

1
2
3 **INTEGRATION OF THE NATIONAL LONG DISTANCE PASSENGER TRAVEL**
4 **DEMAND MODEL WITH THE TENNESSEE STATEWIDE MODEL**
5 **AND CALIBRATION TO BIG DATA**
6
7

8
9 **Vincent L. Bernardin, Jr., PhD**

10 Director, RSG

11 2709 Washington Ave, Ste. 9

12 Evansville, IN 47714

13 Vince.Bernardin@RSGinc.com

14 Ph: 812-200-2351
15

16 **Nazneen Ferdous, PhD**

17 Travel Demand Modeler, CH2M

18 2411 Dulles Corner Park, Suite 500

19 Herndon, VA 20171

20 Nazneen.Ferdous@ch2m.com

21 Ph: 703-376-5000
22

23 **Hadi Sadrsadat, PhD**

24 Consultant, RSG

25 2200 Wilson Blvd., Suite 205

26 Arlington, VA 22201

27 Hadi.Sadrsadat@RSGinc.com

28 Ph: 888-774-5986
29

30 **Steven Trevino**

31 Analyst, RSG

32 2709 Washington Ave, Ste. 9

33 Evansville, IN 47714

34 Steven.Trevino@RSGinc.com

35 Ph: 812-200-2352
36

37 **Chin-Cheng Chen**

38 Forecasting Office Supervisor, Tennessee Department of Transportation

39 Suite 900, James K. Polk Building

40 Nashville, TN 37243-0334

41 chin-cheng.chen@tn.gov

42 Ph: 615-253-6301
43
44

45 *Submitted: November 14, 2016*

46 *Word Count: 6,466 words + 4 tables/figures @ 250 words = 7,466 equivalent words*

47 **ABSTRACT**

48 The Tennessee Department of Transportation chose to replace their quick-response-based long
49 distance component in their statewide model by integrating FHWA's new national long distance
50 passenger travel demand model into their new statewide model and calibrating it to long distance
51 trips observed in cell-phone based origin-destination data from AirSage. The new national long
52 distance model is a national scale, tour-based simulation model developed from FHWA research
53 on long distance travel behavior and patterns. The tool allows the evaluation of many different
54 policy scenarios including fare or service changes for various modes including commercial air
55 travel, intercity bus, and Amtrak as well as highway travel. The availability of this new tool
56 represents a new opportunity for state DOTs developing statewide models. Commercial cell-
57 phone based big data on long distance trips also represents a new opportunity and a new data
58 source on long distance travel patterns which have previously been the subject of very limited
59 data collection in the form of surveys. This project is the first to seize on both of these new
60 opportunities by integrating the new national long distance model with the new Tennessee
61 statewide model and by processing big data for use as a calibration target for long distance travel
62 in a statewide model. The paper demonstrates the feasibility of integrating the new national
63 model with statewide models, the ability of the national model to be calibrated to new data
64 sources, the ability to combine multiple big data sources, the value of big data on long distance
65 travel as well as important lessons on its expansion.

66

67 *Keywords:* Long Distance Travel, Statewide Model, Travel Demand Forecasting, AirSage,
68 rJourney

69 INTRODUCTION

70 The Tennessee Department of Transportation (TDOT) chose to implement an innovative
71 approach to forecasting long distance passenger travel in their new statewide model. The
72 standard practice for handling long distance passenger travel in statewide travel models is to add
73 one or more special long distance trip purposes in a three- or four-step model structure, often
74 borrowing parameters from studies such as NCHRP 735 (1). Instead, TDOT chose to integrate
75 FHWA's new national long distance passenger travel model, rJourney, into their statewide model
76 and calibrate it to long distance trips observed in cell-phone based big data.

77 Although long distance trips are much less common than short distance trips, because
78 each trip has the potential to contribute so many vehicle miles of travel (VMT), these trips have a
79 large and disproportionate effect on congestion and traffic on major intercity corridors such as I-
80 40, I-75, I-24, and I-65 in Tennessee. A significant portion of long distance trips related to
81 business travel also have notably higher value of time than most other trips, so reductions in
82 delays for these trips can produce comparatively large economic benefits.

83 The availability of FHWA's new national rJourney model (2, 3, 4), together with the
84 availability of new big data sources such as cell phone derived big origin-destination (OD) data,
85 presented a new and exciting opportunity to dramatically improve the representation of long
86 distance travel in Tennessee and allow new types of scenario analysis. For instance, the
87 inclusion of a robust long distance mode choice modeling in rJourney allows the evaluation of
88 scenarios such as increased air fares, expanded Amtrak service or new intercity bus services and
89 the impact of such assumptions on highway volumes.

90 This is the first application of the national long distance model to support statewide
91 modeling and forecasting, and is believed to also be the first use of big cell-phone based OD data
92 to support development of a statewide travel model, although it is known that work to update the
93 Virginia statewide model with similar data began at close to the same time.

94 The work to incorporate the new national model within TDOT's statewide model was
95 part of a larger update to the model. TDOT originally developed a simple statewide model for
96 Tennessee in 2003. TDOT developed a new, version 2, statewide model in 2014 to support
97 development of their statewide long range plan. Although the version 2 model was also limited
98 in sophistication due to the project schedule, it included three times as many zones (5) and road
99 miles and offered much finer resolution in the representation of projects and their impacts. The
100 version 2 model also included a truck model supported by the purchase of an eight week dataset
101 of truck GPS based OD data from the American Transportation Research Institute (ATRI). The
102 data included information from over 234,000 individual trucks on over 6.5 million truck trips
103 representing roughly 11% of the trucks on the road for 56 days. The version 2 model also used
104 Longitudinal Employer-Household Dynamics (LEHD) Origin-Destination Employment Statistics
105 developed by the Census Bureau in cooperation with the Bureau of Labor Statistics which
106 provides big OD data on commuting patterns based on administrative tax records.

