# CVAE-H: Conditionalizing Variational Autoencoders via Hypernetworks and Trajectory Forecasting for Autonomous Driving

Geunseob (GS) Oh, Huei Peng
University of Michigan
gsoh@umich.edu, hpeng@umich.edu

## Abstract

*The task of predicting stochastic behaviors of road agents in diverse environments is a challenging problem for autonomous driving. To best understand scene contexts and produce diverse possible future states of the road agents adaptively in different environments, a prediction model should be probabilistic, multi-modal, context-driven, and general. We present Conditionalizing Variational AutoEncoders via Hypernetworks (CVAE-H); a conditional VAE that extensively leverages hypernetwork and performs generative tasks for high-dimensional problems like the prediction task. We first evaluate CVAE-H on simple generative experiments to show that CVAE-H is probabilistic, multi-modal, context-driven, and general. Then, we demonstrate that the proposed model effectively solves a self-driving prediction problem by producing accurate predictions of road agents in various environments.*

## 1. Introduction

Over the last few years, companies have made major progress towards the real-world deployment of autonomous driving technologies. In late 2020, Waymo, Google's autonomous vehicle (AV) company, became the first company that offered commercial autonomous taxi rides to the public in Pheonix, Arizona. With the initial milestone for the deployment made, the challenge of developing effective AV technologies that can operate in a wide variety of complex urban environments has never been more important. As pointed out by a recent publication from Waymo, generalization within and between operating regions is crucial to the overall viability of AV technologies [1].

Autonomous vehicles interact with stochastic road users such as human drivers in different driving scenarios and road topologies. The presence of multiple stochastic agents creates uncertainties that grow over time, contributing to the increase in the number of variations of the environments and scenarios that an autonomy stack (perception, prediction,
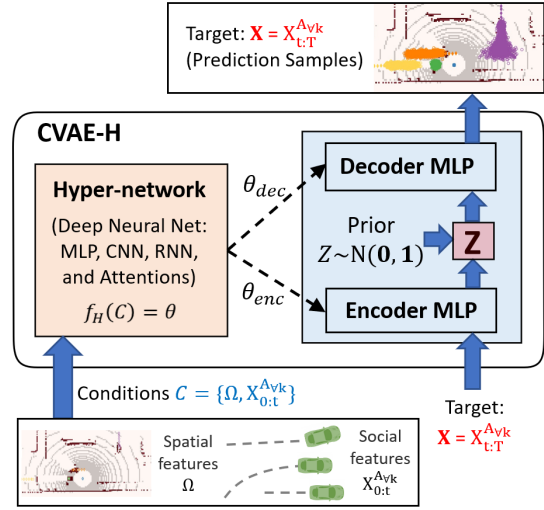


Figure 1: CVAE-H is a conditional VAE that integrates a hypernetwork into a VAE. In a trajectory forecasting problem, CVAE-H takes a high-dimensional features (e.g., lidar and labelled trajectories) as conditions and produces diverse predictions for dynamic road agents in the scene.

planning, and control) of the AV needs to handle. Among the modules of the autonomy stack, the prediction is where the stochasticity and diversity culminates, as it is responsible of understanding scene contexts of various complex environments and predicting a set of possible future trajectories of road users.

There are four key attributes that a prediction module should possess: *probabilistic*, *multi-modal*, *context-driven*, and *general*. Autonomous driving often involves interactions with stochastic road users. Upon the presence of road users, AV needs to predict their states to avoid collisions. In addition to bare predictions, a prediction model should also output the quantified uncertainties of the predictions, which are essential in the downstream modules such as planning.

Another characteristic of road agents is the diversity of possible trajectories. For example, a vehicle approaching a 4-way intersection has four possible *modes*; go straight,

turn left, turn right, or yield if there is a lead vehicle waiting in the queue. A vehicle passing through a crosswalk may slow down or maintain its speed.

There are various factors that contribute to the diversity of the trajectories. Such factors include the relative states of dynamic agents, road geometries, and static obstacles in the scene. To consider these factors, it is important to capture the impact of both social features (e.g., interactions among road-agents) and spatial features (e.g., lanes, stop signs, and crossroads) in the prediction.

