

---

# Efficient Location Sampling Algorithms for Road Networks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Many geographic information systems applications rely on data provided by user  
2 devices in the road network, including traffic monitoring, driving navigation, and  
3 road closure detection. The underlying signal is generally collected by sampling  
4 locations from user trajectories. The sampling process, though critical for various  
5 applications, has not been studied sufficiently in the literature. While the most  
6 natural way to sample a trajectory may be to use a frequency based algorithm, e.g.,  
7 sampling locations every  $x$  seconds, such a sampling strategy can be quite wasteful  
8 in resources (e.g., server-side processing, user battery) as well as stored user data.  
9 In this work, we conduct a horizontal study of various location sampling algorithms  
10 (based on frequency, road geography, reservoir sampling, etc.) on the road network  
11 of New York City and assess their trade-offs in terms of various metrics of interest,  
12 such as the size of the stored data and the induced quality of training for prediction  
13 tasks (e.g., predicting speeds).

## 14 1 Introduction

15 For many geographic information systems (GISs) that operate on road networks, the input received  
16 from user devices is vital for the purposes of providing services back to the users. Specifically,  
17 trajectories obtained from user devices can help in a wide range of downstream tasks. Examples  
18 of such tasks include identifying new roads and correcting the locations of nodes and edges in the  
19 graph [3], mobility prediction and next point of interest recommendation [1], and of course the central  
20 application of monitoring the road network and estimating delays on road segments [4], which in turn  
21 enables other use cases such as routing with dynamic information [2, 10] and global optimization  
22 in the road network [5]. The trajectories that enable these applications typically come in the form  
23 of timestamped location (e.g., GPS) samples. Various problems have been studied with respect to  
24 these timestamped location collection, such as completing long gaps in the trajectory [6] or efficiently  
25 processing these trajectories [9]. However designing exactly how the device samples locations and  
26 how the server decides what to use and what to discard has received little attention.

27 In this work we will focus precisely on this problem and we will explore different algorithms and  
28 strategies for sampling location data from devices. We begin with what is the most natural and  
29 most widely used method: periodic sampling. This method relies on a fixed frequency at which the  
30 device sends a location sample to the server. Our study is motivated by the fact that this method is  
31 inherently wasteful. First, think of a congested highway where devices stuck in the same location  
32 are periodically transmitting the same information over and over. This results in a lot of data being  
33 transferred and processed, wasting user battery and server processing resources for little information  
34 in terms of the state of the road network. Second, this method ignores the unit that forms the road  
35 network: the *road segment*. Road segments are the real-world counterpart of the edges of the graph  
36 on which GIS algorithms typically operate and direct measurements of their delays are important.

37 Periodic sampling will instead use interpolation to identify the timestamps during which a car entered  
38 and exited a segment, assuming a constant speed between location samples.

39 This interpolation process is precisely what makes delay measurements sensitive to the exact sampling  
40 policy used. Consider a simple example of two consecutive 50 meter road segments with an average  
41 traversal time of 10 seconds for the first segment and 20 seconds for the second segment. Consider a  
42 device that passes through these segments and gives a location sample at the start of the first segment  
43 and another one 20 seconds later. If this device travels at the average speed of each segment, the  
44 second sample will be in the middle of the second segment. Assuming a uniform speed between the  
45 samples, we would get that this uniform speed will be 3.75 m/s (since the car would have driven  
46 75 meters in 20 seconds) and we would get that the perceived traversal time for the first segment  
47 is  $50/3.75 = 13.333$  seconds. Similarly a device that gives a location sample at the start of first  
48 segment and another at the end of the second segment would give a perceived traversal time of 15  
49 seconds for the first segment. Finally, a device that provides a location sample exactly every 10  
50 seconds (starting from the beginning of the first segment) would correctly identify the end of the first  
51 segment and give a (correct) perceived traversal time of 10 seconds.

52 The above discussion suggests that it is an interesting question to (i) study the locations at which  
53 devices should provide location samples and (ii) determine which of these samples can be dropped  
54 using sub-sampling strategies. To this end, we study several classes of algorithms for this problem.  
55 In all classes of algorithms, tweaking a corresponding parameter changes the number of location  
56 samples given by the devices. Typically, more samples will improve the signal given by the algorithm,  
57 meaning the performance of downstream tasks such as speed/traversal time prediction is improved.  
58 To understand the trade-off for each class of algorithms and their relative performance, we conduct  
59 an experimental analysis. We create synthetic user trajectories in the network of New York City and  
60 use the induced congestion to define road segment delays. We then apply the different sampling  
61 algorithms on the synthetic trajectories to generate training data for speed estimation. We present  
62 results on the performance of the different algorithms and parameter choices with respect to the  
63 number of locations sampled versus the quality of the predictions.

