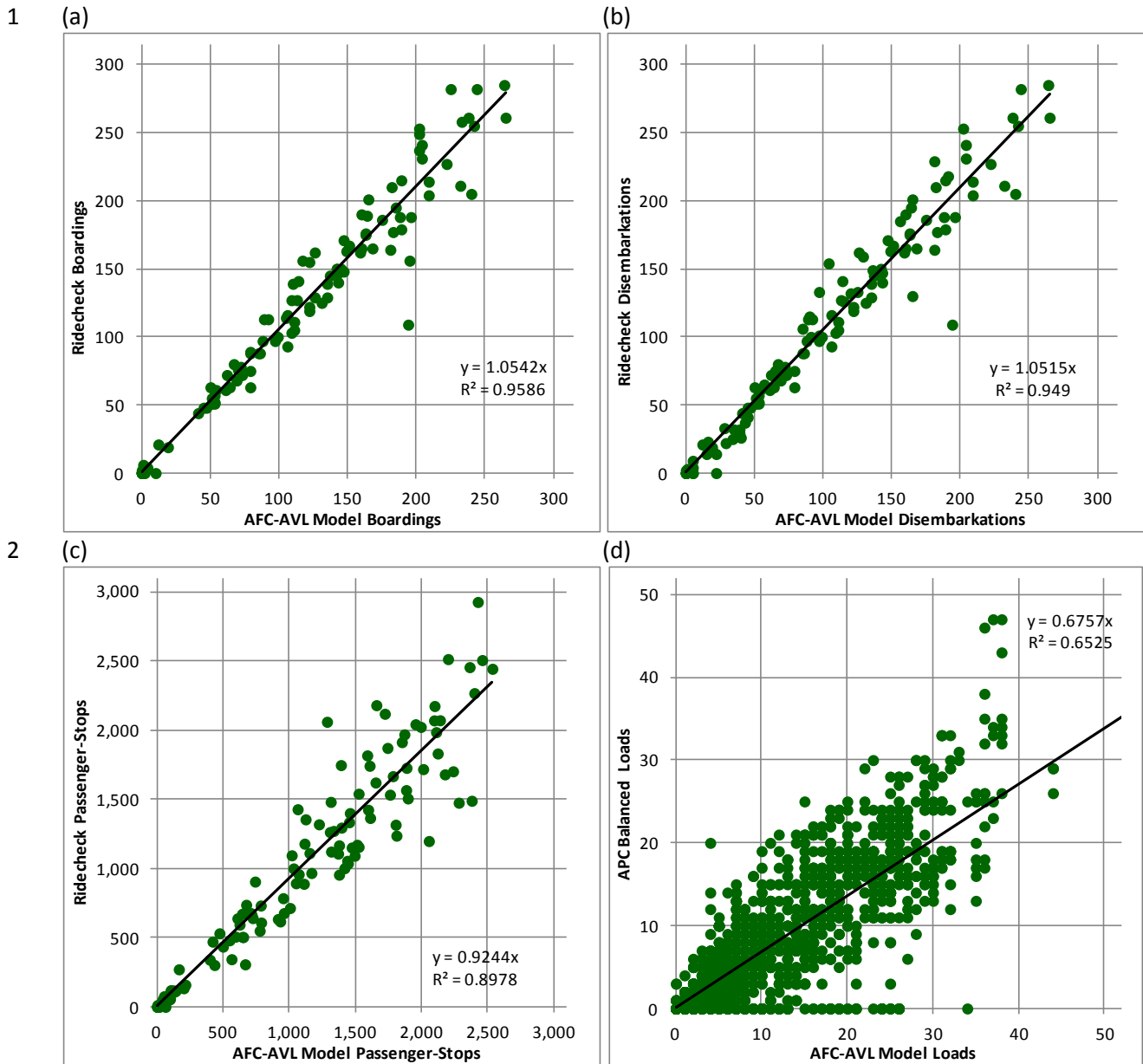
- 16 • Most comparisons yield very satisfying results in terms of load profile shape and values.
- 17 • Monthly averages usually generate better results than individual dates. AFC transaction time is  
18 rounded to nearest six-minutes. For an individual date, model may not accurately pin-point  
19 every stop. As monthly average, correct bus stops have higher probabilities of occurrence.
- 20 • Some routes (e.g. M116) were found to have less total ridership than ridecheck, which suggests  
21 that model underestimated fare evader numbers and/or students failing to swipe their  
22 MetroCard when boarding.
- 23 • Weekday results are usually more consistent than weekend. On weekends, services are more  
24 irregular due to events and G.O.s. Passengers are also more likely to combine bus with other  
25 transportation modes.