Last attribute of autonomous driving is that numerous instances of driving environments exist. Road configuration alone is diverse; examples include un-signalized intersections with single-lane, two-lanes, 2-way stops, or u-turns, one-way roads, straight roads, curved roads, and so on. The presence of static and dynamic agents diversify the urban environments as well. This means that an ad-hoc prediction approach customized to a specific instance of an environment is impractical. A scalable prediction model should adaptively produce predictions in various environments.

It is not trivial for a prediction model to satisfy all four attributes. First of all, not all neural-network models are probabilistic unless they build on explicit density models such as Mixture Density Network (e.g., Gaussian Mixture Models) [2, 3], Normalizing-flow models [4, 5], Variational Autoencoder (VAE) [6]. Secondly, not all probabilistic models are multi-modal, for instance, a uni-modal Gaussian. Third, to be fully context-driven, the model should leverage both the social and spatial information. Some models [7, 8] are based exclusively on Convolutional Neural Network (CNN) and may not best capture social features. Similarly, models that rely solely on social features [9–11] may not consider spatial features of the environment that come from HD Maps, images, and/or lidar point clouds. Finally, generality of a prediction model is hard to achieve. Models that leverage scenario-specific heuristics or knowledge may sacrifice generality over the precision for the targeted scenarios.

To build a model that possesses all four attributes of a ideal prediction model, we leverage the following two architectures: VAE [12], which is an explicit generative model, and hypernetwork [13], which may consist of diverse kinds of neural networks for the context-driven prediction. By integrating a hypernetwork into a VAE, we propose CVAE-H. This facilitates the informational flow of the conditions to the VAE network across all layers and allows the VAE to effectively scale up to work with large conditional inputs. We first evaluate CVAE-H on simple Gaussian experiments to show that CVAE-H is probabilistic, multi-modal, context-driven, and general. Then, we demonstrate that the proposed model effectively solves a self-driving prediction problem by utilizing both social and spatial features of the environments and producing accurate predictions of road agents in various environments.

## 2. Related Works

In this section, we review existing models for the prediction task and categorize them. The first categorization of the prediction approaches is based on input features the prediction model utilizes. Models that focus on the social features include Social-LSTM [9], Social-GAN [10], and Trafficpredict [14]. Some models focused on exploiting the spatial features, for example, IntentNet [7], Fast and Furious [8], that effectively capture underlying traffic rules from the road geometry; however, interactions among road agents are not explicitly modeled. There are also models that utilize both social and spatial features such as DE-SIRE [6], SoPhie [15], R2P2 [4], and MTP [16].

The second categorization is based on the core mechanics of the prediction models; either deterministic or probabilistic. Many earlier deterministic models leveraged either physics-based modeling techniques [17] (e.g., constant velocity model and car-following model) or relatively light-weight neural networks such as multi-layer perceptron (MLP), RNN, and/or CNN [9], while a few recent deterministic models work with larger deep neural networks [7, 8, 16]. Generative models are examples of probabilistic models that model a joint distribution $p(x_i)$ using $z$ that is a variable that encodes the underlying dynamics of the data distribution. Furthermore, a group of generative models that are used to model a conditional distribution $p(x_o|x_i)$ are called *conditional generative models*. Examples include conditional VAE [18], conditional GAN [19], and conditional flow [5].

Modern generative models, including GAN, VAE, Flow, and Mixture Density Networks (MDN) [2] are popular across several domains of artificial intelligence and machine learning due to their capacities to approximate complex probability distributions. GAN is well-known for producing samples of excellent qualities. However, it does not explicitly model the probability [20]. This means that GAN has to rely on expensive Monte-Carlo sampling to approximate the probability distribution of the generated samples to reason the uncertainties of the predictions.

Unlike GAN, VAE is a type of *explicit* generative model that models probability density $p(x)$ explicitly. This property allows us to obtain the uncertainties of the sample predictions directly. Besides, several recent VAE-based models have shown promising results at predicting trajectories of human-driven vehicles in urban environments [6, 16, 21].

