



Trajectron++: Dynamically-Feasible Trajectory Forecasting with Heterogeneous Data

Tim Salzmann¹, Boris Ivanovic^{1(✉)}, Punarjay Chakravarty²,
and Marco Pavone¹

¹ Autonomous Systems Lab, Stanford University, Stanford, USA
{timsal,borisi,pavone}@stanford.edu

² Ford Greenfield Labs, Palo Alto, USA
pchakra5@ford.com

Abstract. Reasoning about human motion is an important prerequisite to safe and socially-aware robotic navigation. As a result, multi-agent behavior prediction has become a core component of modern human-robot interactive systems, such as self-driving cars. While there exist many methods for trajectory forecasting, most do not enforce dynamic constraints and do not account for environmental information (e.g., maps). Towards this end, we present *Trajectron++*, a modular, graph-structured recurrent model that forecasts the trajectories of a general number of diverse agents while incorporating agent dynamics and heterogeneous data (e.g., semantic maps). *Trajectron++* is designed to be tightly integrated with robotic planning and control frameworks; for example, it can produce predictions that are optionally conditioned on ego-agent motion plans. We demonstrate its performance on several challenging real-world trajectory forecasting datasets, outperforming a wide array of state-of-the-art deterministic and generative methods.

Keywords: Trajectory forecasting · Spatiotemporal graph modeling · Human-robot interaction · Autonomous driving

1 Introduction

Predicting the future behavior of humans is a necessary part of developing safe human-interactive autonomous systems. Humans can naturally navigate through many social interaction scenarios because they have an intrinsic “theory of

T. Salzmann and B. Ivanovic—Equal contribution.

T. Salzmann—Work done as a visiting student in the Autonomous Systems Lab.

Electronic supplementary material The online version of this chapter (https://doi.org/10.1007/978-3-030-58523-5_40) contains supplementary material, which is available to authorized users.

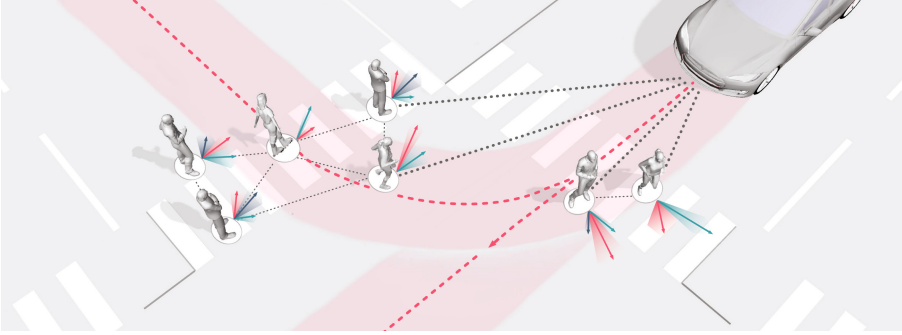


Fig. 1. Exemplary road scene depicting pedestrians crossing a road in front of a vehicle which may continue straight or turn right. The graph representation of the scene is shown on the ground, where each agent and their interactions are represented as nodes and edges, visualized as white circles and dashed black lines, respectively. Arrows depict potential future agent velocities, with colors representing different high-level future behavior modes.

mind,” which is the capacity to reason about other people’s actions in terms of their mental states [14]. As a result, imbuing autonomous systems with this capability could enable more informed decision making and proactive actions to be taken in the presence of other intelligent agents, e.g., in human-robot interaction scenarios. Figure 1 illustrates a scenario where predicting the intent of other agents may inform an autonomous vehicle’s path planning and decision making. Indeed, multi-agent behavior prediction has already become a core component of modern robotic systems, especially in safety-critical applications like self-driving vehicles which are currently being tested in the real world and targeting widespread deployment in the near future [49].

There are many existing methods for multi-agent behavior prediction, ranging from deterministic regressors to generative, probabilistic models. However, many of them were developed without directly accounting for real-world robotic use cases; in particular, they ignore agents’ dynamics constraints, the ego-agent’s own motion (important to capture the interactive aspect in human-robot interaction), and a plethora of environmental information (e.g., camera images, lidar, maps) to which modern robotic systems have access. Table 1 provides a summary of recent state-of-the-art approaches and their consideration of such desiderata.

Accordingly, in this work we are interested in developing a multi-agent behavior prediction model that (1) accounts for the dynamics of the agents, and in particular of ground vehicles [27,35]; (2) produces predictions possibly conditioned on potential future robot trajectories, useful for intelligent planning taking into account human responses; and (3) provides a generally-applicable, open, and extensible approach which can effectively use heterogeneous data about the surrounding environment. Importantly, making use of such data would allow for the incorporation of environmental information, e.g., maps, which would enable producing predictions that differ depending on the structure of the scene (e.g.,

Table 1. A summary of recent state-of-the-art pedestrian (left) and vehicle (right) trajectory forecasting methods, indicating the desiderata addressed by each approach.

| Method | GNA | CD | HD | FCP | OS | Method | GNA | CD | HD | FCP | OS |
|-----------------|-----|----|----|-----|----|---------------|-----|----|----|-----|----|
| DESIRE [31] | ✓ | | ✓ | | | IntentNet [8] | ✓ | | ✓ | | |
| Trajectron [20] | ✓ | | | | ✓ | PRECOG [39] | | | ✓ | ✓ | ✓ |
| S-BiGAT [28] | ✓ | | ✓ | | | MFP [18] | ✓ | | ✓ | | ✓ |
| DRF-Net [23] | ✓ | | ✓ | | | NMP [51] | ✓ | | ✓ | | |
| MATF [53] | ✓ | | ✓ | | ✓ | SpAGNN [7] | ✓ | | ✓ | | |
| Our work | ✓ | ✓ | ✓ | ✓ | ✓ | Our work | ✓ | ✓ | ✓ | ✓ | ✓ |

Legend: GNA = General Number of Agents, CD = Considers Dynamics, HD = Heterogeneous Data, FCP = Future-Conditional Predictions, OS = Open Source.

interactions at an urban intersection are very different from those in an open sports field!). One method that comes close is the Trajectron [20], a multi-agent behavior model which can handle a time-varying number of agents, accounts for multimodality in human behavior (i.e., the potential for many high-level futures), and maintains a sense of interpretability in its outputs. However, the Trajectron only reasons about relatively simple vehicle models (i.e., cascaded integrators) and past trajectory data (i.e., no considerations are made for added environmental information, if available).

In this work we present *Trajectron++*, an open and extensible approach built upon the Trajectron [20] framework which produces dynamically-feasible trajectory forecasts from heterogeneous input data for multiple interacting agents of distinct semantic types. Our key contributions are twofold: First, we show how to effectively incorporate high-dimensional data through the lens of encoding semantic maps. Second, we propose a general method of incorporating dynamics constraints into learning-based methods for multi-agent trajectory forecasting. *Trajectron++* is designed to be tightly integrated with downstream robotic modules, with the ability to produce trajectories that are optionally conditioned on future ego-agent motion plans. We present experimental results on a variety of datasets, which collectively demonstrate that *Trajectron++* outperforms an extensive selection of state-of-the-art deterministic and generative trajectory prediction methods, in some cases achieving 60% lower average prediction error.

