

1 **EM-EVAL: AN EVALUATION PROCEDURE FOR SMARTPHONE-BASED TRAVEL**  
2 **DIARIES**

3

4

5

6 **K. Shankari**

7 shankari@eecs.berkeley.edu

8

9 **Hari Venugopalan**

10 hvenugopalan@ucdavis.edu

11

12 **David E. Culler**

13 culler@berkeley.edu

14

15 **Randy H. Katz**

16 randy@cs.berkeley.edu

17

18

19 Word Count: 6476 words + 1 table(s)  $\times$  250 = 6726 words

20

21

22

23

24

25

26 Submission Date: August 2, 2019

1

## Abstract

2 As smartphone-based Human Mobility Systems (HMSes) shift from informing individuals to influ-  
3 encing public asset allocation, it is critical to assess how well they perform. Evaluation techniques  
4 of HMSes have typically been ad-hoc, limited to single measurements and ignored power-accuracy  
5 tradeoffs. These tradeoffs are important due to built-in and context-sensitive sensing in the phones  
6 - e.g. "Find my iPhone" or "Doze mode".

7 Treating HMSes as physical instruments with inherent physical noise, we propose a novel  
8 evaluation procedure that uses artificial trips and multiple parallel phones to provide controlled,  
9 repeated inputs to the HMSes under test. Artificial trips mitigate privacy concerns and allow re-  
10 peatability while efficiently exploring a wide variety of trip contexts. Parallel sensing with control  
11 phones mitigates the effects of context sensitive power consumption and inherent sensing error.

12 We use this procedure to create three artificial timelines with 15 different modes, including  
13 ebike and scooter. We travel these three timelines three times each to collect and publish a dataset  
14 with over 500 hours of ground truthed data. We use this data to explore the tradeoffs for power  
15 versus spatial and temporal accuracy.

16 Our results show the benefits of thinking statistically about HMSes. They control for out-  
17 liers, revealing meaningful signals about the behavior of smartphone virtual sensors that are rele-  
18 vant to instrumenting human travel behavior. If adopted widely by the community, the resulting  
19 ground truthed, tradeoff-aware, public datasets can form the basis for additional HMS optimiza-  
20 tions.

1	<b>Contents</b>	
2	<b>1 Introduction</b>	<b>2</b>
3	<b>2 Controlled Evaluation of context-sensitive behavior</b>	<b>3</b>
4	2.1 Artificial timeline . . . . .	3
5	2.2 Control phones . . . . .	4
6	<b>3 Discussion of alternative procedures</b>	<b>4</b>
7	3.1 No artificial trips . . . . .	6
8	3.2 No control . . . . .	6
9	<b>4 Related work</b>	<b>8</b>
10	4.1 Context sensitive sensing algorithms: power without accuracy . . . . .	8
11	4.2 Travel diary systems: compare to manual surveys . . . . .	8
12	4.3 Mode inference: accuracy without power, non-uniform data . . . . .	8
13	<b>5 Evaluation system and experiment design</b>	<b>8</b>
14	5.1 System overview . . . . .	9
15	5.2 System iterations and lessons learned . . . . .	11
16	<b>6 Evaluation</b>	<b>11</b>
17	6.1 Experiment design . . . . .	11
18	6.2 Spatial accuracy . . . . .	14
19	6.3 Motion activity accuracy . . . . .	14
20	<b>7 Conclusion and call to action</b>	<b>14</b>

## 1 INTRODUCTION

2 Inspired by the popularity of smartphone-based personal fitness tracking, the transportation com-  
3 munity aims to build Human Mobility Systems (HMSes) that can automatically track and classify  
4 multi-modal travel patterns. Such systems can replace expensive and infrequent travel surveys with  
5 long-term, largely passive data collection augmented with intermittent surveys focused on percep-  
6 tual data. Such data collection can capture changes in travel behavior as they occur and open new  
7 avenues for responsive urban planning.

8 While there has been much work on building HMSes, both in academia and in industry, the  
9 procedure to *evaluate* them has largely been an afterthought. Careful evaluations are critical as we  
10 move from the personal to the societal domain. Users who make decisions based on self-tracking  
11 have an intuition of its accuracy based on their experienced ground truth. The decisions are low-  
12 stakes lifestyle changes, which may be personally meaningful, but are not societally contentious.  
13 However, a Metropolitan Transportation Agency picking projects and allocating millions of dollars  
14 in funding needs to know the accuracy of the data before making its decisions (18).

15 The typical HMS evaluation procedure (e.g. Quantified Traveler (7), prior versions of our  
16 work (8, 17)) is ad-hoc and also functions as a pilot - a small ( $\approx 3$ -12) set of the author's friends  
17 and family are recruited to install the app component of the HMS on their phones, and go about  
18 their daily life for a few days or weeks while annotating the trips with "ground truth". The ground  
19 truth annotation can either directly happen on the app, or through a recap at the end of the day.  
20 Conscientious researchers may ensure that the set of evaluators are demographically diverse, in an  
21 attempt to evaluate against a richer set of travel patterns.

22 While this procedure imposes little additional researcher burden, it conflates the *experi-*  
23 *mental* procedure (understanding human travel behavior in the wild) with the *evaluation* procedure  
24 (evaluating the instrument that will measure the human travel behavior). The first is trying to  
25 understand *behavior*, so it needs human diversity. The second is trying to understand *sensing*  
26 parameters, so it needs diversity of *trip types*. The human functions as a phone transportation  
27 mechanism during evaluation and could be profitably replaced with a self-navigating robot if one  
28 was available.

29 An analogy with classic physical measurements may be useful. Consider the situation in  
30 which a researcher wants to collect data on the weight distribution of the population in a particular  
31 region. Since there are currently no certifying bodies for travel diaries, let us pretend that she  
32 cannot purchase a pre-certified scale. How would she evaluate the available scales before starting  
33 her experiment?

34 The analog to ad-hoc evaluation procedure would involve recruiting several of her friends  
35 and family to weigh themselves on the scales and compare the reported weight with their true  
36 weight. This analogy clearly reveals some of the limitations of the ad-hoc procedure: (i) How does  
37 she trust that the self-reported weights are "true"? (ii) If all her friends are adults weighing 55kg  
38 - 75kg, how does she know how the scales perform outside that range? She can overcome the  
39 range limitations by recruiting a broader set of testers, e.g. through an intercept survey. However,  
40 that modification makes the ground truth limitation worse, since it is less likely that strangers will  
41 reveal their true weight. A further modification might pay contributors to improve the self-reported  
42 accuracy, but at this point, she is essentially running the experiment.

43 A more robust evaluation procedure would involve choosing known weights across a broad  
44 range (e.g. 0kg to 300kg in 10kg increments) and comparing them to the reported weights. Since  
45 no instrument is perfect, there is likely to be some variation in the values reported. She would

1 likely repeat the experiments multiple times in order to establish error bounds.