Normalizing flow, or Flow, is also a type of explicit generative model. While flow-based models [4, 5, 22] have access to the exact probability density $p(x)$, flow-based models are typically better suited for density estimations than generative tasks (e.g., trajectory forecasting).

MDN is another generative model that models explicit probability densities. A common choice for the density model is a mixture of Gaussians where the neural network

estimates the parameters of the Gaussian mixture. MDN with Gaussian mixtures have been used for prediction tasks [3, 14, 23]. Since MDN works with a pre-determined number of modes, a MDN(GMM)-based model may sacrifice the generality.

# 3. Preliminaries

In this section, we review the deep learning concepts that are the two main keys of CVAE-H; Variational autoencoder [12] and Hypernetwork [13].

## 3.1. Variational AutoEncoder (VAE)

Variational Autoencoder [12] is a generative model that learns a data distribution $p(X)$ in an unsupervised way by using an abstract (i.e., latent) variable $Z$ that captures the underlying dynamics of the data distribution $X$.

As the name suggest, VAE is an *autoencoder* model [24] trained using variational inference [25]. Autoencoder is an unsupervised model that learns an identity function using an encoder $f_{enc}$ and decoder $f_{dec}$. The encoder $f_{enc}(X) = Z$ takes the input data $X$ that are high dimensional and compresses it into a lower dimensional data $Z$. The decoder $f_{dec}(Z) = X$ then takes $Z$ and recovers the higher dimensional input $X$. During the reconstruction of the original inputs, the autoencoder aims to discover a more efficient lower dimensional representation of the inputs.

VAE combines the autoencoder with variational inference. Instead of mapping the input into a fixed vector, VAE maps the input into probability distributions; the encoder and decoder each models $p(Z|X)$ and $p(X|Z)$, which are both typically parameterized using neural networks. As the compression $X \rightarrow Z$ and generation $Z \rightarrow X$ steps become probabilistic, VAE can be used for generative tasks by sampling $Z$ and passing it through the decoder to obtain $X$.

Let us denote $p(Z|X)$ and $p_m(Z|X)$ as the true posterior and model posterior distributions. Likewise, $p(X|Z)$ and $p_m(X|Z)$ are denoted as the true likelihood and model likelihood. Finally, we denote $p(Z)$ as the prior distribution. The learning objective is to maximize the likelihood of the training data to find the optimal parameters $\theta^* = argmax(\prod_{k=1}^{N} p_m(X_k))$, where $p(X_k) = \int_Z p(X_k|Z)p(Z)dZ$ via the use of law of total probability.

The computation of $p(X_k)$ is intractable, but can be approximated using variational inference [25]. Instead of maximizing the likelihood, a computable term called *Evidence Lower BOund (ELBO)* of the likelihood can be maximized. As the name indicates, ELBO is a lower bound of the original likelihood, i.e., $p(X) \geq ELBO$. Precisely, the following relationship holds: $log(p(X)) = ELBO + KL(p_m(Z|X_k)||p(Z|X_k))$, where ELBO is defined as $\mathbb{E}_Z \left[ log(p(X_k|Z)) \right] - KL(p_m(Z|X_k)||p(Z))$. The ELBO is

easily computable as long as the posterior, prior, and generative distributions are modeled using explicit density models and through the use of reparameterization trick [12].

VAE approximates the target distribution and is typically used to generate new data or estimate the probability density of the target approximately. In practice, VAE-based models have demonstrated competitive results for trajectory prediction tasks [6, 21] as well as other machine learning tasks such as sequence modeling and density estimation [26–28].