26

27 So far, 90% of express bus routes and over 40% of Bronx and Manhattan local routes have been  
28 validated against historic ridecheck data. Express bus was given priority due to ongoing NYCTA efforts  
29 to restructure and rationalize routes. Most routes generated satisfactory results; special exception  
30 processing continues to be added where necessary as model is fine-tuned and GTFS data is cleaned up.

31

32



3

(e)

Origin	Destination	AFC-AVL Model	Market Survey
Co-op City	Co-op City	8%	12%
Co-op City	Bronx (outside Co-op City)	55%	59%
Co-op City	Manhattan	27%	20%
Co-op City	Brooklyn	3%	3%
Co-op City	Queens	5%	4%
Co-op City	Elsewhere	2%	2%

4  
 5 **Figure 4** Advanced Algorithm Validation: (a) Time-Period Level Same-Day Model-Ridecheck  
 6 Correlation—Ons; (b) Time-Period Level Same-Day Model-Ridecheck Correlation—Offs; (c) Time-Period  
 7 Level Same-Day Model-Ridecheck Correlation—Sum of Link Loads; (d) Trip Level Same-Day Model-APC  
 8 Correlation—Link Loads; (e) Model-Survey Comparison in Co-op City O/D Study.

## 1 CONCLUSIONS

2 New York City surface transit boarding and alighting production application’s success with big data  
3 shows promise of estimating bus boarding and alighting counts with AVL data. Once limited to  
4 academic research, a daily production application has been implemented with high quality of results.  
5 This being a well-researched area, great benefits are leveraged when transit authorities discover  
6 production applications.

7  
8 This application validates results with ridecheck, APC system, and market research surveys, with  
9 encouraging results. Critical passenger information like origin-destination matrix and travel patterns  
10 needed for surface planning and scheduling, traditionally obtained through labor-intensive, time and  
11 budget constrained survey, can now be acquired through applications taking advantage of big transit  
12 data. Besides underlying significant economic benefits, availability of data in term of quantity,  
13 timeliness, type detail, and location precision makes it a clearly better option than surveys.

14  
15 The heart of this innovation, the “dynamic transfer link” network generation method in shortest path,  
16 proved to be effective and feasible in solving large scale and highly complex realistic transit networks.  
17 Even though most effort spent in customizing methodologies to complicated New York City surface  
18 transit system with various subsystems and exceptions, basic methodologies developed is easily adapted  
19 to other transit systems.

### 20 **Future of Production Algorithms**

21  
22 As faster computational tools become available to work with even larger volumes of data, real-time  
23 analytics have become the norm in transportation. Specialized solutions like massively-parallel graphics  
24 co-processors (42) have been adapted for transportation analytics in research settings. Thus, needs for  
25 pure algorithmically-based optimization may decline in future. Nonetheless, core ideas demonstrated in  
26 this paper are useful for practitioners in the following areas:

- 27  
28 1. **Exploratory research**—taking advantage of new data streams and to produce sample analyses  
29 (on commodity desktops) that could justify larger expenditures in specialized analytical  
30 equipment  
31 2. **Short-term production**—public authorities have long procurement cycles for technological  
32 equipment; lean algorithms can be used in production on existing equipment (e.g. general  
33 purpose virtual servers), allowing analytics to be more responsive to operational needs  
34 3. **Scope & capability expansion**—even with specialized equipment, explicit trade-offs to minimize  
35 execution time, problem simplification, and task-specific optimization would allow same  
36 equipment to handle larger problems.

37  
38 With algorithmic innovations to optimize computation time, highly complex tasks can be manageably  
39 adapted for daily production. This successful effort shows it is feasible to introduce automation in  
40 transportation planning—taking advantage of rapid advancement of technology in separate data  
41 sources, and integrates them into a powerful information resource. With efficient algorithm designs,  
42 modern computer hardware can handle large-scale transportation modelling at detail levels, even on  
43 standard desktop PCs.

