

Forecasting Trajectory and Behavior of Road-Agents Using Spectral Clustering in Graph-LSTMs

Rohan Chandra Tianrui Guan Srujan Panuganti
 Trisha Mittal Uttaran Bhattacharya Aniket Bera Dinesh Manocha
 (University of Maryland, College Park, USA)

Supplemental version including Code, Video, Datasets at <https://gamma.umd.edu/spectralcows/>

Abstract— We present a novel approach for traffic forecasting in urban traffic scenarios using a combination of spectral graph analysis and deep learning. We predict both the low-level information (future trajectories) as well as the high-level information (road-agent behavior) from the extracted trajectory of each road-agent. Our formulation represents the proximity between the road agents using a weighted dynamic geometric graph (DGG). We use a two-stream graph-LSTM network to perform traffic forecasting using these weighted DGGs. The first stream predicts the spatial coordinates of road-agents, while the second stream predicts whether a road-agent is going to exhibit overspeeding, underspeeding, or neutral behavior by modeling spatial interactions between road-agents. Additionally, we propose a new regularization algorithm based on spectral clustering to reduce the error margin in long-term prediction (3-5 seconds) and improve the accuracy of the predicted trajectories. Moreover, we prove a theoretical upper bound on the regularized prediction error. We evaluate our approach on the Argoverse, Lyft, Apolloscape, and NGSIM datasets and highlight the benefits over prior trajectory prediction methods. In practice, our approach reduces the average prediction error by more than 75% over prior algorithms and achieves a weighted average accuracy of 91.2% for behavior prediction. Additionally, our spectral regularization improves long-term prediction by up to 70%.

I. INTRODUCTION

Autonomous driving is an active area of research and includes many issues related to navigation [1], trajectory prediction [2], and behavior understanding [3], [4]. Trajectory prediction is the problem of predicting the short-term (1-3 seconds) and long-term (3-5 seconds) spatial coordinates of various road-agents such as cars, buses, pedestrians, rickshaws, and even animals, etc. Accurate trajectory prediction is crucial for safe navigation. Furthermore, road-agents have different dynamic behaviors that may correspond to aggressive or conservative driving styles [5], [6], [7]. While humans can very quickly predict different road-agent behaviors commonly observed in traffic, current autonomous vehicles (AVs) are unable to perform efficient navigation in dense and heterogeneous traffic due to their inability to recognize road-agent behaviors.

While there has been extensive progress in trajectory prediction [2], [8], [9], there has been significantly less research in behavior prediction. The advantage of knowing if a neighboring road-agent is going to overtake another agent or if a road-agent in front is going to brake suddenly is useful for safe navigation. Furthermore, behavior prediction is crucial for making autonomous vehicles socially aware, as opposed to their inherent conservative behavior [10], [11], [12] that

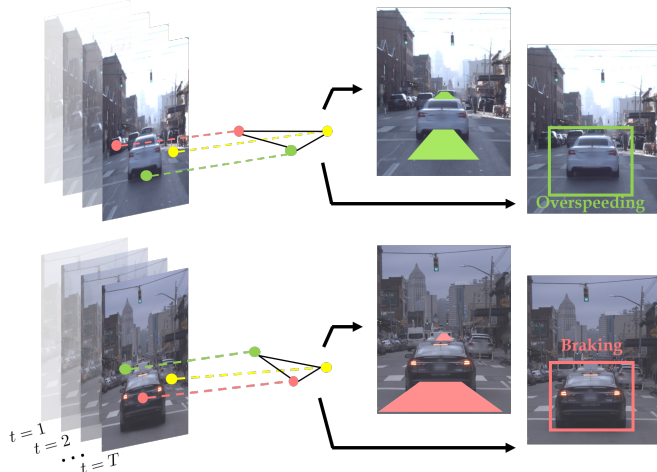


Fig. 1: **Trajectory and Behavior Prediction:** We predict the long-term (3-5 seconds) trajectories of road-agents, as well as their behavior (e.g. overspeeding, underspeeding, etc.), in urban traffic scenes. Our approach represents the spatial coordinates of road-agents (colored points in the image) as vertices of a DGG to improve long-term prediction using a new regularization method.

poses new risks in terms of low efficiency and uncomfortable traveling experiences [13].

Furthermore, a major challenge in traffic forecasting is ensuring accurate long-term prediction (3-5 seconds). As the prediction horizon increases, the temporal correlations in the data between current and previous time-steps grow weaker, which increases the error-margin of long-term prediction ([14], cf. Figure 4 in [8], [2], Figure 3 in [15]). Some approaches have been developed to reduce the long-term error-margin for trajectory forecasting [14], but they assume knowledge of high-order, non-linear traffic dynamics.

a) Main Contributions: We present an algorithm for traffic forecasting that disjointedly predicts trajectories as well as road-agent behavior using two separate streams. We represent the inter-road-agent interactions in the traffic using weighted dynamic geometric graphs (DGGs) [16], where the vertices represent the road-agents, and the weighted edges are a function of the proximity between the road-agents. Our approach makes no assumptions about the size and shape of the road-agents. Our main contributions include:

- 1) A two-stream graph-LSTM network for traffic forecasting in urban traffic. The first stream is a conventional LSTM encoder-decoder network that does not account for neighbor vehicles. It is used to predict the spatial co-

ordinates of the future trajectory. We propose a second stream that predicts the eigenvectors of future DGGs, which serve the dual purpose of behavior prediction as well as regularizing the first stream.

- 2) To reduce the error of long-term predictions, we propose a new regularization algorithm for sequence prediction models called spectral cluster regularization.
- 3) We derive a theoretical upper bound on the prediction error of the regularized forecasting algorithm in the order of $\mathcal{O}(\sqrt{N}\delta_{max})$, where N is the number of road-agents and δ_{max} value corresponds to the distance between the two closest road-agents.
- 4) We present a rule-based behavior prediction algorithm to forecast whether a road-agent is overspeeding (aggressive), underspeeding (conservative), or neutral, based on the traffic behavior classification in psychology literature [17], [18].

We evaluate our approach on four large-scale urban driving datasets – NGSIM, Argoverse, Lyft, and Apolloscape. We also perform an exhaustive comparison with the SOTA trajectory prediction methods and report an average RMSE (root mean square error) reduction of at least 75% with respect to the next best method. We also achieved a weighted average accuracy of 91.2% for behavior prediction. Our regularization algorithm improves long-term prediction by up to 70%.

II. RELATED WORK

Here, we discuss prior work in trajectory prediction, road-agent behavior prediction, and traffic forecasting.

A. Trajectory Prediction

Trajectory prediction is a well-known problem in statistics [19], signal processing [20], and controls and systems engineering [21]. These approaches, however, rely on the knowledge of certain parameters that may not be readily available in traffic videos. In such instances, data-driven methods such as deep learning have become the SOTA for designing trajectory prediction algorithms.

There is some research on trajectory prediction. Deo et al. [8] combined LSTMs with Convolutional Neural Networks (CNNs) to predict the trajectories of vehicles on sparse U.S. highways. Chandra et al. [2], [9] proposed algorithms to predict trajectories in urban traffic with high density and heterogeneity. For traffic scenarios with moderate density and heterogeneity, Ma et al. [22] proposed a method based on reciprocal velocity obstacles. Some additional deep learning-based trajectory prediction methods include [23], [24]. However, these methods only capture road-agent interactions inside a local grid, whereas graph-based approaches such as GRIP [15] for trajectory prediction of road-agents and [25], [26], [27], [28] for traffic density prediction consider all interactions independent of local neighborhood restrictions. Our graph representation differs from that of GRIP by storing the graphs of previous time-steps (III-B). Using our representations, we propose a novel behavior prediction algorithm (IV-C). Additionally, unlike

other trajectory prediction methods in the literature, we propose a new Spectral Regularization-based loss function (IV-D) that automatically corrects and reduces long-term errors. This is a novel improvement over *all* prior prediction methods that do not handle long-term errors.

B. Socially-Aware Autonomous Driving

Current autonomous vehicles lack social awareness due to their inherent conservative behavior [10], [11], [12]. Overly conservative behavior present new risks in terms of low efficiency and uncomfortable traveling experiences [13]. Real-world examples of problems caused by AVs that are not socially adaptable can be seen in this video*. The notion of using driver behavior prediction to make the AVs socially aware is receiving attention [11].

Current driving behavior modeling methods are limited to traffic psychology studies where predictions for driving behavior are made offline, based on either driver responses to questionnaires or data collected over a period of time. Such approaches are not suitable for online behavior prediction. In the following section, we review some of these approaches. In contrast, our behavior prediction algorithm is the first computationally online approach that does not depend on offline data and manually tunable parameters.

C. Road-Agent Behavior Prediction

Many studies have been performed behavior modeling by identifying factors that contribute to different driver behaviors classes such as aggressive, conservative, or moderate driving. These factors can be broadly categorized into four categories. The first category of factors that indicate road-agent behavior is driver-related. These include characteristics of drivers such as age, gender, blood pressure, personality, occupation, hearing, and so on [17], [29], [30]. The second category corresponds to environmental factors such as weather or traffic conditions [31], [32]. The third category refers to psychological aspects that affect driving styles. These could include drunk driving, driving under the influence, state of fatigue, and so on [33], [34]. The final category of factors contributing to driving behavior are vehicular factors such as positions, acceleration, speed, throttle responses, steering wheel measurements, lane changes, and brake pressure [35], [36], [37], [38], [39], [40], [3].