2 HMS evaluation procedures need to be more sophisticated than simple physical measure-  
3 ments since: (i) their operation is based on prior behavior (e.g. HMS duty cycling, android doze  
4 mode) and the potential for feedback loops makes it important to control the *sequence* of evalua-  
5 tion operations, (ii) unlike a physical scale, which has a fixed one-time cost, they have an ongoing,  
6 variable cost in terms of battery drain, so the evaluation must assess the power/accuracy tradeoff,  
7 and (iii) unlike scalar weight data, HMSes generate strongly correlated timeseries data, which is  
8 extremely hard to deanonymize.

9 The main contributions in this paper are:

- 10 1. We propose an **evaluation procedure** for HMSes based on pre-defined, ground truthed,  
11 artificial trips and outline how it addresses the above challenges
- 12 2. We describe the design of a cross-platform **evaluation system** that can be used to per-  
13 form such evaluations reproducibly and publish the results.
- 14 3. We use this system to evaluate the **power/accuracy tradeoffs** of the android and iOS  
15 location, motion accuracy and visit detection virtual sensors in the San Francisco Bay  
16 Area.

17 The rest of this paper is structured as follows. In Section 2 we outline an experiment proce-  
18 dure that can control for data collection inconsistencies, and discuss some alternative approaches  
19 in Section 3. We place the procedure in the context of prior work in Section 4, and describe the  
20 reference implementation in Section 5. In Section 6 we outline a specific experiment involving  
21 15 different modes, and use it to evaluate the power/accuracy tradeoff results for various virtual  
22 sensors, concluding with Section 7.

## 23 CONTROLLED EVALUATION OF CONTEXT-SENSITIVE BEHAVIOR

24 As discussed in Section 1, instruments are typically evaluated by repeatedly exposing them to  
25 controlled inputs to determine their error characteristics. In the case of complex systems such as  
26 HMSes, the evaluation needs additional controls for feedback loops, cost/accuracy tradeoffs and  
27 privacy considerations. In this section, we outline EM-EVAL a procedure for HMS evaluation that  
28 addresses these concerns with two techniques that are novel in this domain:

- 29 (i) pre-defined, **artificial trips** that support spatial ground truth, preserve privacy, in-  
30 crease the breadth of *trip types* and support repetitions for establishing error bounds,  
31 and
- 32 (ii) power and accuracy **control phones** carried at the same time as the experimental  
33 phones, that can cancel out context-sensitive variations in power and accuracy.

### 34 Artificial timeline

35 The core of the experimental procedure is the pre-defined specification of a sequence of artificial  
36 *trips*, potentially with multiple *legs* or *sections* per trip. The trajectory and mode of travel is  
37 also pre-defined. The *data collector* completes the timeline trips by strictly following the specified  
38 trajectory and mode while carrying multiple phones that collect data simultaneously using different  
39 configurations.

40 The specified pre-defined trajectories provide spatial ground truth. We do not pre-define  
41 temporal ground truth since it is extremely hard to control for differences in walking speeds, de-  
42 lays due to traffic conditions, etc. We use manual input from the data collector to collect *coarse*  
43 temporal “ground truth” of the transitions along the timeline. We do not use manual input for *fine-*

1 *grained* temporal ground truth along the trajectory because: (i) human response times are too slow  
2 for fine-grained temporal ground truth during motorized transportation, and (ii) distracting the data  
3 collector during active transportation can be risky.

4 Using an artificial timeline addresses several of the unique challenges associated with HMS  
5 evaluation.

6 **Privacy** Since the trips are artificial, they preserve the data collector’s privacy. Even if his  
7 adversaries would download the trips, they would not be able to learn anything about his normal  
8 travel patterns.

9 **Spatial ground truth** Since even high accuracy (GPS-based) data collection has errors,  
10 pre-defining spatial ground truth allows us to resolve discrepancies (Section 3) and compute the  
11 true accuracy.

12 **Breadth and variety of trips** Artificial trips allow efficient exploration of the breadth of  
13 the trip space. For example, the trips could include novel modes such as e-scooters and e-bikes, or  
14 specify different contexts for conventional modes, such as express bus versus city bus.

15 **Repetitions** Since the trips are pre-defined, they can be repeated exactly. This allows us to  
16 use standard variance and outlier detection to estimate error bounds on the measured values.

## 17 **Control phones**

18 The artificial trips give us spatial ground truth, but they do not give us cost (power consumption)  
19 or temporal ground truth.

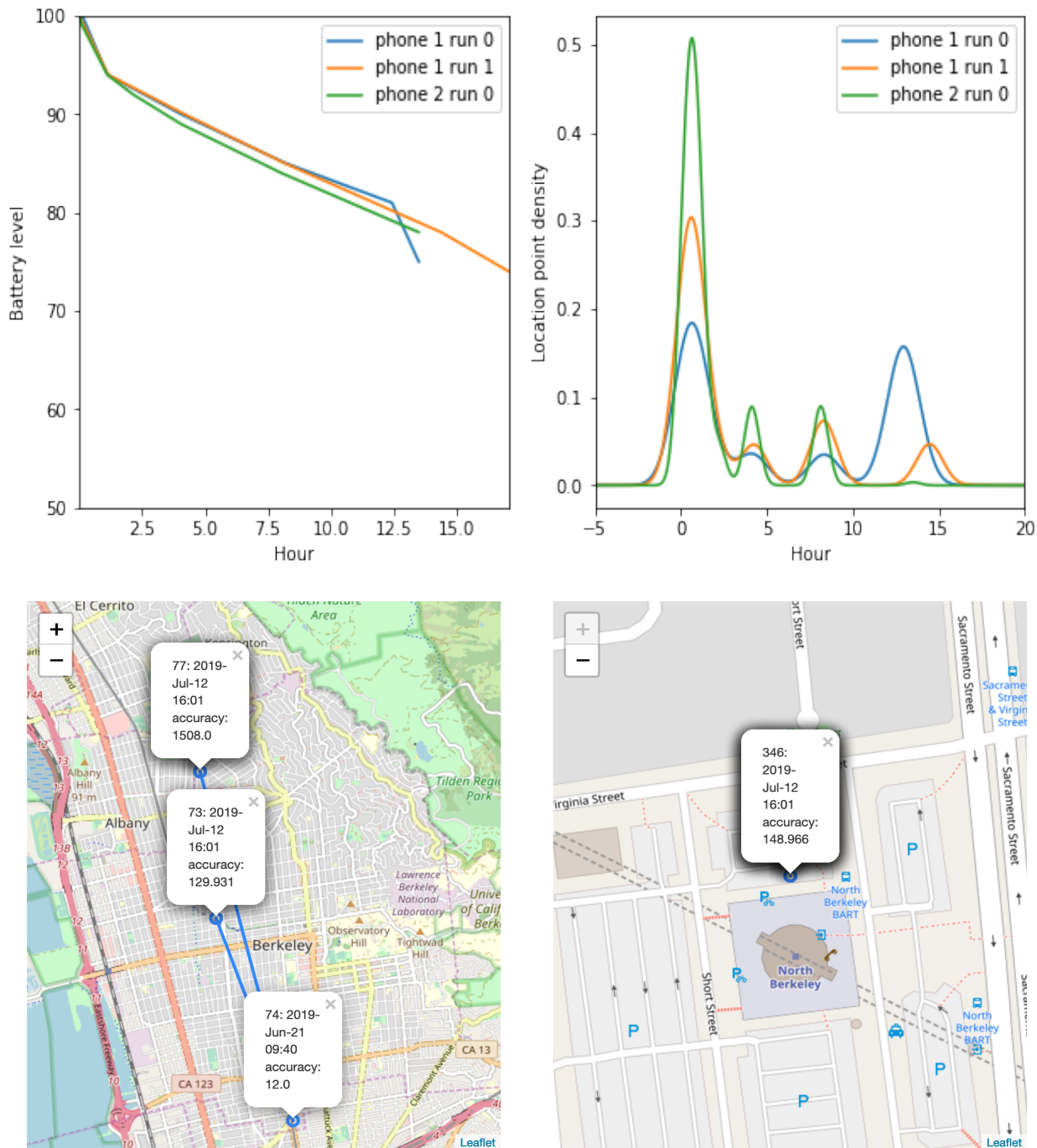
20 We control for the cost through the use of the use of multiple phones, carried at the same  
21 time by the data collector. The phones carried by the data collector are divided into control phones  
22 and experiment phones. The control phones represent the baseline along each of the axes in our  
23 tradeoff and the experiment phones implement a custom sensing regime that is at some intermedi-  
24 ate point. The evaluation procedure allows us to determine those points.

