# ScePT: Scene-consistent, Policy-based Trajectory Predictions for Planning

Yuxiao Chen[1]      Boris Ivanovic[1]      Marco Pavone[1,2]

[1]NVIDIA Research      [2]Stanford University

{yuxiaoc, bivanovic, mpavone}@nvidia.com, pavone@stanford.edu

## Abstract

*Trajectory prediction is a critical functionality of autonomous systems that share environments with uncontrolled agents, one prominent example being self-driving vehicles. Currently, most prediction methods do not enforce scene consistency, i.e., there are a substantial amount of self-collisions between predicted trajectories of different agents in the scene. Moreover, many approaches generate individual trajectory predictions per agent instead of joint trajectory predictions of the whole scene, which makes downstream planning difficult. In this work, we present ScePT, a policy planning-based trajectory prediction model that generates accurate, scene-consistent trajectory predictions suitable for autonomous system motion planning. It explicitly enforces scene consistency and learns an agent interaction policy that can be used for conditional prediction. Experiments on multiple real-world pedestrians and autonomous vehicle datasets show that ScePT matches current state-of-the-art prediction accuracy with significantly improved scene consistency. We also demonstrate ScePT's ability to work with a downstream contingency planner.*

## 1. Introduction

Predicting the future motion of uncontrolled agents is critical to the safety of autonomous systems that interact with them. A prominent example is self-driving cars, where the ego-vehicle shares the road with other road users such as vehicles, pedestrians, and cyclists. The prediction task is difficult as humans are notoriously uncertain and inconsistent. For example, it is well-known that humans demonstrate multimodal behaviors in the context of driving, such as being simultaneously able to maintain their current lane, change lanes, yield, or overtake in the future. As a result, early works on human driving behavior prediction [37] were not accurate enough to be used in an autonomous vehicle's motion planning stack. To remedy this, many researchers have been developing phenomenological methods, i.e., methods that learn the behavior of agents from a wealth of data (e.g., [2, 28, 29, 44]), to great effect.

In a typical autonomy stack, the trajectory prediction module is followed by a planning module which takes the
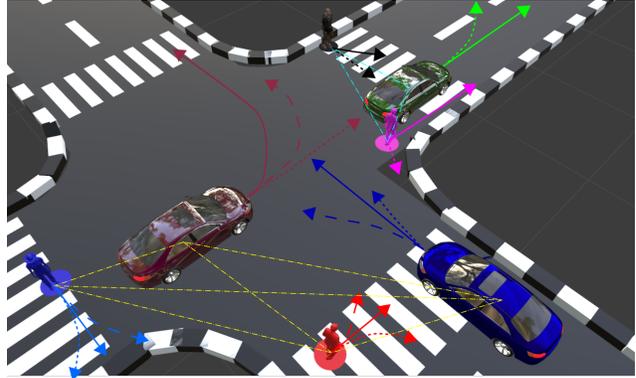


Figure 1. An illustration of ScePT's output, comprised of multimodal trajectory predictions for each agent. Different stroke type (solid, dashed, dotted) represents the different modes of the scene-consistent joint trajectory prediction. Agents in a scene are partitioned into highly-interacting cliques, an example of which is visualized with yellow dashed lines.

predicted trajectories of surrounding agents and plans the ego-motion accordingly. With this downstream planner in mind, several requirements, in addition to prediction accuracy, emerge, and are discussed in detail below.

### 1.1. Desiderata

Typical features desired for a trajectory prediction model include high prediction accuracy, fast inference speed, and calibrated uncertainties. When predictions are subsequently consumed by a downstream planner, the following features are also critical for overall system performance:

**Compatibility:** Trajectory predictions for different agents in a scene should be compatible with each other within a single joint prediction. In particular, collisions among predicted trajectories should be rare, as collisions are themselves rare in reality.

**Tractable Joint Trajectory Prediction:** As mentioned previously, the future motion of the agents can be multimodal. In a scene consisting of multiple agents, if the multimodal predictions are generated for individual nodes, a downstream motion planner needs to consider all combinations of these trajectory predictions. Since the number of modes grows exponentially with the number of agents, the planner

is quickly overwhelmed. Alternatively, the motion planner could take a conservative approach and avoid all predicted trajectories, yet often at the price of compromised planning performance (e.g., bringing the robot to a standstill if all plans seem to collide). As a result, we desire multimodal, joint predictions of all agents with a limited, but fully-representative number of modes such that a downstream planner can perform contingency planning.

**Time Consistency:** With a downstream planner, the motion plan heavily depends on the prediction results. To ensure smooth motion plans, predictions should not change significantly between subsequent time steps if the scene itself has not changed drastically in the meantime. As a result, sampling should be avoided as predictions may change significantly between time steps, causing discontinuity in the resulting motion plans which may hurt planning performance and safety.

**Conditioning:** Conditioning is the operation of fixing one or multiple agents' future trajectories and predicting the resulting distribution of other agents' future trajectories. Conditional prediction is useful for motion planning (conditioned on the ego agent's motion plan) and for understanding interactions between agents. Conditioning is available in several existing works such as [42], yet requires explicit modeling. Ideally, conditional distributions would be generated without requiring structural changes to the model.

## 1.2. Contributions

In this work we present ScePT, a trajectory forecasting method that generates joint trajectory predictions for multiple interacting agents. Our contributions are threefold: First, we propose predicting the futures of *cliques* of agents rather than individuals or the scene graph as a whole (Sec. 3.1), and present a neural network architecture to do so (Sec. 3.2). Second, we leverage insights from motion planning and propose a policy network that autoregressively rolls out closed-loop trajectory predictions via a GNN that models agent-to-agent interactions and maps them to control inputs (Sec. 3.3). Finally, we improve output sample diversity by augmenting our loss function with a tunable risk measure that determines weights between trajectory samples during training (Sec. 3.6).

When evaluated on large-scale, real-world pedestrian and driving datasets, ScePT reduces the dimensionality necessary to capture scene-level multimodality (Sec. 4.1); achieves significant improvements in the scene consistency of its predictions, as measured by collision rate (Sec. 4.2); and easily enables counterfactual analyses (Sec. 4.3); all of which are critical for simulation (Sec. 4.3), downstream planning (Sec. 4.4), and verification of autonomous vehicle performance.

## 2. Related Work

Early trajectory prediction works were predominantly ontological, positing structure about an agent's decision-making process, exemplified by social force models [22, 33], hidden Markov models [27], and the intelligent driver model [45]. However, limited by their expressibility, these models cannot scale to complicated scenarios despite extensive tuning. To remedy this, many researchers have been developing phenomenological methods, i.e., methods that focus on learning the behavior of agents from a wealth of data. Several notable works include the Social LSTM [2], GAIL [28], MFP [44], and DESIRE [29]. Since the trajectory prediction problem is intrinsically a sequence-to-sequence modeling task, recent works commonly apply Recurrent Neural Networks (RNNs) [10, 21, 29, 40, 42, 44] and Transformers [31, 34, 50], achieving strong results.

Another core facet of trajectory prediction is accounting for the interactions between agents and the scene geometry. Two common choices to model agent interactions are Graph Neural Networks (GNNs) and Convolutional Neural Networks (CNNs). GNN-based methods [10, 42, 48, 52] construct scene graphs with agents as nodes and their interaction as edges, performing message passing to congregate information. CNN-based methods [4, 16, 19, 26, 34, 36, 51] typically rasterize scene information into layers of images such as Birds Eye View (BEV) images and velocity images to encode information. In general, CNN-based methods have a fixed computation complexity (usually being faster than GNN-based methods) and preserve geometric information about the scene. GNN-based methods can be viewed as a "sparse" representation of the scene and allow for more complicated features, yet their computation complexity scales at least linearly with the number of agents, and geometric information may be lost without the use of special structures [9].

