

# Beijing City Lab

Wang J, Mao Y, Li J, Wang W, Li C, Xiong Z, 2014, Predictability of road traffic and congestion in urban areas. Beijing City Lab. Working paper #55

# Predictability of road traffic and congestion in urban areas

Jingyuan Wang,<sup>1</sup> Yu Mao,<sup>1</sup> Jing Li,<sup>1</sup> Wen-Xu Wang,<sup>2</sup> Chao Li,<sup>1,3</sup> and Zhang Xiong<sup>1,3</sup>

*<sup>1</sup>School of Computer Science and Engineering,  
Beihang University, Beijing, 100191, China*

*<sup>2</sup>School of Systems Science, Beijing Normal University, Beijing, 100875, China*

*<sup>3</sup>Research Institute of Beihang University in ShenZhen, Sehnzhen, 518057 China*

## Abstract

Mitigating traffic congestion on urban roads, with paramount importance in the urban development, reduction of energy consumption and air pollution, depends on our ability to foresee road usage and traffic condition pertaining to the collective behavior of drivers, raising a significant question: to what degree is road traffic predictable in urban areas? Here we rely on the precise record of daily vehicle mobility based on GPS devices installed in taxis to uncover the potential daily predictability of urban traffic patterns. Using the mapping from the degree of congestion on roads into a time series of symbols and measuring its entropy, we find a relatively high daily predictability of traffic condition despite the absence of any a priori knowledge of drivers' origins and destinations and quite different travel patterns between weekdays and weekends. Moreover, we find a counterintuitive dependence of the predictability on travel speed: the road segment associated with intermediate average travel speed is most difficult to be predicted. We also explore the possibility of recovering the traffic condition of an inaccessible segment from its adjacent segments with respect to limited observability. The highly predictable traffic patterns in spite of the heterogeneity of drivers' behaviors and the variability of their origins and destinations enable the development of accurately predictive models for eventually devising practical strategies to mitigate urban road congestion.

The past decades have witnessed a rapid development of modern society accompanied with an increasing demand for mobility in metropolises [1–4], accounting for conflict between the limits of road capacity and the increment of traffic demand reflected by severe traffic congestions [5, 6]. Induced by such problems, citizens suffer from reduction of travel efficiency, increase of both fuel consumption [7] and air pollution [8] related with vehicle emission. For instance, in recent years, a number of major cities in China have frequently experienced persistent haze, raising the need of better traffic management to mitigate congestion that is likely one of the main factors for the pollution [9, 10]. Despite much effort dedicated to addressing the problems of traffic jam [11], urban planning [12, 13] and traffic prediction [14–16], we still lack a comprehensive understanding of the dynamical behaviors of urban traffic. The difficulty stems from two factors: the lack of systematic and accurate data in conventional researches based on travel surveys and the diversity of drivers’ complex self-adaptive behaviors in making routing choice decision [17]. Fortunately, “big data” as the inevitable outcome in the information era opens new routes to reinvent urban traffic systems and offer solutions for increasingly serious traffic jams [18]. In this light, mobile phone data have been employed to explore road usage patterns in urban areas [19–21]. However, to eventually implement control on road traffic, predict traffic condition is the prerequisite, which prompts us to wonder, to what degree traffic flow on complex road network is predictable with respect to high self-adaptivity of drivers and without a priori knowledge of their origins and destinations.

In this paper, we for the first time explore the predictability of urban traffic and congestion by using comprehensive records of Global Position System (GPS) devices installed in vehicles. The data provide the velocity and locations of a large number of taxis in real time, enabling investigation and quantification of the predictability of segments in main roads in an urban road network. In particular, we establish a mapping from the degree of congestion on a segment of road into a time series of symbols, which allows us to exploit tools in the information theory, such as entropy [22] and Fano’s inequality [23] to measure the predictability of traffic condition on a segment of road. Our methodology is inspired by the seminal work of Song et al. who incorporate information theory into time series analysis to measure the limited predictability of individual mobility [24]. Our main contribution is that we extend the tools of time series analysis to the collective dynamics of road traffic rather than at the individual level, by mapping the vehicle records from GPS into road usage so as to offer the predictability of traffic condition at different locations. In contrast to the traditional way based on origin-destination analysis [25], our approach relies only on short-time historical record of traffic condition without the need of a priori knowledge of drivers’ origins and destinations and their associated navigation strategies. Our accessibility of such individual-level information is inherently limited by the diversity in population, job switching, moving and urbanization. Our research gives rise to a number of interesting findings, including relatively high daily predictability of traffic condition in the three Ring Roads in Beijing [26] despite quite different travel patterns at the weekends compared to working days, the non-monotonic dependence of the predictability on vehicle velocity and the recoverability of the traffic condition of an inaccessible segment by the information of its adjacent observable segments. Thus we present a general and practical approach for understanding the predictability of real time urban road traffic and for devising effective control strategies to improve the roads’ level of service.

