

Spatial and Temporal Analysis of Planet Scale Vehicular Imagery Data

Gautam S. Thakur^{*†}, Pan Hui[†], Hamed Ketabdar[†] and Ahmed Helmy^{*}

^{*}*CISE, University of Florida, Gainesville, FL 32611-6120*

[†]*Deutsche Telekom Laboratories, Berlin*

Email: gsthakur@cise.ufl.edu, pan.hui@telekom.de, hamed.ketabdar@telekom.de, helmy@cise.ufl.edu

Abstract—Vehicular traffic congestion is becoming a major problem in metropolitan cities throughout the world. Looking into the future, this becomes particularly more challenging with the emergent nature combining population explosion, number of vehicles and the organic growth of cities’ infrastructure. In order to study this problem, we need the traffic data and cities’ physical infrastructure and the application of robust data mining and knowledge discovery techniques on this data to identify potential bottlenecks. In this work, we propose a novel method of collecting city-wide traffic information from online vehicular traffic camera. Our resulting dataset is a several months collection of vehicular mobility traces captured from 2709 traffic webcams in 10 different cities across the world, with 7.5 Terabytes of data with 125 million vehicular images. We also collect driving distance and time between geo-coordinate pairs of street intersections for these cities. We apply spatio-temporal data mining techniques to profile these global cities and reason about their geographical backbone and provide an insight into their vehicular traffic density distribution. Our results show that: (i) High correlation between driving time and distance indicate congestion-free traffic, (ii) Traffic follow certain patterns that are stable for a long time (42 days). (iii) Traffic Congestion show high Correlation (80%) for 1-2 hour lag then decrease significantly to 25-30% for four hours lag. We believe our study help to shed light on causes of contention in the present day traffic-jams and provide an insight into the planning and development of future cities and resolution to traffic congestion.

Keywords-Traffic Camera; Level of Congestion; City Dynamics Profiling;

I. INTRODUCTION

Vehicular traffic congestion is becoming an ever increasing problem around the world. In the latest (2010) urban mobility report[1], congestion caused urban Americans to travel 4.8 billion hours and to purchase an extra 3.9 billion gallons of fuel for a cost of \$115 billion. On average, yearly peak period delay caused by the traffic congestion for the average commuter was 34 hours and the cost to the average commuter has increased by 230% in two decades[1]. Congestions not only affect people during the peak period, but also at other hours, approximately half of total delays occur at midday and overnight.

In attempting to identify the causes for this problem, most of the current approaches focus on isolated efforts to improve conditions as and when it arose at few locations. As of now the transportation and engineering sciences have proposed to get as much service as possible. First, by timing

the traffic signals so that more vehicles see green lights, improving road and intersection designs, or adding a short section of roadways. Second, by adding more capacity in critical corridors new streets and highways, new or expanded public transportation facilities, and larger bus and rail fleets. Third, by changing the usage patterns like flexible work hours, avoid traveling during the rush hours. We believe these approaches are currently insufficient unless coupled with a comprehensive picture of city structure and the traffic distribution across its key intersection in a collective manner.

In doing so, we propose a systematic approach to first analyze the structural dynamics of these cities and a temporal analysis of the vehicular traffic flow in them. One way to pursue structural dynamics is to look into driving distances and corresponding time across several intersections. A deviation from the general notion of small driving distance to corresponding small driving time and similarly for large distances require long time can be used to identify the critical sections of the cities that are prone to congestions. It is eminent that small and less number of lanes, closely situated intersection, merging and bi-furcations will only result in longer time to travel despite relatively shorter distances. Furthermore, a temporal analysis of the traffic distribution across several hours of day, during weekdays and weekends give an important insight into the distribution and correlation of traffic in that city. Finally, combining the distribution of traffic with structural dynamics of cities will provide a comprehensive knowledge of locations and time and insight into the efficacy of prior measures.