Once scene information is gathered and encoded, generative models such as GANs [21, 40] and CVAEs [18, 29, 42] are typically used to produce multimodal trajectory predictions. Of these, CVAE models are the most common choice due to their performance and ease of training. CVAEs with continuous latent spaces [10, 48, 52] enjoy stronger expressivity, but require sampling to obtain predictions, removing any time consistency in sequential outputs. On the other hand, a discrete latent space [42, 44] does not require sampling, but is less expressive and more likely to suffer from mode collapse.

Since trajectory prediction usually concerns multiple agents in a scene, the issue of scene consistency emerges, i.e., predictions of different agents should not collide with each other or with static obstacles. For models that pursue only high prediction accuracy, especially agent-centric models, scene consistency is usually poor. To remedy this, [17, 21] use pooling to model interaction between agents in the encoder, [52] introduces a group-level encoder with clustered agents sharing a group mode, and [10] learns a scene interaction model via message passing. Comparing to the encoder, coupled decoding is more difficult as the future trajectories are not known. [32] uses fictitious play to

roll out predictions and gradually improve prediction quality, [44] performs autoregressive decoding, building subsequent predictions on prior steps of prediction. [9] uses a message-passing procedure to search for most likely joint trajectories, yet only predicts a single mode. To the best of our knowledge, most existing methods do not model the agents as policy planners with observations, cost functions, and actuation inputs, all while explicitly enforcing scene consistency.

For autonomous vehicles which involve a downstream planner, it is also important to consider multimodality, i.e., the possibility for multiple distinct futures. However, the number of modes cannot be too large so as not to overwhelm the downstream planner. [12] generates a set of motion primitives that are deemed possible and constrain the planner to avoid them. [24] represents multimodal predictions as mixtures of linear dynamics (as opposed to tracklets), simplifying incorporation in downstream planning. To achieve a wide coverage of modes with a limited computational budget, several diversity sampling techniques have been proposed. [49] proposes using determinantal point processes for diverse latent sampling, [23] uses the farthest point sampling algorithm, and [15] apply a reparameterization trick with coefficients given by a GNN. Overall, most existing diversity sampling techniques are designed for continuous latent spaces. In this work, we introduce a new risk-based loss modification that yields output diversity for a discrete latent space with only a small number of samples.

## 3. ScePT

ScePT[1] is a discrete CVAE model that outputs joint trajectory predictions for multiple agents in a scene, ensuring high scene consistency in its predictions by reasoning about each agent's motion policy and the influence of their neighbors.

**Nomenclature.** Throughout the paper, we use the terms node and agent interchangeably, which may be a vehicle, a pedestrian, a cyclist, or other kinds of road users. We use $s$ to denote an agent's state and $e$ an edge between two nodes. Since our model is a CVAE, we follow standard terminology in the CVAE literature, i.e., $x$ denotes the conditioning variable, $y$ the observed variable, and $z$ the hidden latent variable. We use bold font to denote variables associated with a group of nodes, e.g., a clique. For example, for a clique consisting of nodes 1 through $N$, $z_1, ...z_N$ are the latent variables of each of the nodes and $\mathbf{z} = [z_1, ...z_N]$ is the latent variable of the clique.

### 3.1. Preprocessing

To maintain scene consistency, ScePT is a scene-centric model, i.e., its output predictions are the joint trajectories of multiple nodes in a scene. Given a scene with multiple nodes, a spatiotemporal scene graph is generated where

nodes represent agents and edges represent their interactions. We use agents' closest future distance as a proxy for interaction, propagating forward each node according to a constant velocity model $\Phi_{a_i=0}^{0:T}(s_i)$, where $\Phi_{a_i=0}^{t}$ is the flow operator that maps the initial state to the future states $t$ time steps ahead with the action $a_i$ of node $i$ set to 0 and $T$ is the prediction horizon. The closest future distance between two agents is defined as

$$d_{ij} = \min_{t \in [0,T]} \text{Dis}(\Phi_{a_i=0}^{t}(s_i), \Phi_{a_i=0}^{t}(s_j)), \quad (1)$$

where Dis is the Euclidean distance between the two agents. We then define the scene graph adjacency matrix as

$$A_{ij} = \left\{ \begin{array}{ll} 0, & d_{ij} > d_0(\eta_i, \eta_j) \\ \frac{d_0(\eta_i, \eta_j)}{d_{ij}} & d_{ij} \leq d_0(\eta_i, \eta_j) \end{array} \right. ,$$

where $\eta_i$ is the type of node $i$ (e.g., vehicle, cyclist, pedestrian) and $d_0$ is a distance threshold which is fixed for each edge type.

With the scene graph determined by the adjacency matrix, in contrast to models that keep all nodes in a single graph [10,44], we partition the scene graph into cliques with a maximum size (fixed as a parameter). We do this to reduce the dimensionality of the product latent space, which scales exponentially with the size of the graph, causing prediction accuracy to deteriorate if too large (see Sec. 4.5 for further discussion). While weighted graph partitioning is NP-hard, there are many off-the-shelf algorithms, and we use the well-known Louvain algorithm due to its strong performance [6]. After partitioning, every pair of nodes within a clique is connected (despite the distance threshold) to form a clique. Node histories are then collected and fed to ScePT. Node states and dynamics are explained in detail in Appendix A.2. When available, we also utilize the map information and the relative position to the closest lane.

### 3.2. Encoder

With cliques in hand, agents' state and edge (relative states between agents) histories are encoded into feature vectors via LSTMs. Instead of associating each node with a latent variable distribution that is independent of its neighbors, our encoder models the joint latent distribution. Specifically, each agent is equipped with a discrete latent variable $z_i$ with cardinality $N$, making the joint latent variable of the clique simply $\mathbf{z} = [z_1, z_2, ...z_n]$. This means that the cardinality of the joint latent space grows exponentially with the number of nodes in the clique, and is the reason why we limit clique size.

ScePT represents the distribution of the joint latent variable as a Gibbs distribution consisting of node factors and edge factors,

$$\log \mathbb{P}(\mathbf{z}) \propto \sum_i f_i(x_i, z_i) + \sum_{e_{ij} \in \mathcal{E}} f_{ij}(x_i, x_j, z_i, z_j), \quad (2)$$

---

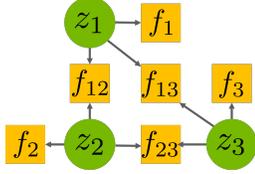[1]Code available at https://github.com/nvr-avg/ScePT

3

Figure 2. Factor graph with individual agent latent variables as variable nodes and factor nodes which are functions of the connected variable nodes. Factor nodes are comprised of individual agent and agent-agent interaction factors, e.g., $f_1$ is a function of $z_1$ while $f_{12}$ is a function of $z_1$ and $z_2$. All factor nodes are summed to obtain the log likelihood.

where $x_i$ is the state history of node $i$ and $f_i$ is the node factor of node $i$, which is a feedforward neural network mapping $x_i$ and $z_i$ to a real number. $f_{ij}$ is the edge factor of the node pair $i, j$, also a feedforward network, and $\mathcal{E}$ is the set of edges. The log-likelihood can be computed by constructing a factor graph [1], which is a bipartite graph with variable nodes and factor nodes. An example factor graph is shown in Fig. 2. Normalization is done by summing up all possible valuations of $\mathbf{z}$ (since $\mathcal{Z}$ is discrete).

While the joint latent space's cardinality scales exponentially with clique size, we found that probability mass typically concentrates on only a few ($< 10$) modes.

### 3.3. Decoder

Our decoder design is inspired by the motion planning process, i.e., we view each agent as a motion planner and emulate their planning process to output trajectory predictions. A typical motion planner takes a reference trajectory, i.e., the desired motion, and adjusts it to satisfy constraints (e.g., collision avoidance) and minimize a specified cost function. Inspired by this process, the structure of our policy net is visualized in Fig. 3.