25 **Power** The power control phone captures the baseline power consumption of a phone that  
26 is not being used for tracking by a HMS. This does not mean that the phone is idle - phone OSes  
27 (e.g. iOS or android) are complex, context-sensitive systems that perform their own location track-  
28 ing (e.g. “Find my iPhone”) and their own duty cycling (Figure 1). Using a power control allows  
29 us to identify the **additional power** consumed by the HMS, even if it is context sensitive.

30 **Accuracy** The accuracy control captures the upper bound on the accuracy of a particular  
31 class of smartphones given sensor and OS limitations. While we would like to compare the exper-  
32 imental accuracy to ground truth, (i) all sensors have errors, so ground truth is not achievable in  
33 practice, (ii) artificial trips give us spatial but not temporal ground truth, and (iii) GIS-based trajec-  
34 tory specifications do not have an associated power tradeoff. Using an accuracy control allows us  
35 to compare the experimental data collection against the best achievable data collection, in addition  
36 to the ground truth.

## 37 **DISCUSSION OF ALTERNATIVE PROCEDURES**

38 While EM-EVAL (Section 2) addresses the complexities of HMS evaluation, it also imposes a much  
39 higher researcher burden than the ad-hoc method. This raises the question of whether all these  
40 controls are necessary or merely sufficient. In this section, we discuss some alternative approaches  
41 and highlight the unexpected behavior that they would miss. This list is not comprehensive but  
42 provides a flavor of the arguments without tedious repetition.



**FIGURE 1: Top:** Power variation illustrated by duty cycling on android. All the phones were configured identically, and placed in the same environment. The built-in duty cycling on android switches all phones to low power mode at around 1 hr. However, phone 1, on run 1 alone, switches back to high power mode at around 12.5 hours. Repeating experiments allows us to distinguish the first consistent duty cycle and the second outlier. **Bottom:** Accuracy variation illustrated by mismatched timestamps during trajectory data collection. Both trajectories are collected from identical phones during a subway trip. Point 74 has an accuracy radius of only 12, but its timestamp is in June instead of July! Spatial ground truth allows us to sort out the varying accuracies here.

## 1 **No artificial trips**

2 Creating pre-determined trips requires an upfront investment in effort, and requires the data collec-  
3 tor to take trips just for data collection. An alternative would use multiple phones, but allow data  
4 collectors to go about their regular routines and tag the modes only. We could use the accuracy  
5 control phones to determine the ground truth trajectory.

6 **No privacy** Capturing the data collector's regular routines compromises their privacy. Even  
7 if the data does not include their name or phone number, a list of their commonly visited places  
8 and trips can form a unique fingerprint that can uniquely identify them (2, 19). This sensitivity  
9 precludes evaluation data from being published and used for reproducible research.

10 **No repetition** The behavior of the same phone with the same configuration can vary over  
11 time, both for power and for accuracy (Figure 1). Repeating the same trip multiple times allows us  
12 to detect and remove outliers. With ad-hoc trips, it is unclear whether any difference in behavior  
13 is real or caused by context-sensitive variation. And without pre-determined trips, it is challenging  
14 to repeat the same trips and trajectories over time.

15 **No spatial ground truth** No sensor is perfect and even the accuracy control phones can  
16 have sensing errors. If we see a divergence between an experiment phone and the accuracy control,  
17 it is unclear which one has the error (Figure 1).

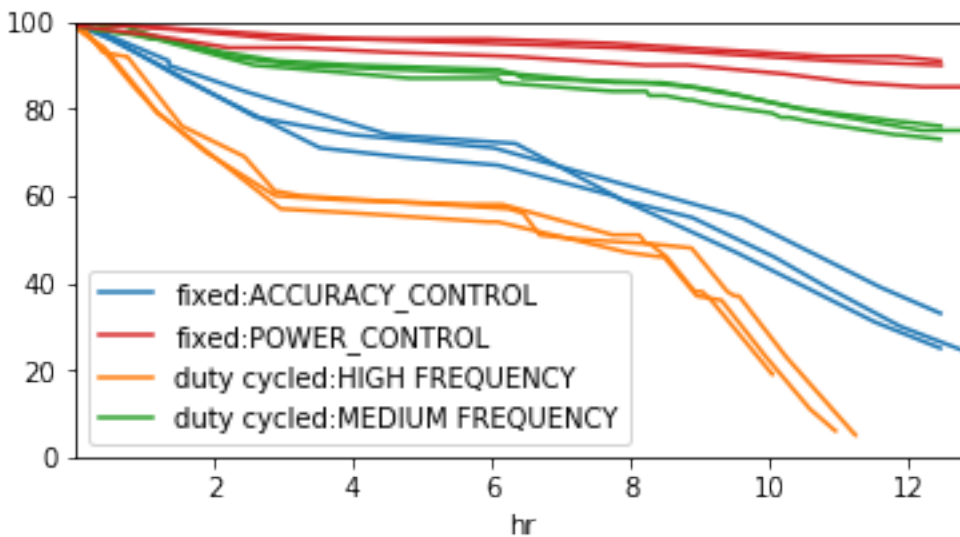
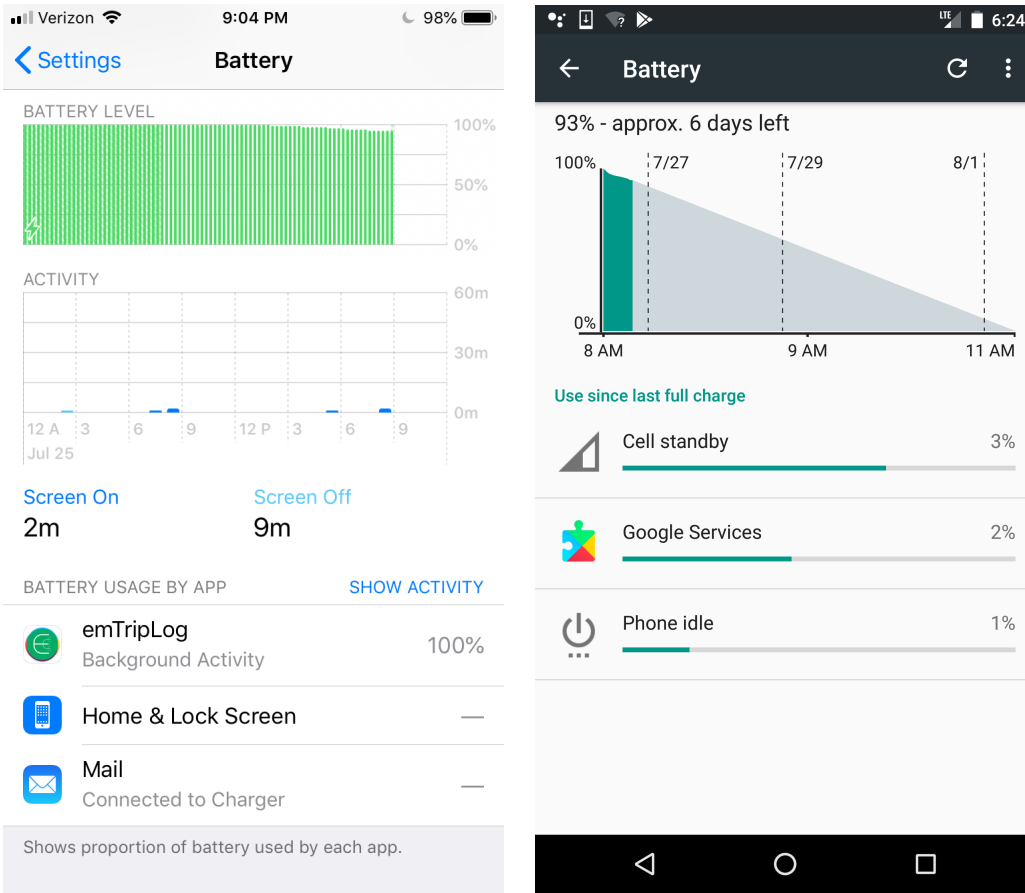
## 18 **No control**