107 This data-driven approach, albeit incomplete and supplemented by traditional synthetic
108 quick response methods in version 2, lead to very good model performance. The model's
109 highway assignment achieved impressive validation statistics versus traffic counts for a
110 statewide model including a 37% root mean squared error (RMSE) and a correlation coefficient
111 of 0.97.

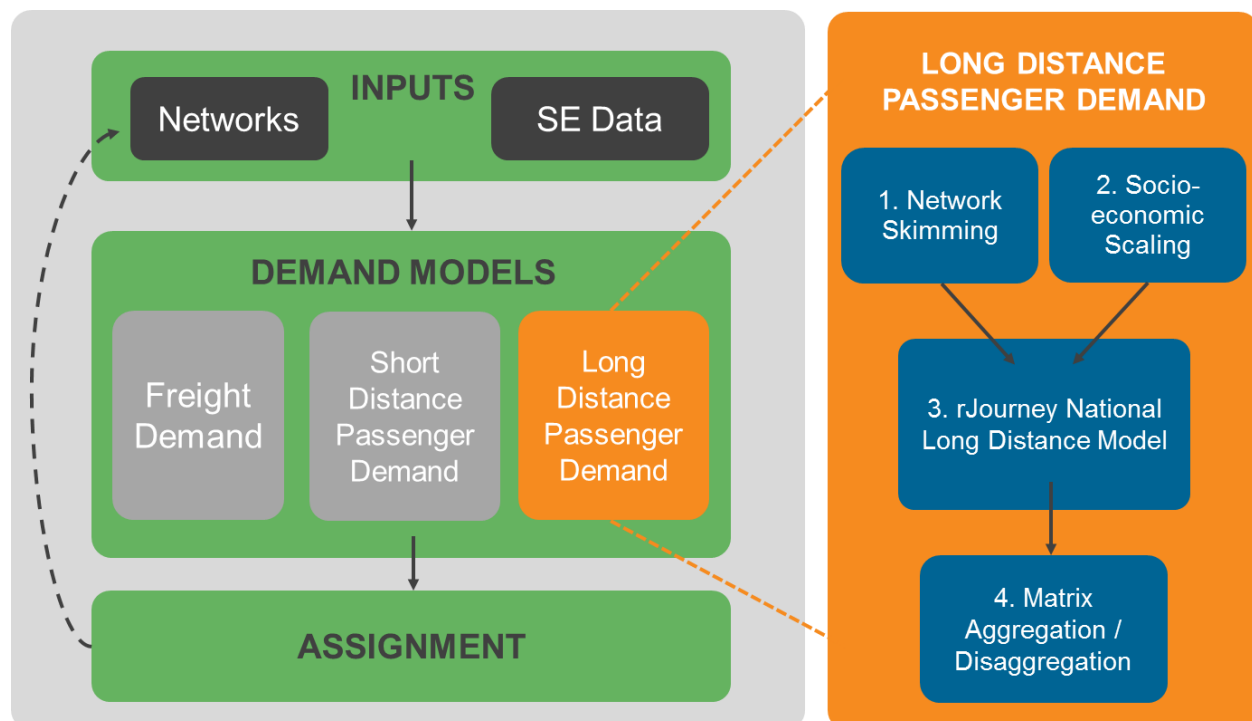
112 Since the project schedule had precluded the incorporation of all the desired functionality
113 and features in the version 2 model, TDOT embarked upon an update to develop the version 3
114 Tennessee Statewide Travel Model (TSTM3). The new model incorporates a new commodity

115 flow freight model, an advanced trip-based model for short distance passenger trips including
 116 mode and destination choice models with non-home-based trips linked to home-based trips (6),
 117 as well as the integration of the new national model for long distance passenger trips. Building
 118 on the success of the data driven approach with the ATRI and LEHD data in version 2, TDOT
 119 purchased cell-phone based data to support the development of the version 3 model, with special
 120 filtering for long distance trips to support the calibration of the national long distance model in
 121 particular.

122 The paper below begins by describing the data and its processing before turning to
 123 describe the integration of the national model with the Tennessee statewide model designed to
 124 produce reasonable runtimes. The paper then documents the success of the project at calibrating
 125 the national model to the big OD data and its support of overall impressive validation statistics
 126 before offering some concluding thoughts.

127 INTEGRATED MODELING METHODOLOGY

128 The demand forecasting components of the Tennessee Statewide Travel Model (TSTM3) can be
 129 grouped into three sets of models: (i) freight demand models, (ii) short distance daily passenger
 130 demand models, and (iii) long distance passenger demand models. The latter, long distance
 131 models are the focus of this paper and include the national long distance passenger demand
 132 model together with several data manipulation and processing steps to achieve integration of and
 133 translation between the networks and zone systems for Tennessee and the national model. (See
 134 Figure 1.)
 135



136
 137 **FIGURE 1 Tennessee Statewide Travel Model Integration with the National Long Distance**
 138 **Passenger Travel Demand Model**

139 The first process involved in the integrated model run is the production of the highway
 140 skim for the national model using the Tennessee model network. The national model natively

141 uses a highway network developed from the Oak Ridge National Highway Planning Network
142 and a zone system of 4,570 zones or national use modeling areas (NUMAs) comprised of
143 counties in rural areas and Census Public Use Microdata Areas (PUMAs) in urban areas. To
144 facilitate the integration of the TSTM3 with the national model, the national model's highway
145 network was added to the TSTM's model network outside of Tennessee, while the more detailed
146 TSTM network was retained within (and immediately adjacent to) Tennessee. Outside of
147 Tennessee the national model's original centroids and centroid connectors were retained. Within
148 Tennessee, the TSTM centroid nearest the center of each NUMA was designated as the centroid
149 for that NUMA. Highway skims for the national model were thus created using the national
150 model's zone system but the TSTM's more detailed highway network within Tennessee. (A
151 separate highway skim is used for the daily short distance passenger trips using the more detailed
152 zone system within Tennessee, while a third highway skim with yet another zone system is used
153 for the freight demand model.) This approach has the advantage of limiting the size and
154 associated run time of the skim and not requiring processing to convert skim matrices between
155 zone systems. Moreover, it allows the national model to be sensitive to network changes in
156 Tennessee without requiring improvements to be coded in two separate networks (one for short
157 distance and one for long distance models).