2 Related Work

Deterministic Regressors. Many earlier works in human trajectory forecasting were deterministic regression models. One of the earliest, the Social Forces model [16], models humans as physical objects affected by Newtonian forces (e.g., with attractors at goals and repulsors at other agents). Since then, many approaches have been applied to the problem of trajectory forecasting, formulating it as a time-series regression problem and applying methods like Gaussian Process Regression (GPR) [38, 48], Inverse Reinforcement Learning (IRL) [32],

and Recurrent Neural Networks (RNNs) [1, 34, 47] to good effect. An excellent review of such methods can be found in [40].

Generative, Probabilistic Approaches. Recently, generative approaches have emerged as state-of-the-art trajectory forecasting methods due to recent advancements in deep generative models [12, 44]. Notably, they have caused a shift from focusing on predicting the single best trajectory to producing a *distribution* of potential future trajectories. This is advantageous in autonomous systems as full distribution information is more useful for downstream tasks, e.g., motion planning and decision making where information such as variance can be used to make safer decisions. Most works in this category use a deep recurrent backbone architecture with a latent variable model, such as a Conditional Variational Autoencoder (CVAE) [44], to explicitly encode multimodality [11, 20, 21, 31, 39, 42], or a Generative Adversarial Network (GAN) [12] to implicitly do so [13, 28, 41, 53]. Common to both approach styles is the need to produce position distributions. GAN-based models can directly produce these and CVAE-based recurrent models usually rely on a bivariate Gaussian Mixture Model (GMM) to output position distributions. However, both of these output structures make it difficult to enforce dynamics constraints, e.g., non-holonomic constraints such as those arising from no side-slip conditions. Of these, the Trajectron [20] and MATF [53] are the best-performing CVAE-based and GAN-based models, respectively, on standard pedestrian trajectory forecasting benchmarks [33, 37].

Accounting for Dynamics and Heterogeneous Data. There are few works that account for dynamics or make use of data modalities outside of prior trajectory information. This is mainly because standard trajectory forecasting benchmarks seldom include any other information, a fact that will surely change following the recent release of autonomous vehicle-based datasets with rich multi-sensor data [6, 9, 26, 50]. As for dynamics, current methods almost exclusively reason about positional information. This does not capture dynamical constraints, however, which might lead to predictions in position space that are unrealizable by the underlying control variables (e.g., a car moving sideways). Table 1 provides a detailed breakdown of recent state-of-the-art approaches and their consideration of these desiderata.

3 Problem Formulation

We aim to generate plausible trajectory distributions for a time-varying number $N(t)$ of interacting agents $A_1, \dots, A_{N(t)}$. Each agent A_i has a semantic class S_i , e.g., Car, Bus, or Pedestrian. At time t , given the state $\mathbf{s} \in \mathbb{R}^D$ of each agent and all of their histories for the previous H timesteps, which we denote as \mathbf{x} , $\mathbf{x} = \mathbf{s}_{1, \dots, N(t)}^{(t-H:t)} \in \mathbb{R}^{(H+1) \times N(t) \times D}$, as well as additional information available to each agent $I_{1, \dots, N(t)}^{(t)}$, we seek a distribution over all agents' future states for the next T timesteps $\mathbf{y} = \mathbf{s}_{1, \dots, N(t)}^{(t+1:t+T)} \in \mathbb{R}^{T \times N(t) \times D}$, which we denote as $p(\mathbf{y} \mid \mathbf{x}, I)$.

We also assume that geometric semantic maps are available around A_i 's position, $M_i^{(t)} \in \mathbb{R}^{\lceil C/r \rceil \times \lceil C/r \rceil \times L}$, with context size $C \times C$, spatial resolution r ,

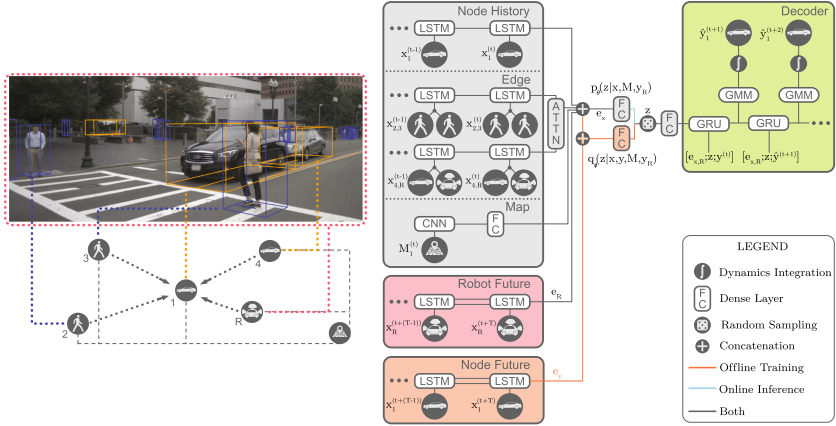


Fig. 2. Left: Our approach represents a scene as a directed spatiotemporal graph. Nodes and edges represent agents and their interactions, respectively. **Right:** The corresponding network architecture for Node 1.

and L semantic channels. Depending on the dataset, these maps can range in sophistication from simple obstacle occupancy grids to multiple layers of human-annotated semantic information (e.g., marking out sidewalks, road boundaries, and crosswalks).

We also consider the setting where we condition on an ego-agent’s future motion plan, for example when evaluating responses to a set of motion primitives. In this setting, we additionally assume that we know the ego-agent’s future motion plan for the next T timesteps, $\mathbf{y}_R = \mathbf{s}_R^{(t+1:t+T)}$.

4 Trajectron++

Our approach¹ is visualized in Fig. 2. At a high level, a spatiotemporal graph representation of the scene in question is created from its topology. Then, a similarly-structured deep learning architecture is generated that forecasts the evolution of node attributes, producing agent trajectories.

Scene Representation. The current scene is abstracted as a spatiotemporal graph $G = (V, E)$. Nodes represent agents and edges represent their interactions. As a result, in the rest of the paper we will use the terms “node” and “agent” interchangeably. Each node also has a semantic class matching the class of its agent (e.g., Car, Bus, Pedestrian). An edge (A_i, A_j) is present in E if A_i influences A_j . In this work, the ℓ_2 distance is used as a proxy for whether agents are influencing each other or not. Formally, an edge is directed from A_i to A_j if $\|\mathbf{p}_i - \mathbf{p}_j\|_2 \leq d_{S_j}$ where $\mathbf{p}_i, \mathbf{p}_j \in \mathbb{R}^2$ are the 2D world positions of agents A_i, A_j ,

¹ All of our source code, trained models, and data can be found online at <https://github.com/StanfordASL/Trajectron-plus-plus>.

respectively, and d_{S_j} is a distance that encodes the perception range of agents of semantic class S_j . While more sophisticated methods can be used to construct edges (e.g., [47]), they usually incur extra computational overhead by requiring a complete scene graph. Figure 2 shows an example of this scene abstraction.