A recent data-driven behavior prediction approach [3] also models traffic through graphs. The method predicts the driving behavior by training a neural network on the eigenvectors of the DGGs using supervised machine learning. The proposed behavior prediction algorithm in this paper extends the approach in [3] by predicting sequences of eigenvectors for future time-steps. Apart from behavior modeling, several methods have used machine learning to predict the intent of road-agents [41], [42].

*<https://www.youtube.com/watch?v=Rm8aPR0aMDE>

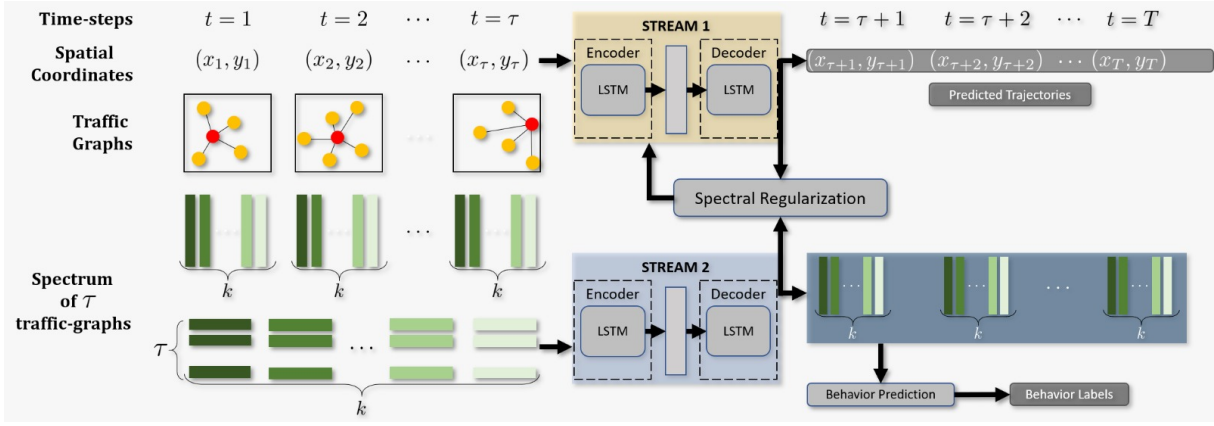


Fig. 2: **Network Architecture:** We show the trajectory and behavior prediction for the i^{th} road-agent (red circle in the DGGs). The input consists of the spatial coordinates over the past τ seconds as well as the eigenvectors (green rectangles, each shade of green represents the index of the eigenvectors) of the DGGs corresponding to the first τ DGGs. We perform spectral clustering on the predicted eigenvectors from the second stream to regularize the original loss function and perform back-propagation on the new loss function to improve long-term prediction.

III. BACKGROUND AND OVERVIEW

In this section, we define the problem statement and give a brief overview of spectral Dynamic Geometric Graphs (DGGs) in the context of road-agents.

A. Problem Statement

We first present a definition of a vehicle trajectory:

Definition III.1. Trajectory: The trajectory for the i^{th} road agent is defined as a sequence $\Psi_i(a, b) \in \{\mathbb{R}^2\}$, where $\Psi_i(a, b) = \{[x_t, y_t]^\top \mid t \in [a, b]\}$. $[x, y] \in \mathbb{R}^2$ denotes the spatial coordinates of the road-agent in meters according to the world coordinate frame and t denotes the time instance.

We define *traffic forecasting* as solving the following two problem statements, simultaneously, but separately using two separate streams.

Problem III.1. Trajectory Prediction: In a traffic video with N road agents, given the trajectory $\Psi_i(0, \tau)$, predict $\Psi_i(\tau^+, T)$ for each road-agent $v_i, i \in [0, N]$.

Problem III.2. Behavior Prediction: In a traffic video with N road agents, given the trajectory, $\Psi_i(0, \tau)$, predict a label from the following set, $\{\text{Overspeeding, Neutral, Underspeeding}\}$ for each road-agent $v_i, i \in [0, N]$.

B. Weighted Dynamic Geometric Graphs (DGGs)

We assume that the trajectories of all the vehicles in the video are provided to us as the input. Given this input, we first construct a DGG [16] at each time-step. In a DGG, the vehicles represent the vertices and the edge weights are a function of the euclidean distance between the vertices. This function [43] is given by,

$$f(v_i, v_j) = e^{-d(v_i, v_j)} \quad (1)$$

where v_i and v_j are the i^{th} and j^{th} vertices and d is the euclidean distance function.

We represent traffic at each time instance using a DGG \mathcal{G} of size $N \times N$, with the spatial coordinates of the road-agent

representing the set of vertices $\mathcal{V} = \{v_1, v_2, \dots, v_n\}$ and a set of undirected, weighted edges, \mathcal{E} . Two road-agents are said to be connected through an edge if $d(v_i, v_j) < \mu$, where $d(v_i, v_j)$ represents the Euclidean distance between the road-agents and μ is a heuristically chosen threshold parameter. In our experiments, we choose $\mu = 10$ meters, taking into account the typical size of road-agents and the width of the road.

For a DGG, \mathcal{G} , we define the symmetrical adjacency matrix, $A \in \mathbb{R}^{N \times N}$ as,

$$A(i, j) = \begin{cases} e^{-d(v_i, v_j)} & \text{if } d(v_i, v_j) < \mu, i \neq j, \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

Equation 1 denotes the interactions between any two road-agents, v_i and v_j . This implies that road-agents at a greater distance are assigned a lower weight, while road-agents in close proximity are assigned a higher weight. This follows the intuition that each road-agent needs to pay more attention to nearby agents than those farther away to avoid collisions.

For the adjacency matrix A at each time instance, the corresponding degree matrix $D \in \mathbb{R}^{N \times N}$ is a diagonal matrix with main diagonal $D(i, i) = \sum_{j=1}^N A(i, j)$ and 0 otherwise. The unnormalized Laplacian matrix $L = D - A$ of the graph is defined as the symmetric matrix,

$$L(i, j) = \begin{cases} D(i, i) & \text{if } i = j, \\ -e^{-d(v_i, v_j)} & \text{if } d(v_i, v_j) < \mu, \\ 0 & \text{otherwise.} \end{cases} \quad (3)$$

The Laplacian matrix for each time-step is correlated with the Laplacian matrices for all previous time-steps. Let the Laplacian matrix at a time instance t be denoted as L_t . Then, the laplacian matrix for the next time-step, L_{t+1} is given by the following update,

$$L_{t+1} = \left[\begin{array}{c|c} L_t & 0 \\ \hline 0 & 1 \end{array} \right] + \delta\delta^\top, \quad (4)$$

where $\delta\delta^\top$ is a perturbation matrix represented by an outer

product of rank 2. Here, $\delta \in \mathbb{R}^{(N+1) \times 2}$ is a sparse matrix $\|\delta\|_0 \ll N$, where N represents the total number of road-agents at time-step t . The presence of a non-zero entry in the j^{th} row of δ implies that the j^{th} road-agent has observed a new neighbor, that has now been added to the current DGG. During training time, the size of L_t is fixed for all time t and is initialized as a zero matrix of size $N \times N$, where N is max number of agents (different N is used for different datasets). For instance, $N = 270$ is used for Lyft Level 5 dataset. At current time t , if the $N < 270$, the zeros in L_t will simply be updated with new values. Once $N = 270$, L_t is reset to zero and the process repeats. During test time, trained models for stream 1 predict trajectories based only on past trajectories; these models for stream 1 do not use graphs. Trained model for stream 2, however, generate traffic-graphs in realtime for behavior prediction at test time. The matrix $U \in \mathbb{R}^{n \times k} := \{u_j \in \mathbb{R}^n | j = 1 \dots k\}$ of eigenvectors of L is called the *spectrum* of L , and can be efficiently computed using eigenvalue algorithms.

IV. TRAJECTORY AND BEHAVIOR FORECASTING

The overall flow of the approach is as follows:

- 1) Our input consists of the spatial coordinates over the past τ seconds as well as the eigenvectors of the DGGs corresponding to the first τ DGGs.
- 2) *Solving Problem III.1*: The first stream accepts the spatial coordinates and uses an LSTM-based sequence model [44] to predict $\Psi_i(\tau^+, T)$ for each $v_i, i \in [0, N]$, where $\tau^+ = \tau + 1$.
- 3) *Solving Problem III.2*: The second stream accepts the eigenvectors of the input DGGs and predicts the eigenvectors corresponding to the DGGs for the next τ seconds. The predicted eigenvectors form the input to the behavior prediction algorithm in Section IV-C to assign a behavior label to the road-agent.
- 4) Stream 2 is used to regularize stream 1 using a new regularization algorithm presented in Section IV-D. We derive the upper bound on the prediction error of the regularized forecasting algorithm in Section V.

A. Network Overview

We present an overview of our approach in Figure 2 and defer the technical implementation details of our network to the supplementary material. Our approach consists of two parallel LSTM networks (or streams) that operate separately.

Stream 1: The first stream is an LSTM-based encoder-decoder network [44] (yellow layer in Figure 2). The input consists of the trajectory history, $\Psi_i(0, \tau)$ and output consists of $\Psi_i(\tau^+, T)$ for each road-agent $v_i, i \in [0, N]$.