## 3.2. Hypernetwork

Assume a function that takes two inputs $X, C$ and outputs $Y$, i.e., $f : (X, C) \rightarrow Y$. Given that $f$ is a probabilistic mapping, $p(Y|X, C)$ represents a probability distribution of the output $Y$ conditioned on the two inputs. Learning such $f$ and the corresponding conditional distribution can be achieved using two different approaches; (1) the embedding and (2) the hypernetwork approach, as depicted in Figure 2.
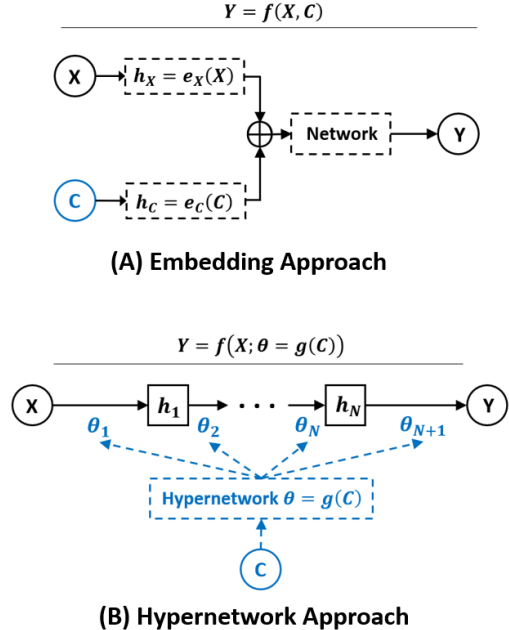


**(A) Embedding Approach**



**(B) Hypernetwork Approach**

Figure 2: Two approaches to build conditional neural nets. $X, Y, C$ each indicates input, output, and conditions.

As commonly used in machine learning, the embedding approach is a straightforward way of learning $f$. This approach implies a model that consists of the embedding modules $e_1(X) = h(X)$ and $e_2(C) = h(C)$ where the two outputs are concatenated $h := h(X) \oplus h(C)$ and passed to another network $f(h) = Y$.

The hypernetwork approach [13] is a less intuitive solution that relies on a hierarchical structure between two models; one called a primary network $g$ produces weights of the

other separate network $f$. This corresponds to $f(X;\theta) = Y$ where $\theta = g(C)$. Observe that the embedding approach uses an embedding network $e$ and treats $C$ equally to $X$ as part of the input. On the other hand, the hypernetwork approach models the conditional mapping $f_C(X) = Y$, i.e., learning the contributions of $X$ to $Y$ via a parameterized function $f(X;\theta)$ where the parameters themselves are the outputs of another network $g$, which is independent of $f$.

The ability to effectively learn conditional functions (i.e., model a function that transforms into different functions depending on the condition $C$) corresponding to the property of modularity [29] is a major benefit of the hypernetwork approach. A recent study showed that hypernetwork exhibits the property of modularity for a large $g$ and small $f$ and the hypernetwork approach uses much smaller networks compared to the embedding approach to achieve the same performance in modeling conditional distributions [29]. Furthermore, hypernetworks offer excellent *generalization* capability when it is combined with deep neural networks. This has been demonstrated across diverse tasks in machine learning [5, 29–35].

Hypernetwork provides a way to make the unconditional generative models (e.g., vanilla VAE and normalizing flow) conditional without a use of RNN. Hypernetwork is a neural-network that can be modeled using various types of neural networks modules such as CNN, RNN, Residual blocks [36], and Attention [37]. The flexibility of the hypernetwork allows the generative models to even condition on bigger models such ResNet [38], BERT [39], or FPN [40] by including such models as part of the hypernetwork. In this sense, hypernetworks take diverse types of data including social and spatial data. The effectiveness of the hypernetworks grows with respect to the capacity of its backbone neural network. Upon use of CNN blocks (spatial features) and RNN/Attention modules (social features), the main neural network can become *context-driven*.

This becomes especially useful for the prediction task for autonomous driving. A challenge of the prediction task is that it may involve very large inputs whose dimensions are in the magnitude of millions. This is because the inputs include social features (e.g., labeled states for all road-agents in the vicinity of the AV) and spatial features (e.g., lidar scans, camera images, and/or HD maps) from past to the current moment. Since our prediction model extensively leverages hypernetworks, the high-dimensional information about the driving environment can be effectively encoded and used for contextual prediction.

# 4. CVAE-H

We propose CVAE-H (Conditionalized Variational AutoEncoder via Hypernetwork), a conditional VAE that integrates a hyper-network into a VAE. The hypernetwork encodes various conditioning information (e.g., social and spatial information) and the VAE utilizes the encoding as the network parameters to perform auto-encoding $X > Z > X$ and generative $Z > X$ tasks.