### 44 **Further Applications**

1 With production-driven generation of extensive O/D matrices on a daily basis, new concept of “micro  
2 transportation planning” emerges, where planning occurs at individual passenger level. Traditionally,  
3 transportation planning utilizes scarce and highly aggregated sample data. With daily processing of  
4 individual trips, farecard histories can be built over longer periods and more useful information could be  
5 inferred.

6  
7 Compared to subway stations, bus stops are much closer to passenger actual residence and work  
8 locations. Over 80% of MetroCard trips start at identical stations each weekday (37). This methodology  
9 was used in Co-op City study to identify residents. Most frequent stops other than residence might be  
10 work location. It could be possible to synthesize journey-to-work data. Based on MetroCard attributes  
11 (student, senior, and other status), together with Census socioeconomic data associated with home  
12 stops, valuable market research could be generated without unreliable and costly manual surveys.

13  
14 Compared to traditional ridecheck, AFC-AVL model has many advantages:

- 15 • Ridecheck has limited time frame due to cost considerations, hence effect of day-to-day  
16 variation is unknown. Using average loads from AFC-AVL model can avoid over- or under-  
17 design caused by perturbations on sample day.
- 18 • Ridecheck is done few routes at a time. AFC model provides system wide information everyday,  
19 which assists planners in optimizing entire system by resource reallocation.
- 20 • Daily data can help planner to identify repeating patterns of operations or ridership disturbance.

21  
22 Potentially, similar applications taking advantage of new data sources, algorithms, and hardware could  
23 contribute to a new era in micro-level transportation planning:

- 24 • Farecards could be linked to demographic (income, race, etc.) groups in aggregate, to observe  
25 travel patterns and economic geography.
- 26 • Observe impacts of non-recurring versus recurring passenger congestion.
- 27 • With suitable assumptions, trip purpose can be derived, and customer activity habits can be  
28 studied and marketing efforts could be coordinated based on this data.

## 30 31 **ACKNOWLEDGEMENTS**

32 Authors gratefully acknowledge support and assistance of Wan Chen, Stella Levin, and Rob Hickey  
33 during algorithm development, and Cheri Tan and Tuan Huynh for assisting with validation and  
34 development of the production system. Responsibility for errors or omissions remains with the authors.  
35 Opinions expressed or implied are the authors’ and do not necessarily reflect official policies of any  
36 organization.

## 1 REFERENCES

- 2 (1) Lee, Sang Gu, M. Hickman, D.Q. Tong. Stop Aggregation Model: Development and Application.  
3 *Transportation Research Record: Journal of the Transportation Research Board 2276* (2012), pp. 38–47.  
4 TRB #12-1287.  
5
- 6 (2) Barry, James Jerome, R. Newhouser, A. Rahbee, and S. Sayeda. Origin and Destination Estimation in  
7 New York City with Automated Fare System Data. *Transportation Research Record 1817* (2002), pp. 183-  
8 187. TRB #02-1045.  
9
- 10 (3) Zhao, Jin Hwa, A. Rahbee, and N.H.M. Wilson. Estimating a Rail Passenger Trip Origin-Destination  
11 Matrix Using Automatic Data Collection Systems. In *Computer-Aided Civil and Infrastructure*  
12 *Engineering*, No. 22, pp. 376-387 (2007).  
13
- 14 (4) Guptill, Robert. Data from MBTA's Automated Fare Collection (AFC). Presented at *TRB National*  
15 *Transportation Planning Applications Conference*, May 2009. Retrieved from [http://www.trb-](http://www.trb-appcon.org/TRB2009presentations/s19/07_impact.ppt)  
16 [appcon.org/TRB2009presentations/s19/07\\_impact.ppt](http://www.trb-appcon.org/TRB2009presentations/s19/07_impact.ppt) on June 27, 2010.  
17
- 18 (5) Sun, Yanshuo and P.M. Schonfeld. Schedule-Based Route Choice Estimation with Automatic Fare  
19 Collection Data for Rail Transit Passengers. Presented at *Transportation Research Board 93rd Annual*  
20 *Meeting*. Paper #14-0834.  
21
- 22 (6) Frumin, Michael, J.H. Zhao, N.H.M. Wilson, and Z. Zhao. Automatic Data for Applied Railway  
23 Management: Case Study on the London Overground. In *Transportation Research Record: Journal of the*  
24 *Transportation Research Board 2353*, pp. 47-56 (2013). TRB #13-2987.  
25
- 26 (7) Wong, S.C., and C.O. Tong. Estimation of Time-Dependent Origin-Destination Matrices for Transit  
27 Networks. In *Transportation Research Part B: Methodological*, Volume 32, Issue 1, pp. 35-48, January,  
28 1998.  
29
- 30 (8) Barry, James Jerome, R. Freimer, and H.L. Slavin. Use of Entry-Only Automatic Fare Collection Data  
31 to Estimate Linked Transit Trips in New York City. In *Transportation Research Record: Journal of the*  
32 *Transportation Research Board 2112*, pp. 53-61 (2009). TRB #09-2812.  
33
- 34 (9) Liao, Chen-Fu, and H. Liu. Mining Bus Location, Passenger Count and Fare Collection Database for  
35 Intelligent Transit Applications. Presented at the *21st Annual Transportation Research Conference*, April  
36 27-28, 2010, St. Paul, Minn. Retrieved from  
37 <http://www.cts.umn.edu/Events/ResearchConf/2010/presentations/24-liao.pdf> on June 27, 2010.  
38
- 39 (10) Ro, Wei-Yuan (羅惟元). Using Taipei EasyCard Transaction Data to Explore the O-D Table of Bus  
40 Passengers. Masters Thesis, Graduate Institute of Transportation Management, Tamkang University,  
41 Damshui, Taiwan, June 2008. Retrieved from  
42 <http://tkuir.lib.tku.edu.tw:8080/dspace/handle/987654321/33825> on June 27, 2010.  
43
- 44 (11) Seaborn, Catherine, J. Attanucci, and N.H.M. Wilson. Analyzing Multimodal Public Transport  
45 Journeys in London with Smart Card Fare Payment Data. In *Transportation Research Record: Journal of*  
46 *the Transportation Research Board 2121*, pp. 55-62 (2009). TRB #09-3419.