We specifically choose to model the scene as a directed graph, in contrast to an undirected one as in previous approaches [1, 13, 20–22, 47], because a directed graph can represent a more general set of scenes and interaction types, e.g., asymmetric influence. This provides the additional benefit of being able to simultaneously model agents with different perception ranges, e.g., the driver of a car looks much farther ahead on the road than a pedestrian does while walking on the sidewalk.

Modeling Agent History. Once a graph of the scene is constructed, the model needs to encode a node’s current state, its history, and how it is influenced by its neighboring nodes. To encode the observed history of the modeled agent, their current and previous states are fed into a Long Short-Term Memory (LSTM) network [19] with 32 hidden dimensions. Since we are interested in modeling trajectories, the inputs $\mathbf{x} = \mathbf{s}_{1,\dots,N(t)}^{(t-H:t)} \in \mathbb{R}^{(H+1) \times N(t) \times D}$ are the current and previous D -dimensional states of the modeled agents. These are typically positions and velocities, which can be easily estimated online.

Ideally, agent models should be chosen to best match their semantic class S_i . For example, one would usually model vehicles on the road using a bicycle model [27, 35]. However, estimating the bicycle model parameters of another vehicle from online observations is very difficult as it requires estimation of the vehicle’s center of mass, wheelbase, and front wheel steer angle. As a result, in this work pedestrians are modeled as single integrators and wheeled vehicles are modeled as dynamically-extended unicycles [29], enabling us to account for key non-holonomic constraints (e.g., no side-slip constraints) [35] without requiring complex online parameter estimation procedures – we will show through experiments that such a simplified model is already quite impactful on improving prediction accuracy. While the dynamically-extended unicycle model serves as an important representative example, we note that our approach can also be generalized to other dynamics models, provided its parameters can either be assumed or quickly estimated online.

Encoding Agent Interactions. To model neighboring agents’ influence on the modeled agent, *Trajectron++* encodes graph edges in two steps. First, edge information is aggregated from neighboring agents of the same semantic class. In this work, an element-wise sum is used as the aggregation operation. We choose to combine features in this way rather than with concatenation or an average to handle a variable number of neighboring nodes with a fixed architecture while preserving count information [3, 21, 22]. These aggregated states are then fed into an LSTM with 8 hidden dimensions whose weights are shared across all edge instances of the same type, e.g., all Pedestrian-Bus edge LSTMs share the same weights. Then, the encodings from all edge types that connect to the modeled node are aggregated to obtain one “influence” representation

vector, representing the effect that all neighboring nodes have. For this, an additive attention module is used [2]. Finally, the node history and edge influence encodings are concatenated to produce a single node representation vector, $e_{\mathbf{x}}$.

Incorporating Heterogeneous Data. Modern sensor suites are able to produce much more information than just tracked trajectories of other agents. Notably, HD maps are used by many real-world systems to aid localization as well as inform navigation. Depending on sensor availability and sophistication, maps can range in fidelity from simple binary obstacle maps, i.e., $M \in \{0, 1\}^{H \times W \times 1}$, to HD semantic maps, e.g., $M \in \{0, 1\}^{H \times W \times L}$ where each layer $1 \leq \ell \leq L$ corresponds to an area with semantic type (e.g., “driveable area,” “road block,” “walkway,” “pedestrian crossing”). To make use of this information, for each modeled agent, *Trajectron++* encodes a local map, rotated to match the agent’s heading, with a Convolutional Neural Network (CNN). The CNN has 4 layers, with filters $\{5, 5, 5, 3\}$ and respective strides of $\{2, 2, 1, 1\}$. These are followed by a dense layer with 32 hidden dimensions, the output of which is concatenated with the node history and edge influence representation vectors.

More generally, one can include further additional information (e.g., raw LIDAR data, camera images, pedestrian skeleton or gaze direction estimates) in this framework by encoding it as a vector and adding it to this backbone of representation vectors, $e_{\mathbf{x}}$.

Encoding Future Ego-Agent Motion Plans. Producing predictions which take into account future ego-agent motion is an important capability for robotic decision making and control. Specifically, it allows for the evaluation of a set of motion primitives with respect to possible responses from other agents. *Trajectron++* can encode the future T timesteps of the ego-agent’s motion plan \mathbf{y}_R using a bi-directional LSTM with 32 hidden dimensions. A bi-directional LSTM is used due to its strong performance on other sequence summarization tasks [5]. The final hidden states are then concatenated into the backbone of representation vectors, $e_{\mathbf{x}}$.

Explicitly Accounting for Multimodality. *Trajectron++* explicitly handles multimodality by leveraging the CVAE latent variable framework [44]. It produces the target $p(\mathbf{y} \mid \mathbf{x})$ distribution by introducing a discrete Categorical latent variable $z \in Z$ which encodes high-level latent behavior and allows for $p(\mathbf{y} \mid \mathbf{x})$ to be expressed as $p(\mathbf{y} \mid \mathbf{x}) = \sum_{z \in Z} p_{\psi}(\mathbf{y} \mid \mathbf{x}, z) p_{\theta}(z \mid \mathbf{x})$, where $|Z| = 25$ and ψ, θ are deep neural network weights that parameterize their respective distributions. z being discrete also aids in interpretability, as one can visualize which high-level behaviors belong to each z by sampling trajectories.

During training, a bi-directional LSTM with 32 hidden dimensions is used to encode a node’s ground truth future trajectory, producing $q_{\phi}(z \mid \mathbf{x}, \mathbf{y})$ [44].

Producing Dynamically-Feasible Trajectories. After obtaining a latent variable z , it and the backbone representation vector $e_{\mathbf{x}}$ are fed into the decoder, a 128-dimensional Gated Recurrent Unit (GRU) [10]. Each GRU cell outputs the parameters of a bivariate Gaussian distribution over control actions $\mathbf{u}^{(t)}$ (e.g., acceleration and steering rate). The agent’s system dynamics are

then integrated with the produced control actions $\mathbf{u}^{(t)}$ to obtain trajectories in position space [25, 46]. The only uncertainty at prediction time stems from *Trajectron++*'s output. Thus, in the case of linear dynamics (e.g., single integrators, used in this work to model pedestrians), the system dynamics are linear Gaussian. Explicitly, for a single integrator with control actions $\mathbf{u}^{(t)} = \dot{\mathbf{p}}^{(t)}$, the position mean at $t + 1$ is $\mu_{\mathbf{p}}^{(t+1)} = \mu_{\mathbf{p}}^{(t)} + \mu_{\mathbf{u}}^{(t)} \Delta t$, where $\mu_{\mathbf{u}}^{(t)}$ is produced by *Trajectron++*. In the case of nonlinear dynamics (e.g., unicycle models, used in this work to model vehicles), one can still (approximately) use this uncertainty propagation scheme by linearizing the dynamics about the agent's current state and control. Full mean and covariance equations for the single integrator and dynamically-extended unicycle models are in the appendix. In contrast to existing methods which directly output positions, our approach is uniquely able to guarantee that its trajectory samples are dynamically feasible by integrating an agent's dynamics with the predicted controls.