As elaborated in Section 3.2, one advantage of hypernetwork over the embedding approach is the flexibility of hypernetwork (i.e., any differentiable neural network can be used as a hypernetwork). In addition, the hypernetwork passes the information from the conditions to the VAE network across all VAE layers and thus allows the VAE to scale up to large conditional inputs of the prediction problem. Another advantage is that hyper-network can encode conditioning information with a smaller number of network parameters compared to the embedding approach [29].
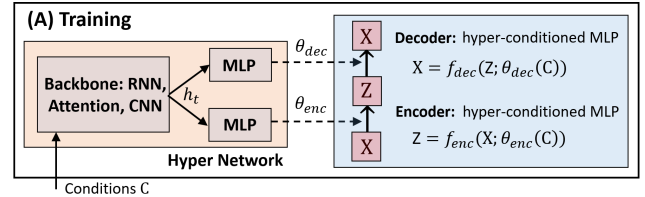


Figure 3: Schematics of the CVAE-H architecture. During the training, CVAE-H takes conditions $C$ and a target variable $X$ and learns the encoder and decoder to reconstruct $X$. For the prediction problem, we use the following conditions $C := [X_{-t:0}^{A_{\forall k}}, \Omega]$ and target variable $X_{0:T} := [x_{0:T}^{A_{\forall k}}, y_{0:T}^{A_{\forall k}}]$.

Figure 3 illustrates CVAE-H, which largely consists of two modules: 1) a VAE and 2) a hypernetwork that computes the network parameters of 1). In the following subsections, the details of the two modules are presented.

## 4.1. Hyper-conditioning VAE

Conditionalizing a VAE via a hypernetwork involves the following process of parameterizing the encoder and decoder of the VAE:

$$
\begin{aligned}
(Encoder) \quad & Z = f_{enc}(X; \theta_{enc}(C)) \\
(Decoder) \quad & X = f_{dec}(Z; \theta_{dec}(C))
\end{aligned}
\tag{1}
$$

where $\theta_{enc}$ and $\theta_{dec}$ are obtained via the hypernetwork $f_H(C)$ as in $f_H(C) = [\theta_{enc}, \theta_{enc}]$. If a VAE is conditionalized using the embedding approach, the resulting equations would instead be $Z = f_{enc}(X, g_1(C))$ and $X = f_{dec}(Z, g_2(C))$. For the reasons described in Section 3.2, we conditionalize the VAE using a hypernetwork.

The VAE architecture can consist any type of neural network. For a 1D target variable such as a sequence of coordinates or a collection of variables, MLP is a good choice. While we expect CNN to be an effective encoder & decoder structure for a 2D target variable such as an image, we leave the search of proper network for 2D target variables to the future work.

Similar to how $p(X)$ is modeled in an unconditional VAE, a latent variable $Z$ is introduced to map the inputs $X$ and condition $C$ to the the conditional distribution $p(X|C)$ as follows: $p(X|C) = \int_Z p(X|Z,C)p(Z|C)dZ$. Here, $p(Z|C)$ and $p(X|Z,C)$ represent the prior and the generative decoder conditioned on $C$, respectively. Ideally, $Z$ should encode abstract features of $p(X|C)$.

Since the integral to compute $p(X|C)$ is intractable, variational inference [25] is used to approximate the integral with a parameterized posterior $q_{\theta_{enc}}(Z|X,C)$. This leads to the construction of ELBO on the log conditional likelihood:

$$\log p(X|C) \geq \mathbb{E}_{q_{\theta_{enc}}(Z|X,C)}[\log p_{\theta_{dec}}(X|Z,C)]$$
$$-\text{KL}(q_{\theta_{enc}}(Z|X,C)||p(Z|C)). \quad (2)$$

We set the conditional prior $P(Z|C)$ to be a zero-mean standard multi-variate Gaussian (i.e., diagonal covariance with all entries equal to one). Different choices of the prior include parameterizing priors using neural networks and/or domain knowledge [41].