The inputs to the policy net are the current states of the clique nodes, their reference trajectories $\mathbf{s}_{\text{des}}$, and the clique latent $\mathbf{z}$. Reference trajectories are generated via Gated Recurrent Unit (GRU) networks that take the state history encoding, map encoding, and latent variable $z$ as input. The current node states are then compared with the reference trajectories to obtain the tracking error $\Delta\mathbf{s}$ and the next waypoint in the local coordinate frame $\Delta\mathbf{s}^+$.

To model an edge, we pair its two node states together and feed the state pair into a pre-encoding network (fully connected) and then an LSTM cell. For each node, depending on the graph structure, there may be a varying number of neighbors. To encode a variable number of neighbors, all of a node's edges are condensed into a single observation encoding via an attention network [14]. The observation encoding, latent variable, and tracking error are then concatenated and fed to a fully connected action network to obtain the node's control action prediction $a$. Here, we assume that the node's dynamics are differentiable functions of the state and control input, which is true for common agent types
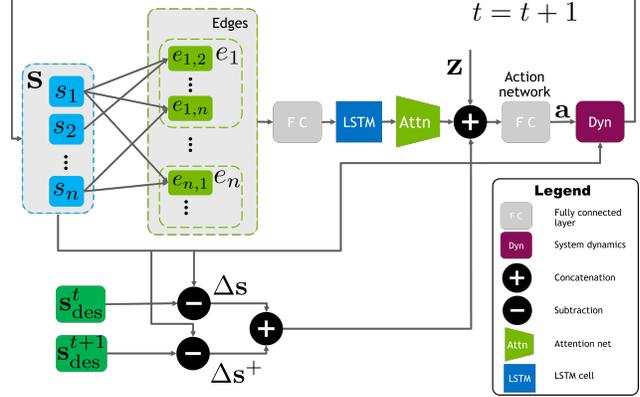


Figure 3. Our autoregressive policy network architecture. For each node, neighboring node states are pooled with an attention mechanism, the resulting encoding then generates control inputs together with the reference trajectory. The control inputs pass through the agent dynamics to produce position predictions and the process repeats for subsequent timesteps.
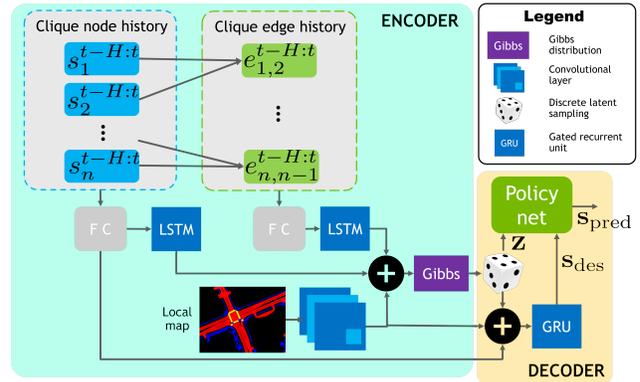


Figure 4. Overview of ScePT: Node history and map information are collected for all nodes within a clique and passed through a Gibbs distribution to generate the discrete joint latent distribution. The policy net (decoder) then generates closed-loop trajectory predictions given latent samples.

such as vehicles (e.g., Dubin's car model [20]) and pedestrians (single or double integrators). The state prediction is then fed back to the state vector and this process repeats.

The overall structure of ScePT is shown in Fig. 4. The encoder takes the LSTM-encoded state and edge history as well as the CNN-encoded local map, and generates a discrete Gibbs distribution over the clique latent variable. The latent variable, together with the state history and map encodings, is used to generate the desired trajectory for each node via GRUs. The desired trajectories and latent variable are then passed to the policy net to obtain closed-loop trajectory predictions.

### 3.4. Conditioning via policy learning

As mentioned in Sec. 1.1, conditioned prediction is an important capability. Conditioning was performed in previous works [25, 42] by explicitly encoding the ego future

4

trajectory in the encoder. However, assuming that only one agent can be conditioned on makes use cases such as driving simulation difficult, as one would need to train explicit conditioning models for each pair of agents. By comparison, PRECOG [38] simply needs to set the latent variable of the robot to produce future-conditional predictions. Similarly, ScePT does not require any structural change to produce conditional predictions since it learns agents' interaction policies. Conditional predictions are generated by simply fixing the trajectory roll-outs of conditioned agents and outputting the trajectory predictions of the rest of the agents in the clique. Since a fixed future trajectory does not fall into any latent mode, we remove any factors concerning conditioned nodes from the Gibbs distribution factor graph.

### 3.5. Training

Following standard CVAE training [43], our objective is the Evidence Lower Bound (ELBO) loss:

$$
\begin{aligned}
ELBO = \ &\mathbb{E}_{\mathbf{z}\sim Q(\mathbf{z}|\mathbf{x},\mathbf{y})}[\log(P(\mathbf{y}|\mathbf{x},\mathbf{z}))] \\
&- \beta D_{KL}(Q(\mathbf{z}|\mathbf{x},\mathbf{y})||P(\mathbf{z}|\mathbf{x})),
\end{aligned}
\tag{3}
$$

where $\mathbf{z}$ is the clique latent variable, $\mathbf{y}$ is the future trajectories of all nodes, and $\mathbf{x}$ is the conditional variable, consisting of node and edge history, map encoding, and lane information for all nodes in the clique. For the likelihood cost, we assume Gaussian noise around the predicted trajectory for each mode, resulting in a 2-norm loss,

$$
\begin{aligned}
&\mathbb{E}_{\mathbf{z}\sim Q(\mathbf{z}|\mathbf{x},\mathbf{y})}[\log(P(\mathbf{y}|\mathbf{x},\mathbf{z}))] \\
&= \sum_{\mathbf{z}\in\mathcal{Z}} Q(\mathbf{z}|\mathbf{x},\mathbf{y})||f_{\mathbf{y}}(\mathbf{x},\mathbf{z}) - \mathbf{y}_{GT}(\mathbf{x})||^2,
\end{aligned}
\tag{4}
$$

where $f_{\mathbf{y}}(\mathbf{x},\mathbf{z})$ is the trajectory prediction from the decoder and $\mathbf{y}_{GT}(\mathbf{x})$ is the ground truth.

We also add a collision penalty, specified in detail in Appendix A.3, as a regularization term to penalize incompatible predictions, the influence of which will be further discussed in Sec. 4.5. Other types of regularization, e.g., ride comfort, can also be added since the node dynamics are explicitly included in the policy net.

**Sampling the Latent Space.** While our discrete latent space is enumerable, the cardinality of $\mathcal{Z}$ grows exponentially with the clique size. Thus, it is sometimes not tractable to decode all modes. To remedy this, we apply diversity sampling. Specifically, we take the $N_g$ highest probability modes and randomly sample $N_r$ modes from the rest. When the total cardinality of $\mathcal{Z}$ is less than $N_g + N_r$, all modes are selected. The sample probabilities are then normalized so that the expected loss does not collapse to 0.

### 3.6. Mode Collapse and Diverse Sampling

Discrete CVAEs for trajectory prediction are prone to mode collapse, i.e., the decoder tends to predict similar trajectories under different modes since the likelihood cost is a

weighted sum of 2-norm errors and the average prediction is likely to be a local minimum. Mode collapse has been discussed in previous works and tackled by methods such as Multiple-Trajectory Prediction (MTP) loss [16], using prior knowledge [11, 19], and assigning modes by classifying the ground truth into categories [31]. Our approach maintains the expected loss function, but introduces CVaR as a new way to avoid mode collapse.