Recently several transportation departments(DoTs) have installed online traffic web cameras at key intersections to know current trends in the traffic flow. At regular intervals of time, these cameras capture still pictures of on-going road traffic and send them in form of feeds to media server. We develop an automatic script to acquire images at a finer interval of around 30 seconds per image. We develop a fast background image subtraction algorithm to extract the traffic densities from these images for the purpose of analyzing the traffic on roads. We also use location information(geo-coordinates) of these cameras to calculate driving distances and driving times(using Google services) between all pairs of locations. We perform k-mediod clustering on the camera pairs based on the distances and compare the clusters formed using driving time for the same pairs. The discrepancy

between driving distance and time clusters provide a good reasoning for the study to identify locations that are prone to traffic congestions. Later, we use the longitudinal densities values extracted from images for these cities to show the distribution of traffic across several hours of the day and across several weeks. In summary, our contributions include :

- 1) A novel approach to collect vehicular traffic flow and driving information using publicly available traffic cameras.
- 2) To the best of our knowledge, we provide by far the largest and most extensive library of vehicular density data, based on processing of millions of images. This addresses a severe shortage of such data sets in the community. The library will be made available to the research community in the future.
- 3) Application of correlation and spatio-temporal analysis effectively in identifying patterns and predictions to current traffic problems.

This paper is organized as: Section-II we provide background to this work and in Section-III we describe the data set and the algorithm to extract traffic densities from the images. In Section-IV we explain the structural dynamics of the city and apply clustering to cities' structural data sets. In Section-V, we perform temporal Analysis of Traffic Densities and discuss the results in a elaborate way. Finally, we conclude our paper in Section-VI detailing future directions.

II. RELATED WORK

Several measures have been proposed recently to counter the traffic congestion and provide better management for the traffic throughput. For transportation and civil engineering point of view segregate efforts have been made to improve junctions, bus and express lanes, and car pooling to name a few. City planning and urban design[2] practices can have a huge impact on levels of future traffic congestion. In [3], clustering have been applied to study vehicular traffic through a sequence of traffic lights on a highway, where all signals turn on and off synchronously is studied and the dynamical behaviors of vehicles are clarified by analyzing trafrc patterns. The results show that clustering of vehicles varies with the cycle time of signals and are controlled by varying both split and cycle time of signals. Along the same lines in [4], authors studied simple aggregation model that mimics the clustering of traffic on a one-lane roadway and derived derive an analytical solution for the probability of a single car and an asymptotically exact expression for the joint mass-velocity distribution function.

CORRSIM[5] and VISSIM[6] are two commonly used simulators used for micro-modeling of traffic. In other studies, researcher attempted to model the congestion and traffic using mathematical models and derive a close form of expression. In [7], Bando et. al. proposed dynamical model of traffic congestion based on the equation of motion

Table I
GLOBAL WEBCAM DATASETS

City	# of Cameras	Duration	Interval	Records	Database Size
Bangalore	160	30/Nov/10 - 01/Mar/11	180 sec	2.8 million	357 GB
Beaufort	70	30/Nov/10 - 01/Mar/11	30 sec	24.2 million	1150 GB
Connecticut	120	21/Nov/10- 20/Jan/11	20 sec	7.2 million	435 GB
Georgia	777	30/Nov/10 - 02/Feb/11	60 sec	32 million	1400 GB
London	182	11/Oct/10 - 22/Nov/10	60 sec	1 million	201 GB
London(BBC)	723	30/Nov/10 - 01/Mar/11	60 sec	20 million	1050 GB
New york	160	20/Oct/10 - 13/Jan/11	15 sec	26 million	1200 GB
Seattle	121	30/Nov/10 - 01/Mar/11	60 sec	8.2 million	600 GB
Sydney	67	11/Oct/10 - 05/Dec/10	30 sec	2.0 million	350 GB
Toronto	89	21/Nov/10 - 20/Jan/11	30 sec	1.8 million	325 GB
Washington	240	30/Nov/10 - 01/Mar/11	60 sec	5 million	400 GB
Total	2,709	-	-	125.2 million	7,468 GB

of each vehicle by analyzing the stability of traffic flow and the evolution of traffic congestion is observed with the development of time. The implications of empirical time headway distributions of traffic flow[8] and underlying stochastic process has shown to model fluctuations of traffic flow. However these approaches also lack a comprehensive study for traffic congestions and its effect on other non-congested segments of the city.

However, to best of our knowledge ours approach is the first of its kind to apply data mining on such a large data set to profile cities and model the traffic densities. In next few sections we introduce our data set and explain algorithm to extract relevant information from the imagery data.

III. MEASUREMENTS AND PROCESSING

In this section, we introduce our vehicular imagery data set and detail the process to gather structural information and then explore the vehicular imagery data that is used to represent traffic densities. Later on we explain the algorithm to process these images.