Stream 2: The second stream is also an LSTM-based encoder-decoder network (blue layer in Figure 2). To prepare the input to this stream, we first form a sequence of DGGs, $\{\mathcal{G}_t | t \in [0, \tau]\}$ for each time instance of traffic until time τ . For each DGG, \mathcal{G}_t , we first compute its corresponding Laplacian matrix, L_t and use SOTA eigenvalue algorithms to obtain the spectrum, U_t consisting of the top k eigenvectors of length n . We form k different sequences, $\{\mathcal{S}_j | j \in [0, k]\}$,

where each $\mathcal{S}_j = \{u_j\}$ is the set containing the j^{th} eigenvector from each U_t corresponding to the t^{th} time-step, with $|\mathcal{S}_j| = \tau$.

The second stream then accepts a sequence, \mathcal{S}_j , as input to predict the j^{th} eigenvectors for the next $T - \tau$ seconds. This is repeated for each \mathcal{S}_j . The resulting sequence of spectrums, $\{\mathcal{U}_t | t \in [\tau^+, T]\}$ are used to reconstruct the sequence, $\{\mathcal{L}_t | t \in [\tau^+, T]\}$, which is then used to assign a behavior label to a road-agent, as explained below.

B. Trajectory Prediction

The first stream is used to solve Problem III.1. We clarify at this point that stream 1 does not take into account road-agent interactions. We use spectral clustering (discussed later in Section IV-D) to model these interactions. It is important to further clarify that the trajectories predicted from stream 1 are not affected by the behavior prediction algorithm (explained in the next Section).

C. Behavior Prediction Algorithm

We define a rule-based behavior algorithm (blue block in Figure 2) to solve Problem III.2. This is largely due to the fact that most data-driven behavior prediction approaches require large, well-annotated datasets that contain behavior labels. Our algorithm is based on the predicted eigenvectors of the DGGs of the next τ seconds.

The degree of i^{th} road-agent, ($\theta_i \leq n$), can be computed from the diagonal elements of the Laplacian matrix L_t . θ_i measures the total number of distinct neighbors with which road-agent v_i has shared an edge connection until time t . As L_t is formed by simply adding a row and column to L_{t-1} , the degree of each road-agent monotonically increases. Let the rate of increase of θ_i be denoted as θ'_i . Intuitively, an aggressively overspeeding vehicle will observe new neighbors at a faster rate as compared to a road-agent driving at a uniform speed. Conversely, a conservative road-agent that is often underspeeding at unconventional spots such as green light intersections (Figure 1) will observe new neighbors very slowly. This intuition is formalized by noting the change in θ_i across time-steps. In order to make sure that slower vehicles (conservative) did not mistakenly mark faster vehicles as new agents, we set a condition where an observed vehicle is marked as ‘new’ if and only if the speed of the observed vehicle is less than the active vehicle (or ego-vehicle). To predict the behavior of the i^{th} road-agent, we follow the following steps:

- 1) Form the set of predicted spectrums from stream 2, $\{\mathcal{U}_t | t \in [\tau^+, T]\}$. We compute the eigenvalue matrix, Λ , of L_t by applying theorem 5.6 of [45] to L_{t-1} . We explain the exact procedure in the supplemental version.
- 2) For each $U_t \in \mathcal{U}$, compute $L_t = U_t \Lambda U_t^T$.
- 3) $\theta_i = i^{\text{th}}$ element of $\text{diag}(L_t)$, where ‘diag’ is the diagonal matrix operator.
- 4) $\theta'_i = \frac{\Delta \theta_i}{\Delta t}$.

where Λ is the eigenvalue matrix of L_t . Based on heuristically pre-determined threshold parameters λ_1 and λ_2 , we define the following rules to assign the final behavior label:

Overspeeding ($\theta' > \lambda_1$), Neutral ($\lambda_2 \leq \theta' \leq \lambda_1$), and Underspeeding ($\theta' < \lambda_2$).

Note that since human behavior does not change instantly at each time-step, our approach predicts the behavior over time periods spanning several frames.

D. Spectral Clustering Regularization

The original loss function of stream 1 for the i^{th} road-agent in an LSTM network is given by,

$$F_i = -\sum_{t=1}^T \log Pr(x_{t+1} | \mu_t, \sigma_t, \rho_t) \quad (5)$$

Our goal is to optimize the parameters, μ_t^*, σ_t^* , that minimize equation 5. Then, the next spatial coordinate is sampled from a search space defined by $\mathcal{N}(\mu_t^*, \sigma_t^*)$. The resulting optimization forces μ_t, σ_t to stay close to the next spatial coordinate. However, in general trajectory prediction models, the predicted trajectory diverges gradually from the ground-truth, causing the error-margin to monotonically increase as the length of the prediction horizon increases ([14], cf. Figure 4 in [8], [2], Figure 3 in [15]). The reason for this may be that while equation 5 ensures that μ_t, σ_t stays close to the next spatial coordinate, it does not, however, guarantee the same for $\hat{x}_{t+1} \sim \mathcal{N}(\mu_t, \sigma_t)$. Our solution to this problem involves regularizing equation 5 by adding appropriate constraints on the parameters, μ_t, σ_t , such that sampled coordinates from $\mathcal{N}(\mu_t^*, \sigma_t^*)$ are close to the ground-truth trajectory.

We assume the ground-truth trajectory of a road-agent to be equivalent to their ‘‘preferred’’ trajectory, which is defined as the trajectory a road-agent would have taken in the absence of other dynamic road-agents. Preferred trajectories can be obtained by minimizing the Dirichlet energy of the DGG, which in turn can be achieved through spectral clustering on the road-agents [46]. Our regularization algorithm (shown in the yellow arrow in Figure 2) is summarized below. For each road-agent, v_i :

- 1) The second stream computes the spectrum sequence, $\{U_{T+1}, \dots, U_{T+\tau}\}$.
- 2) For each U , perform spectral clustering [47] on the eigenvector corresponding to the second smallest eigenvalue.
- 3) Compute cluster centers from the clusters obtained in the previous step.
- 4) Identify the cluster to which v_i belongs and retrieve the cluster center, μ_c and deviation, σ_c .

Then for each road-agent, v_i , the regularized loss function, F_i^{reg} , for stream 1 is given by,

$$\sum_{t=1}^T \left(-\log Pr(\hat{y}_{t+1} | \mu_t, \sigma_t, \rho_t) \right) + b_1 \|\mu_t - \mu_c\|_2 + b_2 \|\sigma_t - \sigma_c\|_2 \quad (6)$$

where $b_1 = b_2 = 0.5$ are regularization constants. The regularized loss function is used to backpropagate the weights corresponding to μ_t in stream 1. Note that F_i^{reg} resembles a Gaussian kernel. This makes sense as the Gaussian kernel models the Euclidean distance non-linearly – greater the Euclidean distance, smaller the Gaussian kernel value and vice versa. This behavior is similarly captured by Equation 1).

Furthermore, we can use Equation 6 to predict multiple modes[8] by computing maneuver probabilities using μ, σ following the approach in Section 4.3 of [8].

V. UPPER BOUND FOR PREDICTION ERROR

In this section, we derive an upper bound on the prediction error, ϕ_j , of the first stream as a consequence of spectral regularization. We present our main result as follows,

Theorem V.1. $\phi_j \leq \frac{\|\delta_t \delta_t^\top\|_2}{\min(\lambda_j, \Lambda)}$, where $\min(\lambda_j, \Lambda)$ denotes the minimum distance between λ_j and $\lambda_k \in \Lambda \setminus \lambda_j$.

Proof. At time instance t , the Laplacian matrix, L_t , its block form, $\begin{bmatrix} L_t & 0 \\ 0 & 1 \end{bmatrix}$, denoted as $\text{block}(L_t)$, and the laplacian matrix for the next time-step, L_{t+1} are described by Equation 4. We compute the eigenvalue matrix, Λ , of L_t by applying theorem 5.6 of [45] to L_{t-1} .

LSTMs make accurate sequence predictions if elements of the sequence are correlated across time, as opposed to being generated randomly. In a general sequence of eigenvectors, the eigenvectors may not be correlated across time. Consequently, it is difficult for LSTM networks to predict the sequence of eigenvectors, \mathcal{U} accurately. This may adversely affect the behavior prediction algorithm described in Section IV-C. Our goal is now to show there exist a correlation between Laplacian matrices across time-steps and that this correlation is lower-bounded, that is, there exist sufficient correlation for accurate sequence modeling of eigenvectors.

Proving a lower-bound for the correlation is equivalent to proving an upper-bound for the noise, or error distance, between the j^{th} eigenvectors of L_t and L_{t+1} . We denote this error distance through the angle ϕ_j . From Theorem 5.4 of [45], the numerator of bound corresponds to the frobenius norm of the error between L_t and L_{t+1} . In our case, the update to the Laplacian matrix is given by Equation 4 where the error matrix is $\delta \delta^\top$. \square

In Theorem V.1, $\phi_j \ll 1$ and δ is defined in equation 4. λ_j represents the j^{th} eigenvalue and Λ represents all the eigenvalues of L_t . If the maximum component of δ_t is δ_{max} , then $\phi_j = \mathcal{O}(\sqrt{N} \delta_{max})$. Theorem V.1 shows that in a sequence of j^{th} eigenvectors, the maximum angular difference between successive eigenvectors is bounded by $\mathcal{O}(\sqrt{N} \delta_{max})$. By setting $N = 270$ (number of road-agents in Lyft), and $\delta_{max} := e^{-3} = 0.049$ (width of a lane), we observe a theoretical upper bound of 0.8 meters. A smaller value of ϕ_j indicates a greater similarity between successive eigenvectors, thereby implying a greater correlation in the sequence of eigenvectors. This allows sequence prediction models to learn future eigenvectors efficiently.