19 Using control phones requires the researcher to purchase multiple phones of the same make, model  
20 and approximate age. While used smartphones are relatively cheap (USD 50 - USD 100), 4 android  
21 phone and 4 iPhones combined will still cost USD 400 - USD 800. An alternative would be to  
22 use one phone each for each OS, perform the timeline trips, and look at the app-based power  
23 consumption reported by the phone OS?

24 **Sensor access attribution and the meaning of the %** Sensor access in modern phone OSes  
25 (android and iOS) is also context-sensitive, making it unclear how it is counted for per-app con-  
26 sumption. For example, if multiple apps request a sensor reading, the OS delays returning a result  
27 until it can batch related requests and serve all of them with a single sensor access. This is why  
28 the OSes treat the sensing frequency as a hint instead of a guarantee. Second, if sensor access  
29 is mediated by a service (e.g. fused location in Google Play Services), it is unclear whether the  
30 sensor access is counted for the service or the app (Figure 2). And finally, although android reports  
31 per app consumption as a % of the battery *capacity*, iOS does so as a % of the battery *consump-*  
32 *tion*. This indicates that on dedicated phones, the HMS under test will always show close to 100%,  
33 whether it is the power control or the accuracy control (Figure 2). Using a control phone for the  
34 power will cancel out these context sensitive effects and estimate the difference in power drain  
35 with and without the HMS app component installed.

36 **Custom duty cycling increases power drain** Sensing is not the only source of power con-  
37 sumption - CPU usage can also have a significant impact on power usage. HMSes can use smart  
38 local processing to reduce local sensing, but the increased power consumption from the CPU can  
39 cancel out the savings from the sensing. Including an accuracy control showed that the basic duty  
40 cycling algorithm in our experiment paid for itself in low frequency sensing but actually increased  
41 power usage for high frequency sensing (Figure 2).





**FIGURE 2: Top left:** iOS power control phone with the sensing app consuming 100% but of a power drain of only 2% over the entire day. **Top right:** android phone showing Google Play services as a separate power consumer. **Bottom:** Explicit duty cycling causes increased power drain at high frequencies, possibly due to greater CPU power consumption. Note that the battery drain flattens out on all curves during the middle, stationary part. The main difference is in the rate of power drain while moving: low frequency sensing consumes the least power, and the checks for android’s built-in duty cycling appear to be optimized to be more efficient than our simple implementation.

## 1 RELATED WORK

2 Human Mobility Systems (HMSes) are complex to evaluate (Section 3). Other researchers have  
3 identified similar challenges as part of survey papers (e.g. phone context, privacy, learning, scal-  
4 ing (9), varying metrics and time scales across research areas (13)). However, to the best of our  
5 knowledge, there is no proposed solution that addresses all of them.

6 Papers related to instrumenting travel behavior fall into three main categories; we list some  
7 work from each as an example. A comprehensive classification of papers into categories is beyond  
8 the scope of this paper.

### 9 **Context sensitive sensing algorithms: power without accuracy**

10 This research area focuses on context sensitive, adaptive power management of sensors. Papers  
11 such as ACE (11) and Jigsaw (10) compare their power requirements to naive sensing techniques.  
12 However, their accuracy evaluations focus on the localization error (Jigsaw), or comparison to  
13 naive inferred results<sup>1</sup> (ACE).

### 14 **Travel diary systems: compare to manual surveys**

15 There is a vast variety of one-off travel diary systems that combine smartphone based sensing with  
16 cloud-based processing to generate travel diaries. Systems such as such as Data Mobile (12), Fu-  
17 ture Mobility Study (FMS) (1, 6) and rMove (5) aim to replace the paper and telephone based  
18 Household Travel Surveys with smartphone and cloud based systems. So they evaluate the ac-  
19 curacy of their systems against the traditional methods, not against ground truth. This can show  
20 that smartphone based methods are significantly better than traditional methods, but not provide a  
21 quantitative estimate of the accuracy of their system. Similarly, they do not include quantitative  
22 power evaluations - preferring statements like "Among the three types of discrepancies, the second  
23 type, data gap due to battery drainage, was most frequently observed." (6) or "The battery con-  
24 sumption test was simply whether, under regular usage, the phone could make it through the day  
25 without having to be charged." (12). So they do not rigorously evaluate either the power or the  
26 accuracy side of the tradeoff.

### 27 **Mode inference: accuracy without power, non-uniform data**

28 Mode inference of travel mode based on sensor data is an extremely popular subject in the litera-  
29 ture<sup>2</sup>. Researchers have used decision trees (14, 20), Hidden Markov Models (15, 21), and neural  
30 networks (3, 4) to distinguish between various subsets of travel modes. However, although the  
31 inference algorithms are different, most such papers use similar methods for evaluating their accu-  
32 racy. They typically recruit a small sample of their friends (e.g. 16 users over one day (14), 4 users  
33 over two weeks (21)) to collect naturalistic data along with annotations of the ground truth. The  
34 data collection focuses on the sensors used for analysis and omits the battery. This kind of evalua-  
35 tion does not meet the any of the requirements outlined above, except privacy, which is addressed  
36 by not publishing the dataset.