A. Online Traffic Web Cameras

To curb traffic issues and disseminate any congestion in real time, several governments have started to install online traffic web cameras that provide qualitative trends about road traffic. These web cameras are installed at strategies locations includes critical intersection and highways to visually monitor traffic. For the purpose of this study, we choose 10 cities with large number of webcam coverage and collect vehicular imagery data for several months, which help analyzing the traffic trends across the city. The details of this dataset are shown in Table-I.

B. Geo-Coordinates Data

We collect the physical information of these cameras as reference points to study the structural analysis. The physical information of these cameras includes latitude and longitude coordinates, zipcode and state, directional view, camera locations. These cameras are installed at critical intersections and highways, so a study involving such locations will eventually give a good estimate of its structural significance.

C. Algorithm for Traffic Density Estimation

We aim to estimate traffic density(d) on roads considering the number of vehicles or pedestrians crossing the road. We have a sequence of images captured by webcams. Considering our problem, we have to be able to separate information we need, e.g. number of vehicles and pedestrians from the back ground image, which is normally road and buildings around. The main factor that can distinguish between vehicles and background image (road, buildings) is the fact that the vehicles are not in a stationary situation for a long period of time, however the back ground is stationary. The solution for the problem then seems to be applying a sort of high pass filtering over a sequence of images captured by a webcam over time. The high pass filter removes the stationary part of the images (road, buildings, etc.), and keeps the moving components (mainly vehicles). In order to implement such a high pass filter, we subtract result of a low pass filter over a sequence of images, from each still image. This is practically equivalent to implementing a high pass filter over sequence of images. In order to obtain low pass filtering effect, we run a moving average filter over a time sequence of images obtained from one webcam. The duration of the moving average filter can be adjusted in an adhoc way. The moving average filter is simply implemented by averaging over the intensity map for several images in a certain duration. At the output of the moving average filter, the intensity of each pixel is obtained by averaging intensity of corresponding pixels in the interval. The output of the moving average filter (low pass filter) is normally the required background image, which is still part of the image. Therefore, subtracting each image from the output of the low pass filter, gives us the moving components (e.g. vehicles). Having the high pass component of the image, the vehicles are highlighted from background. One could then use regular object detection techniques to identify and count number of vehicles in the high pass filtered image. However, this is computationally expensive and unnecessary. As an alternative, we simply count the number of active pixels (pixels with a value higher than a certain threshold). This is much faster than detecting and counting objects in an image. At the same time, it is more effective, because we are looking at the traffic densities (d) i.e., percentage of the street (road) which is covered by vehicles (as an indicator of how crowded is the street), rather than number of vehicles. The number of vehicles is not a good indicator of crowd, as a long vehicle may introduce more traffic than a small one. Second, our method overcomes the issues that object detection face in case of severe congestion. Counting number of active pixels can indicate what percentage of the road is covered, no matter how many vehicles are in the road. In many instances, images are duplicate, corrupted with zero sized or with extraneous bytes (noise). We use semi-supervised learning and hierarchical clustering to overcome the challenges of

Table II
DETAILS OF GEO-COORDINATE MEASUREMENTS

City	# Pairs	Avg. Dis.(km)	Avg. Time(m)	Area
Connecticut	74801	6.4	42	60%
London	32580	1.56	26	53%
Sydney	4422	3.2	33	67%

Table III
CLUSTERING ANALYSIS

City	D.I	# Clusters	ρ
Connecticut	0.001	6	6%
London	0.0002	10	52%
Sydney	0.005	4	5%

outliers' detection and removal. Due to limited space, we are omitting the details, however the results of outliers detection and removal is shown in here[11].

D. Geo-Coordinates Pairing

After gathering the physical information of cameras, we calculate the driving distances and corresponding driving time for all camera pairs within each city. We use the Google Maps API to gather information about all such pairs of locations. The details show in the Table-II. Thus, average time to distance ratio is very high to travel in the city of London.

IV. SPATIAL ANALYSIS OF CITIES

In this section, we study the geographical access and city dynamics in terms of driving distance and respective time between a set of cameras for each city separately. By default, the driving time and corresponding distance are obviously correlated, but organic growth of cities, number of traffic signals, lanes and width of road to name a few can affect this correlation. Without any exceptions, a high correlation value reflect a congestion-free traffic system and normal movement of vehicles. On the other hand, a low to negative correlation is prone to traffic congestions and uneven distribution of traffic speeds. The knowledge of this kind is important for example, in hospital emergency vehicles, where patients are much more sensitive to time than to distance differentials[9].