## I. RESULTS

We explore the predictability of traffic condition by using the GPS records of more than 20000 taxis in Beijing, China, (see Methods for data description and processing). We focus on the three Ring Roads, the 2nd, 3rd and 4th Rings in Beijing by mapping the states of vehicles into the traffic condition on the roads. The three rings bear the most heavy traffic burden in Beijing and the data records pertaining to them with high frequency are sufficient for quantifying their traffic condition. In particular, we divide each ring into a number of segments with given *segment length*  $\Delta L$ , and measure the traffic condition of each segment by the average velocity of vehicles. To simplify our study, we discretize the average velocity of the segments in the range from 0km/h to the speed limit 100km/h with a certain *speed level interval*  $\Delta V$ , e.g., 10km/h. Thus, the mapping gives rise to a time series of discrete states of speed for each road segment, which allows us to do some analysis of discrete time series to reveal intrinsic traffic patterns. The dynamical behavior of a whole ring can then be quantified by that of all segments of it. Figure 1 shows the transition probabilities between different ranges of speed, namely, speed states. We find that on average, a speed state is more likely to remain unchanged or shift to its nearby states rather than change to a distant state. These observations imply the existence of a potentially stable transition pattern that may facilitate the prediction of traffic condition and congestion from historical records.

We exploit information entropy [22] to quantify the uncertainty of speed transition and the degree of predictability characterizing the time series of the speed at each segment. By following Ref. [24], we assign three entropy measures to each road segment's traffic pattern: (i) *Random Entropy*  $S_i^{\text{rand}}$ . Random entropy is defined as  $S_i^{\text{rand}} = \log_2 N_i$  where  $N_i$  is the number of distinct states, or speed levels, reached by road segment  $i$ . (ii) *Temporal-uncorrelated Entropy*  $S_i^{\text{unc}}$ . Temporal-uncorrelated entropy is defined as  $S_i^{\text{unc}} = -\sum_{j=1}^{N_i} p_i(j) \log_2 p_i(j)$ , where  $p_i(j)$  is the probability that the state  $j$  is reached by the road segment  $i$ . (iii) *Actual Entropy*  $S_i$ . Actual entropy is defined as  $-\sum_{T'_i \subset T_i} P(T'_i) \log_2 [P(T'_i)]$ , where  $T_i = \{X_1, X_2, \dots, X_L\}$  denotes the sequence of states road segment  $i$  reached in observation.  $P(T'_i)$  is the probability of finding the time-ordered subsequence  $T'_i$  in the state transition sequence of the segment  $i$ . It is noteworthy that the random entropy  $S_i^{\text{rand}}$  reflects the degree of predictability of a road segment's state transition based on the assumption that each state is visited with equal probability. For the temporal-uncorrelated entropy  $S_i^{\text{unc}}$ , it takes the heterogeneity in the probability into account, but omits the order of the transition. In contrast, the actual entropy  $S_i$  by considering both heterogeneous probability and temporal correlation offers more realistic characterization of the traffic pattern.

The sufficient data with high record frequency on the three ring roads allow us to calculate the actual entropy  $S_i$  that in principle requires a continuous record of a road segment's momentary state. As shown in Fig. 2(a), we can see remarkable difference between  $P(S)$  and  $P(S^{\text{rand}})$ . To be concrete,  $S^{\text{rand}}$  peaks at about 2.6, indicating that on average each update of the speed state represents 2.6 bits per hour new information. In other words, the new speed level could be found in average  $2^{2.6} \approx 6$  states. In contrast, the fact that  $P(S)$  of the actual entropy peaks at  $S = 0.9$  demonstrates that the real uncertainty in a segment's speed state is  $2^{0.9} \approx 1.87$  rather than 6.

The entropy of a segment's speed allows us to measure the predictability  $\Pi$  that a suitable