## 37 EVALUATION SYSTEM AND EXPERIMENT DESIGN

38 The EM-EVAL procedure (Section 2) allows us to estimate the power/accuracy tradeoff of various  
39 sensing configurations used in Human Mobility Systems (HMSes). One of the novel components

---

<sup>1</sup>e.g. based on speed

<sup>2</sup>Probably because it is hard, and nobody has really solved it well yet for modes other than walking

1 of the procedure involves the specification of pre-defined, artificial trips with ground truthed tra-  
2 jectories and modes.

3 This section explores the nuances of implementing such a procedure. We first describe a  
4 publicly available reference implementation of a system - EM-EVAL-ZEPHYR - that can be used  
5 perform this procedure. We then discuss challenges encountered while using the system to perform  
6 an experiment in the San Francisco Bay Area (Section 6). Some of these challenges were addressed  
7 by system improvements, while others are documented as best practices for future data collectors.

## 8 **System overview**

9 EM-EVAL is a generic *procedure* for HMS evaluation - it does not actually collect any data. To  
10 use it, we need a concrete system that configures data collection based on the spec configurations,  
11 collects coarse temporal ground truth, periodically reads battery levels and stores data for future  
12 analysis.

13 As part of our evaluation (Section 6), we built a system EM-EVAL-ZEPHYR that combines  
14 our prior work on power evaluations (16) with our existing HMS platform (8) and supports per-  
15 forming the EM-EVAL procedure. The system consists of three main parts:

16 **Evaluation Specification** The *spec* describes an evaluation that has been performed or  
17 will be performed in the future. In addition to mode and trajectory ground truth, it includes the  
18 app configurations to be compared and the mapping from phones to evaluation *roles*. The spec  
19 automatically configures both the data collection app and the standard analysis modules.

20 To reduce evaluator burden, we provide pre-processing functions to fill in trajectory in-  
21 formation based on route waypoints for road trips and OSM relations for public transit. We also  
22 provide sample notebooks to verify timelines and their components before finalizing and uploading  
23 the evaluation spec.

24 **Auto-configured Smartphone App** We have generated a custom UI skin for our E-MISSION  
25 platform (8) that is focused on evaluation. It allows evaluators to select the current spec from the  
26 public datastore, and automatically downloads the potential comparisons to be evaluated, the role  
27 mappings and the timeline.

28 Since the e-mission platform data collection settings are configurable through the UI, the  
29 sensing configurations defined in the spec are automatically applied based on the phone role when  
30 the data collector starts an experiment. For example, when starting an experiment to compare high  
31 accuracy (HAHFDC) versus medium accuracy (MAHFDC) data collection, the second experiment  
32 phone will automatically be set to MAHFDC settings. Finally, when the data collector performs  
33 the trips, he marks the transition ground truth in the UI, and the app automatically displays the next  
34 step in the timeline (Figure 3).

35 **Public Data + Sample Access Modules** Since there are no privacy constraints, EM-EVAL-  
36 ZEPHYR uploads all collected data to a public instance of the E-MISSION server. The associated  
37 repository contains sample notebooks that can download, visualize and evaluate the data associated  
38 with a particular spec. All the data used in this paper is publicly available, and the notebooks can  
39 be manipulated interactively using binder.

40 Note that although the EM-EVAL **procedure** is general, the current implementation of the  
41 EM-EVAL-ZEPHYR **system** is integrated only with the E-MISSION platform. Using the procedure  
42 with other HMSes will require re-implementing the EM-EVAL procedure with the other HMS, or  
43 using a combination of systems for the evaluation. For example, EM-EVAL-ZEPHYR can still read  
44 the battery level periodically, display the trip sequence to the data collector, and be used to mark



**FIGURE 3:** Top: Spec components in EM-EVAL-ZEPHYR include configuration, timeline and trip details. Bottom: Sample spec for a multi-modal trip, including transfers and waits for public transit.

1 the transition ground truths. However, the evaluator needs to configure the settings for the app  
2 being tested manually, and to download, clean and analyse the resulting data.

### 3 **System iterations and lessons learned**

4 As we started collecting data, we had to resolve some ambiguities around exactly when the transi-  
5 tion ground truth should be collected. We also discovered best practices that increased the likeli-  
6 hood of successful data collection. This section outlines these lessons learned.

7 **System change: capture transition complexity** One of the big promises of using HMSes  
8 for instrumenting human travel is that we don't have to focus only on the primary mode. Instead,  
9 with fine-grained data collection, we can understand the full complexity of end to end travel.

10 In fact, the only true unimodal trips are walking trips. Everything else is multi-modal.  
11 Thus, a significant change to the system was to restore the hidden complexity that is elided from  
12 user descriptions of travel diaries. For example, consider the trip description "Drive from Mountain  
13 View Library to Los Altos Library". Although that appears to be a unimodal trip, it is actually a  
14 multi-modal trip which involves implicit walk access sections to and from the car at the source and  
15 destination respectively.

16 It is not possible to pre-determine the ground truth for these walk access sections since  
17 we cannot control which parking spaces are available when we perform the trip. We address such  
18 issues by adding *shim sections*, and expanding the start and end from points to 100m polygons. We  
19 can then relax the constraints around ground truth within the polygon by only using the reference  
20 dataset, but still check the accuracy of the mode inference (Figure 3).

21 **Best practice: Pilots are critical** In spite of reviewing the pre-determined trajectories ahead  
22 of time as part of the validation process, and also having them displayed on the EM-EVAL-ZEPHYR  
23 UI, we found that we frequently made small mistakes, during the first round of data collection for  
24 a new timeline. Sometimes, we found that the route suggested by OSRM felt unsafe to bicycle  
25 on and we had to accordingly change our specification. The second repetition generally resolved  
26 these issues. In order to avoid a stressful data collection experience, we suggest running through a  
27 new timeline with a trial run before starting full-featured data collection.