Analysis

We separately apply k-medoid clustering in two-dimensions on geo-coordinates pairs of cameras, first clustering them based on driving distance and then on corresponding driving time between them for each city. To measure the variation we compare deviation in the cluster membership. This helps: (i) to discover the irregularity caused by spatial feature around these locations, (ii) to know the set of cameras that are similar according to specific metrics. In order to identify correct number of clusters, we used Dunn index[10] since k-medoid clustering is hard partitioning.

We sample the distance and time distribution of all camera pairs to investigate if there is a linear mapping in

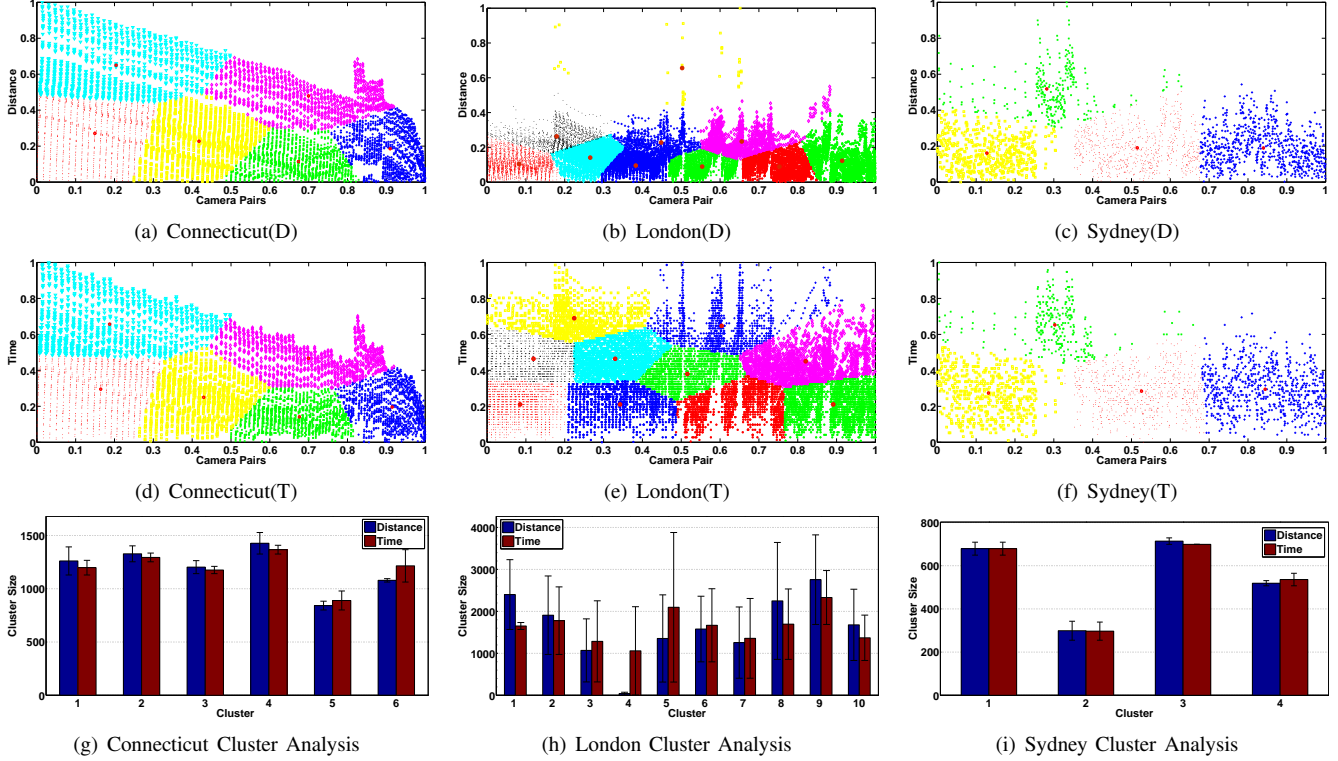


Figure 1. **K-medoid Clustering for the distance (D) and time (T) pairing between the cameras for the cities with centroids (in red). The size and boundaries of clusters show a non-linear distance and time correspondence. The errorbars show the deviation in the membership of objects for that cluster. Quantitative results of clustering show connecticut geo-spatial distribution is more streamlined.**