An alternative approach to computing the spectrums $\{U_{T+1}, \dots, U_{T+\tau}\}$ is to first form traffic-graphs from the predicted trajectory given as the output from the stream 1. After obtaining the corresponding Laplacian matrices for these traffic-graphs, standard eigenvalue algorithms can be used to compute the spectrum sequence. This is, however,

a relatively sub-optimal approach as in this case, $\phi = \mathcal{O}(NL_{max})$, with $L_{max} \gg \delta_{max}$.

VI. EXPERIMENTS AND RESULTS

We begin by listing the datasets used in our approach in Section VI-A. We list the evaluation metrics used and methods compared within Section VI-B. We analyze the main results and discuss the results of comparison methods and ablation studies of our approach in Section VI-C. In Section VI-D, we analyse the theoretical upper bound. We present an ablation analysis of the radius parameter μ in Section VI-F. We make all the implementation and training details available in the supplementary material.

A. Datasets

We use the NGSIM [48], Lyft Level 5 [49], Argoverse Motion Forecasting [50], and the ApolloScape Trajectory [22] datasets for evaluation. These are large-scale urban trajectory prediction datasets for autonomous driving. We give a brief description of all the datasets in the supplemental version.

Level 5 Lyft: The LYFT Level 5 dataset contains 180 videos of traffic in Palo Alto, California, U.S. Each video consists of 126 frames covering a duration of 20 seconds. A single traffic video consists of around 300 distinct road-agents. The data format is similar to the nuScenes format [51].

Argoverse Motion Forecasting: Argoverse motion forecasting data consists of 324,557 video segments of 5 seconds each. The total video length is 320 hours. The dataset contains traffic videos recorded in Miami (204 kilometers) and Pittsburgh (6 kilometers). The format of the data includes the timestamp, road-agent I.D., road-agent type, the spatial coordinates, and the location.

ApolloScape Trajectory: The ApolloScape trajectory dataset consists of 53 minutes of training sequences and 50 minutes of testing sequences captured at two fps. The dataset has been collected in Beijing, China. The format of the data includes the frame I.D., road-agent I.D., road-agent type, 3D spatial coordinates, heading angle and height, length, and width of the object.

NGSIM: While the previous three datasets are moderately dense datasets of urban traffic videos, the NGSIM dataset contains videos of sparse highway traffic. Each segment in the dataset is 45 minutes. In addition to trajectory information, the dataset also contains lane annotations for three lanes (left, center, right).

It is worth mentioning that there are several other datasets related to autonomous driving, for instance, the TRAF [2] and the Honda Driving Dataset [52]. However, TRAF is not relevant for our current work as the number of traffic videos in TRAF is less than 50, and the annotations consist of pixel coordinates instead of meters in the world coordinate system.

B. Evaluation Metrics and Methods

1) *Metrics:* For trajectory prediction, we use the standard metrics followed by prior trajectory prediction approaches [53], [54], [2], [8], [9].

- 1) Average Displacement Error (ADE): The root mean square error (RMSE) of all the predicted positions and real positions during the prediction window.
- 2) Final Displacement Error (FDE): The RMSE distance between the final predicted positions at the end of the predicted trajectory and the corresponding true location.

For behavior prediction, we report a weighted classification accuracy (W.A.) over the 3 class labels: {overspeeding, neutral, underspeeding}.

2) *Methods:* We compare our approach with SOTA trajectory prediction approaches for road-agents. Our definition of SOTA is not limited to ADE/FDE values. We consider SOTA additionally with respect to the deep learning architecture used in a different approach. Combined, our basis for selecting SOTA methods not only evaluates the ADE/FDE scores but also evaluates the benefits of using the two-stream network versus other deep learning-based architectures.

- Enc-Dec: This method is based on a standard encoder-decoder architecture similar to the Seq2Seq model [55].
- Deo et al. [8] (CS-LSTM): This method combines CNNs with LSTMs to perform trajectory prediction on U.S. highways.
- Chandra et al. [2] (TraPHic): This approach also uses a CNN + LSTM approach along with spatial attention-based pooling to perform trajectory prediction of road-agents in dense and heterogeneous traffic.
- Gupta et al. [54] (Social-GAN): This GAN-based trajectory prediction approach is originally trained on pedestrian crowd datasets. The method uses the encoder-decoder architecture to act as the generator and trains an additional encoder as the discriminator.
- Li et al. [15] (GRIP): This is a graph-based trajectory prediction approach that replaces standard CNNs with graph convolutions and combines GCNs with an encoder-decoder framework.

We use the publicly available implementations for CS-LSTM, TraPHic, and Social-GAN, and train the entire model on all three datasets. We performed hyper-parameter tuning on all three methods and reported the best results. Moreover, we compare with the officially published results for GRIP as reported on the NGSIM [15] and the ApolloScape datasets[†].

C. Analysis and Discussion

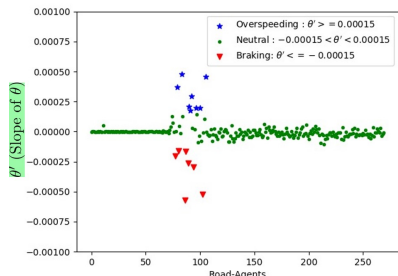
We compare the ADE and FDE scores of our predicted trajectories with prior methods in Table I and show qualitative results in the supplementary material. We compare with several SOTA trajectory prediction methods and reduce the average RMSE by at least 75% with respect to the next best method (GRIP).

Ablation Study of Stream 1 vs. Both Streams: To highlight the benefit of the spectral cluster regularization on long-term prediction, we remove the second stream and only train the LSTM encoder-decoder model (Stream 1) with the original loss function (equation 5). Our results (Table I, last four columns) show that regularizing stream 1 reduces

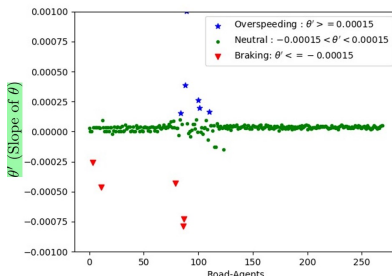
[†]http://apolloscape.auto/leader_board.html

TABLE I: **Main Results:** We report the Average Displacement Error (ADE) and Final Displacement Error (FDE) for prior road-agent trajectory prediction methods in meters (m). Lower scores are better and **bold** indicates the SOTA. **Revisions:** Green colored cells indicate revised results for *our* method while blue colored cells indicates revised results for the original implementations *GRIP* and *Social-GAN* method.

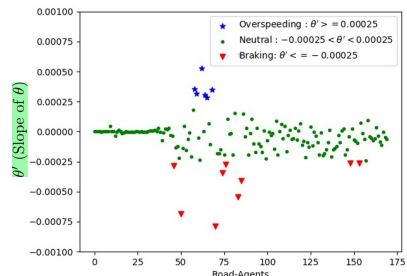
Dataset (Pred. Len.)	Comaprison Methods										Ablation		Our Approach	
	Enc-Dec		CS-LSTM		TraPHic		Social-GAN		GRIP		Stream 1		Both Streams	
	ADE	FDE	ADE	FDE	ADE	FDE	ADE	FDE	ADE	FDE	ADE	FDE	ADE	FDE
Lyft (5 sec.)	-	-	4.423	8.640	5.031	9.882	7.860	14.340	-	-	5.77	11.20	2.65	2.99
Argoverse (5 sec.)	-	-	1.050	3.085	1.039	3.079	3.610	5.390	-	-	2.40	3.09	0.99	1.87
ApolloScape (3 sec.)	2.24	8.25	2.144	11.699	1.283	11.674	3.980	6.750	1.25	2.34	2.14	9.19	1.12	2.05
NGSIM (5 sec.)	6.86	10.02	7.250	10.050	5.630	9.910	5.650	10.290	1.61	3.16	1.31	2.98	0.40	1.08



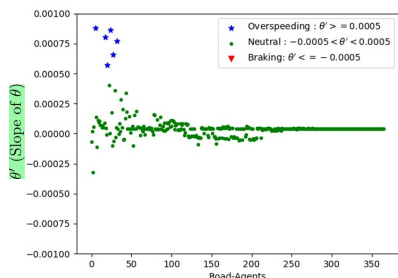
(a) Lyft Ground-Truthwith $\lambda=0.00015$.



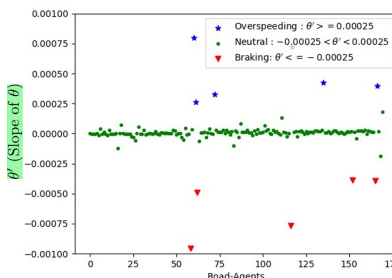
(b) Lyft Behavior Predictions.



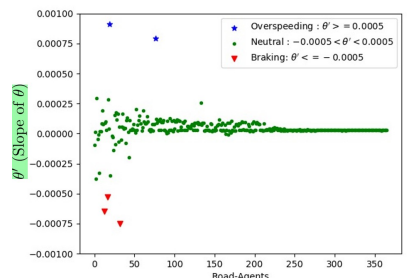
(c) Argoverse Ground-Truth with $\lambda=0.00025$.