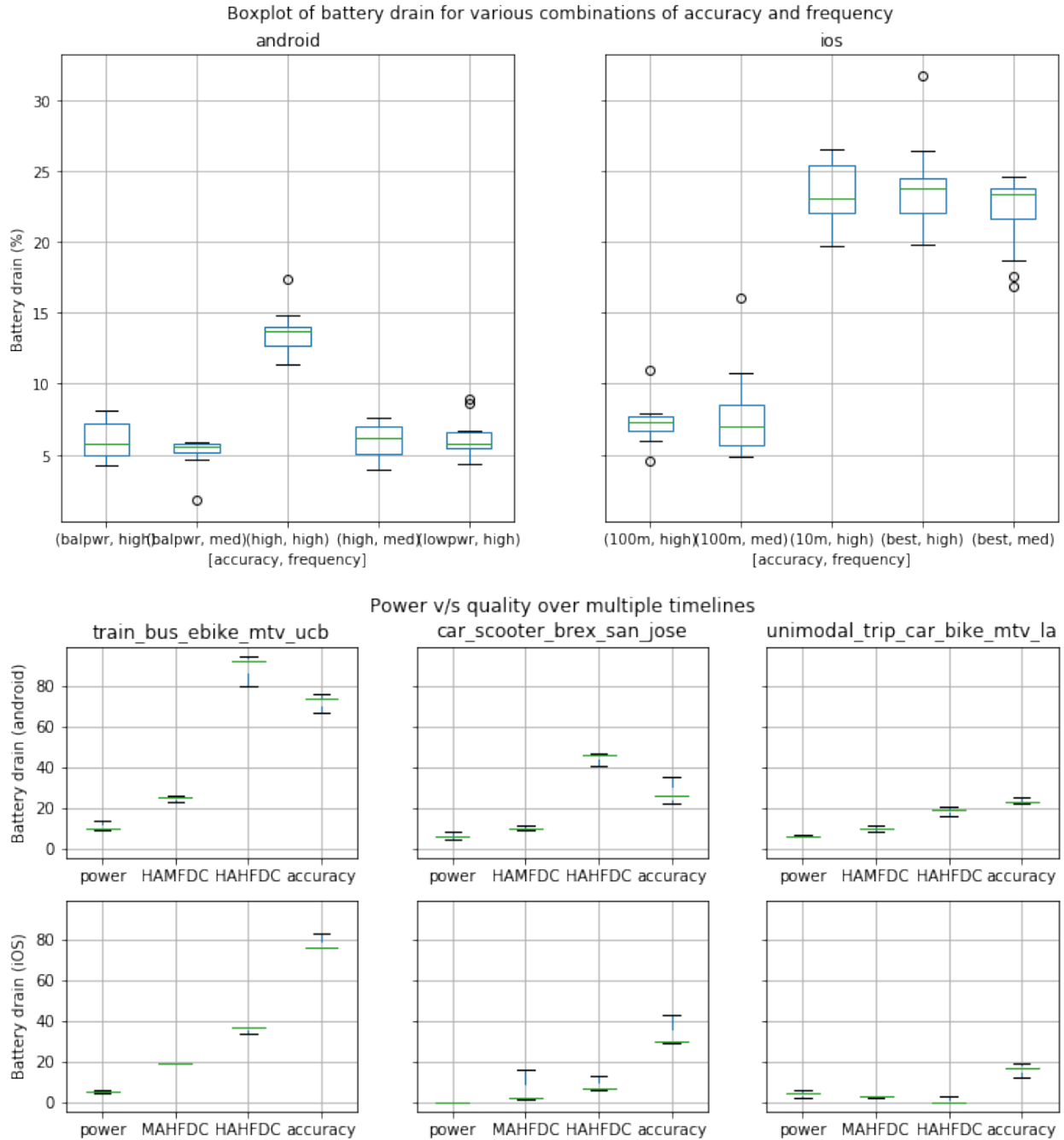
28 **Best practice: Mindfulness** Remembering to mark the transition ground truths was one  
29 of the hardest parts of the ongoing data collection and really highlights the challenges of ground  
30 truth collection. In spite of the fact that she was performing artificial trips to collect data for her  
31 own project, one of the authors forgot to mark wait -> move transitions during the pilot for the long  
32 multi-modal timeline because she had started checking her email while waiting. It is important to  
33 be present in the moment and pay attention to the context while collecting data.

## 34 **EVALUATION**

### 35 **Experiment design**

36 EM-EVAL is a generic evaluation procedure and can be used with any kind of HMS. EM-EVAL-  
37 ZEPHYR is a reference implementation of EM-EVAL that used to evaluate many experimental set-  
38 tings relevant to HMSes over any set of trajectories. The evaluator can pick her settings based on  
39 her research goals. In this section, we outline our goals for this evaluation, and use them to define  
40 three timelines that cover 15 separate modes, including recently popular modes such as scooter  
41 and ebike.

42 **Dwell time** Instead of focusing only on trips, we wanted to evaluate a timeline that in-  
43 cluded significant dwell time. We could see from our calibration runs that android appears to have



**FIGURE 4:** Power calibration and evaluation. **Calibration:** All phones were stationary, configured identically and had the app-based sensing turned always on. Although there are outliers, repeated experiments allow us to see clear behavior differences. High accuracy, high frequency (HAHF) is the only high power configuration on android, lowering either the accuracy or the frequency leads to essentially the same power drain. On iOS, the frequency does not matter at all, the only way to reduce the power drain is to lower the accuracy. **Evaluation:** Based on these calibration results, our experiment design compares high v/s med frequency on android and high v/s low accuracy on iOS. Given the calibration results, we would expect the drain to increase with quality, and we do see that behavior on the longest timeline for iOS and the shortest timeline for android. On shorter timelines for iOS, the differences are small and fall within our current error bounds. On longer timelines for android, the CPU drain of our basic duty cycling overwhelms any improvements from duty cycling, so the HAHFDC performs worse than the accuracy control.

id	Description	Outgoing trip modes	Incoming trip modes	Dwell time	Overall time
unimodal trip car bike mtv la	Suburban round trip	car (city streets)	bike	1 hrs	3 hrs
car scooter brex san jose	Downtown library	car (freeway)	escooter Bus Rapid Transit	3 hrs	5.5 hrs
train bus ebike mtv ucb	Multi-modal trip across the bay	commuter train subway city bus	ebike (shared) express bus downtown walk light rail commuter rail	6 hrs	12.5 hrs

**TABLE 1:** Brief description of timelines, covered modes, dwell times and overall times

1 built-in duty cycling and including significant dwell time would allow us to capture the impact  
 2 of this context sensitive behavior. Therefore, we structured our timeline trips as round trips to  
 3 libraries with an intermediate dwell time  $\approx 3 \times$  the mean travel time to the location.

4 **Broad range of modes** HMS evaluations should cover a broad spectrum of trip types, and  
 5 since we are creating artificial trips, we can structure them to maximize mode variety. In order to  
 6 efficiently cover this space, we tried to ensure that no mode was repeated. We only had to include  
 7 commuter rail twice since there were few other transit options to reach the starting point chosen.

8 **Multi-modal transfers** Detecting multi-modal transfers in a HMS is tricky because there  
 9 isn't a clear signal similar to a trip end. We ensure that there are many transition examples by  
 10 emphasizing multi-modal transfers.

11 With those goals in mind, we decided on three artificial timelines of varying lengths that  
 12 cover a total of 15 separate modes. We chose each timeline to be round trips to libraries so as to  
 13 not include identifiable location data (e.g. home location) in our experiments. A description of  
 14 each timeline with the associated modes and dwell times is given in Table 1.