these metrics. Surprisingly, we find a deviation in this statistics as show in Fig.-1. The bar graphs show the size of clusters generated after running 1000 iterations of k-medoid algorithm. The error bars are the deviation in the data points (camera pairs) that are present in distance cluster but not in time clusters and vice versa for time. In the Table-III we show this deviation using ρ , which is a function of membership of a pair of coordinates that belong to one cluster. A correlation exists, when a pair of coordinates fall in the same cluster number in case of both time and distance. In the city of London, we find that road intersections are very close together and hence tendency to have a slow traffic in inherent in the city infrastructure. In the state of Connecticut the clusters 1, 4, 5 and 6 have large deviation and have a comparatively larger number of small distance intersections. While for the sydney, it seems the clusters have very low deviation and an analysis of clusters show an even distribution of traffic across all its roads. This further shows that Sydney is more well planned and has less organic growth as compared to other cities. We find similar analysis for other cities as well, however brevity we only discuss three of the ten cities. In Table-III the average deviation among clusters membership shows the inconsistency among distance and time. Our analysis suggests that distance and time variation indicate one the main features of congestion

and slow traffic in the cities.

V. TEMPORAL ANALYSIS OF TRAFFIC DENSITIES

In the previous section on spatial analysis, we discussed the clustering variation for geodesic driving distance and time among a set of critical intersection points of cities. At any given time these statistics remain the same unless cities are structurally modified. In this section, we focus more on the temporal aspect of traffic distribution in these cities. We use traffic densities extracted from imagery data set of cameras to perform this analysis. In this activity we ask the following questions.

- Q.1: What is the nature of traffic distribution for these cameras? Do all cameras have same distribution?*
- Q.2: Is the nature of distribution predictable over a long period of time?*
- Q.3: Which events cause deviation in the traffic distribution? How to identify such events ?*

In order to get an impression of the traffic distribution, we start with a qualitative analysis of densities in three cities. This step help to get a high level representation and reasoning methods about the behavior of traffic. Later on, we make it a logical base to investigate precise quantitative information.

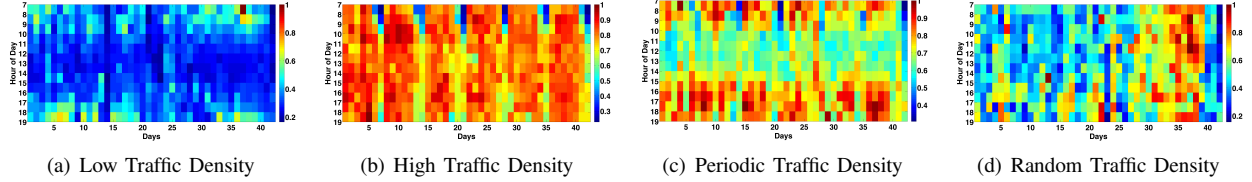


Figure 2. A 42 days traffic distribution snapshot from four different cameras is shown. Fig-(a) shows relatively mild traffic during various hours of the day, while (b) shows high traffic recording for the full trace periods. In Fig-(c) we find a regularity patterns during the morning and evening hours when the traffic is relatively higher than afternoon intervals. A random traffic characterization is recorded in the last.

A. Traffic Distribution

To answer the first question, we sample the data set into hourly represent for a period of 42 days and use spectral method to generate the relationship between various days. A sampled view taken from three cities show that cameras have varying traffic distribution against the popular notion of 'rush hours'. We find that it is difficult to estimate an aggregate statistical parameter that define all cameras. For example, in Fig.-2(a), we find a very low instance of traffic for a period of 42 days, while its opposite in the Fig.-2(b) with consistently high traffic, indicating a busy market area. One interesting pattern is the reduction during the weekends (day 7,14 etc.). In Fig.-2(c) the temporal activity riches its height during the morning and evening hours while relative smoothness during the afternoon hours. Finally in Fig.-2(d), we found that the patterns are not very regular and we find it really hard to estimate a correct traffic model. This kind of variation in traffic give challenges in developing an aggregate forecasting mechanism for a city wide traffic. It also reject the a popular notion of 'rush hours' concept since different cameras seems to posses their own distribution of traffic during varying hours of the day.