158 The second process required to integrate the TSTM3 with the national model is to ensure
159 that the socioeconomic growth assumptions from the TSTM3 are incorporated into the national
160 model's inputs. There are two parts to achieving this. The first is the scaling of the synthetic
161 population for the national model to reflect population growth assumptions and the second is the
162 update of the national model's destination choice size variables to reflect employment growth
163 assumptions. The national model uses a detailed synthetic population covering the whole
164 country. Control variables for the synthetic population come from the Census's American
165 Community Survey (ACS) in the base year. Creating future year synthetic populations for the
166 whole country poses several challenges. The process is both data and runtime intensive.
167 Detailed demographic control total forecasts are required covering the whole country. Moreover,
168 these control totals must be carefully crafted or checked to ensure that they are internally
169 consistent (e.g., the distribution of households by size and the distribution of households by
170 workers must have the same number of total households for each zone and the number of
171 households with X workers cannot exceed the number of households with at least X people, etc.).
172 Given both the data and computational challenges of synthesizing new populations for
173 alternative future scenarios, an alternative approach was developed to simply rescale the base
174 year synthetic population based on forecast growth in households. This approach has the
175 limitation of not being able to reflect future changes in the characteristics (such as income) of
176 households in an area, but allows household growth forecasts from the TSTM3 zones to be used
177 to automatically and reliably update the national model's inputs with very limited runtime,
178 without requiring the development of detailed socioeconomic forecasts for the whole country or
179 complex data reconciliation.

180 The other part of the socioeconomic updating is the recalculation of the national model's
181 destination choice size variables. Fortunately, this process is substantially simpler, and only
182 involves the recalculation of formulas using updated employment data forecasts taken from the
183 TSTM3. Thus, the national model reflects the household and employment growth scenarios
184 from the TSTM3 zonal data within Tennessee without requiring the duplication of this data in
185 another dataset. Employment and household growth outside Tennessee is taken from a simple
186 input table with household and employment totals provided at the county level.

187 After creating the necessary inputs for the national model using the TSTM3's highway
188 network and zonal data, the national model is run. The national model is a household level
189 disaggregate tour-based simulation model. In some regards, it can be considered akin to a
190 simplified activity-based model. (For details of the national model, see 4.) The national model's
191 components begin with tour generation, scheduling, duration, and party-size models by purpose
192 followed by mode and destination choice models similarly segmented by purposes including
193 leisure/vacation, visit friends or relatives, personal business, commute, and employer's business.
194 The national model includes four modes: highway, intercity bus, intercity rail, and commercial
195 air travel. As described above, the national model uses the TSTM3's highway network for its
196 highway travel times. It also uses the TSTM3 highway travel times to update intercity bus travel
197 times. Tennessee implementation uses the national model's original networks for intercity
198 passenger rail and commercial flights, but the user can adjust these to create alternative scenarios
199 such as increased or decreased commercial air service or fares or new intercity rail service.

200 The final process in the integrated modeling system is matrix manipulation to convert the
201 trip list from the national model, using its NUMA zone system, into a trip table matrix using the
202 TSTM3's assignment zone system. This involves both the disaggregation of national model
203 zones to TSTM3 zones within Tennessee and immediately surrounding areas and the aggregation
204 of national model zones farther away from Tennessee. The demand within Tennessee (and
205 nearby) is disaggregated based on a simple function of the socioeconomic characteristics of the
206 TSTM3 zones within each NUMA, designed to approximate the number of long distance trips
207 produced by and attracted to each zone. Demand farther from Tennessee is simply aggregated
208 into a larger zone system (at the level of states for much of the country far from Tennessee) to
209 keep the number of zones limited for assignment to help manage runtimes.

210 The resulting integrated system provided an efficient approach, allowing the national
211 model to be run as part of the TSTM3 modeling system, using the information in the TSTM3's
212 highway network and zone system. The long distance components of the TSTM3 including both
213 the national model itself together with the ancillary pre- and post-processing procedures
214 described runs in close to one hour on a machine with 12 physical cores and 32 GB of RAM.

215 **CELL-PHONE DATA**

216 As noted previously, this is believed to be the first use of big cell-phone based OD data to
217 support development of a statewide travel model, although it is known that work to update the
218 Virginia statewide model with similar data began at close to the same time. Large scale,
219 aggregated, anonymous, passively-collected cell-phone OD data such as used in this study has a
220 history of use for various purposes including the estimation of travel times (7) and origin-
221 destination patterns (8) and resulting data has been compared to and incorporated in metropolitan
222 area travel demand models (9, 10, 11, 12, 13, 14).