Output Configurations. Based on the desired use case, *Trajectron++* can produce many different outputs. The main four are outlined below.

1. *Most Likely (ML)*: The model's deterministic and most-likely single output. The high-level latent behavior mode and output trajectory are the modes of their respective distributions, where

$$z_{\text{mode}} = \arg \max_z p_{\theta}(z \mid \mathbf{x}), \quad \mathbf{y} = \arg \max_{\mathbf{y}} p_{\psi}(\mathbf{y} \mid \mathbf{x}, z_{\text{mode}}). \quad (1)$$

2. *z_{mode}* : Predictions from the model's most-likely high-level latent behavior mode, where

$$z_{\text{mode}} = \arg \max_z p_{\theta}(z \mid \mathbf{x}), \quad \mathbf{y} \sim p_{\psi}(\mathbf{y} \mid \mathbf{x}, z_{\text{mode}}). \quad (2)$$

3. *Full*: The model's full sampled output, where z and \mathbf{y} are sampled sequentially according to

$$z \sim p_{\theta}(z \mid \mathbf{x}), \quad \mathbf{y} \sim p_{\psi}(\mathbf{y} \mid \mathbf{x}, z). \quad (3)$$

4. *Distribution*: Due to the use of a discrete latent variable and Gaussian output structure, the model can provide an analytic output distribution by directly computing $p(\mathbf{y} \mid \mathbf{x}) = \sum_{z \in Z} p_{\psi}(\mathbf{y} \mid \mathbf{x}, z) p_{\theta}(z \mid \mathbf{x})$.

Training the Model. We adopt the InfoVAE [52] objective function, and modify it to use discrete latent states in a conditional formulation (since the model uses a CVAE). Formally, we aim to solve

$$\begin{aligned} \max_{\phi, \theta, \psi} \sum_{i=1}^N \mathbb{E}_{z \sim q_{\phi}(\cdot \mid \mathbf{x}_i, \mathbf{y}_i)} [\log p_{\psi}(\mathbf{y}_i \mid \mathbf{x}_i, z)] \\ - \beta D_{KL}(q_{\phi}(z \mid \mathbf{x}_i, \mathbf{y}_i) \parallel p_{\theta}(z \mid \mathbf{x}_i)) + \alpha I_q(\mathbf{x}; z), \end{aligned} \quad (4)$$

where I_q is the mutual information between \mathbf{x} and z under the distribution $q_{\phi}(\mathbf{x}, z)$. To compute I_q , we follow [52] and approximate $q_{\phi}(z \mid \mathbf{x}_i, \mathbf{y}_i)$ with

$p_{\theta}(z \mid \mathbf{x}_i)$, obtaining the unconditioned latent distribution by summing out \mathbf{x}_i over the batch. Notably, the Gumbel-Softmax reparameterization [24] is not used to backpropagate through the Categorical latent variable z because it is not sampled during training time. Instead, the first term of Eq. (4) is directly computed since the latent space has only $|Z| = 25$ discrete elements. Additional training details can be found in the appendix.

5 Experiments

Our method is evaluated on three publicly-available datasets: The ETH [37], UCY [33], and nuScenes [6] datasets. The ETH and UCY datasets consist of real pedestrian trajectories with rich multi-human interaction scenarios captured at 2.5 Hz ($\Delta t = 0.4$ s). In total, there are 5 sets of data, 4 unique scenes, and 1536 unique pedestrians. They are a standard benchmark in the field, containing challenging behaviors such as couples walking together, groups crossing each other, and groups forming and dispersing. However, they only contain pedestrians, so we also evaluate on the recently-released nuScenes dataset. It is a large-scale dataset for autonomous driving with 1000 scenes in Boston and Singapore. Each scene is annotated 2 Hz ($\Delta t = 0.5$ s) and is 20 s long, containing up to 23 semantic object classes as well as HD semantic maps with 11 annotated layers.

Trajectron++ was implemented in PyTorch [36] on a desktop computer running Ubuntu 18.04 containing an AMD Ryzen 1800X CPU and two NVIDIA GTX 1080 Ti GPUs. We trained the model for 100 epochs (~ 3 h) on the pedestrian datasets and 12 epochs (~ 8 h) on the nuScenes dataset.

Evaluation Metrics. As in prior work [1, 13, 20, 28, 41, 53], our method for trajectory forecasting is evaluated with the following four error metrics:

1. *Average Displacement Error (ADE)*: Mean ℓ_2 distance between the ground truth and predicted trajectories.
2. *Final Displacement Error (FDE)*: ℓ_2 distance between the predicted final position and the ground truth final position at the prediction horizon T .
3. *Kernel Density Estimate-based Negative Log Likelihood (KDE NLL)*: Mean NLL of the ground truth trajectory under a distribution created by fitting a kernel density estimate on trajectory samples [20, 45].
4. *Best-of- N (BoN)*: The minimum ADE and FDE from N randomly-sampled trajectories. We compare our method to an exhaustive set of state-of-the-art deterministic and generative approaches.

Deterministic Baselines. Our method is compared against the following deterministic baselines: (1) *Linear*: A linear regressor with parameters estimated by minimizing least square error. (2) *LSTM*: An LSTM network with only agent history information. (3) *Social LSTM* [1]: Each agent is modeled with an LSTM and nearby agents’ hidden states are pooled at each timestep using a proposed social pooling operation. (4) *Social Attention* [47]: Same as [1], but all other agents’ hidden states are incorporated via a proposed social attention operation.

Generative Baselines. On the ETH and UCY datasets, our method is compared against the following generative baselines: (1) *S-GAN* [13]: Each agent is modeled with an LSTM-GAN, which is an LSTM encoder-decoder whose outputs are the generator of a GAN. The generated trajectories are then evaluated against the ground truth trajectories with a discriminator. (2) *SoPhie* [41]: An LSTM-GAN with the addition of a proposed physical and social attention module. (3) *MATF* [53]: An LSTM-GAN model that leverages CNNs to fuse agent relationships and encode environmental information. (4) *Trajectron* [20]: An LSTM-CVAE encoder-decoder which is explicitly constructed to match the spatiotemporal structure of the scene. Its scene abstraction is similar to ours, but uses undirected edges.