B. Traffic Congestion Correlation

We utilize the power of correlation coefficients to measure the degree to which traffic congestion values and direction of a linear relationship are associated with itself. In our case we using this technique to correlate the traffic change during several hours of day. We are analyzing the correlations for 1-4 hour lags for each camera with itself. For example, we investigate what is the correlation between the traffic at 7 AM and 8 AM, 7 AM and 9 AM etc. We use the following to calculate the value of correlation coefficient ρ . For two cameras C_1 and C_2 at hour h of the day k , their with traffic densities will be $d(C_{1h}^k)$ and $d(C_{2h}^k)$, respectively. The correlation of their densities is shown in Eq.-1,

The value of ρ is varying: $-1 \leq \rho \leq +1$. If two hour lags have a strong positive linear correlation in the traffic values, the value of ρ is reaches to $+1$ else a strong negative linear correlation the value of ρ is close to -1 . If there is no linear correlation or a weak linear correlation, the value of ρ is close to zero. A value near zero means that there is a random, nonlinear relationship between the two

hours traffic values. In order to analyze the traffic change we investigate the correlation of traffic in four different time lags from one to four hours. To do this, we sample the traffic densities of each camera into hours from 7 A.M to 7 P.M. Then we find correlations between consecutive hours for 42 days. For example, we find the correlation coefficient of traffic densities between 42 days of 7 A.M and 8 A.M, 7 A.M and 9 A.M. We accumulate the results for each camera separately and compile them in form a CDF that shown in Fig.-3. We see that for the city of Connecticut and Sydney in Fig.-3(a) and 3(c), the hourly traffic change is highly correlated, almost 80% of cameras' next hour traffic is 70% correlated to its current hour. More surprising is the two hours difference of traffic densities for cameras of these two cities, where 80% of next two hours traffic is only 50% or less correlated to the current hour. And around 60% cameras have 30% correlation for a time lag 3-4 hours. While in case of the city of London the next hour traffic density for 80% cameras is close to 60% correlated to the current hour. It goes further down to 30% for next two hours and around 15-20% for a 3-4 hour difference. These observations tell us that London traffic has high fluctuations than Connecticut and Sydney, which are rather very smooth. Any prediction model build on this requires more insight into the transitioning traffic for London and Connecticut because of their high variability. After looking into a high level picture of traffic change for several consecutive hours, we move to analyze the variability contributed to the CDFs of Fig.-3 by the individual hours. This helps to know which hour of the day's traffic change on a large scale compared to its previous hours. We start by analyzing the city of Connecticut, where the correlation for next hour is low at 7 AM and and 12 noon and it is relatively low also for 2-4 hours lag. The city of London normally has a relatively low correlation with around 50% during 7 AM, but it suddenly drops to less than 1% for 2-4 hours lag. However, afternoon traffic remain relatively 40% correlated for 1-2 hour lag and very high during the later hours of the day with 60%-80% correlation for 1-3 hour lags. The London hourly traffic fluctuates more during the morning times and early evening, but relatively stable during the noon and evening hour, which is opposite to connecticut, although the average correlation for London is nearly 50%-60% only. As expected for the city of Sydney,

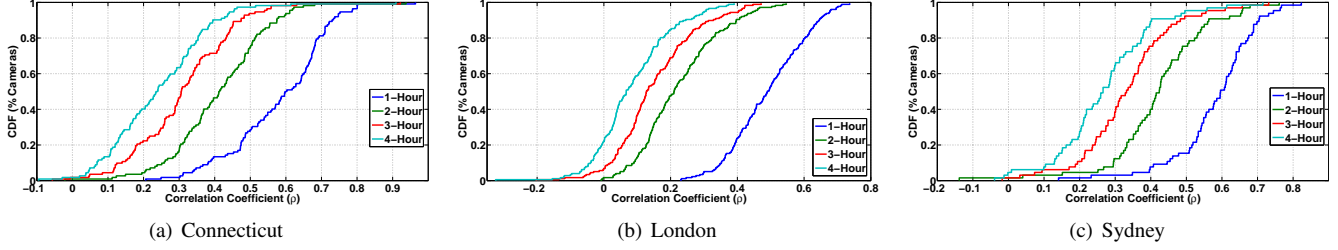


Figure 3. Cumulative distribution function show camera auto-correlation of traffic densities between hour differences of the day for three cities.