## 64 **2 Sampling Methods**

65 In this section, we give an overview of various algorithms that we consider for sampling location  
66 data. Our algorithms are tailored toward our special use case of road networks and they are based on  
67 well-known sampling algorithms in the literature.

### 68 **2.1 Uniform Sampling**

69 The uniform sampling algorithm is perhaps the most natural one and the one most widely used in  
70 practice, due to its simplicity. In uniform sampling, we want to sample with a preset frequency.  
71 This frequency is set in advance and it is communicated to each device so as the device travels their  
72 assigned path, they will send related data at the requested frequency. It is easy to see that in cases  
73 that there is congestion or device is traveling on a high-traffic road, we collect too many unnecessary  
74 samples. This class of algorithms is parameterized by the sampling frequency.

### 75 **2.2 Randomized Segment-Based Sampling**

76 In this approach, roads are broken down into smaller segments. Each device provides a location  
77 sample as close as possible to the start of a new segment. This is typically not precise, due to GPS  
78 noise, or due to limitations on the device's ability to ping frequently. We parameterize this class of  
79 algorithms with a probability of measurement  $p$ . Each device provides a measurement for any given  
80 road segment with probability  $p$ . Note that to provide such a measurement the device will need to  
81 ping at the beginning of the segment and at the beginning of the next segment.

### 82 **2.3 Congestion-Based Sampling**

83 This is another class of segment-based sampling algorithms. The difference from randomized  
84 segment-based sampling is in how the probability of pinging is determined. In order to avoid having  
85 unnecessary samples, we propose a probability of pinging that correlates with the inverse of the

86 congestion. Now one challenge here is that, we want to have independent (decentralized) sampling  
 87 on each device but in order to get the congestion, we need to consider all devices traveling on  
 88 the road. To circumvent this difficulty, we rely on additional data on the road such as expected  
 89 traveling speed or speed limit, if we observe that our traveling speed is much lower then we lower  
 90 our sampling probability. The parameter of this class of algorithms is a scaling factor that is applied  
 91 to the congestion-based probability. We can think of this scaling factor as the expected number of  
 92 samples we wish to extract per road segment.

## 93 2.4 Reservoir Sampling

94 This class of algorithms also relies on segment-based sampling. The idea of this sampling algorithm,  
 95 is that more samples is not necessarily better. By the law of large numbers, the larger the sample size  
 96 the more the sample will look like the population, so we want to make sure we get enough samples  
 97 that can help us estimate the average speed correctly, however the gain beyond certain point becomes  
 98 marginal. To this end, we consider Reservoir sampling which chooses  $k$  samples from a list of  $n$   
 99 objects where  $n$  is not required to be known in advance. Not requiring the number of objects in  
 100 advance is a positive point, that both helps with running sampling methods in streaming settings and  
 101 also not collecting too many samples for large  $n$ . More precisely, in this sampling, we put the first  $k$   
 102 items in the reservoir list and iteratively go through the remaining items by drawing a random number  
 103  $x$  between 1 and  $t$  for the  $t$ -th item, and then swapping this item with  $x$ -th item in the reservoir if  
 104  $x \leq k$ . The caveat about this method is that it requires each item to know where in the list they are<sup>1</sup>  
 105 which means that we need to have a centralized algorithm. Since routes are communicated through a  
 106 centralized platform, we can ask the platform to communicate an index for each road to each device  
 107 and then the device can decide to send their data if their generated number falls below  $k$  for their  
 108 road segment. Note that this approach works with uniform  $k$  or specialized  $k$  for each road, however  
 109 that needs to be communicated to the device traveling on the road. The obvious parameter for this  
 110 class of algorithms is the size of the reservoir,  $k$ .