On the nuScenes dataset, the following methods are also compared against: (5) *Convolutional Social Pooling (CSP)* [11]: An LSTM-based approach which explicitly considers a fixed number of movement classes and predicts which of those the modeled agent is likely to take. (6) *CAR-Net* [42]: An LSTM-based approach which encodes scene context with visual attention. (7) *SpAGNN* [7]: A CNN encodes raw LIDAR and semantic map data to produce object detections, from which a Graph Neural Network (GNN) produces probabilistic, interaction-aware trajectories.

Evaluation Methodology. For the ETH and UCY datasets, a leave-one-out strategy is used for evaluation, similar to previous works [1, 13, 20, 28, 41, 53], where the model is trained on four datasets and evaluated on the held-out fifth. An observation length of 8 timesteps (3.2s) and a prediction horizon of 12 timesteps (4.8s) is used for evaluation. For the nuScenes dataset, we split off 15% of the train set for hyperparameter tuning and test on the provided validation set.

Throughout the following, we report the performance of *Trajectron++* in multiple configurations. Specifically, *Ours* refers to the base model using only node and edge encoding, trained to predict agent velocities and Euler integrating velocity to produce positions; *Ours*+ \int is the base model with dynamics integration, trained to predict control actions and integrating the agent’s dynamics with the control actions to produce positions; *Ours*+ \int, M additionally includes the map encoding CNN; and *Ours*+ \int, M, \mathbf{y}_R adds the robot future encoder.

5.1 ETH and UCY Datasets

Our approach is first evaluated on the ETH [37] and UCY [33] Pedestrian Datasets, against deterministic methods on standard trajectory forecasting metrics. It is difficult to determine the current state-of-the-art in deterministic methods as there are contradictions between the results reported by the same authors in [13] and [1]. In Table 1 of [1], Social LSTM *convincingly* outperforms a baseline LSTM without pooling. However, in Table 1 of [13], Social LSTM is actually *worse* than the same baseline on average. Thus, when comparing against Social LSTM we report the results summarized in Table 1 of [13] as it is the most recent work by the same authors. Further, the values reported by Social Attention in [47] seem to have unusually high ratios of FDE to ADE. Nearly every

Table 2. (a) Our model’s deterministic Most Likely output outperforms other deterministic methods on displacement error metrics, even if it was not originally trained to do so. (b) Our model’s probabilistic Full output significantly outperforms other methods, yielding accurate predictions even in a small number of samples. Lower is better. Bold indicates best.

| Dataset | (a) ADE/FDE (m) | | | | | |
|---------|-----------------|-----------|-------------|------------------|-------------------|-------------------|
| | Linear | LSTM | S-LSTM [13] | S-ATTN [47] | Ours (ML) | Ours+ \int (ML) |
| ETH | 1.33/2.94 | 1.09/2.41 | 1.09/2.35 | 0.39/3.74 | 0.71/ 1.66 | 0.71/1.68 |
| Hotel | 0.39/0.72 | 0.86/1.91 | 0.79/1.76 | 0.29/2.64 | 0.22/0.46 | 0.22/0.46 |
| Univ | 0.82/1.59 | 0.61/1.31 | 0.67/1.40 | 0.33/3.92 | 0.44/1.17 | 0.41/ 1.07 |
| Zara 1 | 0.62/1.21 | 0.41/0.88 | 0.47/1.00 | 0.20/0.52 | 0.30/0.79 | 0.30/0.77 |
| Zara 2 | 0.77/1.48 | 0.52/1.11 | 0.56/1.17 | 0.30/2.13 | 0.23/0.59 | 0.23/0.59 |
| Average | 0.79/1.59 | 0.70/1.52 | 0.72/1.54 | 0.30/2.59 | 0.38/0.93 | 0.37/ 0.91 |

| Dataset | (b) ADE/FDE, Best of 20 Samples (m) | | | | | |
|---------|-------------------------------------|-------------|-----------------|-----------|------------------|---------------------|
| | S-GAN [13] | SoPhie [41] | Trajectron [20] | MATF [53] | Ours (Full) | Ours+ \int (Full) |
| ETH | 0.81/1.52 | 0.70/1.43 | 0.59/1.14 | 1.01/1.75 | 0.39/0.83 | 0.43/0.86 |
| Hotel | 0.72/1.61 | 0.76/1.67 | 0.35/0.66 | 0.43/0.80 | 0.12/0.21 | 0.12/0.19 |
| Univ | 0.60/1.26 | 0.54/1.24 | 0.54/1.13 | 0.44/0.91 | 0.20/0.44 | 0.22/ 0.43 |
| Zara 1 | 0.34/0.69 | 0.30/0.63 | 0.43/0.83 | 0.26/0.45 | 0.15/0.33 | 0.17/ 0.32 |
| Zara 2 | 0.42/0.84 | 0.38/0.78 | 0.43/0.85 | 0.26/0.57 | 0.11/0.25 | 0.12/ 0.25 |
| Average | 0.58/1.18 | 0.54/1.15 | 0.47/0.92 | 0.48/0.90 | 0.19/0.41 | 0.21/ 0.41 |

Legend: \int = Integration via Dynamics, M = Map Encoding, \mathbf{y}_R = Robot Future Encoding.

other method (including ours) has FDE/ADE ratios around 2–3 \times whereas Social Attention’s are around 3–12 \times . Social Attention’s errors on the Univ dataset are especially striking, as its FDE of 3.92 is 12 \times its ADE of 0.33, meaning its prediction error on the other 11 timesteps is essentially zero. We still compare against the values reported in [47] as there is no publicly-released code, but this raises doubts of their validity. To fairly compare against prior work, neither map encoding nor future motion plan encoding is used. Only the node history and edge encoders are used in the model’s encoder. Additionally, the model’s deterministic ML output scheme is employed, which produces the model’s most likely single trajectory. Table 2(a) summarizes these results and shows that our approach is competitive with state-of-the-art deterministic regressors on displacement error metrics (outperforming existing approaches by 33% on mean FDE), even though our method was not originally trained to minimize this. It makes sense that the model performs similarly with and without dynamics integration for pedestrians, since they are modeled as single integrators. Thus, their control actions are velocities which matches the base model’s output structure.