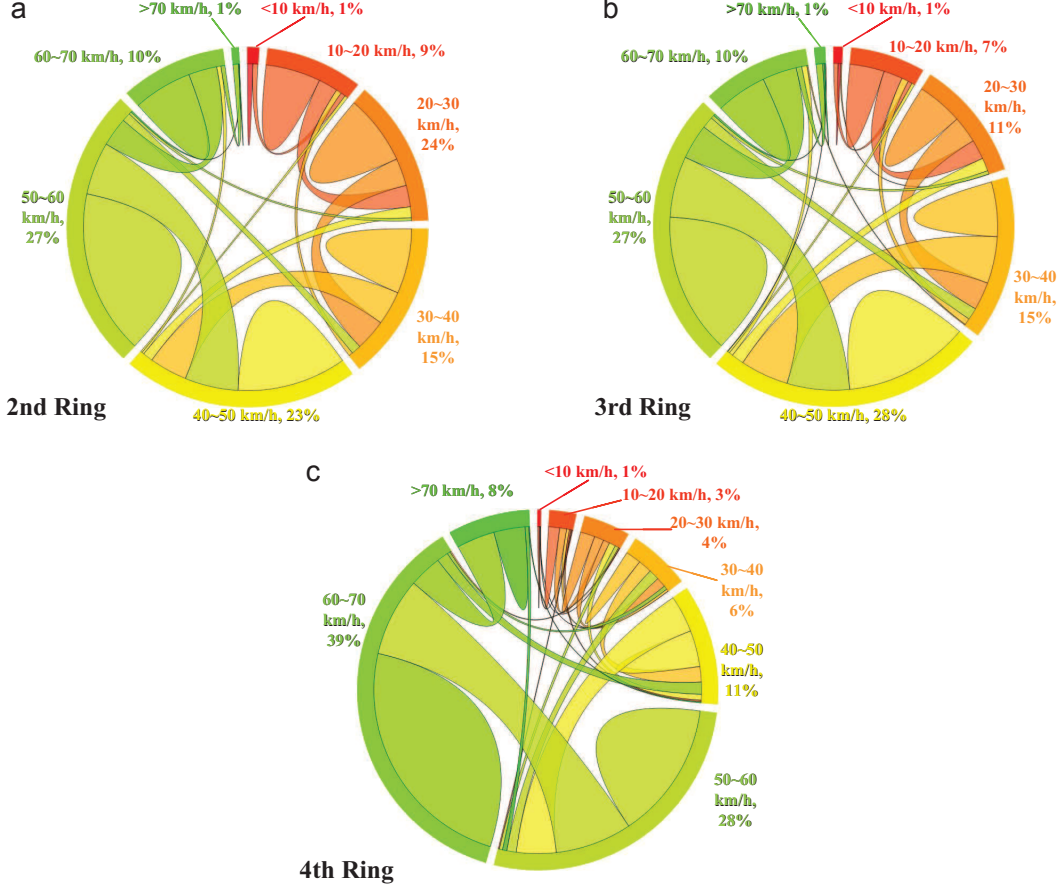


FIG. 1: **Transition probability of speed states.** (a)-(c) Transition probability between different speed states in the 2nd (a), 3rd (b) and 4th Ring Roads of Beijing. The speed  $V$  between 10km/h and 70km/h is divided into 6 states with equal speed interval  $\Delta V = 10\text{km/h}$ . Due to rare observations for  $V > 70\text{km/h}$ , they are set to be two states respectively, without any further partitions. For each Ring Road, the result is obtained by averaging over all road segments with equal length  $\Delta L = 1\text{km}$ . We see that for each state, remaining unchanged and shifting to its adjacent states constitute a very large proportion, implying a potential stable regulation in the traffic patterns.

predictive algorithm can correctly predict the segment's future speed state. In analogy with Ref.[24], the predictability measure is subject to Fano's inequality. Specifically, if the speed level of a single road segment is updated in  $N$  states with the time, then its predicability  $\Pi \leq \Pi^{\max}(S, N)$ , where  $\Pi^{\max}$  could be acquired by solving

$$S = H(\Pi^{\max}) + (1 - \Pi^{\max}) \log_2(N - 1),$$

where  $H(\Pi^{\max})$  represents the binary entropy function, namely

$$H(\Pi^{\max}) = -\Pi^{\max} \log_2(\Pi^{\max}) - (1 - \Pi^{\max}) \log_2(1 - \Pi^{\max}).$$

For a road segment with  $\Pi^{\max} = 0.1$ , we could predict its state transition accurately only in 10% of the cases. An equivalent statement is that 10% is the upper bound of probability

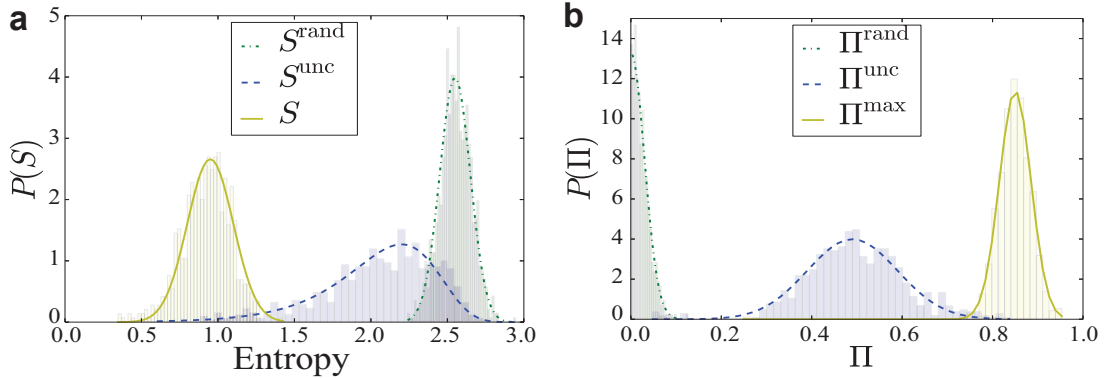


FIG. 2: **Distributions of entropy and probability  $\Pi$ .** (a) The distribution of the random entropy  $S^{\text{rand}}$ , the uncorrelated entropy  $S^{\text{unc}}$  and the entropy  $S_i$  of road segments in the 2nd Ring Road in Beijing. (b) The distribution of the  $\Pi^{\text{rand}}$ , the  $\Pi^{\text{unc}}$  and the  $\Pi^{\text{max}}$  across all road segments. The road segments are of identical length  $\Delta L = 1\text{km}$  and the interval of speed state is  $\Delta V = 10\text{km/h}$ . The 3rd and 2th Ring Roads show similar results of  $P(S)$  and  $P(\Pi)$  to that of the 2nd Ring Roads.