15 Since this paper focuses on the evaluation procedure, to avoid bias, we do not use it to eval-  
 16 uate any particular HMS. Instead, we use it to evaluate the virtual sensors provided by smartphones  
 17 themselves. These virtual sensors are exposed by smartphone operating systems (OSes) by com-  
 18 bining underlying physical sensors using proprietary algorithms. iOS already restricts developer  
 19 access to the GPS, instead supporting a virtual location sensor that chooses underlying sources  
 20 based on developer-supplied accuracy constants. As phone OSes impose greater background re-  
 21 strictions on apps, we can only expect the use of such sensors to grow. In this section, we analyze  
 22 spatial accuracy and motion activity accuracy from our collected data, and draw inferences from  
 23 them to demonstrate the benefits that can be gained by using a procedure such as EM-EVAL.

## 1 **Spatial accuracy**

2 Box plots of the distribution of spatial error of the measured locations against the groundtruth in the  
3 different timelines are shown on the top row of Figure 5. It is interesting to note the skewed nature  
4 of the plots, and the presence of different outliers, whose values seem to vary across phones across  
5 timelines. Ad-hoc evaluation schemes would probably identify them by comparing error values  
6 against a threshold. Choosing this threshold without the knowledge of the distribution of error  
7 values would be difficult and introduce ambiguity to the evaluation. Our statistical and systematic  
8 approach eliminates these difficulties.

## 9 **Motion activity accuracy**

10 The distribution of temporal error between the ground truth section transitions and sensed activity  
11 transitions is shown in the bottom row of Figure 5. The plot shows that despite the presence  
12 of outliers, the quality of sensing does not seem to impact any difference in accuracy. Thus, a  
13 study focusing on places and modes but not trajectories could choose a lower quality sensing  
14 configuration and reduce power drain.

## 15 **CONCLUSION AND CALL TO ACTION**

16 Human Mobility Systems (HMSes) are complex software systems that run on equally complex  
17 smartphone operating systems (OSes). This complexity implies that there is rarely a simple linear  
18 relationship between their inputs and outputs, which complicates their evaluation.

19 We outline a procedure, based on repeated travel over pre-defined artificial *timelines* car-  
20 rying experiment and control phones, to control this complexity. We show that it can control for  
21 outliers, and also reveal meaningful signals about the behavior of smartphone virtual sensors that  
22 are relevant to instrumenting human travel data.

23 The procedure is privacy-preserving, so it does not need human subjects approval. It fo-  
24 cuses on trip diversity, not demographic diversity, so it can be undertaken by a small research  
25 group, or even a single researcher, as a pre-pilot before recruiting study participants. It uses pre-  
26 determined trips and modes, so it can efficiently explore complex or newly emerging travel pat-  
27 terns and modes, such as scooters. The procedure, and the associated reference implementation  
28 can simplify the testing required before a study is launched.

29 The procedure is not a replacement for a pilot, since it does not capture user experience  
30 or other behavioral factors. However, it can shorten and simplify the pilot, since the pilot does  
31 not need evaluate quantitative “hard” factors. As a general principle, this controlled procedure  
32 evaluates computational features, such as power drain or accuracy, while the pilot assesses “soft”  
33 features such as usability and motivation.

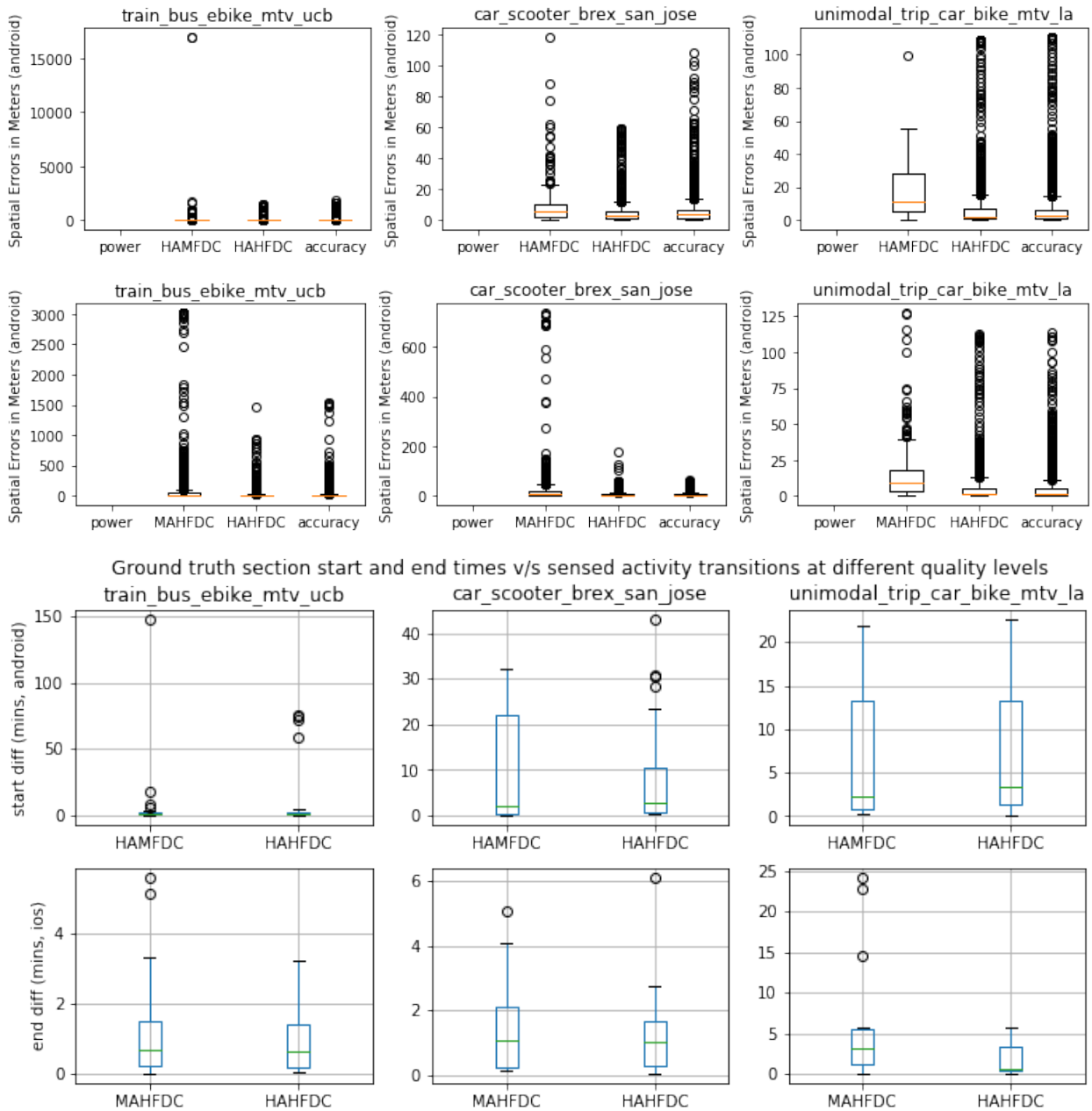
34 The procedure primarily supports direct comparisons between configurations or systems  
35 for the tested contexts. Based on the procedure, we may be able to model individual features of the  
36 HMS or the underlying system (e.g. the rate of power drain while moving). However, since both  
37 hardware specifications and phone OS implementations can and do change<sup>3</sup>, researchers will need  
38 to periodically re-tune their models, which will entail re-running the procedure.