## 111 3 Experiments

112 We conduct experiments on synthetic user trajectories built on the real network of New York City.  
 113 The use of synthetic trajectories is due to the fact that real trajectories are not available on the public  
 114 data sets due to their sensitive nature. We now describe how we obtained these trajectories. We  
 115 extract the map of New York City from OpenStreetMap [8] and generate a congestion function for  
 116 each road segment using the functional form of the Bureau of Public Roads [7]. Specifically, for an  
 117 edge (road segment)  $e$  with traffic demand  $x$ , we set:

$$l_e(x) = \frac{0.6t_e^f}{c_e^4}x^4 + t_e^f,$$

118 where  $t_e^f$  is the time needed to cross the edge when the road is empty, i.e., the free-flow travel time,  
 119 and  $c_e$  is the capacity of the street, defined as the number of lanes multiplied by the free-flow speed.  
 120 We then partition the map into 8 areas and pick a random travel demand for each pair or areas. For  
 121 each demand, we compute a set of (at most 10) candidate routes between the origin and destination  
 122 using the alternative route computation algorithm from [10]. We split the demand equally over these  
 123 candidate routes. This process induces a certain congestion on each road segment, which in turn  
 124 determines the average travel time of all cars passing through it by means of applying the congestion  
 125 function  $l_e(x)$ . For each car and each segment that it passes through, we add a noise parameter, which  
 126 finally gives rise to our ground truth user trajectories. This noise models how different cars traverse  
 127 the same road segment with different delays and is the cause of downstream prediction errors.

128 We apply each one of the sampling algorithms to each one of these “ground truth” trajectories. Each  
 129 one of our four algorithms types uses a parameter (the sampling frequency for uniform sampling,  
 130 the sampling probability for random sampling, the expected number of samples per segment for  
 131 congestion-based sampling, and the size of the reservoir for reservoir sampling) and we select 4  
 132 values of each parameter, giving rise to 16 total algorithms. Using the location samples selected by  
 133 the algorithms, we recompute the perceived road segment entry and exit timestamps by interpolating

<sup>1</sup>Note that the order does not matter, however we need to assign a unique number from 1 to  $n$  to each of the items

Sampling Method	Parameter	MSE	MAE	MBE	Sample Size
Uniform Sampling	5	6.90	1.35	0.40	1252746
Uniform Sampling	10	11.34	2.03	0.38	630995
Uniform Sampling	15	17.63	2.68	0.36	423639
Uniform Sampling	20	24.25	3.17	0.34	320088
Congestion-Based Sampling	10	5.68	0.77	0.46	35104
Congestion-Based Sampling	20	4.90	0.73	0.40	66116
Congestion-Based Sampling	30	5.38	0.74	0.41	92717
Congestion-Based Sampling	40	4.59	0.73	0.41	115688
Reservoir Sampling	10	4.59	0.72	0.43	33954
Reservoir Sampling	20	4.42	0.72	0.42	57404
Reservoir Sampling	30	4.29	0.71	0.41	76407
Reservoir Sampling	40	4.51	0.72	0.41	93329
Randomized Segment-Based Sampling	25	4.47	0.72	0.41	129616
Randomized Segment-Based Sampling	50	4.37	0.71	0.41	223111
Randomized Segment-Based Sampling	75	4.41	0.71	0.40	279424
Randomized Segment-Based Sampling	100	4.42	0.71	0.40	299244

Table 1: Performance of various sampling algorithm on test data set.

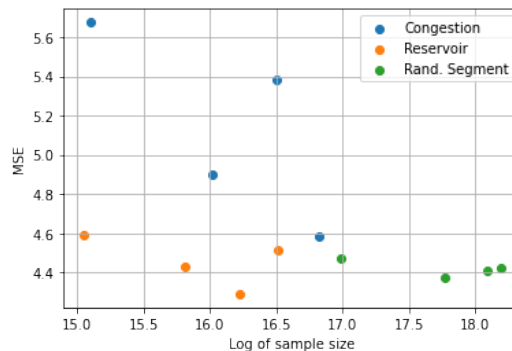


Figure 1: The mean square error of sampling algorithms on test data. X-Axis is on log scale and the Uniform sampling for more focused view.

134 between location samples and assuming a constant speed. This process induces a different perceived  
135 road segment travel time for each location sampling strategy. These perceived travel times are in fact  
136 our training data for road segment delay prediction. We have 16 training data sets, one per algorithm.