**Conditional Value at Risk** (CVaR) [39] is a risk measure commonly used in finance and optimization, defined as

$$
\begin{aligned}
\text{CVaR}_{1-\alpha}(X) &= \inf_{\eta\in\mathbb{R}}\{\eta + \frac{1}{\alpha}\int_{-\infty}^{\infty}[x-\eta]_+ P(x)\} \\
&= \min_{0\leq P'(x)\leq\frac{1}{\alpha}P(x),\int P'(x)dx=1} \mathbb{E}_{P'}[X],
\end{aligned}
\tag{5}
$$

where $P$ is the probability distribution of $X$ and $\alpha$ tunes the level of risk-averseness. CVaR is the mean of the lowest $\alpha$-percentile values of $x$ under $P$. At the limits of $\alpha$, $\alpha \to 0$ yields the essential infimum of $X$ and $\alpha = 1$ yields $\mathbb{E}[X]$.

The second line in (5) is the dual form of CVaR, which can be understood as shifting the distribution $P$ to $P'$ under the constraint that $P'$ has to be a proper distribution and for all $x$, $P'(x) \leq \frac{1}{\alpha}P(x)$. Inspired by the dual form, we modify the expectation loss in (4) to

$$
\min_{0\leq Q'(z)\leq\frac{1}{\alpha}Q(z|x,y),\sum Q'(z)=1} \mathbb{E}_{z\sim Q'}[||f_y(x,z)-y_{GT}(x)||^2],
\tag{6}
$$

which is the best $\alpha$-percentile loss value among the discrete modes. This CVaR loss does not force all modes to match the ground truth, only those that are already close, directly preventing mode collapse. Compared to common usages of risk measures which typically focus on the worst outcomes, we use CVaR to focus on the best predictions to maintain output diversity. During training, $\alpha$ is used to trade-off the model's focus on encoder accuracy vs diversity, see Appendix A.4 for details. In addition to incorporating CVaR, we also use a greedy algorithm to diversely sample the product latent space, see Appendix A.5 for further details.

## 4. Experiments

We evaluate the performance of ScePT on the tasks of pedestrian and vehicle motion prediction. In particular, we make use of the well-known ETH [35], UCY [30], and nuScenes [7] datasets.

**Metrics.** We use the common Average and Final Displacement Error (ADE/FDE) metrics to measure the quality of trajectory predictions. Since our outputs are multimodal, we adopt the Best-of-N (BoN) extension and take the $N$ highest probability modes from the encoder to compute BoN ADE/FDE values. The sampling process in ScePT is different from prior works with continuous latent spaces as we do not randomly sample from the latent distribution, but rather pick the $N$ modes deemed most likely by the encoder. As a result, when the number of samples is larger than the
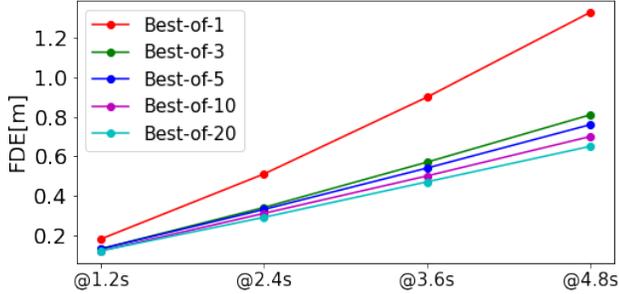
Figure 5. FDE of the ETH dataset under varying numbers of samples, the most significant improvement is seen at 3 samples.

number of possible clique modes (i.e., $N > |\mathcal{Z}|$), we take only $|\mathcal{Z}|$ samples.

## 4.1. Pedestrian Motion Prediction

The ETH (containing ETH and Hotel scenes) [35] and UCY (containing Univ, Zara1 and Zara2 scenes) [30] datasets are widely-used benchmarks for pedestrian motion prediction. Together, they contain 9,514 unique pedestrians in many challenging, interactive real-world scenarios.

For all pedestrian datasets, the maximum clique size is 5 and each node's latent space cardinality is 6. The ADE/FDE results are shown in Tables 1 (deterministic) and 2 (multi-modal). While ScePT is designed mainly for autonomous driving, it performs remarkably well on pedestrian datasets, achieving the best or second-best performance among state-of-the-art models in the field. In particular, ScePT outperforms prior methods on most datasets in ADE, yet performs slightly worse in UCY scenes in terms of FDE. The likely reason is that the UCY datasets are much denser than ETH, forcing ScePT to partition the large scene graph into small cliques. Once partitioned, interactions between cliques are ignored, hurting prediction accuracy. Using a larger maximum clique size, however, would cause the joint latent cardinality to be too large and further deteriorate performance. Even in these cases, ScePT still performs on par with state-of-the-art prediction methods.

Fig. 5 shows the FDE of the ETH dataset under different numbers of samples, Due to our use of CVaR in the loss function, ScePT is able to generate diverse modes. After only 3 samples, our method's predictions are already very accurate. We observe this phenomenon across all datasets, and find that sampling 3 to 5 modes achieves a good balance of prediction quality and runtime complexity.

## 4.2. Vehicle Motion Prediction

The nuScenes dataset [8] consists of 1000 driving scenes, each 20 seconds long and containing up to 23 object classes. To match the nuScenes prediction challenge set, only vehicles and pedestrians are predicted during training and evaluation. Tab. 3 summarizes the prediction accuracy of our method alongside a set of state-of-the-art approaches. Despite ScePT discarding edges between
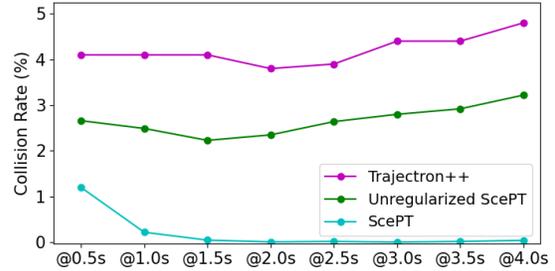


Figure 6. ScePT achieves a better collision rate than Trajectron++ [42] without collision cost regularization. With collision cost regularization, the collision rate becomes virtually zero.

cliques, the prediction accuracy is near the state-of-the-art, especially in later timesteps. Furthermore, thanks to our diversity-promoting design, prediction accuracy can be significantly improved by adding only 1 or 2 extra modes.

As mentioned in the introduction, scene consistency in trajectory predictions is critical for planning, even more so when simulating and verifying the performance of autonomous vehicles. Fig. 6 compares the collision rate achieved by Trajectron++ [42] and ScePT at different prediction horizons. The collision rate is averaged over all vehicles whose trajectories are predicted by the two models, including vehicle-to-vehicle and vehicle-to-pedestrian collisions. The result shows that ScePT's predictions have much less collisions than Trajectron++, even without collision cost regularization; becoming virtually zero with regularization. Moreover, the collision rate tends to decrease with the horizon, implying that collisions at earlier timesteps may be due to poor initial conditions and that the learned policy net is able to resolve conflicts in later timesteps.

## 4.3. Conditioning and Counterfactual Analyses

Conditioning is an important capability, enabling a downstream planner to obtain trajectory predictions conditioned on ego motion. It is also useful for performing counterfactual, "what if?" analyses. Fig. 7 demonstrates ScePT's ability to perform conditional prediction, where the left and right figures show unconditioned and conditioned predictions, respectively. In the top right, we condition on vehicle A and C braking at $-4m/s^2$ to a full stop. Accordingly, our method predicts that (1) vehicle B will brake to avoid a collision, and (2) vehicle D will perform a lane change to avoid vehicle C. In the bottom right, we condition on vehicle C accelerating at $4m/s^2$. This causes a chain reaction, making vehicle B's prediction swerve to avoid a collision, which then subsequently affects vehicle A and the two pedestrians on the bottom left, showing that nodes within a clique are influenced by all other nodes in that clique.