for any algorithms attempting to predict the segment's speed state transition. Since we calculate  $\Pi^{\text{max}}$  based on  $S^{\text{rand}}$ ,  $S^{\text{unc}}$  and  $S$ , the result is encouraging. We found that under the condition where  $\Delta L = 1\text{km}$  and  $\Delta V = 20\text{km/h}$ , the predictability of the 2nd Ring Road segments is narrowly peaked approximately at 0.83, indicating that it is theoretically possible to predict the transition of speed status in 83% of the cases. This high predictability with bounded distribution indicates that, despite the diversity of drivers' origins, destinations, their routing decisions and adaptive behaviors, strikingly the traffic pattern as a collective behavior of a large number of drivers is of high degree of potential predictability exclusively from the historical records of daily traffic patterns in the absence of any individual level information. We have also explored the maximum predictability  $\Pi^{\text{unc}}$  and  $\Pi^{\text{rand}}$  based on  $S^{\text{unc}}$  and  $S^{\text{rand}}$ , as shown in Fig. 2(b). We see that both maxima in  $P(\Pi^{\text{unc}})$  and  $P(\Pi^{\text{rand}})$  are much lower than that of  $P(\Pi^{\text{max}})$ , manifesting that  $\Pi^{\text{max}}$  is a much better predictive tool than the other two and the temporal order of traffic pattern contains significant information for precisely predicting future patterns.

We further explore how the settings of the road segment length  $\Delta L$  and speed level interval  $\Delta V$  affect the predictability. As shown in Fig. 3, except very small  $\Delta V$  and very short  $\Delta L$ , quite high average predictability is observed. This provides strong evidence for the generally high predictability of traffic condition of the three ring roads. Figure. 3 also shows that the predictability decreases with decrease of road segment length and speed level interval. Because the shorter road segment length and smaller speed level interval mean higher prediction granularity, the phenomenon that higher prediction accuracy corresponds lower predictability limit meets the intuition. The relatively low predictability for extreme cases is ascribed to the relatively big fluctuations in the average speed resulted from insufficient records. For example, for a road segment with very short length, the probability of finding a taxi in it within a certain time interval will be low. In other words, in this scenario, the data record of taxis will become insufficient to capture the actual average speed in the segment,

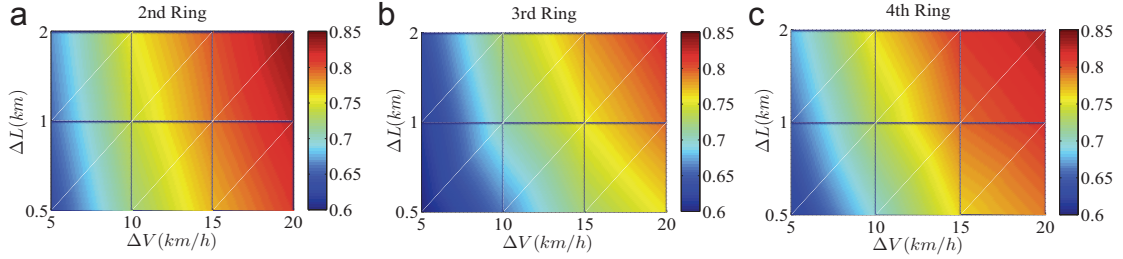


FIG. 3: **Predictability of the three Ring Roads.** (a)-(c) The dependence of the maximum value  $\Pi^{\max}$  on  $\Delta L$  and  $\Delta V$  for the 2nd (a), the 3rd (b) and the 4th (c) Ring Road. The color bars represent the values of  $\Pi^{\max}$ . The results for each Ring Road are the average over all road segments in the Ring Road.

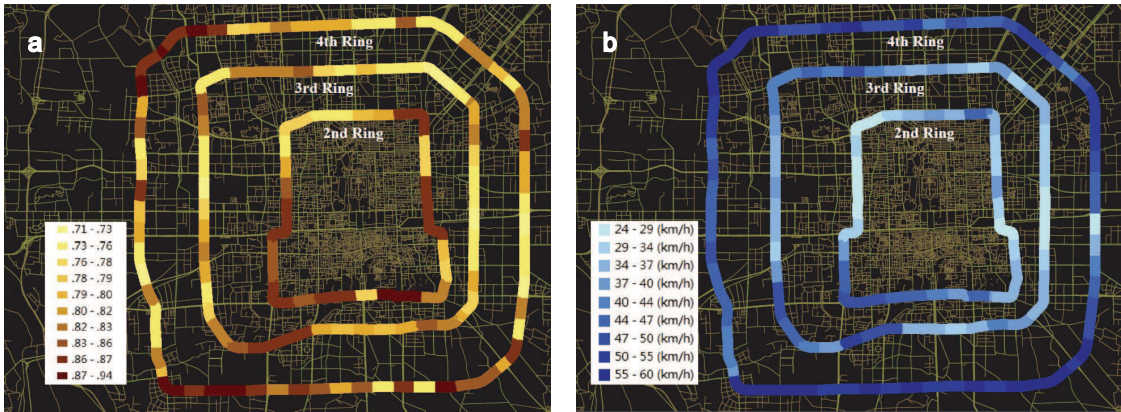


FIG. 4: **Local predictability and average speed.** (a) The local predictability of road segments in the three Ring Roads. (b) The local average speed of road segments in the three Ring Roads. In (a), the color bar represents the maximum value  $\Pi^{\max}$  of road segments and In (b), the color bar represents the average speed of road segments.