(d) ApolloScape Ground-Truth with $\lambda=0.0005$.



(e) Argoverse Behavior Predictions.



(f) ApolloScape Behavior Predictions.

Fig. 3: **Behavior Prediction Results:** We classify the three behaviors— overspeeding(blue), neutral(green), and underspeeding(red), for all road-agents in the Lyft, Argoverse, and ApolloScape datasets, respectively. The y-axis shows θ' and the x-axis denotes the road-agents. We follow the behavior prediction protocol described in Section IV-C. Each figure in the top row represents the ground-truth labels, while the bottom row shows the predicted labels. In our experiments, we set $\lambda = \lambda_1 = -\lambda_2$.

the FDE by up to 70%. This is as expected since stream 1 does not take into account neighbor information. Therefore, it should also be noted that stream 1 performs poorly in dense scenarios but rather well in sparse scenarios. This is evident from Table I where stream 1 outperforms comparison methods on the sparse NGSIM dataset.

Additionally, Figure 5 shows that in the presence of regularization, the RMSE for our spectrally regularized approach (“both streams”, purple curve) is much lower than that of stream 1 (red curve) across the entire prediction window.

RMSE depends on traffic density: The upper bound for the increase in RMSE error is a function of the density of the traffic since $\phi = \mathcal{O}(\sqrt{N}\delta_{max})$, where N is the total number of agents in the traffic video and $\delta_{max} = 0.049$ meters for a three-lane wide road system. The NGSIM dataset contains the sparsest traffic with the lowest value for N and therefore the RMSE values are lower for the NGSIM (0.40/1.08) compared to the other three datasets that contain dense urban

traffic.

Comparison with other methods: Our method learns weight parameters for a spectral regularized LSTM network (Figure 2), while GRIP learns parameters for a graph-convolutional network (GCN). We outperform GRIP on the NGSIM and ApolloScape datasets, while comparisons on the remaining two datasets are unavailable.

TraPHic and CS-LSTM are similar approaches. Both methods require convolutions in a heuristic local neighborhood. The size of the neighborhood is specifically adjusted to the dataset that each method is trained on. We use the default neighborhood parameters provided in the publicly available implementations, and apply them to the NGSIM, Lyft, Argoverse, and ApolloScape datasets. We outperform both methods on all benchmark datasets.

Lastly, Social-GAN is trained on the scale of pedestrian trajectories, which differs significantly from the scale of vehicle trajectories. This is primarily the reason behind

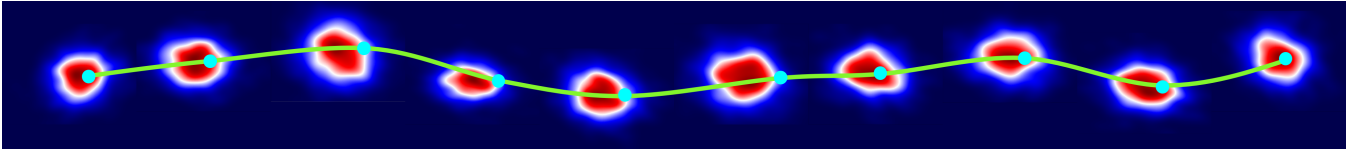


Fig. 4: **Qualitative Analysis:** We compare the predicted trajectory with the ground truth trajectory (green line with cyan coordinates). The prediction time is 5 seconds. Each red blob in the figure represents a predicted bi-variate Gaussian distribution, $\mathcal{N}(\mu^*, \sigma^*, \rho^*)$. The prediction is more accurate when the cyan points are closer to the center of the red blobs.

Social-GAN placing last among all methods.

Comparison with other methods: Our method learns weight parameters for a spectral regularized LSTM network (Figure 2), while GRIP learns parameters for a graph-convolutional network (GCN). We outperform GRIP on the NGSIM and Apolloscape datasets, while comparisons on the remaining two datasets are unavailable. TraPHic and CS-LSTM are similar approaches. Both methods require convolutions in a heuristic local neighborhood. The size of the neighborhood is specifically adjusted to the dataset that each method is trained on. We use the default neighborhood parameters provided in the publicly available implementations, and apply them to the NGSIM, Lyft, Argoverse, and Apolloscape datasets. We outperform both methods on all benchmark datasets. Lastly, Social-GAN is trained on the scale of pedestrian trajectories, which differs significantly from the scale of vehicle trajectories. This is primarily the reason behind Social-GAN placing last among all methods.

D. Long-Term Prediction Analysis

The goal of improved long-term prediction is to achieve a lower FDE, which can be clearly observed in our results in Table I. We achieve this goal by successfully upper-bounding the worst-case maximum FDE that can be obtained. These upper bounds are a consequence of the theoretical results in Section V. We denote the worst-case theoretical FDE by T-FDE. This measure represents the maximum FDE that can be obtained using Theorem V.1 under fixed assumptions. In Table II, we compare the T-FDE with the empirical FDE results obtained in Table I. The T-FDE is computed by $T\text{-FDE} = \frac{\phi}{n} \times (T - \tau)$.

The formula for T-FDE is derived as follows. The error incurred by all vehicles at a current time-step during spectral clustering is bounded by ϕ (Theorem V.1). Since the cluster centers update the positions for each vehicle in constant time, the increase in overall RMSE values at each time-step is bounded by ϕ as well. The increase in RMSE for a single agent, therefore, is bounded by $\frac{\phi}{n}$, where $n = \frac{N}{T}$ is the average number of vehicles per frame in each dataset. $T - \tau$ is the length of the prediction window and $n = 10$ on average for the Lyft, Argoverse, and Apolloscape datasets as indicated by the datasets. We do not have the data needed to compute ϕ for the NGSIM dataset as the total number of lanes are not known.

We note a 73%, 82%, 100% agreement between the theoretical FDE and the empirical FDE on the Apolloscape, Lyft, and Argoverse datasets, respectively. The main cause for disagreements in the first two datasets is the choice for

TABLE II: **Upper Bound Analysis:** ϕ is the upper bound on the RMSE for all agents at a time-step. $T - \tau$ is the length of the prediction window. T-FDE is the theoretical FDE that should be achieved by using spectral regularization. E-FDE is the empirical FDE observed in Table I. The % agreement is the agreement between the T-FDE and E-FDE computed using $\frac{T\text{-FDE}}{E\text{-FDE}}$, if $T\text{-FDE} < E\text{-FDE}$, else 100%. *Conclusion:* Theorem V.1 is empirically verified with at least 73% guarantee.

Dataset	ϕ	$(T - \tau)$	T-FDE	E-FDE	% Agreement
Lyft Level 5	0.80	30	2.46	2.99	82%
Apolloscape	1.50	10	1.50	2.05	73%
Argoverse	0.64	30	1.95	1.87	100%

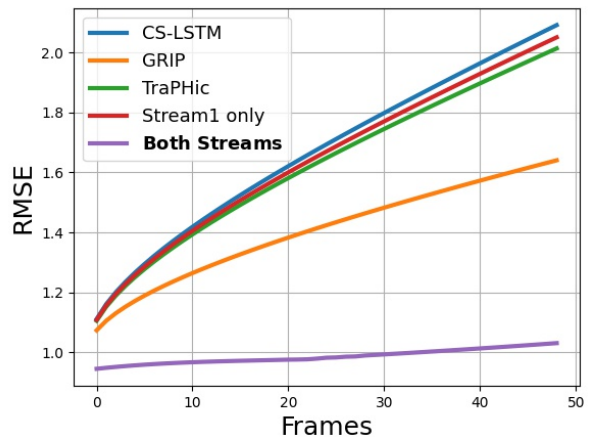


Fig. 5: **RMSE Curves:** We plot the RMSE values for all methods. The prediction window is 5 seconds corresponding to a frame length of 50 for the NGSIM dataset.

the value of $\delta_{\max} = 0.049$. This value is assumed for a three-lane wide road system that was observed in majority of the videos in both datasets. However, it may be the case that several videos contain one- or two-lane traffic. In such cases, the values for δ_{\max} changes to 0.36 and 0.13, respectively, thereby increasing the upper bound for increase in RMSE.

Note, in Figure 5, the increase in RMSE for our approach (purple curve) is much lower than that of other methods, which is due to the upper bound induced by spectral regularization.

E. Behavior Prediction Results

We follow the behavior prediction algorithm described in Section IV-C. The values for λ_1 and λ_2 are based on the ground truth labels and are hidden from the test set. We observe a weighted accuracy of 92.96% on the Lyft dataset, 84.11% on the Argoverse dataset, and 96.72% on the Apolloscape dataset. In the case of Lyft, Figures 3a and 3b

TABLE III: Ablation experiments of the radius parameter μ . Each column contains averaged RMSE values over the corresponding range interval. *Conclusion:* The optimum results are obtained by setting $0 < \mu \leq 10$ meters.