Interestingly, the policy learned from the nuScenes dataset is quite robust to clique size. Our model is currently trained with a maximum clique size of 4, but we evaluated it with sizes as large as 8. Even in such cases, the model

6

| Dataset | S-LSTM [2] | S-ATTN [46] | Trajectron++ [42] | ScePT |
|---------|-----------|-------------|-------------------|-------|
| ETH     | 1.09/2.35 | *0.39*/3.74 | 0.71/*1.66*       | **0.19/1.33** |
| Hotel   | 0.79/1.76 | 0.29/2.64   | *0.22*/**0.46**   | **0.18**/*1.12* |
| Univ    | 0.67/1.40 | *0.33*/3.92 | 0.44/**1.17**     | **0.19**/*1.19* |
| Zara1   | 0.47/1.00 | *0.20*/**0.52** | 0.30/*0.79*   | **0.18**/1.10 |
| Zara2   | 0.56/*1.17* | 0.30/2.13 | 0.23/**0.59**     | **0.19**/1.20 |
| Average | 0.71/1.54 | 0.30/2.59   | *0.38*/**0.93**   | **0.19**/*1.19* |

Table 1. ADE/FDE in meters, using the most likely mode. **Bold**/*italic* font indicates the best/second-best value. Lower is better.

| Dataset | S-GAN [21] | SoPhie [40] | MATF [51] | Trajectron++ [42] | ScePT |
|---------|-----------|-------------|-----------|-------------------|-------|
| ETH     | 0.81/1.52 | 0.70/1.43   | 1.01/1.75 | *0.39/0.83*       | **0.10/0.65** |
| Hotel   | 0.72/1.61 | 0.76/1.67   | 0.43/0.80 | **0.12/0.21**     | *0.13/0.77* |
| Univ    | 0.60/1.26 | 0.54/1.24   | 0.44/0.91 | *0.20*/**0.44**   | **0.12**/*0.65* |
| Zara1   | 0.34/0.69 | 0.30/0.63   | 0.26/*0.45* | *0.15*/**0.33** | **0.13**/0.77 |
| Zara2   | 0.42/0.84 | 0.38/0.78   | 0.26/*0.57* | **0.11/0.25**   | *0.14*/0.81 |
| Average | 0.58/1.18 | 0.54/1.15   | 0.48/0.90 | *0.19*/**0.41**   | **0.12**/*0.73* |

Table 2. Best-of-20 ADE/FDE in meters. **Bold**/*italic* font indicates the best/second-best value. Lower is better.

| Method | @1s | @2s | @3s | @4s |
|--------|-----|-----|-----|-----|
| S-LSTM [2, 9]       | 0.47 | -    | 1.61 | -    |
| CSP [9, 18]         | 0.46 | -    | 1.50 | -    |
| CAR-Net [9, 41]     | 0.38 | -    | 1.35 | -    |
| SpAGNN [9]          | 0.35 | -    | 1.23 | -    |
| Trajectron++ [42]   | 0.07 | 0.45 | 1.14 | 2.20 |
| Ours                | 0.44 | 0.93 | 1.63 | 2.58 |
| Ours (Best-of-2)    | 0.41 | 0.83 | 1.44 | 2.20 |
| Ours (Best-of-3)    | 0.40 | 0.80 | 1.36 | 2.14 |

Table 3. ScePT's prediction accuracy on nuScenes vehicles nearly matches state-of-the-art methods. Furthermore, thanks to its diversity-promoting design, ScePT's prediction accuracy can be drastically improved by adding only 1 or 2 extra modes.
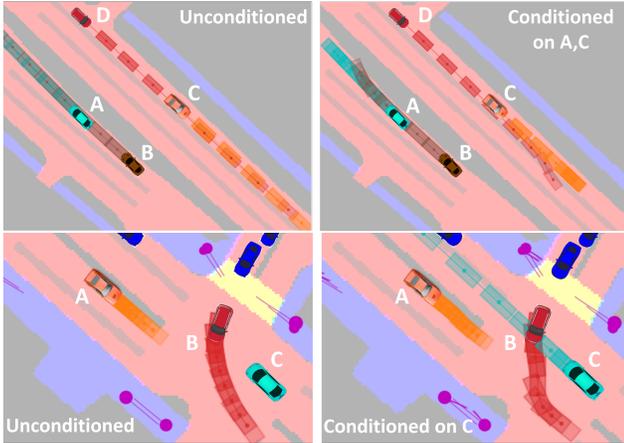


Figure 7. By design, ScePT can produce predictions that are conditioned on any number of agents. **Left:** Original, unconditioned predictions. **Top Right:** Predictions conditioned on vehicles A and C braking. **Bottom Right:** Predictions conditioned on vehicle C accelerating.

is still able to produce reasonable conditioned predictions, verifying the efficacy of the policy net's attention network.

## 4.4. Integration with a Downstream Planner

To demonstrate ScePT's performance when integrated with a downstream planner, we feed its predictions to a downstream model predictive control (MPC)-based planner. The MPC planner takes the multimodal trajectory predictions into account and performs contingency planning via branching [3, 13]. Given $M$ joint trajectory predictions, the MPC plans $M$ corresponding ego trajectories with the additional constraint that the first control inputs for all $M$ ego trajectories must be the same. Formally,

$$
\min_{\{a_{0:T-1}^i, s_{0:T}^i\}_{i=1}^M} \sum_i \pi_i \mathcal{J}(a_{0:T-1}^i, s_{1:T}^i)
$$
$$
\text{s.t. } s_{t+1}^i = \text{Dyn}(s_t^i, a_t^i), \ \forall i, t, a_t^i \in \mathcal{A}, \ x_t^i \in \mathcal{S},
$$
$$
\mathcal{C}(s_{1:T}^i, y_{1:T}^i) \leq 0, i = 1..., M,
$$
$$
s_0^i = s_0, \ a_0^1 = a_0^2 = ...a_0^M,
$$
$$
(7)
$$

where $\pi_i$ is the probability of prediction mode $i$, $s^i$ and $u^i$ are the planned state and input sequences of the ego vehicle under the $i$-th mode, $\mathcal{J}$ is the cost function, and $\mathcal{C}$ is a constraint (e.g. collision avoidance). Eq. (7) is a nonlinear optimization problem and is solved with IPOPT [47]. As an example of runtime, when $M = 3$, our unoptimized PyTorch prediction code executes in less than 240ms and MPC planning takes less than 60 ms, all on a CPU. Fig. 8 shows the result of combining ScePT's predictions with this downstream MPC planner, visualizing prediction modes and their resulting ego motion plans.

## 4.5. Ablation study

**Collision Cost Regularization and Accuracy.** As we have seen in Sec. 4.2, the collision rate of ScePT without the collision penalty already outperforms prior works. When collision penalty is added, the collision rate of ScePT drops
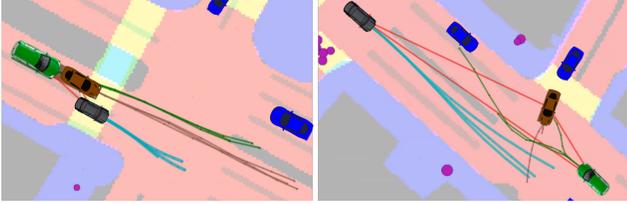
Figure 8. The integration of prediction and planning. **Black vehicle**: ego vehicle; blue vehicles: adjacent vehicles outside the ego clique; cyan trajectory: planned trajectory (3 modes); green and brown vehicles: adjacent vehicles within the ego clique; green and brown trajectories: predicted trajectories (top 3 modes); magenta circle: pedestrians; red lines: connected nodes within the ego's clique.

|  | @1s | @2s | @3s | @4s |
|---|---|---|---|---|
| With collision penalty | 0.44 | 0.93 | 1.63 | 2.58 |
| Without collision penalty | 0.40 | 0.92 | 1.70 | 2.71 |

Table 4. Influence of including collision penalty regularization on nuScenes vehicle prediction accuracy.

| Loss Ablation | @1s | @2s | @3s | @4s |
|---|---|---|---|---|
| With CVaR | 0.44 | 0.93 | 1.63 | 2.58 |
| Without CVaR | 0.45 | 0.97 | 1.74 | 2.79 |

Table 5. Influence of CVaR loss on nuScenes prediction FDE.

to virtually zero. Tab. 4 summarizes the effect of including the collision penalty on prediction accuracy. We can see that prediction accuracies are either the same or better after the collision penalty is added, indicating that avoiding collisions also produces more accurate outputs.