accounting for the big fluctuation of speed and inaccurate reflection of the traffic pattern in the segment. Similarly, for small  $\Delta V$ , the insufficient data subject to each speed state is incapable of characterizing the real situation, leading to the specious low predictability. Nevertheless, based on our findings, insofar as the records are adequate to measure traffic conditions, the traffic pattern is highly predictable, regardless of the settings of the road segment length and speed interval.

Although the traffic pattern of the three ring roads on average is highly predictable, there are certain variations between different segments. Figure 4(a) shows the local predictability of each segment on the map. We find that the local predictability is correlated with the average local speed (Fig. 4(b)), prompting us to investigate the correlation between them. Interestingly, we observe a non-monotonic correlation between the local predictability and average speed with the lowest predictability arising at intermediate speed, as shown in Fig. 5(a) and 5(b). As a result, we also find that it is most difficult to predict the traffic condition of the 3rd ring road, due to its intermediate average speed compared to the 2nd and 4th ring roads. A heuristic explanation for this phenomenon can be provided with

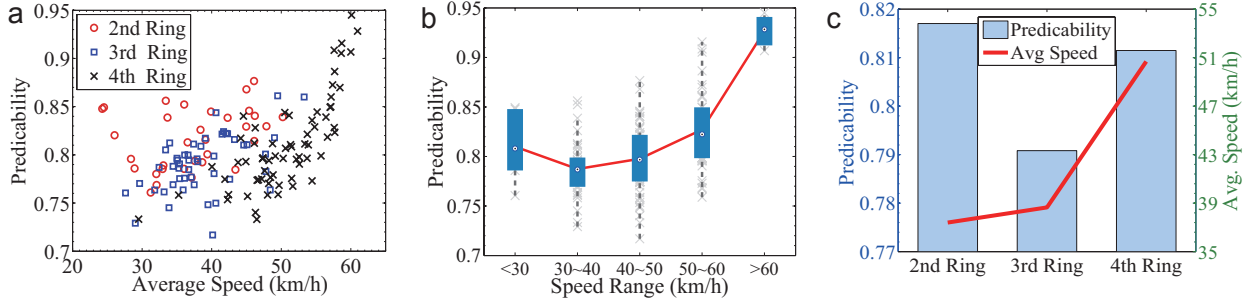


FIG. 5: **Relationship between predictability and average speed.** (a) Predictability as measured by  $\Pi^{\max}$  as a function of the average speed for the three Ring Roads. (b) Box plots of the predictability in different ranges of the average speed. (c) The predictability and the average speed of each entire Ring Road. The results are obtained for  $\Delta L = 1\text{Km}$  and  $\Delta V = 10\text{km/h}$ . The bars represent mean predictability of 2nd, 3rd and 4th ring roads.

respect to the variational direction of speed. Suppose that in a segment all the vehicles is fully stopped because of heavy congestion. One minute later, remaining stopped or starting to pull away are the only two possible situations. Let's consider another extreme case in which all vehicles are moving along the speed limit of a road without any congestion. One minute later, there are also only two possible scenarios, i.e., their speeds remain unchanged or reduce because of some suddenly emerged congestions. In contrast to the extreme cases, for a car with intermediate speed, the car may accelerate, decelerate or keep its current speed some time later, relying on what happens in the near future. Therefore, due to more variant possibilities of intermediate speed compared to that of low and high speed, the traffic condition of a segment with intermediate average speed is relatively most difficult to be predicted.

To gain a deeper understanding of the predictability of traffic patterns, we explore the effect of commuter demand on daily traffic predictability in terms of the comparison between weekdays and weekends. It is intuitive that the commuter demand during weekdays may induce quite different traffic patterns and congestion distribution compared to that at weekends. However, to our surprise, despite these obvious difference, we find that the daily traffic patterns in a week are of very similar predictability, nearly regardless of the commuter demand, as shown in Fig. 6. These striking results suggest that both weekdays and weekends have their specific inherent patterns encoded in the historical records, accounting for the relatively high and similar predictability.