- 1  
2 (12) Morency, Catherine, M. Trepanier, and B. Agard. Measuring Transit Use Variability with Smartcard  
3 Data. In *Transport Policy* **14**:3, pp. 193-203.  
4
- 5 (12a) Wang, Wei, J.P. Attanucci, and N.H.M. Wilson. Bus passenger O-D estimation and related analyses  
6 using automated data collection systems. In *Journal of Public Transport* **14**:4, pp. 132-150 (2011).  
7
- 8 (14) Munizaga, Marcela, C. Palma, and D. Fischer. Estimation of a Disaggregate Public Transport OD  
9 Matrix from Passive SmartCard Data from Santiago, Chile. Presented at *Transportation Research Board*  
10 *90th Annual Meeting*. TRB #11-0430.  
11
- 12 (15) Farzin, Janine M. Constructing an Automated Bus Origin-Destination Matrix Using Farecard and  
13 Global Positioning System Data in São Paulo, Brazil. In *Transportation Research Record 2072*,  
14 *Transportation Research Board of the National Academies*, 2008.  
15
- 16 (16) McCord, Mark R., R.G. Mishalani, P. Goel, and B. Strohl. Iterative Proportional Fitting Procedure to  
17 Determine Bus Route Passenger Origin–Destination Flows. In *Transportation Research Record: Journal*  
18 *of the Transportation Research Board 2145*, pp. 59-65 (2010). TRB #10-3425.  
19
- 20 (17) Shi, Xin Yi and H.F. Lin. The Analysis of Bus Commuters’ Travel Characteristics Using Smartcard  
21 Data: the Case of Shenzhen, China. Presented at *Transportation Research Board 93rd Annual Meeting*.  
22 TRB #14-2571.  
23
- 24 (18) Zúñiga, Felipe G., J.C.A. Muñoz, R.E. Giesen. Real-Time Prediction and Update of Dynamic Origin-  
25 Destination Matrices on a Transit Corridor. Presented at *TransLog Transportation and Logistics*  
26 *Workshop*, Hamilton, Ont., Canada (2009).  
27
- 28 (19) Furth, Peter G., B. Hemily, T.H.J. Muller, and J.G. Strathman. *Using Archived AVL-APC Data to*  
29 *Improve Transit Performance and Management*. Transit Cooperative Research Program, Report 113  
30 (2006).  
31
- 32 (20) Lu, Alex, and A.V. Reddy. Algorithm to Measure Daily Bus Passenger Miles Using Electronic Farebox  
33 Data for National Transit Database Section 15 Reporting. In *Transportation Research Record: Journal of*  
34 *the Transportation Research Board 2216*, pp. 19-32 (2011). TRB #11-0368.  
35
- 36 (21) Levine, Brian, A. Lu, and A.V. Reddy. Measuring Subway Service Performance at New York City  
37 Transit: A Case Study Using Automated Train Supervision (ATS) Track-Occupancy Data. In *Transportation*  
38 *Research Record: Journal of the Transportation Research Board 2353*, pp. 57-68 (2013). TRB #13-2997.  
39
- 40 (22) Reddy, Alla V., A. Lu, S. Kumar, V. Bashmakov, and S. Rudenko. Entry-Only Automated Fare-  
41 Collection System Data Used to Infer Ridership, Rider Destinations, Unlinked Trips, and Passenger Miles.  
42 In *Transportation Research Record: Journal of the Transportation Research Board 2110*, pp. 128-136  
43 (2009). TRB #09-0809.  
44
- 45 (23) Hickey, Robert L., A. Lu, and A.V. Reddy. Using Quantitative Methods in Equity and Demographic  
46 Analysis to Inform Transit Fare Restructuring Decisions. In *Transportation Research Record: Journal of*  
47 *the Transportation Research Board 2144*, pp. 80-92 (2010). TRB #10-0280.  
48