223 TDOT acquired origin-destination (OD) data from AirSage, Inc., for the state of
224 Tennessee and a halo area surrounding it. AirSage aggregates and processes information from
225 wireless data providers to provide mobility information such as trip tables. While the exact
226 number of unexpanded observed trips is unknown, an extremely conservative lower bound can
227 be established based on the number of OD pairs reported since at least one unexpanded trip of
228 each type must be observed to be expanded. Using this method, the cell phone data was based
229 on a minimum of 3,355,539 observed trips, although the actual number of observed trips is likely
230 significantly higher. In contrast, combined household travel survey prepared for TDOT from an
231 add-on sample to the 2008-2009 National Household Travel Survey (NHTS) and travel surveys

232 from local metropolitan planning organizations in the state contained a total of 81,065 trips by
233 10,344 households in 39,782 OD pairs. Thus, the cell-phone data contains at least 84 times as
234 many observations as the household and likely substantially more. The result is that the big cell-
235 phone data provides a much more complete picture of OD patterns compared to the survey.

236 Another way to understand the difference in the completeness of the new big data versus
237 traditional survey data is to consider the amount of the origin-destination space that the data
238 covers or the percentage of cells within the origin-destination matrix with an observed frequency.
239 TDOT's traditional household survey data included observations of 39,782 origin-destination
240 pairs or 0.3% of the cells in the origin-destination matrix. In contrast, the cell-phone data
241 included observations of 3,355,539 origin-destination pairs or 26.3% of the cells in the origin-
242 destination matrix. This substantially better coverage offered by big data is one major
243 motivation for its use to support travel modeling in general.

244 There is even further motivation for the use of cell phone or similar passively collected
245 big data for studying and modeling long distance trips in particular. It generally takes significant
246 additional effort in travel surveys to collect an adequate sample of long distance trips. For this
247 reason, it has tended to be expensive and rarely done. For example, TDOT's combined travel
248 survey dataset included only 1,076 long distance trips (over 50 miles in length) out of the 81,065
249 total trips, and these were clearly skewed towards long distance commute trips and the shorter
250 end of the spectrum of long distance trips. Across the United States over the past twenty years
251 only five useful attempts to collect a representative sample of long distance trips could be
252 identified for the development of the national long distance model. While its anonymous nature
253 precludes it from supplying the same kind of rich detailed information that surveys can,
254 passively collected big data such as the cell-phone based data used in this study provides a cost
255 effective alternative to at least for understanding long distance OD patterns.

256 In order to use the cell-phone based data to calibrate the national long distance passenger
257 model, it was necessary to first remove commercial travel from the dataset because cell-phone
258 data captures both personal and commercial trips since travelers, including truck drivers, carry
259 their cell phones regardless of their travel purpose. As was noted in the introduction, TDOT had
260 acquired and processed truck GPS data from ATRI. The initial plan was to simply subtract the
261 truck ODs based on the truck GPS data from the total cell-phone based ODs. However, the
262 initial attempt to do so revealed that there were more truck trips than total trips for 11% of the
263 OD pairs observed in the cell-phone data. Although only 0.2% of the total cell-phone trips were
264 involved, given the large number of OD pairs, this was considered problematic.

265 Upon investigation, it became clear that the primary reason for this was a difference in
266 the way the two datasets were processed relative to the definition of trips and long distance trips
267 in particular. The cell-phone data had been purchased with filtering to remove intermediate
268 stops (such as for fuel, meals, etc.) on long distance trips. Based on AirSage's description of
269 their methodology, if a traveler traveled 50 miles from home the criteria for defining a stop
270 changed and rather than being based on the amount of time the traveler spent in the same place,
271 instead, a stop was coded only when the traveler reached the point furthest from home and began
272 traveling back towards home. In this way intermediate stops between home and the assumed
273 destination at the farthest point from home are removed from the dataset. The truck GPS data
274 was originally not processed in an analogous way, so it included intermediate stops on long
275 distance trips. It was therefore necessary to re-process the truck GPS data, filtering out
276 intermediate stops. However, it is less clear how to define home for trucks and in many cases
277 much more difficult to identify than for most residents of an area who return home most nights.

278 Moreover, it was deemed important and desirable to allow for multiple destinations on a long
279 distance tour (e.g., a truck carrying one shipment from Nashville to Knoxville may then pick up
280 another shipment and take this to Chattanooga before returning to Nashville). For both these
281 reasons, a slightly different algorithm was used for removing intermediate stops from the truck
282 trips. When a truck traveled more than 50 miles from one origin, A, to a stop, B, based on dwell
283 time, location B was not immediately logged as a stop. Rather, at the next stop C (based on
284 dwell time), the distance between A and B plus the distance between B and C was compared to
285 the direct distance between A and C. If the direct distance between A and C was more than 95%
286 of the sum of the distance between A and B and between B and C, then B was considered an
287 intermediate stop and removed, otherwise it was retained. Thus, this criterion identified
288 intermediate stops based on whether the truck went out of its way to reach the location. This
289 method allowed the removal of many intermediate stops while still allowing multiple “true”
290 stops on a long distance tour.