$$\rho(C_{1h}^k, C_{2h}^k) = \frac{n * \sum d(C_{1h}^k) * d(C_{2h}^k) - (\sum d(C_{1h}^k)) * (\sum d(C_{2h}^k))}{\sqrt{n * (\sum d(C_{1h}^k)^2) - (\sum d(C_{1h}^k))^2} * \sqrt{n * (\sum d(C_{2h}^k)^2) - (\sum d(C_{2h}^k))^2}} \quad (1)$$

the traffic is relatively 60%-70% correlated for one hour lag and show nearly consistent correlation of 40%-55% for 2-4 hour lag. This analysis shows the trends as seen in previous section where CDFs of Sydney are relatively more stable and highly correlated. This analysis provide a stepping stone to develop a model for predicting next hours traffic based on the current condition. More details are available in [11].

VI. CONCLUSION

The congestion problem has become widespread in recent years. In this work, we introduced a novel data collection method for vehicular traffic using globally available web cams and a novel method to mine vehicular imagery data for the measuring level of traffic congestion. We studied three cities to analyze the spatio temporal distribution of traffic. The spatial gave an impression about the internal road networks of the city and its mapping to the driving time. A variation in such statistics results in city being highly critical to a potential congestion problem that result in a non-uniform traffic. In the second step, we used the traffic estimates from 42 days measurement to analyze its pattern across many dimension of days and weeks. The results of correlation for different hour lags gave the stability in the traffic. The temporal analysis is done to show the traffic congestion are variable in nature and can be selective and different in time scales. We find the London traffic is very uncorrelated and hence difficult to predict. Our results show that: (i) High correlation between driving time and distance indicate congestion-free traffic, (ii) Traffic follow certain patterns that are stable for a long time (42 days). (iii) Traffic Congestion show high Correlation (80%) for 1-2 hour lag then decrease significantly to 25-30% for four hours lag. In future, we look forward to develop a framework for real time analysis of traffic and incorporate them as data mining tool for the community.

REFERENCES

- [1] T. L. David Schrank and S. Turner, "Ttis 2010 urban mobility report powered by inrix traffic data," *Texas Transportation Institute, The Texas A&M University System*, December 2010.
- [2] J. Barnett, *An introduction to urban design*, 1st ed. Harper & Row, New York :, 1982.
- [3] T. Nagatani, "Clustering and maximal flow in vehicular traffic through a sequence of traffic lights," *Physica A: Statistical Mechanics and its Applications*, vol. 377, no. 2, pp. 651 – 660, 2007.
- [4] E. Ben-Naim, P. L. Krapivsky, and S. Redner, "Kinetics of clustering in traffic flows," *Phys. Rev. E*, vol. 50, no. 2, pp. 822–829, Aug 1994.
- [5] A. Halati, H. Lieu, and S. Walker, "CORSIM-corridor traffic simulation model," in *Proceedings of the Traffic Congestion and Traffic Safety in the 21st Century Conference*, 1997, pp. 570–576.
- [6] N. E. Lownes and R. B. Machemehl, "Vissim: a multi-parameter sensitivity analysis," in *Proceedings of the 38th conference on Winter simulation*, ser. WSC '06. Winter Simulation Conference, 2006, pp. 1406–1413.
- [7] M. Bando, K. Hasebe, A. Nakayama, A. Shibata, and Y. Sugiyama, "Dynamical model of traffic congestion and numerical simulation," *Physical Review*, vol. 51, pp. 1035–1042, Feb. 1995.
- [8] P. Wagner, "Modeling traffic flow fluctuations," *Journal of the ICE*, vol. 4, pp. 121–130, 1936.
- [9] C. S. Phibbs and H. S. Luft, "Correlation of travel time on roads versus straight line distance," *Medical Care Research and Review*, vol. 52, no. 4, pp. 532–542, 1995. [Online]. Available: <http://mcr.sagepub.com/content/52/4/532.abstract>
- [10] J. C. Dunn, "Well-separated clusters and optimal fuzzy partitions," *Journal of Cybernetics*, vol. 4, no. 1, p. 95, 1974.
- [11] G. S. Thakur, P. Hui, H. Ketabdar, and A. Helmy, "Towards realistic vehicular network modeling using planet-scale public webcams," *CoRR*, vol. abs/1105.4151, 2011.