Dataset	$\mu = 0$	$0 < \mu \leq 10$	$10 < \mu \leq 20$
ApolloScape	2.14	1.12	2.62
Argoverse	2.40	0.99	3.15
Lyft Level 5	5.77	2.65	3.36
NGSIM	1.31	0.40	2.03

show the ground truth and predictions for Lyft, respectively. We plot the value of θ' on the vertical axis and the road-agent I.D.s on the horizontal axis. More similarity across the two plots indicates higher accuracy. For instance, the red (aggressive) and blue (conservative) dotted regions in **3a** and **3b** are nearly identical indicating a greater number of correct classifications. Similar results follow for the ApolloScape and Argoverse datasets, which we show in the supplementary material due to lack of space. Due to the lack of diverse behaviors in the NGSIM dataset, we do not perform behavior prediction on the NGSIM.

An interesting observation is that road-agents towards the end of the x-axis appear late in the traffic video while road-agents at the beginning of the x-axis appear early in the video. The variation in behavior class labels, therefore, decreases towards the end of the x-axis. This intuitively makes sense as θ' for a road-agent depends on the number of distinct neighbors that it observes. This is difficult for road-agents towards the end of the traffic video.

F. Ablation Study of the Radius Parameter (μ)

We conducted ablation experiments in which we vary the radius parameter μ (See Section III-B for a discussion on μ) from 0 to 20. We obtained results on the ApolloScape, Argoverse, Lyft, and NGSIM datasets which we present in Table III. We measured the average RMSE values over 3 range intervals: $\mu = 0$, $0 < \mu \leq 10$, and $10 < \mu \leq 20$. We use range intervals to clearly and succinctly capture the trend of the RMSE values for $\mu > 10$ meters and $\mu < 10$ meters. We observe that the best performance is achieved from the latter range ($0 < \mu \leq 10$).

It is clear that setting $\mu = 0$ and thus ignoring neighborhood information in dense traffic severely degrades performance. But on the other hand, increasing the radius beyond 10 meters also increases the RMSE error. This is because by increasing the radius beyond a certain threshold, we inadvertently include in our spectral clustering algorithm those road-agents that are too far to interact with the ego-agent. In order to accommodate these “far-away” road-agents, the clusters expand and shift the cluster center from its true center. This phenomenon is common in statistics where an outlier corrupts the data distribution. The far-away agents are outliers in the spectral clustering algorithm, thereby leading to an increase in RMSE. We conclude that our method produces optimum results for $0 < \mu \leq 10$ in dense traffic systems.

G. Training Details

We use 2 Nvidia GeForce RTX 2080 Ti GPUs with 12GB memory each, for all experiments. Initially, we trained both streams together for 20 epochs. However, we observed that by training stream one first for 20 epochs, and then training both streams for another 5 epochs generated the results reported in Table I. All datasets are normalized prior to training as each dataset is recorded at a different scale.

Stream 1: The input to stream one consists of trajectory tensors of size $B \times T \times 2$, with $B = 128$ representing the batch size, $T = 3$ seconds denoting the length of the observed trajectory, and 2-dimensional spatial coordinates. In stream 1, our training converges in 20 epochs in approximately 20 minutes for each data set. The best results from stream one are obtained by using the RMSprop optimizer with a learning rate of 0.001.

Stream 2: The input to stream 2 consists of a sequence of eigenvectors represented as a tensor of size $k \times B \times T \times N$, where k denotes the number of eigenvectors for each Laplacian matrix, N is the number of road-agents in the traffic-video. It takes approximately 42 hours per epoch to process the data for stream 2. Therefore, we pre-compute the input data for stream two offline, which reduces the training time to 2.2, 1.5, and 1.2 hours per epoch for Lyft, Argoverse, and ApolloScape, respectively. The hyper-parameters for stream 2 are identical to those used in stream 1.

VII. CONCLUSION, LIMITATIONS, AND FUTURE WORK

We present a unified algorithm for trajectory prediction and behavior prediction of road-agents. We use a two-stream LSTM network in which the first stream predicts the trajectories, while the second stream predicts the behavior of road-agents. We also present a regularization algorithm to reduce long-term prediction errors.

Our method has some limitations. Currently, we use only one feature to design our behavior prediction model, which may not be able to generalize to new traffic scenarios. In addition, our training is slow and takes several hours due to the number of computations required for computing the traffic-graphs and corresponding Laplacian matrices. We plan to make our behavior prediction model data-driven, rather than rule-based. We will also explore ways to improve trajectory prediction using our behavior prediction algorithm.

REFERENCES

- [1] E. Cheung, A. Bera, and D. Manocha, “Efficient and safe vehicle navigation based on driver behavior classification,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 1024–1031.
- [2] R. Chandra, U. Bhattacharya, A. Bera, and D. Manocha, “Trajectory prediction in dense and heterogeneous traffic using weighted interactions,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8483–8492.
- [3] R. Chandra, U. Bhattacharya, T. Mittal, X. Li, A. Bera, and D. Manocha, “Graphrqi: Classifying driver behaviors using graph spectrums,” *arXiv preprint arXiv:1910.00049*, 2019.
- [4] W. Schwarting, A. Pierson, J. Alonso-Mora, S. Karaman, and D. Rus, “Social behavior for autonomous vehicles,” *Proceedings of the National Academy of Sciences*, vol. 116, no. 50, pp. 24972–24978, 2019.
- [5] R. Hoogendoorn, B. van Arerm, and S. Hoogendoorn, “Automated driving, traffic flow efficiency, and human factors: Literature review,” *Transportation Research Record*, vol. 2422, no. 1, pp. 113–120, 2014.