- 1 (24) Holmes, Mark. HASTUS Implementation at MTA Bus Company. Presented at *American Public*  
2 *Transit Association Multimodal Operations Planning Workshop*, Seattle, Wash. (2011). Retrieved from  
3 [http://www.apta.com/mc/multimodal/previous/2011/Presentations/Session-10-MTA-Bus-Company's-](http://www.apta.com/mc/multimodal/previous/2011/Presentations/Session-10-MTA-Bus-Company's-Transition-into-HASTUS-M-Holmes.pdf)  
4 [Transition-into-HASTUS-M-Holmes.pdf](http://www.apta.com/mc/multimodal/previous/2011/Presentations/Session-10-MTA-Bus-Company's-Transition-into-HASTUS-M-Holmes.pdf) on July 5, 2014.  
5
- 6 (25) Levine, Brian, S. Iyer, and A.V. Reddy. Development of Automated Vehicle Location (AVL) Data  
7 Based System to Improve Bus Service at New York City Transit. Presented at *Transportation Research*  
8 *Board 93rd Annual Meeting*. TRB #14-1711.  
9
- 10 (26) Glikin, Michael, and A.V. Reddy. MTA Bus Time Implementation & New Applications. Presented at  
11 *American Public Transit Association Multimodal Operations Planning Workshop*, San Francisco, Calif.  
12 (2013). Retrieved from  
13 [http://www.apta.com/mc/multimodal/previous/2013/presentations/Presentations/Glikin%20Reddy%2](http://www.apta.com/mc/multimodal/previous/2013/presentations/Presentations/Glikin%20Reddy%20MOPW%202013.pdf)  
14 [0MOPW%202013.pdf](http://www.apta.com/mc/multimodal/previous/2013/presentations/Presentations/Glikin%20Reddy%20MOPW%202013.pdf) on July 5, 2014.  
15
- 16 (27) Barbeau, Sean J., A. Borning, and K. Watkins. OneBusAway Multi-region: Rapidly Expanding Mobile  
17 Transit Apps to New Cities. Presented at *Transportation Research Board 93rd Annual Meeting*. TRB #14-  
18 4285.  
19
- 20 (28) Metropolitan Transportation Authority. Archived B63 Bus Time Data. Retrieved from  
21 <http://bustime.mta.info/wiki/Developers/ArchiveData> on July 5, 2014.  
22
- 23 (29) Thorup, Mikkel (1999). Undirected Single-Source Shortest Paths with Positive Integer Weights in  
24 Linear Time. In *Journal of the Association for Computing Machinery* **46**(3):pp.362–394.  
25 doi:10.1145/316542.316548.  
26
- 27 (30) RSM Services. Ridecheck Plus—Denver, Colo. Case Study. Retrieved from  
28 <http://www.rsm-services.com/rsm-projects/Denver-CO/30.html> on July 12, 2014.  
29
- 30 (31) Tollerson, Ernest. MTA's Open-Data Opportunities in 2013 and Beyond—Developing All-Agencies  
31 Policies, Projects & Public Private Partnerships (P3s). Presented at *APTA TransiTech Conference*,  
32 Phoenix, Ariz., March 20, 2013.  
33
- 34 (32) Reddy, Alla V., J.A. Kuhls, and A. Lu. Measuring and Controlling Subway Fare Evasion: Improving  
35 Safety and Security at New York City Transit. In *Transportation Research Record: Journal of the*  
36 *Transportation Research Board* 2216, pp. 85-99 (2011).  
37
- 38 (33) Beaton, Eric B., T.V. Orosz, R. Ravit, R. Thompson, and D. Tyson. A Limited for the 21st Century:  
39 Applying BRT Principles to Create Select Bus Service on Hylan Boulevard. In *TRB 93rd Annual Meeting*  
40 *Compendium of Papers*. TRB #14-2807.  
41
- 42 (34) Beaton, Eric B., E. Bialostozky, O. Ernhofer, T.V. Orosz, T. Reiss, and D. Yuratovac. Designing Bus  
43 Rapid Transit Facilities for Constrained Urban Arterials: Case Study of the Selection Process for the  
44 Webster Avenue Bus Rapid Transit Running Way Design in New York City. In *Transportation Research*  
45 *Record: Journal of the Transportation Research Board* 2352, pp. 50-60.  
46
- 47 (35a) Sheffi, Yosef. Urban Transportation Networks: Equilibrium Analysis with Mathematical  
48 Programming Methods. 1984. Prentice-Hall, Inc., Englewood Cliffs, New Jersey 07632.