To more concretely compare generative methods, we use the KDE-based NLL metric proposed in [20, 45], an approach that maintains full output distributions and compares the log-likelihood of the ground truth under different methods’ outputs. Table 3 summarizes these results and shows that our method significantly outperforms others. This is also where the performance improve-

Table 3. Mean KDE-based NLL for each dataset. Lower is better. 2000 trajectories were sampled per model at each prediction timestep. Bold indicates the best values.

| Dataset | KDE NLL | | | |
|---------|------------|-----------------|-------------|---------------------|
| | S-GAN [13] | Trajectron [20] | Ours (Full) | Ours+ \int (Full) |
| ETH | 15.70 | 2.99 | 1.80 | 1.31 |
| Hotel | 8.10 | 2.26 | -1.29 | -1.94 |
| Univ | 2.88 | 1.05 | -0.89 | -1.13 |
| Zara 1 | 1.36 | 1.86 | -1.13 | -1.41 |
| Zara 2 | 0.96 | 0.81 | -2.19 | -2.53 |
| Average | 5.80 | 1.79 | -0.74 | -1.14 |

Legend: \int = Integration via Dynamics, M = Map Encoding, \mathbf{y}_R = Robot Future Encoding.

ments brought by the dynamics integration scheme are clear. It yields the best performance because the model is now explicitly trained on the distribution it is seeking to output (the loss function term $p_\psi(\mathbf{y}|\mathbf{x}, z)$ is now directly over positions), whereas the base model is trained on velocity distributions, the integration of which (with no accounting for system dynamics) introduces errors. Unfortunately, at this time there are no publicly-released models for SoPhie [41] or MATF [53], so they cannot be evaluated with the KDE-based NLL metric. Instead, we evaluate *Trajectron++* with the Best-of- N metric used in their works. Table 2(b) summarizes these results, and shows that our method *significantly* outperforms the state-of-the-art [53], achieving 55–60% lower average errors.

Map Encoding. To evaluate the effect of incorporating heterogeneous data, we compare the performance of *Trajectron++* with and without the map encoder. Specifically, we compare the frequency of obstacle violations in 2000 trajectory samples from the Full model output on the ETH - University scene, which provides a simple binary obstacle map. Overall, our approach generates colliding predictions 1.0% of the time with map encoding, compared to 4.6% without map encoding. We also study how much of a reduction there is for pedestrians that are especially close to an obstacle (i.e. they have at least one obstacle-violating trajectory in their Full output), an example of which is shown in the appendix. In this regime, our approach generates colliding predictions 4.9% of the time with map encoding, compared to 21.5% without map encoding.

5.2 nuScenes Dataset

To further evaluate the model’s ability to use heterogeneous data and simultaneously model multiple semantic classes of agents, we evaluate it on the nuScenes dataset [6]. Again, the deterministic ML output scheme is used to fairly compare with other single-trajectory predictors. The trajectories of both Pedestrians and Cars are forecasted, two semantic object classes which account for most of the 23

Table 4. [nuScenes] (a): Vehicle-only FDE across time for *Trajectron++* compared to that of other single-trajectory and probabilistic approaches. Bold indicates best. (b): Pedestrian-only FDE and KDE NLL across time for *Trajectron++*.

| (a) Vehicle-only | | | | | | | |
|----------------------|-------------|-------------|-------------|-------------|--|--|--|
| Method | FDE (m) | | | | | | |
| | @1 s | @2 s | @3 s | @4 s | | | |
| Const. Velocity | 0.32 | 0.89 | 1.70 | 2.73 | | | |
| S-LSTM* [1, 7] | 0.47 | – | 1.61 | – | | | |
| CSP* [7, 11] | 0.46 | – | 1.50 | – | | | |
| CAR-Net* [7, 42] | 0.38 | – | 1.35 | – | | | |
| SpAGNN* [7] | 0.36 | – | 1.23 | – | | | |
| Ours (ML) | 0.18 | 0.57 | 1.25 | 2.24 | | | |
| Ours+ \int, M (ML) | 0.07 | 0.45 | 1.14 | 2.20 | | | |

| (b) Pedestrian-only | | | | | | | |
|----------------------|---------|-------|-------|-------|---------|------|------|
| Method | KDE NLL | | | | FDE (m) | | |
| | @1 s | @2 s | @3 s | @4 s | @1 s | @2 s | @3 s |
| Ours (ML) | –2.69 | –2.46 | –1.76 | –1.09 | 0.03 | 0.17 | 0.37 |
| Ours+ \int, M (ML) | –5.58 | –3.96 | –2.77 | –1.89 | 0.01 | 0.17 | 0.37 |

*We subtracted 22–24 cm from these reported values (their detection/tracking error [7]), as we do not use a detector/tracker. This is done to establish a fair comparison. Legend: \int = Integration via Dynamics, M = Map Encoding, \mathbf{y}_R = Robot Future Encoding.

possible object classes present in the dataset. To obtain an estimate of prediction quality degradation over time, we compute the model’s FDE at $t = \{1, 2, 3, 4\}$ s for all tracked objects with at least 4 s of available future data. We also implement a constant velocity baseline, which simply maintains the agent’s heading and speed for the prediction horizon. Table 4(a) summarizes the model’s performance in comparison with state-of-the-art vehicle trajectory prediction models. Since other methods use a detection/tracking module (whereas ours does not), to establish a fair comparison we subtracted other methods’ detection and tracking error from their reported values. The dynamics integration scheme and map encoding yield a noticeable improvement with vehicles, as their dynamically-extended unicycle dynamics now differ from the single integrator assumption made by the base model. Note that our method was only trained to predict 3 s into the future, thus its performance at 4 s also provides a measure of its capability to generalize beyond its training configuration. Other methods do not report values at 2 s and 4 s. As can be seen, *Trajectron++* outperforms existing approaches without facing a sharp degradation in performance after 3 s. Our approach’s performance on pedestrians is reported in Table 4(b), where the inclusion of HD maps and dynamics integration similarly improve performance as in the pedestrian datasets.

Table 5. [nuScenes] (a): Vehicle-only prediction performance for ablated versions of our model. (b): The same, but excluding the ego-robot from consideration (as it is being conditioned on). This shows that our model’s robot future conditional performance does not arise from merely removing the ego-vehicle.

| (a) Including the Ego-vehicle | | | | | | | | | | | | |
|-------------------------------|---------|-------|-------|-------|------------|------|------|------|--------------|-----|-----|-----|
| Ablation | KDE NLL | | | | FDE ML (m) | | | | B. Viol. (%) | | | |
| \int M \mathbf{y}_R | @1s | @2s | @3s | @4s | @1s | @2s | @3s | @4s | @1s | @2s | @3s | @4s |
| – – – | 0.81 | 0.05 | 0.37 | 0.87 | 0.18 | 0.57 | 1.25 | 2.24 | 0.2 | 0.6 | 2.8 | 6.9 |
| ✓ – – | –4.28 | –2.82 | –1.67 | –0.76 | 0.07 | 0.45 | 1.13 | 2.17 | 0.2 | 0.7 | 3.2 | 8.1 |
| ✓ ✓ – | –4.17 | –2.74 | –1.62 | –0.70 | 0.07 | 0.45 | 1.14 | 2.20 | 0.3 | 0.6 | 2.8 | 7.6 |

| (b) Excluding the Ego-vehicle | | | | | | | | | | | | |
|-------------------------------|---------|-------|-------|-------|------------|------|------|------|--------------|-----|-----|-----|
| Ablation | KDE NLL | | | | FDE ML (m) | | | | B. Viol. (%) | | | |
| \int M \mathbf{y}_R | @1s | @2s | @3s | @4s | @1s | @2s | @3s | @4s | @1s | @2s | @3s | @4s |
| ✓ ✓ – | –4.26 | –2.86 | –1.76 | –0.87 | 0.07 | 0.44 | 1.09 | 2.09 | 0.3 | 0.6 | 2.8 | 7.6 |
| ✓ ✓ ✓ | –3.90 | –2.76 | –1.75 | –0.93 | 0.08 | 0.34 | 0.81 | 1.50 | 0.3 | 0.5 | 1.6 | 4.2 |

Legend: \int = Integration via Dynamics, M = Map Encoding, \mathbf{y}_R = Robot Future Encoding.