- [6] R. Yoshizawa, Y. Shiomi, N. Uno, K. Iida, and M. Yamaguchi, "Analysis of car-following behavior on sag and curve sections at intercity expressways with driving simulator," *International Journal of Intelligent Transportation Systems Research*, pp. 56–65, 2012.
- [7] M. Saifuzzaman and Z. Zheng, "Incorporating human-factors in car-following models: a review of recent developments and research needs," *Transportation Research: emerging technologies*, 2014.
- [8] N. Deo and M. M. Trivedi, "Convolutional social pooling for vehicle trajectory prediction," *arXiv preprint arXiv:1805.06771*, 2018.
- [9] R. Chandra, U. Bhattacharya, C. Roncal, A. Bera, and D. Manocha, "Robusttp: End-to-end trajectory prediction for heterogeneous road-agents in dense traffic with noisy sensor inputs," *arXiv preprint arXiv:1907.08752*, 2019.
- [10] B. Brown and E. Laurier, "The trouble with autopilots: Assisted and autonomous driving on the social road," in *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, 2017, pp. 416–429.
- [11] W. Schwarting, A. Pierson, J. Alonso-Mora, S. Karaman, and D. Rus, "Social behavior for autonomous vehicles," *Proceedings of the National Academy of Sciences*, vol. 116, no. 50, pp. 24972–24978, 2019.
- [12] Z. N. Sunberg, C. J. Ho, and M. J. Kochenderfer, "The value of inferring the internal state of traffic participants for autonomous freeway driving," in *2017 American Control Conference (ACC)*. IEEE, 2017, pp. 3004–3010.
- [13] D. Seth and M. L. Cummings, "Traffic efficiency and safety impacts of autonomous vehicle aggressiveness," *simulation*, vol. 19, p. 20.
- [14] R. Yu, S. Zheng, A. Anandkumar, and Y. Yue, "Long-term forecasting using tensor-train rns," *arXiv preprint arXiv:1711.00073*, 2017.
- [15] X. Li, X. Ying, and M. C. Chuah, "Grip: Graph-based interaction-aware trajectory prediction," *arXiv preprint arXiv:1907.07792*, 2019.
- [16] B. M. Waxman, "Routing of multipoint connections," *IEEE journal on selected areas in communications*, pp. 1617–1622, 1988.
- [17] Z.-X. Feng, J. Liu, Y.-Y. Li, and W.-H. Zhang, "Selected model and sensitivity analysis of aggressive driving behavior," *Zhongguo Gonglu Xuebao(China Journal of Highway and Transport)*, pp. 106–112, 2012.
- [18] E. R. Dahlen, B. D. Edwards, T. Tubré, M. J. Zypur, and C. R. Warren, "Taking a look behind the wheel: An investigation into the personality predictors of aggressive driving," *Accident Analysis & Prevention*, vol. 45, pp. 1–9, 2012.
- [19] G. E. Box, G. M. Jenkins, G. C. Reinsel, and G. M. Ljung, *Time series analysis: forecasting and control*. John Wiley & Sons, 2015.
- [20] C. Hermes, C. Wohler, K. Schenk, and F. Kummert, "Long-term vehicle motion prediction," in *Intelligent vehicles symposium*, 2009.
- [21] L. Ljung, "System identification," *Wiley Encyclopedia of Electrical and Electronics Engineering*, 2001.
- [22] Y. Ma, X. Zhu, S. Zhang, R. Yang, W. Wang, and D. Manocha, "Trafficpredict: Trajectory prediction for heterogeneous traffic-agents," *arXiv preprint arXiv:1811.02146*, 2018.
- [23] F.-C. Chou, T.-H. Lin, H. Cui, V. Radosavljevic, T. Nguyen, T.-K. Huang, M. Niedoba, J. Schneider, and N. Djuric, "Predicting motion of vulnerable road users using high-definition maps and efficient convnets," *arXiv preprint arXiv:1906.08469*, 2019.
- [24] N. Djuric, V. Radosavljevic, H. Cui, T. Nguyen, F.-C. Chou, T.-H. Lin, and J. Schneider, "Short-term Motion Prediction of Traffic Actors for Autonomous Driving using Deep Convolutional Networks," *ArXiv e-prints*, Aug. 2018.
- [25] Z. Cui, K. Henrickson, R. Ke, and Y. Wang, "Traffic graph convolutional recurrent neural network: A deep learning framework for network-scale traffic learning and forecasting," *arXiv preprint arXiv:1802.07007*, 2018.
- [26] S. Guo, Y. Lin, N. Feng, C. Song, and H. Wan, "Attention based spatial-temporal graph convolutional networks for traffic flow forecasting," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 922–929.
- [27] B. Yu, H. Yin, and Z. Zhu, "Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting," *ArXiv*, vol. abs/1709.04875, 2018.
- [28] Z. Diao, X. Wang, D. Zhang, Y. Liu, K. Xie, and S. He, "Dynamic spatial-temporal graph convolutional neural networks for traffic forecasting," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 890–897.
- [29] M. R. Barrick and M. K. Mount, "The big five personality dimensions and job performance: a meta-analysis," *Personnel psychology*, 1991.
- [30] B. Krahé and I. Fenske, "Predicting aggressive driving behavior: The role of macho personality, age, and power of car," *Aggressive Behavior: Official Journal of the International Society for Research on Aggression*, vol. 28, no. 1, pp. 21–29, 2002.
- [31] A. H. Jamson, N. Merat, O. M. Carsten, and F. C. Lai, "Behavioural changes in drivers experiencing highly-automated vehicle control in varying traffic conditions," *Transportation research part C: emerging technologies*, vol. 30, pp. 116–125, 2013.
- [32] P. Rämä, "Effects of weather-controlled variable speed limits and warning signs on driver behavior," *Transportation Research Record*, vol. 1689, no. 1, pp. 53–59, 1999.
- [33] J. Dai, J. Teng, X. Bai, Z. Shen, and D. Xuan, "Mobile phone based drunk driving detection," in *2010 4th International Conference on Pervasive Computing Technologies for Healthcare*, 2010.
- [34] M. Saifuzzaman, M. M. Haque, Z. Zheng, and S. Washington, "Impact of mobile phone use on car-following behaviour of young drivers," *Accident Analysis & Prevention*, vol. 82, pp. 10–19, 2015.
- [35] A. Aljaafreh, N. Alshabat, and M. S. N. Al-Din, "Driving style recognition using fuzzy logic," *2012 IEEE International Conference on Vehicular Electronics and Safety (ICVES 2012)*, pp. 460–463, 2012.
- [36] Y. L. Murphey, R. Milton, and L. Kiliaris, "Driver's style classification using jerk analysis," *2009 IEEE Workshop on Computational Intelligence in Vehicles and Vehicular Systems*, pp. 23–28, 2009.
- [37] I. Mohamad, M. A. M. Ali, and M. Ismail, "Abnormal driving detection using real time global positioning system data," *Proceeding of the 2011 IEEE International Conference on Space Science and Communication (IconSpace)*, pp. 1–6, 2011.
- [38] G. Qi, Y. Du, J. Wu, and M. Xu, "Leveraging longitudinal driving behaviour data with data mining techniques for driving style analysis," *IET intelligent transport systems*, vol. 9, no. 8, pp. 792–801, 2015.
- [39] B. Shi, L. Xu, J. Hu, Y. Tang, H. Jiang, W. Meng, and H. Liu, "Evaluating driving styles by normalizing driving behavior based on personalized driver modeling," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 45, pp. 1502–1508, 2015.
- [40] E. Cheung, A. Bera, E. Kubin, K. Gray, and D. Manocha, "Identifying driver behaviors using trajectory features for vehicle navigation," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 3445–3452.
- [41] A. Zyner, S. Worrall, and E. Nebot, "A recurrent neural network solution for predicting driver intention at unsignalized intersections," *IEEE Robotics and Automation Letters*, 2018.
- [42] A. Zyner, S. Worrall, J. Ward, and E. Nebot, "Long short term memory for driver intent prediction," in *IEEE Intelligent Vehicles Symposium (IV)*, 2017.
- [43] M. Belkin and P. Niyogi, "Laplacian eigenmaps for dimensionality reduction and data representation," *Neural computation*, 2003.
- [44] A. Graves, "Generating sequences with recurrent neural networks," *arXiv preprint arXiv:1308.0850*, 2013.
- [45] J. W. Demmel, *Applied numerical linear algebra*. Siam, 1997, vol. 56.
- [46] B. Osting, C. D. White, and É. Oudet, "Minimal dirichlet energy partitions for graphs," *SIAM Journal on Scientific Computing*, 2014.
- [47] U. Von Luxburg, "A tutorial on spectral clustering," *Statistics and computing*, 2007.
- [48] U. F. H. Administration, "U.s. highway 101 and i-80 dataset," 2005.
- [49] R. Kesten, M. Usman, J. Houston, T. Pandya, K. Nadhamuni, A. Ferreira, M. Yuan, B. Low, A. Jain, P. Ondruska, S. Omari, S. Shah, A. Kulkarni, A. Kazakova, C. Tao, L. Platinisky, W. Jiang, and V. Shet, "Lyft level 5 av dataset 2019," <https://level5.lyft.com/dataset/>, 2019.
- [50] M.-F. Chang, J. W. Lambert, P. Sangkloy, J. Singh, S. Bak, A. Hartnett, D. Wang, P. Carr, S. Lucey, D. Ramanan, and J. Hays, "Argoverse: 3d tracking and forecasting with rich maps," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [51] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nuscenes: A multimodal dataset for autonomous driving," *arXiv preprint arXiv:1903.11027*, 2019.
- [52] V. Ramanishka, Y.-T. Chen, T. Misu, and K. Saenko, "Toward driving scene understanding: A dataset for learning driver behavior and causal reasoning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7699–7707.
- [53] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese, "Social lstm: Human trajectory prediction in crowded spaces," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 961–971.
- [54] A. Gupta, J. Johnson, L. Fei-Fei, S. Savarese, and A. Alahi, "Social GAN: Socially Acceptable Trajectories with Generative Adversarial Networks," *ArXiv e-prints*, 2018.

- [55] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Advances in neural information processing systems*, 2014, pp. 3104–3112.
- [56] G. Hinton, N. Srivastava, and K. Swersky, "Neural networks for machine learning lecture 6a overview of mini-batch gradient descent."

VIII. DATA PREPROCESSING

We include all data and code used for data preprocessing with the supplementary material. Our data structure format (See Section IX) includes the time-stamp, road-agent ID, and the road-agent’s spatial coordinates in the world coordinate frame. The process of obtaining these attributes, and utilizing them, to construct the data structures differs for all three datasets – Lyft, Apolloscape, and Argoverse. We converted the three datasets to one unique representation that includes frame_ID, road-agent_ID, X, Y, dataset_ID.

A. Metadata

To the best of our knowledge, there is very little known prior work using these datasets as they are relatively recent. As such, the raw datasets are not trivial to work with, due to their large size. In particular, to understand the performance, and gain interpretability, of an approach on a dataset, it is essential to study the underlying meta-features of each dataset. In Table IV, we list some descriptive meta-features of the datasets. Additionally in this work, we also release code for efficient implementations for several key operations such as storage, extraction, querying, addition, and deletion on these datasets to facilitate future research using these datasets.

Dataset	Batch	Avg. Density
Lyft	Train	0.80
	Val	0.83
	Test	0.84
Argoverse	Train	0.75
	Val	0.67
	Test	1.67
Apolloscape	Train	3.49
	Val	3.50
	Test	2.56

TABLE IV: **Meta-Feature information** for the Lyft Level 5, Apolloscape, and the Argoverse datasets. The average density is reported by measuring the average number of road-agents per frame.

IX. DATA STRUCTURES

This section describes the different data structures that are used in our approach. The implementations of these data structures are included with the code.

A. Adjacency Matrices

Figure 6 shows a schematic for the data structure used to create the adjacency matrices for the whole dataset. The adjacency matrix corresponding to a traffic graph captures the interactions between all agents within that frame. A python list is used to store adjacency data of each traffic video, where each video is again a list of all frames in that dataset. Each frame is a python dictionary containing the ‘dataset_ID’, ‘frame_ID’ and ‘adj_matrix’. Each adjacency matrix is an array of size $N \times N$, where N is the total

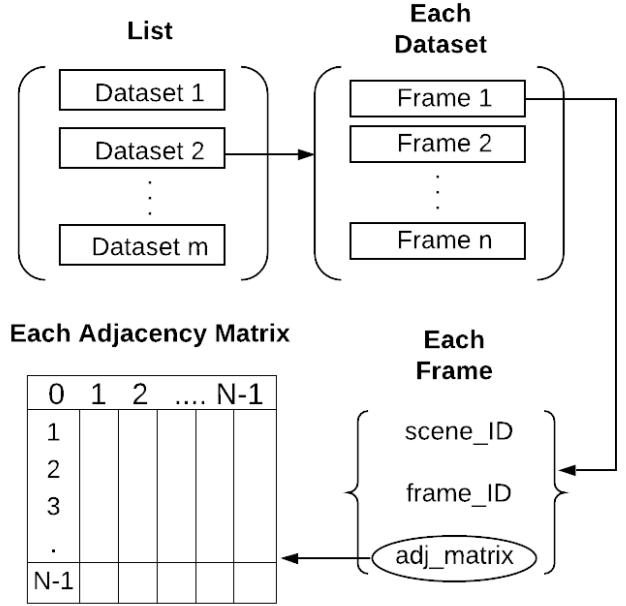


Fig. 6: **Schematic:** Data structure for Adjacency matrices.

number of agents in the dataset. The adjacency matrices are used to form the Laplacian matrices that are updated at every time-step according to equation 3.