### 4.2. Hypernetwork and Training of CVAE-H

The hypernetwork aims to effectively propagate the conditioning information to the VAE. An earlier work that extensively leveraged hypernetwork to conditionalize a normalizing-flow [5] showed that the main-network becomes versatile by customizing the backbone neural networks of the hypernetwork. In this regard, for sequential targets, a hypernetwork can include RNN, attention modules, and/or Transformer [37]. For computer vision tasks, a hypernetwork can use CNN-based models.

The goal of CVAE-H is to learn the target conditional distribution $p(X|C)$ by minimizing the negative log-likelihood of model distribution $p_{model}(X|C)$. This is achieved by maximizing the conditional ELBO presented in Equation 2 (i.e., the ELBO is the loss function). The first term of the ELBO is the reconstruction error of the VAE and can be computed as long as the estimation of $p_{\theta_{dec}}(X|Z,C)$ is tractable. For this reason, we design the decoder to be an explicit density model (e.g., softmax, MDN [2]). In this work, we model the decoder to output mixture parameters for a Gaussian mixture: the mean $\mu$, covariance $\Sigma$, and mixing probabilities of the output Gaussian mixture $\rho$. The second term of the ELBO is a KL divergence and can be easily computed as we use a multi-variate standard Gaussian for the prior $p(Z|C)$ and a multi-variate (uni-modal) Gaussian with a diagonal covariance for the posterior $q_{\theta_{enc}(Z|X,C)}$.

### 4.3. Inference: CVAE-H for Generative Tasks

The training and inference processes for CVAE-H are different. During the training, we have access to the target variable $X$ and we run the encoder $Z = f_{enc}(X; \theta_{enc}(C))$ to
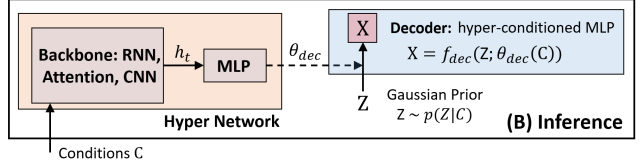


Figure 4: A schematic of the CVAE-H during inference. As opposed to the training, CVAE-H does not take $X$ and only utilizes the conditions $C$. Instead, $Z$ is sampled using the prior and is passed through the decoder to obtain $X$.

obtain the posterior $Z$. The decoder then takes $Z$ to recover $X = f_{dec}(Z; \theta_{dec}(C))$. As illustrated in Figure 4, during the inference, the encoder is not used because we do not have access to $X$. Instead, we leverage the Gaussian prior distribution $p(Z|C)$, sample latent variable $Z \sim p(Z|C)$, and run the decoder to generate $X = f_{dec}(Z; \theta_{dec}(C))$. Similar to the training, $C$ is passed to the hyper-network and $\theta_{dec}$ is computed to construct the decoder $f_{dec}$, but $\theta_{enc}$ is not computed as the we do not use the encoder in the inference. In this regard, CVAE-H can perform both the *auto-encoding* of a target variable like a vanilla VAE ($X > Z > X$) and *generative process* ($Z > X$) via the inference process described above. We note that we use the generative process for the self-driving prediction problem as we do not have access to the target variable $X$ during the inference.

## 5. Experiments: Gaussian Generative Tasks

The next two Sections are dedicated to the experiments we conducted to evaluate CVAE-H. This Section introduces the *Gaussians* experiments that are designed to examine the effectiveness of CVAE-H on generative tasks. Then in the next Section, we evaluate CVAE-H on a prediction task for autonomous driving, which is a more challenging problem that involves high-dimensional spatio-temporal sensor information such as lidar scans of the environments and labelled trajectories.

In this Section, we assess the model in simple generative tasks used in the literature [5, 27, 42, 43]. The task of forecasting a distribution of trajectories is fundamentally a generative task. Therefore, generative tasks help us evaluate the quality of samples generated from CVAE-H.

The Gaussian experiments allow us to evaluate CVAE-H in simpler and more interpretable settings than the self-driving prediction task. They provide us with better insights regarding the strengths and weaknesses of CVAE-H. Most importantly, the experiments examine the capacity of CVAE-H as an ideal prediction model introduced in Section 1. In other words, we validate if the proposed prediction model is probabilistic, multi-modal, context-driven, and general.