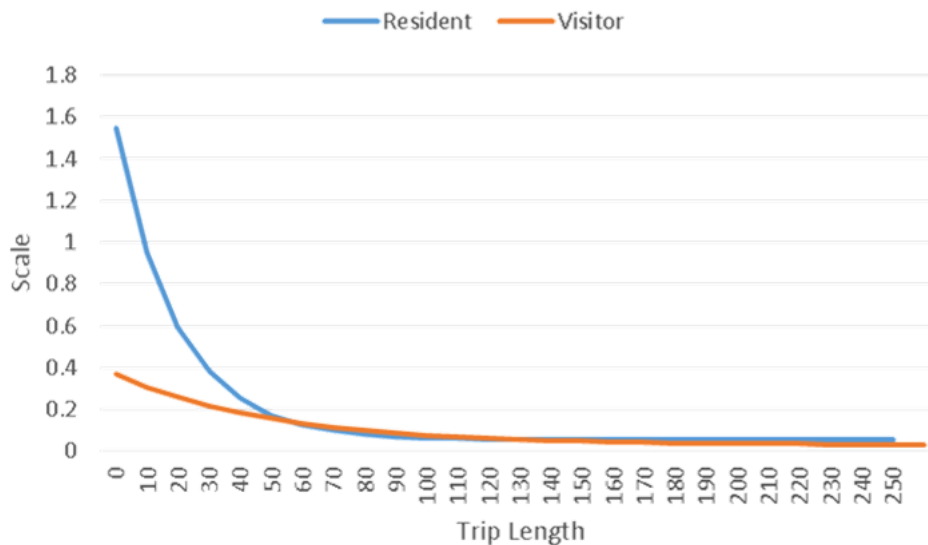
291 After re-processing the truck GPS data using this algorithm to remove intermediate stops,
292 the resulting truck ODs were again subtracted from the total cell-phone ODs. The number of OD
293 pairs with more truck trips than total trips was reduced by 87% from 11% of the ODs to only
294 1.3% involving less than 0.1% of the total trips. Although still not perfect, the output was
295 deemed acceptable because more than 98% of all cells have reasonable auto trips consisting
296 99.9% of total trips. The remaining OD pairs with more truck trips than total trips were reduced
297 to a fraction of a trip (to retain the information that some sort of trip was observed and allow for
298 expansion given some possibility that a passenger trip may have been observed). This
299 experience points to the importance of a common definition of trips (or stops) when combining
300 multiple big OD datasets.

301 The expansion of the cell phone data was tested and ultimately adjusted in a two stage
302 process. While the details of AirSage’s data expansion algorithm are proprietary trade secrets,
303 their documentation indicates that they use methods to expand their data based on the ratio of
304 cell-phones to population data at the inferred residence location. This basic approach has been
305 described and studied in academic literature (15, 16, 17, 18) as well as have more advanced
306 methods that make use of traffic assignment and/or optimization methods (19, 20, 21) to expand
307 cell-phone data based on traffic counts. The authors of this paper believe the latter methods to be
308 generally superior because traffic counts provide unbiased information on the total amount of
309 vehicular traffic on the road in various locations and the former methods based only on cell-
310 phone ownership levels fail to allow for any variety of factors which can affect the detection of
311 trips in particular locations, of particular durations, etc.

312 Initially, the cell-phone data was assigned to the Tennessee model network as a check on
313 its expansion and in order to determine necessary scaling. It is typically necessary to scale cell-
314 phone data to account for the number of cell phones per vehicle (closely related to but slightly
315 different than vehicle occupancy). However, the test revealed relatively poor fit to traffic counts,
316 (somewhat in contrast to experiences with using the data in metropolitan areas where simple
317 scaling usually produces at least reasonable agreement with counts). In particular, when scaled
318 to minimize total loading error, there was substantial underloading in urban areas on the order of
319 -10% versus counts and substantial overloading in rural areas on the order of +15% versus
320 counts. Since it is known that vehicle occupancy is substantially higher on long distance trips
321 than short distance trips, the initial response was to attempt scaling the trips based on distance
322 rather than uniformly as a whole. While this did improve the loading issues, it quickly became
323 apparent that the difference in scaling required to address the loading errors could not be

324 accounted for simply by higher vehicle occupancy on long distance trips. While vehicle
 325 occupancy may be two to three times higher on long distance trips than short distance trips, long
 326 distance trips appeared to be over-represented by a factor of ten or more. Other hypotheses were
 327 therefore also explored, such as that the bias may be related to area type or density rather than
 328 trip length. However, none of these explained or corrected the loading errors better than a
 329 distance-based correction, and generally they did worse. Therefore, a distance-based scaling was
 330 ultimately adopted, but based primarily on the hypothesis of a bias in the cell-phone data rather
 331 than on differences in vehicle occupancy. Upon further reflection, the possibility of a bias in
 332 cell-phone data toward the detection of long distance trips seems plausible. Since cell-phone
 333 mobility data depends on signaling between phones and towers, the likelihood of detection
 334 increases with the likelihood of this signaling and this signaling becomes more likely on longer
 335 trips for a variety of reasons including that people are more likely to use their phone the longer
 336 the trip. The probability of a person using their phone while on a local shopping trip is
 337 presumably much lower than the probability that a person uses their phone a trip to another city
 338 which is likely to take several hours if not a day or more. This line of reasoning provides at least
 339 a plausible explanation for a potential over-representation of long distance trips relative to short
 340 distance ones in a cell phone dataset.

341 While techniques for origin-destination matrix estimation (ODME) from counts without
 342 distance-based scaling could have been applied directly to simultaneously address distance bias
 343 and other potential issues with the expansion, the authors preferred a two-step process, first
 344 scaling trips parametrically based on functions of distance and using non-parametric ODME
 345 methods second. This approach helps avoid large adjustments from ODME without a clear
 346 understanding of the underlying problem or issue. A good review of ODME techniques, their
 347 limits, and effectiveness can be found in the study by Marzano et al. (22).



348

349 **FIGURE 2: Distance-based Scaling Functions for Resident and Visitor Trips**

350 Scaling functions were estimated separately for resident and visitor trips as visitor trips
 351 exhibited much more consistent over-representation independent of trip distance. This is
 352 consistent with the hypothesis of a bias against short trips since in this context visitors are
 353 already by definition (given the structuring of the data) long distance travelers. Parameters were
 354 fitted to the function $scale = c + a \times \text{Exp}(b \times \text{distance})$ using least squared errors. For resident

355 trips, $c = 0.0612$, $a = 1.6404$, and $b = -0.0507$. For visitors, $c = 0.0292$, $a = 0.3376$, and b
356 $= -0.0195$). The curves can be seen in Figure 2. The implication of the resident curve is that a
357 100-mile trip is 12 times as likely to be detected in the cell-phone data as a 10-mile trip. Given
358 that there may be 2 to 3 times as many people on a 100-mile trip as a 10-mile trip, this suggests
359 that a 100-mile trip is 4 to 6 times as likely to be detected than a 10-mile trip for reasons other
360 than vehicle occupancy. The application of these scaling factors did not completely resolve the
361 observed loading errors but significantly improved them reducing urban underloading to roughly
362 -2% and rural overloading to roughly 5%.