We also assess conditional predictions without the collision penalty, and collisions are much more likely when some nodes are assigned errant behavior. For example, the situation in Fig. 7 (top) results in vehicles B and D hitting their lead vehicles when predicted by a model trained without collision penalty regularization. This phenomenon implies that, without collision penalty regularization, ScePT is not able to maintain scene consistency in out-of-distribution scenarios. Overall, this regularization term forces the model to maintain scene consistency, making it more robust to out-of-distribution scenarios.

**Effect of $\alpha$.** To study the utility of using CVaR in Eq. (4), we conduct an ablation study over different values of $\alpha$. Our baseline model uses a varying $\alpha$ (initialized at 0.2 and gradually increasing to 1.0 during training). We also train a version with constant $\alpha = 1.0$, in which case the CVaR loss function is equivalent to the original expectation. Tab. 5 summarizes the results, showing that using CVaR outperforms expectation in the original loss function.

**Effect of Clique Formation.** To assess the sensitivity of ScePT to the clique forming process, we change the distance criteria described in Sec. 3.1 to the Euclidean distance between nodes at the current time. Tab. 6 summarizes the

| Distance Ablation | @1s | @2s | @3s | @4s |
|---|---|---|---|---|
| Flow | 0.44 | 0.93 | 1.63 | 2.58 |
| $t = 0$ Euclidean | 0.43 | 0.95 | 1.70 | 2.71 |

Table 6. Influence of distance criteria for clique forming on nuScenes prediction accuracy (FDE).

| Clique Size Ablation | @1s | @2s | @3s | @4s |
|---|---|---|---|---|
| Max clique size 2 | 0.41 | 0.90 | 1.63 | 2.68 |
| Max clique size 4 | 0.44 | 0.93 | 1.63 | 2.58 |
| Max clique size 6 | 0.64 | 1.27 | 1.90 | 2.90 |

Table 7. Effect of maximum clique size on nuScenes prediction accuracy (FDE).

results, which shows that the flow distance partition leads to better performance, yet using a simple Euclidean distance leads to decent performance.

Another important parameter is the maximum clique size. Since a node only considers neighbors within its clique, a small clique size leads to the nodes overlooking some nearby neighbors that were not partitioned into the same clique. On the other hand, if a clique gets too big, the cardinality of the clique's product latent space becomes too large, causing problems when sampling greedily. After experimenting, we found that a maximum clique size of 4 leads to the best result. Tab. 7 summarizes this ablation.

# 5. Conclusion and Future Work

**Summary.** This paper presents ScePT, a CVAE-based model that generates multimodal joint trajectory predictions with high scene consistency. The encoder uses a Gibbs distribution to capture interactions between agents and outputs prediction mode probabilities for a whole clique instead of individual nodes. The decoder learns agent interaction policies to generate closed-loop trajectory predictions with high scene consistency, thanks to an explicit collision penalty used as regularization during training. Experiments on the ETH, UCY, and nuScenes datasets show that ScePT is able to achieve state-of-the-art prediction accuracy with significantly improved scene consistency. We also demonstrate its capability to integrate with a downstream contingency planner and to generate human-like behaviors via conditioning.

**Limitations and Future Work.** One important limitation (and area of future work) of ScePT is its loss of sparsity. Since the decoder generates predictions concurrently for all nodes in a clique, one cannot utilize the sparsity of the interaction graph to only consider neighbors for each node. This causes the exponential cardinality of the joint latent space and forces the size of cliques to be limited. Subsequently, interactions between nodes in different cliques are ignored, leading to compromises in accuracy especially when there are large numbers of crowded agents, such as in the UCY dataset. Another limitation is the computation time due to the autoregressive policy net, which could be further improved by more efficient coding and parallelization.

# References

[1] Pieter Abbeel, Daphne Koller, and Andrew Y Ng. Learning factor graphs in polynomial time and sample complexity. *The Journal of Machine Learning Research*, 7:1743–1788, 2006. 4

[2] Alexandre Alahi, Kratarth Goel, Vignesh Ramanathan, Alexandre Robicquet, Li Fei-Fei, and Silvio Savarese. Social lstm: Human trajectory prediction in crowded spaces. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 961–971, 2016. 1, 2, 7

[3] John P Alsterda, Matthew Brown, and J Christian Gerdes. Contingency model predictive control for automated vehicles. In *2019 American Control Conference (ACC)*, pages 717–722. IEEE, 2019. 7

[4] Mayank Bansal, Alex Krizhevsky, and Abhijit Ogale. Chauffeurnet: Learning to drive by imitating the best and synthesizing the worst. *arXiv preprint arXiv:1812.03079*, 2018. 2

[5] Dhruv Batra, Payman Yadollahpour, Abner Guzman-Rivera, and Gregory Shakhnarovich. Diverse m-best solutions in markov random fields. In *European Conference on Computer Vision*, pages 1–16. Springer, 2012. 12

[6] Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. Fast unfolding of communities in large networks. *Journal of statistical mechanics: theory and experiment*, 2008(10):P10008, 2008. 3

[7] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *IEEE/CVF conference on computer vision and pattern recognition (CVPR)*, pages 11621–11631, 2020. 5

[8] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11621–11631, 2020. 6

[9] Sergio Casas, Cole Gulino, Renjie Liao, and Raquel Urtasun. Spagnn: Spatially-aware graph neural networks for relational behavior forecasting from sensor data. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9491–9497. IEEE, 2020. 2, 3, 7

[10] Sergio Casas, Cole Gulino, Simon Suo, Katie Luo, Renjie Liao, and Raquel Urtasun. Implicit latent variable model for scene-consistent motion forecasting. In *European Conference on Computer Vision (ECCV)*, pages 624–641. Springer, 2020. 2, 3

[11] Sergio Casas, Cole Gulino, Simon Suo, and Raquel Urtasun. The importance of prior knowledge in precise multimodal prediction. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2295–2302. IEEE, 2020. 5

[12] Yuxiao Chen, Ugo Rosolia, Chuchu Fan, Aaron D Ames, and Richard Murray. Reactive motion planning with probabilistic safety guarantees. *arXiv preprint arXiv:2011.03590*, 2020. 3

[13] Yuxiao Chen, Ugo Rosolia, Wyatt Ubellacker, Noel Csomay-Shanklin, and Aaron D Ames. Interactive multi-modal motion planning with branch model predictive control. *arXiv preprint arXiv:2109.05128*, 2021. 7

[14] Jan Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, and Yoshua Bengio. Attention-based models for speech recognition. *arXiv preprint arXiv:1506.07503*, 2015. 4

[15] Alexander Cui, Sergio Casas, Abbas Sadat, Renjie Liao, and Raquel Urtasun. Lookout: Diverse multi-future prediction and planning for self-driving. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 16107–16116, 2021. 3

[16] Henggang Cui, Vladan Radosavljevic, Fang-Chieh Chou, Tsung-Han Lin, Thi Nguyen, Tzu-Kuo Huang, Jeff Schneider, and Nemanja Djuric. Multimodal trajectory predictions for autonomous driving using deep convolutional networks. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 2090–2096. IEEE, 2019. 2, 5