Ablation Study. To develop an understanding of which model components influence performance, a comprehensive ablation study is performed in Table 5. As can be seen in the first row, even the base model’s deterministic ML output performs strongly relative to current state-of-the-art approaches for vehicle trajectory forecasting [7]. Adding the dynamics integration scheme yields a drastic reduction in NLL as well as FDE at all prediction horizons. There is also an associated slight increase in the frequency of road boundary-violating predictions. This is a consequence of training in position (as opposed to velocity) space, which yields more variability in the corresponding predictions. Additionally including map encoding maintains prediction accuracy while reducing the frequency of boundary-violating predictions.

The effect of conditioning on the ego-vehicle’s future motion plan is also studied, with results summarized in Table 5(b). As one would expect, providing the model with future motion plans of the ego-vehicle yields significant reductions in error and road boundary violations. This use-case is common throughout autonomous driving as the ego-vehicle repeatedly produces future motion plans at every timestep by evaluating motion primitives. Overall, dynamics integration is the dominant performance-improving module.

Qualitative Comparison. Figure 3 shows trajectory predictions from the base model, with dynamics integration, and with dynamics integration + map encoding. In it, one can see that the base model (predicting in velocity space) under-shoots the turn for the red car, predicting that it will end up in oncoming traffic. With the integration of dynamics, the model captures multimodality in the agent’s action, predicting both the possibility of a right turn and continuing

straight. With the addition of map encoding, the predictions are not only more accurate, but nearly all probability mass now lies within the correct side of the road. This is in contrast to versions of the model without map encoding which predict that the red car might move into oncoming traffic.

Online Runtime. A key consideration in robotics is runtime complexity. As a result, we evaluate the time it takes *Trajectron++* to perform forward inference on commodity hardware. The results are summarized in the appendix, and confirm that our model scales well to scenes with many agents and interactions.

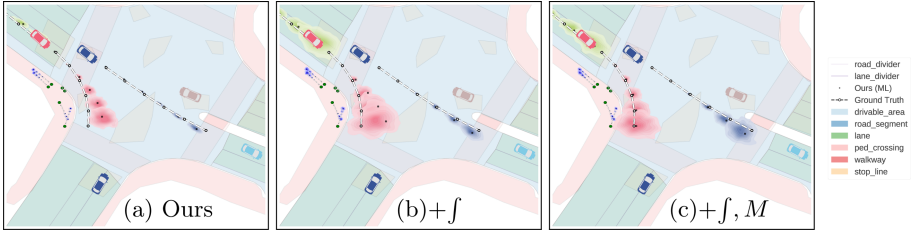


Fig. 3. [nuScenes] The same scene as forecast by three versions of *Trajectron++*. (a) The base model tends to under-shoot turns, and makes overly-confident predictions. (b) Our approach better captures position uncertainty with dynamics integration, producing well-calibrated probabilities. (c) The model is able to leverage the additional information that a map provides, yielding accurate predictions.

6 Conclusion

In this work, we present *Trajectron++*, a generative multi-agent trajectory forecasting approach which uniquely addresses our desiderata for an open, generally-applicable, and extensible framework. It can incorporate heterogeneous data beyond prior trajectory information and is able to produce future-conditional predictions that respect dynamics constraints, all while producing full probability distributions, which are especially useful in downstream robotic tasks such as motion planning, decision making, and control. It achieves state-of-the-art prediction performance in a variety of metrics on standard and new real-world multi-agent human behavior datasets.

Acknowledgment. This work was supported in part by the Ford-Stanford Alliance. This article solely reflects the opinions and conclusions of its authors.

References

1. Alahi, A., Goel, K., Ramanathan, V., Robicquet, A., Fei-Fei, L., Savarese, S.: Social LSTM: human trajectory prediction in crowded spaces. In: IEEE Conference on Computer Vision and Pattern Recognition (2016)

2. Bahdanau, D., Cho, K., Bengio, Y.: Neural machine translation by jointly learning to align and translate. In: International Conference on Learning Representations (2015)
3. Battaglia, P.W., Pascanu, R., Lai, M., Rezende, D., Kavukcuoglu, K.: Interaction networks for learning about objects, relations and physics. In: Conference on Neural Information Processing Systems (2016)
4. Bowman, S.R., Vilnis, L., Vinyals, O., Dai, A.M., Jozefowicz, R., Bengio, S.: Generating sentences from a continuous space. In: Proceedings of the Annual Meeting of the Association for Computational Linguistics (2015)
5. Britz, D., Goldie, A., Luong, M.T., Le, Q.V.: Massive exploration of neural machine translation architectures. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing, pp. 1442–1451 (2017)
6. Caesar, H., et al.: nuScenes: a multimodal dataset for autonomous driving (2019)
7. Casas, S., Gulino, C., Liao, R., Urtasun, R.: SpAGNN: spatially-aware graph neural networks for relational behavior forecasting from sensor data (2019)
8. Casas, S., Luo, W., Urtasun, R.: IntentNet: learning to predict intention from raw sensor data. In: Conference on Robot Learning, pp. 947–956 (2018)
9. Chang, M.F., et al.: Argoverse: 3D tracking and forecasting with rich maps. In: IEEE Conference on Computer Vision and Pattern Recognition (2019)
10. Cho, K., et al.: Learning phrase representations using RNN encoder-decoder for statistical machine translation. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing, pp. 1724–1734 (2014)
11. Deo, M.F., Trivedi, J.: Multi-modal trajectory prediction of surrounding vehicles with maneuver based LSTMs. In: IEEE Intelligent Vehicles Symposium (2018)
12. Goodfellow, I., et al.: Generative adversarial nets. In: Conference on Neural Information Processing Systems (2014)
13. Gupta, A., Johnson, J., Li, F., Savarese, S., Alahi, A.: Social GAN: socially acceptable trajectories with generative adversarial networks. In: IEEE Conference on Computer Vision and Pattern Recognition (2018)
14. Gweon, H., Saxe, R.: Developmental cognitive neuroscience of theory of mind, chap. 20. In: Neural Circuit Development and Function in the Brain, pp. 367–377. Academic Press (2013). <https://doi.org/10.1016/B978-0-12-397267-5.00057-1>. <http://www.sciencedirect.com/science/article/pii/B9780123972675000571>
15. Hallac, D., Leskovec, J., Boyd, S.: Network lasso: clustering and optimization in large graphs. In: ACM International Conference on Knowledge Discovery and Data Mining (2015)
16. Helbing, D., Molnár, P.: Social force model for pedestrian dynamics. *Phys. Rev. E* **51**(5), 4282–4286 (1995)
17. Higgins, I., et al.: β -VAE: learning basic visual concepts with a constrained variational framework. In: International Conference on Learning Representations (2017)
18. Ho, J., Ermon, S.: Multiple futures prediction. In: Conference on Neural Information Processing Systems (2019)
19. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural Comput.* **9**, 1735–1780 (1997)
20. Ivanovic, B., Pavone, M.: The trajectron: probabilistic multi-agent trajectory modeling with dynamic spatiotemporal graphs. In: IEEE International Conference on Computer Vision (2019)
21. Ivanovic, B., Schmerling, E., Leung, K., Pavone, M.: Generative modeling of multimodal multi-human behavior. In: IEEE/RSJ International Conference on Intelligent Robots & Systems (2018)