Next, we explore the probability of inferring the state of a segment from the state series of its adjacent segments. This problem is related to the observability that in the control theory is defined as if a system's state can be fully recovered from a set of observable quantities [27]. To the urban road traffic, inferring traffic condition at some locations from the observation of the other segments has important applications in monitoring and controlling traffic in real time from a limited number of speed detectors. In analogy with the predictability, we calculate the inference probability  $\tilde{\Pi}$  of a segment based on the information entropy and the Fano's inequality. However, different from the predictability, here the information entropy is calculated by  $S'_i = -\sum_{R'_i \subset R_i} P(R'_i) \log_2[P(R'_i)]$ , where  $R_i = \{X_1, X_2, \dots, X_L\}$  denotes the



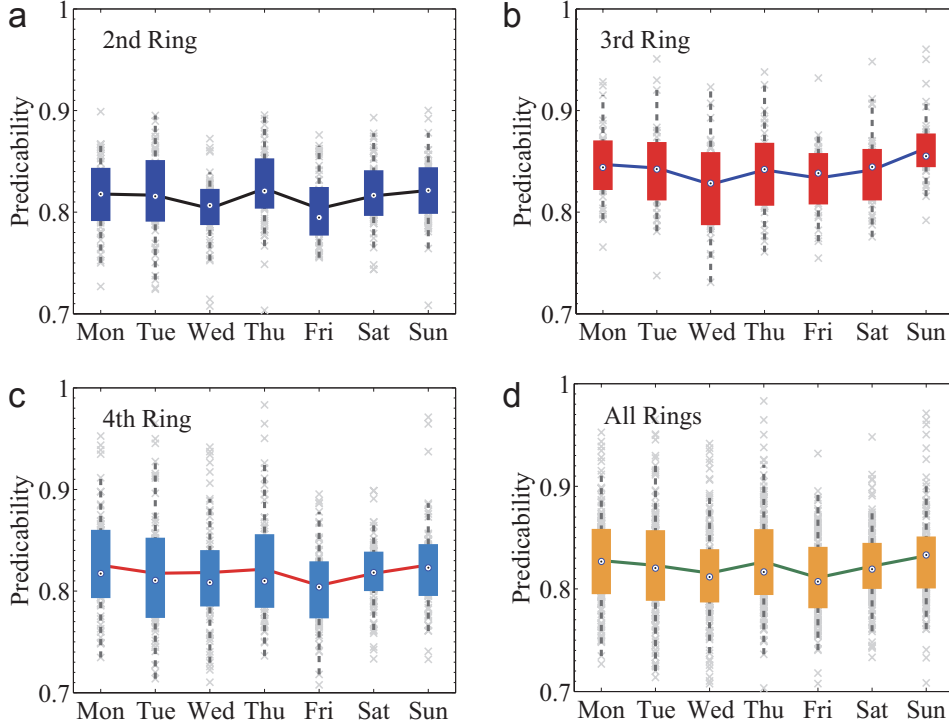


FIG. 6: **Daily predictability.** (a)-(c), The daily predictability during a week of the 2nd (a), 3rd (b) and 4th (c) Ring Road. (d) The daily predictability averaging over all of the three Ring Roads during a week. The parameter values and the box plots are the same as in Fig. 5.

states observed within a single time interval of  $L$  road segments connected in a sequence, and  $P(R'_i)$  is the probability of finding the subsequence  $R'_i$  in this sequence. Similarly, by solving  $S' = H(\tilde{\Pi}^{\max}) + (1 - \tilde{\Pi}^{\max}) \log_2(N - 1)$ , we get an upper bound  $\tilde{\Pi}^{\max}$  which captures the inference probability of the traffic pattern of a road segment from its observable adjacent segments.

As shown in Fig. 7, we see that the inference probability increases as the amount of segments increases for all the three ring roads. This phenomenon can be heuristically explained as follows. For sufficiently short segment lengths (sufficient number of segments), the average vehicle speed in a segment will be sufficiently close to that in its adjacent segments, enabling an accurate inference of the segment's state by trivially using that in its neighborhood. The increment of segment length induces more difference between adjacent segments, rendering the inference more difficult. As a result, the inference probability is an increase function of the amount of segments. More importantly, our results provide a quantitative understanding of the inference probability in terms of number of segments, which is valuable for determining the density of speed detectors installed so as to infer the traffic condition of the entire road in real time in certain accuracy. In addition, we also find that the inference probability of the 3rd ring road exhibits the lowest values compared to the 2nd and 4th ring roads, which is the same as the predictability rank of the three ring roads, i.e., the 3rd ring road is of the lowest predictability. This suggests that the average vehicle speed plays similar role in both predictability and inference probability, which deserves deeper explorations.

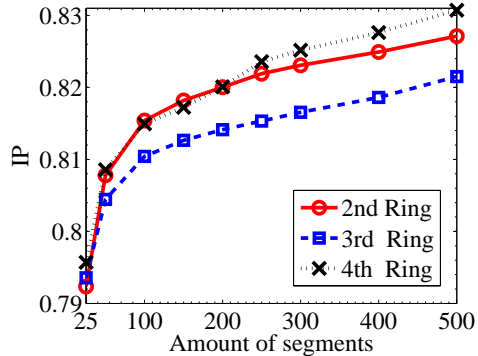


FIG. 7: **Inference probability.** The inference probability (IP) as a function of the amounts of segments for all the three Ring Roads. Here IP is measured by the upper bound of  $\tilde{\Pi}^{\max}$ .