[17] Nachiket Deo and Mohan M Trivedi. Convolutional social pooling for vehicle trajectory prediction. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 1468–1476, 2018. 2

[18] Nachiket Deo and Mohan M Trivedi. Multi-modal trajectory prediction of surrounding vehicles with maneuver based lstms. In *2018 IEEE Intelligent Vehicles Symposium (IV)*, pages 1179–1184. IEEE, 2018. 2, 7

[19] Nemanja Djuric, Henggang Cui, Zhaoen Su, Shangxuan Wu, Huahua Wang, Fang-Chieh Chou, Luisa San Martin, Song Feng, Rui Hu, Yang Xu, et al. Multixnet: Multiclass multistage multimodal motion prediction. *arXiv preprint arXiv:2006.02000*, 2020. 2, 5

[20] Lester E Dubins. On curves of minimal length with a constraint on average curvature, and with prescribed initial and terminal positions and tangents. *American Journal of mathematics*, 79(3):497–516, 1957. 4

[21] Agrim Gupta, Justin Johnson, Li Fei-Fei, Silvio Savarese, and Alexandre Alahi. Social gan: Socially acceptable trajectories with generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2255–2264, 2018. 2, 7

[22] Dirk Helbing and Peter Molnar. Social force model for pedestrian dynamics. *Physical review E*, 51(5):4282, 1995. 2

[23] Xin Huang, Stephen G McGill, Jonathan A DeCastro, Luke Fletcher, John J Leonard, Brian C Williams, and Guy Rosman. Diversitygan: Diversity-aware vehicle motion prediction via latent semantic sampling. *IEEE Robotics and Automation Letters*, 5(4):5089–5096, 2020. 3

[24] Boris Ivanovic, Amine Elhafsi, Guy Rosman, Adrien Gaidon, and Marco Pavone. Mats: An interpretable trajectory forecasting representation for planning and control. *Conference on Robot Learning (CoRL)*, 2020. 3

[25] Boris Ivanovic and Marco Pavone. The trajectron: Probabilistic multi-agent trajectory modeling with dynamic spatiotemporal graphs. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2375–2384, 2019. 4

[26] Alexey Kamenev, Lirui Wang, Ollin Boer Bohan, Ishwar Kulkarni, Bilal Kartal, Artem Molchanov, Stan Birchfield, David Nistér, and Nikolai Smolyanskiy. Predictionnet: Real-time joint probabilistic traffic prediction for planning, control, and simulation. *arXiv preprint arXiv:2109.11094*, 2021. 2

[27] Kris M Kitani, Brian D Ziebart, James Andrew Bagnell, and Martial Hebert. Activity forecasting. In *European conference on computer vision*, pages 201–214. Springer, 2012. 2

[28] Alex Kuefler, Jeremy Morton, Tim Wheeler, and Mykel Kochenderfer. Imitating driver behavior with generative adversarial networks. In *IEEE Intelligent Vehicles Symposium (IV)*, pages 204–211. IEEE, 2017. 1, 2

[29] Namhoon Lee, Wongun Choi, Paul Vernaza, Christopher B Choy, Philip HS Torr, and Manmohan Chandraker. Desire: Distant future prediction in dynamic scenes with interacting agents. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 336–345, 2017. 1, 2

[30] Alon Lerner, Yiorgos Chrysanthou, and Dani Lischinski. Crowds by example. In *Computer graphics forum*, volume 26, pages 655–664. Wiley Online Library, 2007. 5, 6

[31] Yicheng Liu, Jinghuai Zhang, Liangji Fang, Qinhong Jiang, and Bolei Zhou. Multimodal motion prediction with stacked transformers. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7577–7586, 2021. 2, 5

[32] Wei-Chiu Ma, De-An Huang, Namhoon Lee, and Kris M Kitani. Forecasting interactive dynamics of pedestrians with fictitious play. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 774–782, 2017. 2

[33] Ramin Mehran, Alexis Oyama, and Mubarak Shah. Abnormal crowd behavior detection using social force model. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 935–942. IEEE, 2009. 2

[34] Jiquan Ngiam, Benjamin Caine, Vijay Vasudevan, Zhengdong Zhang, Hao-Tien Lewis Chiang, Jeffrey Ling, Rebecca Roelofs, Alex Bewley, Chenxi Liu, Ashish Venugopal, et al. Scene transformer: A unified multi-task model for behavior prediction and planning. *arXiv preprint arXiv:2106.08417*, 2021. 2

[35] Stefano Pellegrini, Andreas Ess, Konrad Schindler, and Luc Van Gool. You'll never walk alone: Modeling social behavior for multi-target tracking. In *2009 IEEE 12th International Conference on Computer Vision*, pages 261–268. IEEE, 2009. 5, 6

[36] Tung Phan-Minh, Elena Corina Grigore, Freddy A Boulton, Oscar Beijbom, and Eric M Wolff. Covernet: Multimodal behavior prediction using trajectory sets. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14074–14083, 2020. 2

[37] Thomas A Ranney. Models of driving behavior: a review of their evolution. *Accident analysis & prevention*, 26(6):733–750, 1994. 1

[38] Nicholas Rhinehart, Rowan McAllister, Kris Kitani, and Sergey Levine. Precog: Prediction conditioned on goals in visual multi-agent settings. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2821–2830, 2019. 5

[39] R Tyrrell Rockafellar, Stanislav Uryasev, et al. Optimization of conditional value-at-risk. *Journal of risk*, 2:21–42, 2000. 5

[40] Amir Sadeghian, Vineet Kosaraju, Ali Sadeghian, Noriaki Hirose, Hamid Rezatofighi, and Silvio Savarese. Sophie: An attentive gan for predicting paths compliant to social and physical constraints. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1349–1358, 2019. 2, 7

[41] Amir Sadeghian, Ferdinand Legros, Maxime Voisin, Ricky Vesel, Alexandre Alahi, and Silvio Savarese. Car-net: Clairvoyant attentive recurrent network. In *European Conference on Computer Vision (ECCV)*, pages 151–167, 2018. 7

[42] Tim Salzmann, Boris Ivanovic, Punarjay Chakravarty, and Marco Pavone. Trajectron++: Dynamically-feasible trajectory forecasting with heterogeneous data. In *European Conference on Computer Vision (ECCV)*, pages 683–700. Springer, 2020. 2, 4, 6, 7

[43] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. *Advances in neural information processing systems*, 28:3483–3491, 2015. 5

[44] Charlie Tang and Russ R Salakhutdinov. Multiple futures prediction. *Advances in Neural Information Processing Systems*, 32:15424–15434, 2019. 1, 2, 3

[45] Martin Treiber, Ansgar Hennecke, and Dirk Helbing. Congested traffic states in empirical observations and microscopic simulations. *Physical review E*, 62(2):1805, 2000. 2

[46] A. Vemula, K. Muelling, and J. Oh. Social attention: Modeling attention in human crowds. In *International Conference on Robotics and Automation (ICRA)*, 2018. 7

[47] Andreas Wächter and Lorenz T Biegler. On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming. *Mathematical programming*, 106(1):25–57, 2006. 7

[48] Xinshuo Weng, Ye Yuan, and Kris Kitani. Ptp: Parallelized tracking and prediction with graph neural networks and diversity sampling. *IEEE Robotics and Automation Letters*, 6(3):4640–4647, 2021. 2

[49] Ye Yuan and Kris Kitani. Diverse trajectory forecasting with determinantal point processes. *arXiv preprint arXiv:1907.04967*, 2019. 3

[50] Ye Yuan, Xinshuo Weng, Yanglan Ou, and Kris Kitani. Agentformer: Agent-aware transformers for socio-temporal multi-agent forecasting. *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 2