B. Input for Stream 1

Figure 7 shows the schematic for the data structure used to prepare the input for stream 1. The implementation consists of a python list of dictionaries. Each dictionary denoted as item_1, item_2...item_n in Figure 7, has three keys- ‘dataset_ID’, ‘agent_ID’, and ‘sequence’. The value of the ‘sequence’ consists of an array of size $n \times 2$, where n is the length of either the observation sequence or prediction sequence. Each row of this sequence array consists of the global X, Y coordinates, respectively, of the road-agent at that observation or prediction time step.

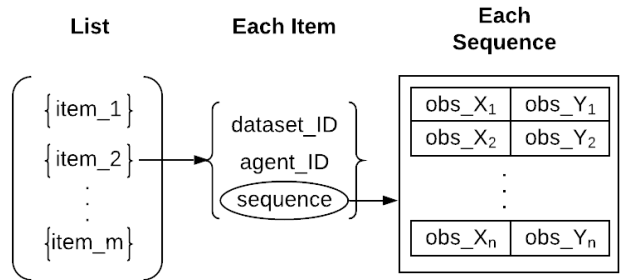


Fig. 7: **Schematic:** Data structure for the input to Stream1.

C. Input for Stream 2

Figure 8 shows the schematic for the data structure used to prepare the input for stream 2. This data structure is similar to the stream1 data structure in that it also consists of a list of dictionaries. Each dictionary has the keys

‘dataset_ID’, ‘agent_ID’, ‘mean_theta_hat’, ‘mean_theta’, and F_1, F_2, \dots, F_n , where n is the length of the observation sequence or prediction sequence, respectively. Each F_i represents the i^{th} frame which is an array of size $2 \times N$, where N denotes the total number of road-agents in that traffic video. The columns of this array stores the global X, Y coordinates of all the road-agents at i^{th} frame. The keys, ‘mean_theta_hat’ and ‘mean_theta’ store information corresponding to the ground-truth behavior labels for that sequence of frames.

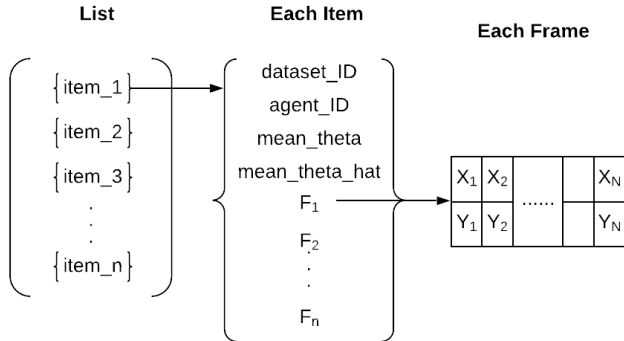


Fig. 8: **Schematic:** Data structure for the input to Stream2.

X. MOTIVATION FOR EQUATION 3

Our intuition is that a road-agent should remember the interactions with not just its current neighbors, but all neighbors it has observed up until current time t . Therefore, the first term on the RHS stores the information of road-agents up to time t , while the second term on the RHS contains the information of the newly observed road-agents at time $t+1$.

XI. ANALYSIS OF STREAM 2 CONTINUED

In Section 4.2, we presented a behavior prediction algorithm that begins by forming the set of predicted spectrums from the second stream, $\mathcal{U} = \{U_{T+1}, U_{T+2}, \dots, U_{t+\tau}\}$. The success of the algorithm depends on the accuracy of these predictions, that further depends on the amount of correlation existing between the sequence of eigenvectors. In Section 4.3, we proved an upper bound for the error distance between the j^{th} eigenvectors of L_t and L_{t+1} , denoted as ϕ_j . We showed that $\phi_j = \mathcal{O}(\sqrt{n}\delta_{max})$, where δ_{max} is the maximum component of δ_t .

An alternative approach to computing the spectrums $\{U_{T+1}, \dots, U_{T+\tau}\}$ is to first form traffic-graphs from the predicted trajectory given as the output from the stream 1. Next, obtain the corresponding Laplacian matrices for these traffic-graphs. Finally, use standard eigenvalue algorithms to compute the spectrum sequence. This is, however, a relatively sub-optimal approach as in this case, $\phi = \mathcal{O}(nL_{max})$, with $L_{max} \gg \delta_{max}$.

XII. ADDITIONAL RELATED WORK

In Section 6, we compared our approach against several prominent SOTA deep learning-based trajectory prediction

algorithms. However, it is worth noting that there are other additional methods in the trajectory prediction literature that we have not considered in this present work. Social-LSTM [53] is a popular approach for trajectory prediction of pedestrians in crowds. However, in the interest of clarity, we provide here an explanation regarding its exclusion from the experimental setup:

- The official implementation of Social-LSTM has been recently withdrawn from public access. Therefore, any results obtained via custom implementations would not offer an unbiased view of the method compared to prior works, that have used the official implementation before its removal.
- Instead of presenting a somewhat biased comparisons with a well-known pedestrian trajectory prediction method, we compare with CS-LSTM [8], which is a slight modification of the Social-LSTM method for road-agent trajectory prediction.

XIII. TRAINING DETAILS OF COMPARISON METHODS

- **TraPHic & CS-LSTM:** Implementations of both methods can be found in [9]. On the Lyft, Argoverse, and Apolloscape datasets, we use NLL loss and MSE loss for pretraining and training. We perform a hyperparameter grid search with 10-40 epochs of batch size 64 and 128. We use the adam, adamax, Rprop, and RMSprop optimizers, with learning rates of 0.01, 0.001, and 0.005, respectively, and dropouts between 0.4–0.6. In TraPHic, we get the best result on Lyft and Argoverse with the adamax optimizer with a dropout of 0.4 and 0.5, respectively, and on Apolloscape with the adam optimizer with a dropout of 0.5. In CS-LSTM, our best result on Lyft is obtained with the adamax optimizer with a dropout of 0.4, and on Argoverse and Apolloscape with the adam optimizer with a dropout of 0.5. However, compared to other methods, these two methods perform poorly as time increases.
- **Social-GAN:** We use the standard implementation in [54]. On each of the three datasets, our hyperparameter grid search ranges from 1500 – 3000 epochs of batch size 32, using adam, adamax, Rprop, and RMSprop optimizers, with a learning rate of 0.001, 0.0001 and 0.0005, respectively, on both generator and discriminator. We achieve the best result with the adamax optimizer on Lyft and Argoverse, and adam optimizer on Apolloscape, all with a learning rate of 0.0005. Due to limited GPU memory, we were not able to experiment on larger batch sizes. Since Social-GAN is a pedestrian trajectory prediction approach, we scaled down the trajectories in all three datasets by a factor of 20 to resemble the trajectories of pedestrians. Even so, we achieve unstable results during testing.
- **Enc-Dec:** We implement this method from scratch, similar to the implementation of stream 1 in our two-stream network. The key difference is that we train Enc-Dec with the MSE loss whereas we train our two stream network with the NLL loss function. We obtain

optimum results for Enc-Dec after 50 epochs with batch size of 40 and learning rate of 0.001.

- **GRIP:** We also implement this method following the approach in [15]. We obtain the best results with batch size of 128 and learning rate of 0.001.

A. Training Details

We use 2 Nvidia GeForce RTX 2080 Ti GPUs with 12GB memory each, for all experiments. Initially, we trained both streams together for 20 epochs. However, we found that by training stream one first for 20 epochs, and then training both streams for another five epochs generated the best results, reported in Table I.

Due to the computations involved in obtaining the Laplacian matrices and their corresponding eigenvectors, data processing for stream two is both time-consuming and expensive in terms of computational resources. Consequently, we choose 6310, 5126, and 5073 valid trajectories to form the training set, and 769, 1678 and 1012 valid trajectories to form the testing set for the Lyft, Argoverse and Apolloscape datasets, respectively. We consider the trajectory for a road-agent to be valid if that road-agent is present for at least 8 seconds (3 observation + 5 prediction) in Lyft and Apolloscape. In Argoverse, we only use the first 2 seconds as an observed trajectory to predict the next 3 seconds since each video in Argoverse is limited to 5 seconds.

Stream 1: The input to stream one consists of trajectory tensors of size $B \times T \times 2$, with $B = 128$ representing the batch size, $T = 3$ seconds denoting the length of the observed trajectory, and 2-dimensional spatial coordinates. In stream 1, our training converges in 20 epochs in approximately 20 minutes for each data set. The best results from stream one are obtained by using the RMSprop optimizer [56] with a learning rate of 0.001.

Stream 2: The input to stream 2 consists of a sequence of eigenvectors represented as a tensor of size $k \times B \times T \times N$, where k denotes the number of eigenvectors for each Laplacian matrix, N is the number of road-agents in the traffic-video. It takes approximately 42 hours per epoch to process the data for stream 2. Therefore, we pre-compute the input data for stream two offline, which reduces the training time to 2.2, 1.5, and 1.2 hours per epoch for Lyft, Argoverse, and Apolloscape, respectively. The hyper-parameters for stream 2 are identical to those used in stream 1.