## II. DISCUSSIONS

In summary, using the GPS records of vehicles to capture the traffic patterns on urban roads in the combination of entropy and Fano’s inequality demonstrates that daily traffic pattern in the three major ring roads in Beijing is high predictable by relying only on short-time historical records, without any a priori knowledge of drivers’ origins and destinations, driving habits, navigation strategies, and adaptive behaviors. We have also found that despite the apparently different traffic patterns in weekdays from that in weekends, where the former is highly affected by the commuter demand, their traffic patterns exhibit similarly high predictability. This result indicates that each day has its specific inherent regularity and traffic pattern encoded in the historical records. Another striking finding is that the local predictability is non-monotonically correlated with the average velocity and the lowest predictability arises at intermediate velocity. Consequently, the traffic condition of the 3rd ring road due to its intermediate average velocity compared to the 2nd and 4rd ring roads, is most difficult to be predicted. We have provided a heuristic explanation for this counterintuitive phenomenon. Furthermore, the probability of inferring the traffic pattern of an inaccessible road segment from the state series of its adjacent segment is explored by using entropy and Fano’s inequality, which is important for monitoring the traffic condition of the entire road network with respect to the limits of our ability to observe every location in real time.

All of these findings are valuable for the development of predictive models and algorithms for achieving actual predictions of traffic condition in real time based solely on short-time historical records, without the need of individual-level information that in principle is impossible to be fully accessed. Relying on the successful prediction of traffic patterns, it is feasible to implement effective control to release and prohibit congestions by exploiting traditional approaches in traffic engineering [28] and the recently developed controllability theory for complex networks [29, 30]. Urban road network as a typical complex networked system exhibits a variety of dynamical behaviors, such as the phantom jam and the diffusion of congestion [11]. Thus, it is imperative to control the road network as a whole in virtue of the controllability framework rather than controlling a single road or area individually. Our approach gains new insight into mitigating increasingly severe congestions in urban areas by

combining “big data” and the tools in information theory and for time series analysis. Further effort, we hope, will be inspired toward predicting traffic pattern and devising effective strategies to alleviate traffic congestion in urban areas.

### III. METHODS

We use OpenStreetMap [31] to extract all roads in the spatial range of Beijing from available database. We then retrieve the trajectories of vehicles. The data set that we used contains the trajectories of 20000 taxis recorded every minute within a month in Beijing. For each record, the location (the latitude and longitude), the direction, the state (whether there are any passengers in the taxi), the time stamp and the velocity updated in every minutes are included. Because of the inevitable error in the GPS locating process, all the records are preprocessed to match the GPS trajectories to the road by exploiting the ST-Matching algorithm [32]. After that, each GPS record is mapped to a road segments of OpenStreetMap.

To be concrete, ST-Matching algorithm of Ref. [32] is implemented via four steps: (i) *Candidate Preparation*. Firstly, for each GPS record point, the ST-Matching algorithm retrieves a set of candidate road segments within a fixed radius  $r$ , which is set to be 20 meters. For the points without any candidates within  $r$ , the algorithm discards them as invalid records. (ii) *Spatial Analysis*. The algorithm next evaluates the given candidate segments by using “observation probability” and “transmission probability” to express the geometric and topological information of each candidate segments and the spatial relationship between them. This step gives rise to the spatial analysis function  $F_s(c_{i-1}^t \rightarrow c_i^s)$ , which is simply the product of the observation probability and transmission probability. In this function,  $c_i^s$  represents the  $s$ th candidate segment of the  $i$ th GPS sampling record. This function measures the probability that the  $i$ th record is on  $c_i^s$ , given an assumed real segment mapping of the  $(i-1)$ th record, that is  $c_{i-1}^t$ . (iii) *Temporal Analysis*. The ST-Matching algorithm exploits the temporal analysis function  $F_t(c_{i-1}^t \rightarrow c_i^s)$  to further incorporate the temporal features into the map-matching process. This step is available for the situation that only spatial analysis could not handle. Specifically, if the trajectory of a vehicle lies between a freeway and a service road, and it moves in a relatively high speed, then more likely it is that the vehicle is on the freeway. (iv) *Result Matching*. Finally, after  $F_s(c_{i-1}^t \rightarrow c_i^s)$  and  $F_t(c_{i-1}^t \rightarrow c_i^s)$  is computed, the algorithm uses the ST-function to evaluate each candidate segments, that is  $F(c_{i-1}^t \rightarrow c_i^s) = F_s(c_{i-1}^t \rightarrow c_i^s) \times F_t(c_{i-1}^t \rightarrow c_i^s)$ ,  $2 \leq i \leq n$ . Thus, the problem is converted to finding a path with the highest ST-function value, given the candidates for all sampling points.

After the map-matching process, each point is assigned with an attribute which represents the road segment that the point is on. Based on the work before, we could generate the time series of each road segment’s speed states.