363 After the distance based scaling, ODME techniques were applied to further improve the
364 expansion of the cell-phone data versus counts. Careful consideration was given to setting
365 appropriate bounds on the ODME adjustments. On the one hand, the most limited adjustments
366 capable of producing good agreement with counts are desirable. At the same time, it is important
367 to acknowledge and allow ODME to factor trips to and from certain areas up and down to
368 account for varying degrees of cell coverage and other factors which can cause necessary
369 expansion factors to vary beyond simply the variance in cell-phone market shares by resident
370 areas. After some experimentation, ultimately, a minimum factor of 0.5 and a maximum factor
371 of 5.0 were chosen to limit ODME scaling of any given OD pair. In addition to these limits, the
372 average amount change in the trip matrix from ODME was closely monitored. The average
373 absolute difference between cells in the final adjusted trip matrix and in the scaled matrix was
374 4.3 trips and the average absolute percentage difference was 1.5%. Together with the limits on
375 minimum and maximum adjustments, these were deemed to be generally reasonable
376 adjustments. The trip length frequency distribution of the adjusted matrix was also compared to
377 the original matrix. The comparison showed that ODME resulted in a modest additional increase
378 in the expansion of short distance trips versus longer trips. This seemed to suggest that the
379 distance based scaling was not excessive but successful in accounting for most of the distance
380 related adjustments. The ODME adjustments improved the fit of the cell-phone based data from
381 55.5% RMSE to 36.6% RMSE versus over 12,000 traffic counts across the state of Tennessee
382 and given the relatively limited adjustments necessary to achieve this improved fit, this was
383 deemed a successful and helpful improvement to the expansion.

384 RESULTS

385 The national model was calibrated to the cell-phone data primarily through the adjustment of
386 constants in its component choice models. In particular, the calibration effort focused on the
387 adjustment of the tour frequency and destination choice models, since the cell-phone data
388 provided information primarily on these dimensions of long distance travel and did not provide
389 information by mode.

390 For purposes of calibration comparisons, Tennessee zones were grouped into eight
391 districts, each named for their largest/best known urban area(s). The total number of long
392 distance trips bound to or from each district in the model and in the cell-phone data are compared
393 in Table 1. As can be seen, the national model was able to be calibrated to closely reproduce the
394 observed long distance trip generation rates observed from the cell-phone data. This was
395 accomplished through the judicious adjustment of existing constants in the national model and
396 without the addition of any special constants specific to these districts or other districts or zones
397 in Tennessee. Most districts are within about 3,000 trips per day and less than 10% of their total.
398 The Knoxville district is somewhat under-predicted, most likely because the Smoky Mountains
399 and associated tourist areas attract more trips than predicted by the model. Trips to and from the

400 Tri-Cities district are over-predicted and the reason for this is less clear, but may be due to model
 401 not understanding the psychological and/or physical barrier posed by the mountainous
 402 topography of this area. Overall, however, the model does an impressive job of reproducing the
 403 number of trips observed for each district.

404 **TABLE 1 COMPARISON OF MODELED AND OBSERVED (CELL PHONE DATA)**
 405 **LONG DISTANCE TRIPS BY TENNESSEE DISTRICTS**

TN Districts	Observed	Modeled	% Difference
Tri-Cities	17,746	23,531	32.6%
Knoxville	59,149	53,239	-10.0%
Chattanooga	32,455	34,212	5.4%
Cookeville	22,486	21,239	-5.5%
Lynchburg	22,038	19,954	-9.5%
Nashville	88,502	85,622	-3.3%
Jackson	37,264	35,409	-5.0%
Memphis	30,340	31,067	2.4%
Total	309,980	304,272	-1.8%

406

407 Calibration of destination choice in the national model was more challenging. The cell-
 408 phone data revealed a significant bias against trips crossing the state border with the total number
 409 of long distance trips within the state slightly higher than trips to and from the state crossing the
 410 state border. The pattern cannot be predicted or explained on the basis of distance alone. For
 411 that reason, the gravity models based on NCHRP 735 in the version 2 TSTM could not
 412 reproduce the pattern, nor could the original national model. In order to reproduce the observed
 413 pattern, a single new term had to be added to the utility function of the national model's
 414 destination choice models to account for a psychological bias against crossing the state border.
 415 Similar psychological boundary effects associated with rivers, railroads, freeways, and
 416 governmental boundaries are commonly observed and incorporated metropolitan destination
 417 choice models. The addition of this term allowed the calibration of the national model to

418 reproduce the pattern observed in the cell-phone data. No other district or zone or other special
419 constants were added to the model specification.