[51] Tianyang Zhao, Yifei Xu, Mathew Monfort, Wongun Choi, Chris Baker, Yibiao Zhao, Yizhou Wang, and Ying Nian Wu. Multi-agent tensor fusion for contextual trajectory prediction. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12126–12134, 2019. 2, 7

[52] Rui Zhou, Hongyu Zhou, Masayoshi Tomizuka, Jiachen Li, and Zhuo Xu. Grouptron: Dynamic multi-scale graph convolutional networks for group-aware dense crowd trajectory forecasting. *arXiv preprint arXiv:2109.14128*, 2021. 2

## A. Preprocessing and Dynamics

### A.1. Pre-encoding of Nodes and Edges

Since all datasets we use provide global coordinates and headings of the nodes, we first transform the global coordinates of the state history and state future to local frames fixed at the nodes' current position. For vehicles with heading angle $\psi$, we use $\cos(\psi)$ and $\sin(\psi)$ as features instead of $\psi$ to prevent the $\pm 2\pi$ issue. For every type of nodes, the raw features are passed through a fully connected layer as pre-encoding, and the network is shared in all modules that require state information, i.e., whenever we use state information, the state vector passes through the pre-encoding layer first.

For edges, we construct a pre-encoding layer for every edge type (e.g., vehicle-vehicle, vehicle-pedestrian, pedestrian-vehicle, and pedestrian-pedestrian). The pre-encoding layer extracts raw features from the states of the two agents then passes them through a fully connected layer. The raw features contain agents' relative position and relative velocity in the local frame as well as their sizes.

### A.2. Dynamic Models and Collision Checking

We use a Dubin's car model for vehicles in the scene with

$$
s = \begin{bmatrix} X \\ Y \\ v \\ \psi \end{bmatrix}, a = \begin{bmatrix} \dot{v} \\ \dot{\psi} \end{bmatrix}, s^+ = \begin{bmatrix} X + v\cos(\psi)\Delta t \\ Y + v\sin(\psi)\Delta t \\ v + \dot{v}\Delta t \\ \psi + \dot{\psi}\Delta t \end{bmatrix} \text{ where}
$$

$v$ and $\dot{v}$ are the longitudinal velocity and acceleration, $\psi$ and $\dot{\psi}$ are the heading angle and yaw rate.

The pedestrians follow a double integrator model with

$$
s = \begin{bmatrix} X \\ Y \\ v_x \\ v_y \end{bmatrix}, a = \begin{bmatrix} \dot{v}_x \\ \dot{v}_y \end{bmatrix}, s^+ = \begin{bmatrix} X + v_x\Delta t \\ Y + v_y\Delta t \\ v_x + \dot{v}_x\Delta t \\ v_y + \dot{v}_y\Delta t \end{bmatrix}.
$$

Indeed, both models consists of basic differentiable functions that can be incorporated in a neural network. We also put bound on the inputs $\dot{v} \in [-5m/s, 5m/s], \dot{\psi} \in [-1m/s^2, 1m/s^2], v_x, v_y \in [-5m/s, 5m/s]$ so that the generated trajectory predictions are dynamically feasible.

### A.3. Collision Check

We model pedestrians as circles with varying radius and vehicles as rectangles. The collision between pedestrians is straightforward to check, simply by taking the Euclidean distance between the two pedestrians. Collisions involving vehicles are checked in the local coordinate frame of the vehicle. Fig. 9 shows the case with a pedestrian and a vehicle, and the collision function is

$$
\text{Col}(\Delta X, \Delta Y, L, W) = \max\{|\Delta X| - \frac{L}{2}, |\Delta Y| - \frac{W}{2}\}.
$$

Vehicle-to-vehicle collision is a bit tricky since it involves two rectangles. As shown in Fig. 10, we use the
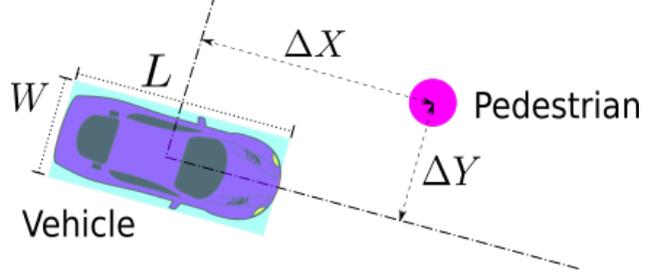


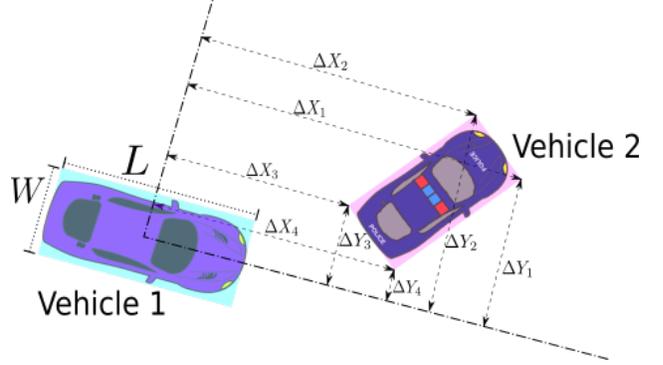Figure 9. Collision check between a vehicle and a pedestrian



Figure 10. Collision check between two vehicles

four corners to calculate the Col function:

$$
\text{Col}(\Delta X_{1:4}, \Delta Y_{1:4}, L, W)
$$
$$
= \max \left\{ \begin{array}{l} |\Delta X_1| - \dfrac{L}{2}, ...|\Delta X_4| - \dfrac{L}{2}, \\ |\Delta Y_1| - \dfrac{W}{2}, ...|\Delta Y_4| - \dfrac{L}{2} \end{array} \right\}.
$$

Note that the collision functions are all differentiable (at least piecewise differentiable), making it convenient to include them in the training process as regularization.

### A.4. Diversity Scheduling during Training

The parameter $\alpha$ serves as a tuning knob to adjust the tradeoff between encoder accuracy and diversity. During training, we start with a small $\alpha$ so that the decoder can learn diverse trajectory patterns without mode collapse, then increase $\alpha$ to improve the encoder's prediction accuracy. When $\alpha$ is above a threshold, we detach the prediction error loss of all modes but the one with the largest $Q$ in Eq. (4) from the gradient graph to avoid mode collapse. This allows us to reduce the mode collapse under a small $\alpha$ while continue to improve the encoder on the mode probability prediction.

### A.5. Diverse Sampling from Product Latent Space

Since $P(\mathbf{z}|\mathbf{x})$ is calculated with factor graphs, two clique modes with similar latent variables, e.g., two $\mathbf{z}$-s that only

differ at one node in the clique, may have similar probabilities,causing the greedy sampling result to lose diversity. This issue is well-known in Markov random fields and we follow the simple diverse sampling scheme in [5]. To be specific, we use a greedy algorithm to pick $\mathbf{z}$ one by one as

$$\mathbf{z}^{k+1} = \arg\max_{\mathbf{z} \in \mathcal{Z}} p(\mathbf{z}|x)$$
$$\text{s.t. } \forall \mathbf{z}^i, i = 1, ..., k, \Delta(\mathbf{z}, \mathbf{z}^i) \geq \beta,$$

where $\Delta$ is a distance function, for our setup, $\Delta(\mathbf{z}^1, \mathbf{z}^2) = |\{j|z_j^1 \neq z_j^2\}|$, i.e., the number of nodes with different latent variables under $\mathbf{z}^1$ and $\mathbf{z}^2$. For example, $\Delta([0, 1, 2], [1, 1, 2]) = 1, \Delta([1, 0, 0], [2, 1, 0]) = 2$.