22. Jain, A., Zamir, A.R., Savarese, S., Saxena, A.: Structural-RNN: deep learning on spatio-temporal graphs. In: IEEE Conference on Computer Vision and Pattern Recognition (2016)
23. Jain, A., et al.: Discrete residual flow for probabilistic pedestrian behavior prediction. In: Conference on Robot Learning (2019)
24. Jang, E., Gu, S., Poole, B.: Categorical reparameterization with Gumbel-Softmax. In: International Conference on Learning Representations (2017)
25. Kalman, R.E.: A new approach to linear filtering and prediction problems. *ASME J. Basic Eng.* **82**, 35–45 (1960)
26. Kesten, R., et al.: Lyft Level 5 AV Dataset 2019 (2019). <https://level5.lyft.com/dataset/>
27. Kong, J., Pfeifer, M., Schildbach, G., Borrelli, F.: Kinematic and dynamic vehicle models for autonomous driving control design. In: IEEE Intelligent Vehicles Symposium (2015)
28. Kosaraju, V., et al.: Social-BiGAT: multimodal trajectory forecasting using bicycle-GAN and graph attention networks. In: Conference on Neural Information Processing Systems (2019)
29. LaValle, S.M.: Better unicycle models. In: Planning Algorithms, p. 743. Cambridge University Press (2006)
30. LaValle, S.M.: A simple unicycle. In: Planning Algorithms, pp. 729–730. Cambridge University Press (2006)
31. Lee, N., et al.: DESIRE: distant future prediction in dynamic scenes with interacting agents. In: IEEE Conference on Computer Vision and Pattern Recognition (2017)
32. Lee, N., Kitani, K.M.: Predicting wide receiver trajectories in American football. In: IEEE Winter Conference on Applications of Computer Vision (2016)
33. Lerner, A., Chrysanthou, Y., Lischinski, D.: Crowds by example. *Comput. Graph. Forum* **26**(3), 655–664 (2007)
34. Morton, J., Wheeler, T.A., Kochenderfer, M.J.: Analysis of recurrent neural networks for probabilistic modeling of driver behavior. *IEEE Trans. Pattern Anal. Mach. Intell.* **18**(5), 1289–1298 (2017)
35. Paden, B., Čáp, M., Yong, S.Z., Yershov, D., Frazzoli, E.: A survey of motion planning and control techniques for self-driving urban vehicles. *IEEE Trans. Intell. Veh.* **1**(1), 33–55 (2016)
36. Paszke, A., et al.: Automatic differentiation in PyTorch. In: Conference on Neural Information Processing Systems - Autodiff Workshop (2017)
37. Pellegrini, S., Ess, A., Schindler, K., Gool, L.: You’ll never walk alone: modeling social behavior for multi-target tracking. In: IEEE International Conference on Computer Vision (2009)
38. Rasmussen, C.E., Williams, C.K.I.: Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning), 1st edn. MIT Press, Cambridge (2006)
39. Rhinehart, N., McAllister, R., Kitani, K., Levine, S.: PRECOG: prediction conditioned on goals in visual multi-agent settings. In: IEEE International Conference on Computer Vision (2019)
40. Rudenko, A., Palmieri, L., Herman, M., Kitani, K.M., Gavrila, D.M., Arras, K.O.: Human motion trajectory prediction: a survey (2019). <https://arxiv.org/abs/1905.06113>

41. Sadeghian, A., Kosaraju, V., Sadeghian, A., Hirose, N., RezaTofighi, S.H., Savarese, S.: SoPhie: an attentive GAN for predicting paths compliant to social and physical constraints. In: IEEE Conference on Computer Vision and Pattern Recognition (2019)
42. Sadeghian, A., Legros, F., Voisin, M., Vesel, R., Alahi, A., Savarese, S.: CAR-Net: Clairvoyant attentive recurrent network. In: Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y. (eds.) ECCV 2018. LNCS, vol. 11215. Springer, Cham (2018). https://doi.org/10.1007/978-3-030-01252-6_10
43. Schöller, C., Aravantinos, V., Lay, F., Knoll, A.: What the constant velocity model can teach us about pedestrian motion prediction. *IEEE Robot. Autom. Lett.* **5**, 1696–1703 (2020)
44. Sohn, K., Lee, H., Yan, X.: Learning structured output representation using deep conditional generative models. In: Conference on Neural Information Processing Systems (2015)
45. Thiede, L.A., Brahma, P.P.: Analyzing the variety loss in the context of probabilistic trajectory prediction. In: IEEE International Conference on Computer Vision (2019)
46. Thrun, S., Burgard, W., Fox, D.: The extended Kalman filter. In: Probabilistic Robotics, pp. 54–64. MIT Press (2005)
47. Vemula, A., Mueller, K., Oh, J.: Social attention: modeling attention in human crowds. In: Proceedings of the IEEE Conference on Robotics and Automation (2018)
48. Wang, J.M., Fleet, D.J., Hertzmann, A.: Gaussian process dynamical models for human motion. *IEEE Trans. Pattern Anal. Mach. Intell.* **30**(2), 283–298 (2008)
49. Waymo: Safety report (2018). <https://waymo.com/safety/>. Accessed 9 Nov 2019
50. Waymo: Waymo Open Dataset: An autonomous driving dataset (2019). <https://waymo.com/open/>
51. Zeng, W., et al.: End-to-end interpretable neural motion planner. In: IEEE Conference on Computer Vision and Pattern Recognition (2019)
52. Zhao, S., Song, J., Ermon, S.: InfoVAE: balancing learning and inference in variational autoencoders. In: Proceedings of the AAAI Conference on Artificial Intelligence (2019)
53. Zhao, T., et al.: Multi-agent tensor fusion for contextual trajectory prediction. In: IEEE Conference on Computer Vision and Pattern Recognition (2019)