- 1  
2 (35) Beaton, Eric B., J.E. Barr, J.V. Chiarmonte, T.V. Orosz, D. McAfee, and A. Sugiura. Select Bus Service  
3 on M15 in New York City: Bus Rapid Transit Partnership. In *Transportation Research Record: Journal of*  
4 *the Transportation Research Board 2277*, pp.1-10 (2012). TRB #12-1809.  
5  
6 (36) Barr, Joseph E., E.B. Beaton, J.V. Chiarmonte, and T.V. Orosz. Select Bus Service on Bx12 in New York  
7 City: Bus Rapid Transit Partnership of New York City DOT and Metropolitan Transit Authority New York  
8 City Transit. In *Transportation Research Record: Journal of the Transportation Research Board 2145*, pp.  
9 40-48 (2010). TRB #10-1693.  
10  
11 (37) Tempest Ludlow, Pete, and T. Stasko. Validate O-D Assumption Using Entry Only Swipe Data.  
12 Working Paper, System Data Research, Operations Planning, New York City Transit (2014).  
13  
14 (38) Gordon, J. B., H. Koutsopoulos, N. Wilson, and J. Attanucci. Automated Inference of Linked Transit  
15 Journeys in London using Fare-Transaction and Vehicle Location Data. In *Transportation Research*  
16 *Record: Journal of the Transportation Research Board 2343*, pp.17-24 (2009).  
17  
18 (39) Cui, A. Bus Passenger Origin-Destination Matrix Estimation Using Automated Data Collection  
19 Systems. MS Thesis. Massachusetts Institute and Technology, Cambridge, 2006. Retrieved from  
20 <http://dspace.mit.edu/handle/1721.1/37970> on July 15, 2014.  
21  
22 (40) Furth, Peter G., J.G. Strathman, and B. Hemily. Making Automatic Passenger Counts Mainstream:  
23 Accuracy, Balancing Algorithms, and Data Structures. In *Transportation Research Record: Journal of the*  
24 *Transportation Research Board 1927*, pp.207-216 (2005).  
25  
26 (41) Li, Da Ming, Y.J. Lin, X.L. Zhao, H.J. Song, and N. Zou. Estimating a transit passenger trip origin-  
27 destination matrix using automatic fare collection system. In *Proceedings of the 16th International*  
28 *Conference on Database Systems for Advanced Applications*, pp. 502-513 (2011).  
29  
30 (42) Zhang, Jianting, S.M. You, and L. Greenwald. High-Performance Spatial Query Processing on Big Taxi  
31 Trip Data using GPGPUs. Retrieved from [http://www-cs.cuny.cuny.edu/~jzhang/papers/taxi\\_p2pdis.pdf](http://www-cs.cuny.cuny.edu/~jzhang/papers/taxi_p2pdis.pdf)  
32 on July 18, 2014.  
33  
34