420 **TABLE 2 PERCENT DIFFERENCE BETWEEN MODELED AND OBSERVED (CELL**
421 **PHONE DATA) LONG DISTANCE TRIPS WITHIN TENNESSEE**

Origin districts	Destination districts								Total
	Tri-Cities	Knoxville	Chattanooga	Cookeville	Lynchburg	Nashville	Jackson	Memphis	
Tri-Cities	0.8%	0.3%	0.0%	0.0%	0.0%	-0.3%	0.0%	0.0%	0.7%
Knoxville	0.6%	2.0%	0.4%	-0.1%	0.0%	-1.5%	-0.3%	-0.3%	0.9%
Chattanooga	0.0%	0.1%	0.7%	0.3%	0.2%	-0.5%	-0.2%	-0.2%	0.4%
Cookeville	0.0%	-0.1%	0.4%	0.1%	0.0%	-0.6%	-0.1%	-0.1%	-0.3%
Lynchburg	0.0%	-0.1%	0.2%	0.0%	0.1%	-1.3%	0.1%	0.0%	-0.9%
Nashville	-0.3%	-1.4%	-0.2%	-0.9%	-1.5%	6.1%	-0.3%	-1.2%	0.4%
Jackson	0.0%	-0.3%	-0.1%	-0.1%	0.1%	-0.3%	-0.1%	0.8%	0.0%
Memphis	0.0%	-0.3%	-0.2%	-0.1%	0.0%	-1.1%	0.3%	0.1%	-1.2%
Total: All	1.0%	0.3%	1.1%	-0.8%	-1.0%	0.5%	-0.5%	-0.7%	0.0%

422
423 The same districts within Tennessee used for trip generation comparisons were also used
424 to help evaluate the distribution of long distance trips within Tennessee. As can be seen in Table
425 2, the model is able to achieve very good agreement with the observed pattern of long distance
426 trips within Tennessee. The total modeled trips to and from each district are within 1.5% of
427 observed trips for the district, and with the exception of long-distance trips within the Nashville
428 district, all the modeled district level OD flows are within 2% of the observed flows. The
429 distribution of too many long distance trips within the Nashville district may be a result of the
430 fact that the long distance destination choice models are more driven by distance than travel
431 times, so congestion within the Nashville region may not be deterring as many trips as it should.
432 Alternatively, it may simply reflect the inability of the national model to reproduce the complex
433 long distance commuting patterns of this region given the limited spatial resolution of the
434 national model. Despite this particular issue, the overall agreement between the modeled and
435 observed data is quite good.

436 **TABLE 3 PERCENT DIFFERENCE BETWEEN MODELED AND OBSERVED LONG**
437 **DISTANCE TRIPS TO AND FROM TENNESSEE**

Internal districts	External districts							Total
	Northwest	North Atlantic	Northcentral	Carolinas	Alabama-Gulf	Southwest	Georgia-Florida	
Tri-Cities	0.4%	0.1%	0.8%	3.6%	0.0%	0.2%	0.3%	5.3%
Knoxville	0.5%	-2.6%	-1.2%	-1.7%	-0.7%	0.3%	-2.0%	-7.3%
Chattanooga	0.0%	-0.1%	-0.5%	-0.4%	-1.1%	0.1%	2.7%	0.8%
Cookeville	0.0%	-0.2%	0.9%	-0.3%	-0.1%	-0.1%	-0.2%	0.0%
Lynchburg	-0.4%	0.1%	0.4%	0.0%	0.7%	-0.1%	-0.4%	0.2%
Nashville	-0.7%	-0.3%	6.6%	-0.8%	-3.6%	-2.3%	-2.0%	-3.1%
Jackson	0.0%	0.1%	0.6%	0.0%	0.0%	-1.9%	0.0%	-1.2%
Memphis	0.5%	0.3%	0.8%	0.1%	-0.1%	3.4%	0.3%	5.2%
Total	0.3%	-2.6%	8.4%	0.5%	-4.9%	-0.4%	-1.3%	0.0%

438
439 Trips to and from the eight internal Tennessee districts and seven external districts
440 covering the rest of the country were also evaluated. As shown in Table 3, as with the internal

441 trips, the national model is able to generally reproduce the observed pattern fairly well. The total
442 modeled trips to and from each district are all within 8.5% of observed trips for the district and
443 most are within about 5%, and with the exception of long-distance trips between the Nashville
444 and Northcentral districts, all the modeled district level OD flows are within 4% of the observed
445 flows. The model under-predicts trips to and from the Knoxville region, most likely
446 underestimating the number of trips attracted to the Smoky Mountains (which is the most visited
447 National Park) and associated tourist areas. The model also over-predicts trips between
448 Tennessee and the Northcentral region, but the reason for this is less clear. Even so, the ability
449 of the national model to reproduce the complex pattern of long distance trips to and from the
450 state is quite good.

451 CONCLUSIONS

452 This paper has described the first integration of the national long distance passenger demand
453 model with a statewide travel model and its calibration to cell-phone based OD data for the state
454 of Tennessee, illustrating one of if not the first application of big OD data to statewide modeling.
455 The case demonstrates the ability of the national model to be calibrated to observed data and of
456 an integrated modeling system to produce reasonable runtimes. The case also illustrates the
457 general importance of the processing of cell-phone OD data and hypothesizes an importance bias
458 in cell-phone data towards the detection of long distance trips over shorter ones based on
459 evidence from the Tennessee application. However, the case also illustrates the value of such
460 data through, for instance, its ability to reveal important aspects of long distance travel patterns
461 such as a psychological boundary effect corresponding to the state border in the case of
462 Tennessee. While each state must evaluate the usefulness of various modeling approaches for
463 their own planning and modeling, the case of Tennessee's new statewide model demonstrates
464 that both the new national long distance model and cell-phone OD data can be successfully used
465 and add value to a statewide model.

466 ACKNOWLEDGEMENTS

467 This research was sponsored by the Tennessee Department of Transportation in part through the
468 use of state planning and research (SPR) funds.