- 
- [1] Brockmann, D., Hufnagel, L. & Geisel, T. The scaling laws of human travel. *Nature* **439**, 462–465 (2006).
  - [2] Belik, V., Geisel, T. & Brockmann, D. Natural human mobility patterns and spatial spread of infectious diseases. *Phys. Rev. X* **1**, 011001 (2011).
  - [3] Kölbl, R. & Helbing, D. Energy laws in human travel behaviour. *New. J. Phys.* **5**, 48 (2003).
  - [4] Gonzalez, M. C., Hidalgo, C. A. & Barabasi, A.-L. Understanding individual human mobility patterns. *Nature* **453**, 779–782 (2008).

- [5] Schrank, D. *Urban Mobility Report (2004)* (DIANE Publishing, 2008).
- [6] Helbing, D. A section-based queueing-theoretical traffic model for congestion and travel time analysis in networks. *J. Phys. A: Math. Gen.* **36**, L593 (2003).
- [7] Chin, A. T. Containing air pollution and traffic congestion: transport policy and the environment in singapore. *Atmos. Environ.* **30**, 787–801 (1996).
- [8] Rosenlund, M. *et al.* Comparison of regression models with land-use and emissions data to predict the spatial distribution of traffic-related air pollution in rome. *J. Expo. Sci. Env. Epid.* **18**, 192–199 (2008).
- [9] Zhang, X. *et al.* Atmospheric aerosol compositions in china: spatial/temporal variability, chemical signature, regional haze distribution and comparisons with global aerosols. *Atmos. Chem. Phys. Discuss.* **12**, 779–799 (2012).
- [10] Zheng, Y., Liu, F. & Hsieh, H.-P. U-air: when urban air quality inference meets big data. In *Proceedings of SIGKDD'13*, 1436–1444 (ACM, 2013).
- [11] Helbing, D. Traffic and related self-driven many-particle systems. *Rev. Mod. Phys.* **73**, 1067 (2001).
- [12] Batty, M. The size, scale, and shape of cities. *Science* **319**, 769–771 (2008).
- [13] Barthélemy, M. Spatial networks. *Phys. Rep.* **499**, 1–101 (2011).
- [14] Herrera, J. C. *et al.* *Dynamic estimation of OD matrices for freeways and arterials* (Institute of Transportation Studies, UC Berkeley, 2007).
- [15] Herrera, J. C. *et al.* Evaluation of traffic data obtained via GPS-enabled mobile phones: The Mobile Century field experiment. *Transport. Res. C-Emer.* **18**, 568–583 (2010).
- [16] Wynter, L. & Shen, W. Real-time traffic prediction using GPS data with low sampling rates: A hybrid approach. In *Transportation Research Board 91st Annual Meeting*, 12-1692 (2012).
- [17] Thomas, J. M. & Darnton, J. Social diversity and economic development in the metropolis. *J. Plan. Liter.* **21**, 153–168 (2006).
- [18] Mayer-Schönberger, V. & Cukier, K. *Big data: A revolution that will transform how we live, work, and think* (Houghton Mifflin Harcourt, 2013).
- [19] Hornsey, R. *'He who Thinks, in Modern Traffic, is Lost': Automation and the Pedestrian Rhythms of Interwar London* (Ashgate, 2010).
- [20] Wang, P., Hunter, T., Bayen, A. M., Schechtner, K. & González, M. C. Understanding road usage patterns in urban areas. *Sci. Rep.* **2** (2012).
- [21] Wang, J., Wei, D., He, K., Gong, H. & Wang, P. Encapsulating urban traffic rhythms into road networks. *Sci. Rep.* **4** (2014).
- [22] Brabazon, A., O'Neill, M. & Maringer, D. *Natural computing in computational finance* (Springer, 2008).
- [23] Cover, T. M. & Thomas, J. A. *Elements of information theory* (John Wiley & Sons, 2012).
- [24] Song, C., Qu, Z., Blumm, N. & Barabási, A.-L. Limits of predictability in human mobility. *Science* **327**, 1018–1021 (2010).
- [25] Jia, T. & Jiang, B. Exploring human activity patterns using taxicab static points. *ISPRS International Journal of Geo-Information* **1**, 89–107 (2012).
- [26] [http://en.wikipedia.org/wiki/ring\\_roads\\_of\\_beijing](http://en.wikipedia.org/wiki/ring_roads_of_beijing).
- [27] Hautus, M. Controllability and observability conditions of linear autonomous systems. *Ned. Akad. Wetenschappen, Proc. Ser. A* **72**, 443 (1969).
- [28] Ortuzar, J. d. & Willumsen, L. G. *Modelling transport* (John Wiley & Sons, 2011).

- [29] Liu, Y.-Y., Slotine, J.-J. & Barabási, A.-L. Controllability of complex networks. *Nature* **473**, 167–173 (2011).
- [30] Yuan, Z., Zhao, C., Di, Z., Wang, W.-X. & Lai, Y.-C. Exact controllability of complex networks. *Nat. Commun.* **4** (2013).
- [31] <http://www.openstreetmap.org>.
- [32] Lou, Y. *et al.* Map-matching for low-sampling-rate GPS trajectories. In *Proceedings of ACM SIGSPATIAL'17*, 352–361 (ACM, 2009).