Two generative experiments, namely *Gaussian 1* and

*Gaussian 2*, are used to verify whether CVAE-H have the four attributes of the ideal prediction model. To evaluate the performance of the model, we utilize both quantitative and qualitative measures. The quantitative measures include NLL (i.e., cross-entropy) and KL divergence ($D_{KL}$), which are defined in the following.

$$\text{NLL} = -E_{X \sim p(X|C)}[\log(p_{model}(X|C))],$$

$$D_{KL}(p||p_{model}) = \sum_{X \sim p(X|C)} p(X|C)\log\left(\frac{p(X|C)}{p_{model}(X|C)}\right). \tag{3}$$

To evaluate CVAE-H qualitatively, we visualize the outputs of the models. Recall that CVAE-H outputs sample coordinates $X \sim p(X|C)$. In this regard, the visualizations include the resulting probability densities of the generated samples from CVAE-H. Specifically, we discretize a 2D grid and estimate the density of each cell by counting the number of generated samples in the cell.
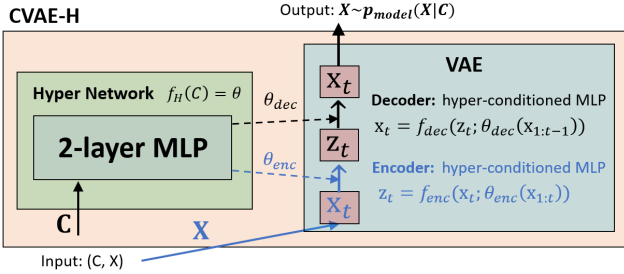


Figure 5: The CVAE-H customized for the Gaussian experiments. The inputs to the CVAE-H are a tuple of target variable $X$ and conditioning information $C$ where $X$ corresponds to the $(x, y)$ coordinates and $C$ corresponds to a set of discrete class variable $C \in \{0, 1, 2\}$ for the Gaussian 1 experiment and $C \in \{(0, 0), (-4, 4), (-4, 4), (4, -4), (4, 4)\}$ for the Gaussian 2 experiment. The parts colored blue are only used for the training. In the inference, Z is sampled from a bivariate standard Gaussian prior.

Figure 5 portrays the CVAE-H network designs customized for the Gaussian experiments. Note that a hypernetwork is a versatile neural network that aims to extract features from the conditioning information $C$. In this regard, a hypernetwork design changes depending on the task and types of $C$.

## 5.1. The Gaussian 1 Experiment

The Gaussian 1 experiment is the experiment used in [5, 43] and aims to show the model's conditioning ability for three distinct multi-modal distributions $p_1(x, y), p_2(x, y), p_3(x, y)$ of 2-by-2, 5-by-5, and 10-by-10 Gaussians depicted in Figure 6. We compare the results

against a conditional generative model, HCNAF [5], which accurately reproduces conditional distributions using a single model.

In this experiment, the hypernetwork is designed using a 2-layer MLP with Relu activation functions and takes a condition $C \in \mathbb{R}^1$ and outputs $\theta_{enc}$ and $\theta_{dec}$. We design the VAE of CVAE-H using 4-layer hyper-conditioned MLPs with Relu activation functions for both encoder and decoder. The last layer of the decoder outputs parameters of a Gaussian mixture (i.e., means, variances, and mixing probabilities). The outputs of CVAE-H are generated by sampling the output Gaussian mixture. All the other parameters were set identically to the setting reported in [5], including those for the Adam optimizer (the learning rate $5e^{-3}$ decays by a factor of 0.5 every 2,000 iterations with no improvement in validation samples). The NLL (i.e., $p_{model}(X_{target})$) values in Table 1 were computed using 10,000 samples.
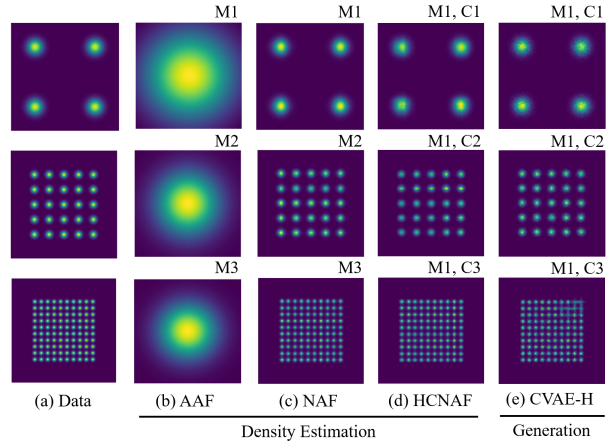