469 REFERENCES

- 470 1. NCHRP Report 735: Long-Distance and Rural Travel Transferable Parameters for Statewide Travel
471 Forecasting Models. Cambridge Systematics. Transportation Research Board of the National
472 Academies, Washington, DC, 2012.
- 473 2. Outwater, M., M. Bradley, N. Ferdous, C. Bhat, R. Pendyala, S. Hess, A. Daly, and J. LaMondia.
474 Tour-Based National Model System to Forecast Long-Distance Passenger Travel in the United States.
475 TRB 94th Annual Meeting Compendium of Papers, 2015.
- 476 3. Outwater, M., M. Bradley, N. Ferdous, R. Pendyala, V. Garikapati, C. Bhat, S. Dubey, J. LaMondia,
477 S. Hess, and A. Daly. *Foundational Knowledge to Support a Long-Distance Passenger Travel
478 Demand Modeling Framework*. FHWA, U.S. Department of Transportation, 2015.
- 479 4. Outwater, M., M. Bradley, N. Ferdous, S. Trevino, and H. Lin. *Foundational Knowledge to Support
480 a Long-Distance Passenger Travel Demand Modeling Framework: Implementation Report*. FHWA,
481 U.S. Department of Transportation, 2015.

- 482 5. Bernardin, V., E. Rentz, and B. Grady. TMIP How-To: Create Travelshed TAZs. Travel Model
483 Improvement Program (TMIP) of the Federal Highway Administration (FHWA), Washington, DC,
484 2015.
- 485 6. Bernardin, V. and J. Chen. New Methods for Improving Non-Home-Based Trips in Trip-Based
486 Models. Presented at the 95th Annual Meeting of the Transportation Research Board, Washington,
487 D.C., 2016.
- 488 7. Liu, H., A. Dancyk, R. Brewer and R. Starr. Evaluation of Cell Phone Traffic Data in Minnesota.
489 *Transportation Research Record: Journal of the Transportation Research Board*. No 2086, 2008, pp.
490 1-7.
- 491 8. Calabrese, F., G. Di Lorenzo, L. Liu and C. Ratti. Estimating Origin-Destination Flows using Mobile
492 Phone Location Data. *IEEE Pervasive Computing*. Vol. 10, No. 4, 2011, pp. 36-44.
- 493 9. Gur, Y. J., S. Bekhor, C. Solomon and L. Kheifits. Intercity Person Trip Tables for Nationwide
494 Transportation Planning in Israel Obtained from Massive Cell Phone Data. *Transportation Research*
495 *Record: Journal of the Transportation Research Board*. No 2121, 2009, pp. 145-151.
- 496 10. Lee, R. J., I. N. Sener and J. A. Mullins. An Evaluation of Emerging Data Collection Technologies
497 for Travel Demand Modelling: from Research to Practice. *Transportation Letters: The International*
498 *Journal of Transportation Research*. Vol. 8, No. 4, 2016, pp. 181-193.
- 499 11. Huntsinger, L. F. and R. Donnelly. Reconciliation of Regional Travel Model and Passive Device
500 Tracking Data. Presented at the 93rd Annual Meeting of the Transportation Research Board,
501 Washington, D.C., 2014.
- 502 12. Bindra, S., B. Grady and J. Deshaies. Using Cellphone Origin-Destination Data for Regional Travel
503 Model Validation. Presented at the 15th National TRB Transportation Planning Applications
504 Conference, Atlantic City, NJ, 2015.
- 505 13. Zhang, W., A. Kuppam, V. Livshits and B. King. Evaluation of Cellular-based Travel Data –
506 Experience from Phoenix Metropolitan Region. Presented at the 15th National TRB Transportation
507 Planning Applications Conference, Atlantic City, NJ, 2015.
- 508 14. Milone, R. Preliminary Evaluation of Cellular Origin-Destination Data as a Basis for Forecasting
509 Non-Resident Travel. Presented at the 15th National TRB Transportation Planning Applications
510 Conference, Atlantic City, NJ, 2015.
- 511 15. Calabrese, F., M. Colonna, P. Lovisolo, D. Parata and C. Ratti. Real-time Urban Monitoring using
512 Cell Phones: A Case Study in Rome. *IEEE Transactions on Intelligent Transportation Systems*. Vol.
513 12. No. 1, 2011, pp. 141-151.
- 514 16. Wang, P., T. Hunter, A. M. Bayen, K. Schechtner and M. C. Gonzalez. Understanding Road Usage
515 Patterns in Urban Areas. *Scientific Reports*. Vol. 2, No. 1001, 2012.
- 516 17. Wang, J., D. Wei, K. He, H. Gong and P. Wang. Encapsulating Urban Traffic Rhythms into Road
517 Networks. *Scientific Reports*. Vol. 4, No. 4141, 2014.
- 518 18. Alexander, L., S. Jiang, M. Murga and M. Gonzalez. Origin-Destination Trips by Purpose and Time
519 of Day Inferred from Mobile Phone Data. *Transportation Research Part C: Emerging Technologies*,
520 Vol. 58, 2015, pp. 240-250.
- 521 19. Ma, J., F. Yuan, C. Joshi, H. Li and T. Bauer. A New Framework for Development of Time-Varying
522 O-D Matrices based on Cellular Phone Data. Presented at the 4th TRB Innovations in Travel
523 Modeling Conference, Tampa, FL, 2012.

- 524 20. Iqbal, M. S., C. F. Choudhury, P. Wang, M. Gonzalez. Development of Origin-Destination Matrices
525 using Mobile Phone Call Data. *Transportation Research Part C: Emerging Technologies*, Vol. 40,
526 2014, pp. 63-74.
- 527 21. Toole, J. L., S. Colak, B. Sturt, L. P. Alexander, A. Evsukoff, M.C. Gonzalez. The Path Most
528 Traveled: Travel Demand Estimation using Big Data Resources. *Transportation Research Part C:
529 Emerging Technologies*, Vol. 58, 2015, pp. 162-177.
- 530 22. Marzano, V., Papola, A., Simonelli, F. Limits and Perspectives of Effective O-D Matrix Correction
531 using Traffic Counts. *Transportation Research Part C: Emerging Technologies*, Vol. 17, 2009, pp.
532 120-132.
- 533