137 Since we only have travel time and not much additional information on segment properties or  
138 trajectories, for square loss minimization, the average speed values for each segment, yields the  
139 optimal prediction. We then use this calculated average per segment and evaluate each algorithm on  
140 the test data set which again is from the trajectories that were sampled uniformly and excluded from  
141 the training. All trajectories are from the same time period. For evaluate our algorithms by reporting  
142 the mean square loss (MSE), mean absolute loss (MAE), and mean biased error (MBE) along side  
143 the number of samples in the training set (see Table 1).

144 We also plot the trade-off between the size of the sample set and the quality of the predictions for  
145 each algorithm in Figure 1. We exclude the Uniform sampling since it requires a lot of samples and  
146 has high error. Since we add noise to our data, we don't necessarily see decrease in the reported  
147 error as we allow more samples. Based on reported numbers, it is easy to see that the Randomized  
148 segment-based and Reservoir sampling have the best performance. The Congestion based sampling  
149 starts to get more competitive results as we allow more samples. Overall, the Reservoir sampling  
150 yields the best result, especially if we consider the Pareto curve where we prefer a point with lower  
151 error and lower sample size, we see that the curve consists of the three points on the bottom left  
152 corresponding to the Reservoir sampling.

## References

- 153
- 154 [1] Q. Chen, R. Jiang, C. Yang, Z. Cai, Z. Fan, K. Tsubouchi, R. Shibasaki, and X. Song. Dualsin:  
155 Dual sequential interaction network for human intentional mobility prediction. In *Proceed-*  
156 *ings of the 28th International Conference on Advances in Geographic Information Systems,*  
157 *SIGSPATIAL '20*, page 283–292, New York, NY, USA, 2020. Association for Computing  
158 Machinery.
- 159 [2] D. Delling, A. Goldberg, T. Pajor, and R. Werneck. Customizable route planning. In *Proceedings*  
160 *of the 10th International Symposium on Experimental Algorithms (SEA'11)*, Lecture Notes in  
161 Computer Science. Springer Verlag, May 2011.
- 162 [3] A. Hendawi, S. S. Sabbineni, J. Shen, Y. Song, P. Cao, Z. Zhang, J. Krumm, and M. Ali. Which  
163 one is correct, the map or the gps trace. *SIGSPATIAL '19*, page 472–475, New York, NY, USA,  
164 2019. Association for Computing Machinery.
- 165 [4] W. Jiang and J. Luo. Graph neural network for traffic forecasting: A survey. *CoRR*,  
166 abs/2101.11174, 2021.
- 167 [5] K. Kollias, A. Chandrashekarapuram, L. Fawcett, S. Gollapudi, and A. K. Sinop. Weighted  
168 stackelberg algorithms for road traffic optimization. In *Proceedings of the 29th International*  
169 *Conference on Advances in Geographic Information Systems, SIGSPATIAL '21*, page 57–68,  
170 New York, NY, USA, 2021. Association for Computing Machinery.
- 171 [6] J. Krumm. Maximum entropy bridgelets for trajectory completion. In *Proceedings of the 30th*  
172 *International Conference on Advances in Geographic Information Systems, SIGSPATIAL '22*,  
173 New York, NY, USA, 2022. Association for Computing Machinery.
- 174 [7] B. of Public Roads. *Traffic assignment manual*. US Department of Commerce, 1964.
- 175 [8] OpenStreetMap contributors. Planet dump retrieved from <https://planet.osm.org> . <https://www.openstreetmap.org>, 2017.
- 176
- 177 [9] S. A. Shaikh, H. Kitagawa, A. Matono, and K.-s. Kim. Tstream: A framework for real-time and  
178 scalable trajectory stream processing and analysis. In *Proceedings of the 30th International*  
179 *Conference on Advances in Geographic Information Systems, SIGSPATIAL '22*, New York,  
180 NY, USA, 2022. Association for Computing Machinery.
- 181 [10] A. K. Sinop, L. Fawcett, S. Gollapudi, and K. Kollias. Robust routing using electrical flows.  
182 In X. Meng, F. Wang, C. Lu, Y. Huang, S. Shekhar, and X. Xie, editors, *SIGSPATIAL '21:*  
183 *29th International Conference on Advances in Geographic Information Systems, Virtual Event /*  
184 *Beijing, China, November 2-5, 2021*, pages 282–292. ACM, 2021.