39 Finally, we have used the procedure to evaluate metrics related to sensing accuracy under  
40 various configurations. Post-processing algorithms can potentially compensate for low sensing  
41 accuracy through integration with other sources, such as GIS systems. Large-scale datasets, such as

---

<sup>3</sup>battery optimization android neural network





**FIGURE 5:** Spatio-temporal accuracy evaluations. Top: Distribution of the *spatial* error across the timelines. The medium sensing accuracy is surprisingly close to the high accuracy, although it has many more outliers, specially in the timeline with the underground segments. **Bottom:** Distribution of the *temporal* error between the ground truth section transitions and sensed activity transitions. There is essentially no difference between high and low quality sensing, even for the longest timelines. The vast majority of the transitions occur within 5 mins, although there are significant outliers.

1 ImageNet, and the resulting toolchains, have greatly accelerated the development of algorithms in  
2 other domains. If HMS reseachers adopt this procedure and publish resulting data, the aggregation  
3 over multiple evaluations can generate a large scale, ground truthed, public dataset. This dataset  
4 can accelerate the development of analysis algorithms, such as mode and purpose inference, for  
5 HMSes.

## 1 REFERENCES

- 2 1. C. D. Cottrill, F. C. Pereira, F. Zhao, I. F. Dias, H. B. Lim, M. Ben-Akiva, and P. C. Zegras.  
3 The Future Mobility Survey: Experiences in developing a smartphone-based travel survey  
4 in Singapore. *Transportation Research Record: Journal of the Transportation Research*  
5 *Board*, (2354):59–67, 2013.
- 6 2. Y.-A. De Montjoye, C. A. Hidalgo, M. Verleysen, and V. D. Blondel. Unique in the crowd:  
7 The privacy bounds of human mobility. *Scientific reports*, 3:1376, 2013.
- 8 3. Fei Yang, Zhenxing Yao, and Peter Jin. Multi-mode Trip Information Recognition Based  
9 on Wavelet Transform Modulus Algorithm by Using GPS and Acceleration Data. In *TRB*  
10 *94th Annual Meeting Compendium of Papers*, Washington, DC, Jan. 2015.
- 11 4. P. Gonzalez, J. Weinstein, S. Barbeau, M. Labrador, P. Winters, N. Georggi, and R. Perez.  
12 Automating mode detection for travel behaviour analysis by using global positioning  
13 systems-enabled mobile phones and neural networks. *IET Intelligent Transport Systems*,  
14 4(1):37, 2010.
- 15 5. E. Greene, L. Flake, K. Hathaway, and M. Geilich. A seven-day smartphone-based gps  
16 household travel survey in indiana. Washington, D.C, Jan. 2016. Transportation Research  
17 Board.
- 18 6. Hongik University, J. S. Lee, P. C. Zegras, F. Zhao, D. Kim, and J. Kang. Testing the  
19 Reliability of a Smartphone-Based Travel Survey: An Experiment in Seoul. *The Journal*  
20 *of The Korea Institute of Intelligent Transport Systems*, 15(2):50–62, Apr. 2016.
- 21 7. J. Jariyasunant, M. Abou-Zeid, A. Carrel, V. Ekambaram, D. Gaker, R. Sengupta, and  
22 J. L. Walker. Quantified Traveler: Travel Feedback Meets the Cloud to Change Behavior.  
23 *Journal of Intelligent Transportation Systems*, 19(2):109–124, Apr. 2015.
- 24 8. K. Shankari, Mohamed Amine Bouzaghrane, Samuel M. Maurer, Paul Waddell, David E.  
25 Culler, and Randy H. Katz. E-mission: An open-source smartphone platform for collect-  
26 ing human travel data. *Transportation Research Record: Journal of the Transportation*  
27 *Research Board*, 2018.
- 28 9. N. D. Lane, E. Miluzzo, H. Lu, D. Peebles, T. Choudhury, and A. T. Campbell. A survey  
29 of mobile phone sensing. *Communications Magazine, IEEE*, 48(9):140–150, 2010.
- 30 10. H. Lu, J. Yang, Z. Liu, N. D. Lane, T. Choudhury, and A. T. Campbell. The jigsaw  
31 continuous sensing engine for mobile phone applications. In *Proceedings of the 8th ACM*  
32 *Conference on Embedded Networked Sensor Systems*, pages 71–84, 2010.
- 33 11. S. Nath. ACE: exploiting correlation for energy-efficient and continuous context sensing.  
34 In *Proceedings of the 10th international conference on Mobile systems, applications, and*  
35 *services*, pages 29–42. ACM, 2012.
- 36 12. Z. Patterson and K. Fitzsimmons. Datamobile. *Transportation Research Record: Journal*  
37 *of the Transportation Research Board*, 2594:35–43, 2016.
- 38 13. A. C. Prelipean, G. Gidófalvi, and Y. O. Susilo. Transportation mode detection – an  
39 in-depth review of applicability and reliability. *Transport Reviews*, 37(4):442–464, July  
40 2017.
- 41 14. S. Reddy, M. Mun, J. Burke, D. Estrin, M. Hansen, and M. Srivastava. Using mo-  
42 bile phones to determine transportation modes. *ACM Transactions on Sensor Networks*,  
43 6(2):1–27, Feb. 2010.
- 44 15. R. C. Shah, C.-y. Wan, H. Lu, and L. Nachman. Classifying the mode of transportation on  
45 mobile phones using GIS information. pages 225–229. ACM Press, 2014.

- 1 16. K. Shankari, J. Fürst, Y. Wang, P. Bonnet, D. E. Culler, and R. H. Katz. Zephyr: Simple,  
2 ready-to-use software-based power evaluation for background sensing smartphone  
3 applications. Technical Report UCB/EECS-2018-168, EECS Department, University of  
4 California, Berkeley, Dec 2018.
- 5 17. K. Shankari, M. Yin, D. Culler, and R. H. Katz. E-Mission: Automated transportation  
6 emission calculation using smartphones. In *Pervasive Computing and Communication  
7 Workshops (PerCom Workshops)*, pages 268–271, Mar. 2015.
- 8 18. A. Vij and K. Shankari. When is big data big enough? Implications of using GPS-based  
9 surveys for travel demand analysis. *Transportation Research Part C: Emerging Technolo-  
10 gies*, 56:446–462, July 2015.
- 11 19. H. Zang and J. Bolot. Anonymization of location data does not work: A large-scale mea-  
12 surement study. In *Proceedings of the 17th annual international conference on Mobile  
13 computing and networking*, pages 145–156. ACM, 2011.
- 14 20. Y. Zheng, Y. Chen, Q. Li, X. Xie, and W.-Y. Ma. Understanding transportation modes  
15 based on GPS data for web applications. *ACM Transactions on the Web*, 4(1):1–36, Jan.  
16 2010.
- 17 21. M. Zhong, J. Wen, P. Hu, and J. Indulska. Advancing Android activity recognition service  
18 with Markov smoother. In *Pervasive Computing and Communication Workshops (PerCom  
19 Workshops), 2015 IEEE International Conference on*, pages 38–43. IEEE, 2015.