Figure 6: Qualitative result of the Gaussian 1 experiment. M and C each denotes model and condition respectably. This results suggest that CVAE-H models conditional distributions effectively, comparable to the HCNAF results.

In order to reproduce the three distributions, AAF and NAF require three different and separately trained models $p_1(x, y; \theta_1), p_2(x, y; \theta_2), p_3(x, y; \theta_3)$. Conversely, a *single* model of HCNAF ($p(x, y|C_k; \theta)$) and CVAE-H ($p(x, y|Z, C_k; \theta)$) can model all three conditions using a single set of network parameters $\theta$. The conditioning information $C_k \in \{0, 1, 2\}$ each value represents the class of the 2-by-2, 5-by-5, and 10-by-10 Gaussians.

Results from Figure 6 and Table 1 show that CVAE-H is able to reproduce the three nonlinear target distributions, confirming that it are *probabilistic* and *multi-modal*. Depending on the condition $C_k$, CVAE-H produced three very different distributions, thus *context-driven* with respect to simple conditions. CVAE-H also achieves comparable results to HCNAF, which is a powerful conditional density estimator that can model complex conditional distributions.

Table 1: NLL for the experiment depicted in Figure 6. Lower values are better. Other than CVAE-H, we reuse the numbers reported in [5].

|          | AAF   | NAF   | HCNAF | CVAE-H (ours) |
|----------|-------|-------|-------|---------------|
| 2 by 2   | 6.056 | 3.775 | 3.896 | 3.895         |
| 5 by 5   | 5.289 | 3.865 | 3.966 | 3.978         |
| 10 by 10 | 5.087 | 4.176 | 4.278 | 4.292         |

It is important to note that our model uses a *single* model to produce the 3 distinct pdfs like HCNAF, whereas AAF and NAF used *3 distinctly trained* models.

We point out that the plots displayed in Figure 6 contains results from both density estimation (for AAF, NAF, and HCNAF) and sample generation (for CVAE-H). In the density estimation task, a model evaluates $p_{model}(x,y|C_k)$ where $(x,y)$ belongs to a cell in a discrete 2-dimensional grid. In other words, $(x,y) \in \text{Grid}([x_{min},y_{min}],[x_{max},y_{max}])$. In the sample generation task, a set of outputs $(x,y) \sim p_{model}(x,y|Z,C_k;\theta), k=1,2,3$ are generated by passing $Z$, which is sampled from the prior, through the decoder. Then, a histogram is used to classify the generated $(x,y)$ into the 2-dimensional grid. While NLLs for the three density estimators and CVAE-H in the same table, the NLLs were evaluated differently. Precisely, the NLLs for the density estimators were estimated by evaluating $p_{model}(X|C_k)$ whereas the NLL for CVAE-H was obtained by computing $p_{model}(X|Z,C_k), Z \sim N(0,I)$.

## 5.2. The Gaussian 2 Experiment

The Gaussian 1 experiment tested the model's ability to learn target distributions. On the other hand, the Gaussian 2 experiment is designed to evaluate how CVAE-H can generalize over *unseen* conditions $C_{unseen}$. By assessing the model with a set of conditions beyond what it was trained with, we can evaluate the model' capacity to interpolate and extrapolate, which is important for the application of CVAE-H on problems such as a self-driving prediction task.

We train a single CVAE-H model to learn five distinct pdfs, where each pdf represents a Gaussian distribution with its mean used as conditions $C := (x_c,y_c) \in \mathbb{R}^2$ and an isotropic standard deviation $\sigma$ of 0.5. We train the models with five different discrete conditions $C_{train} = \{C_1,...,C_5\}$, where $C_i$ represents the mean of an isotropic bivariate Gaussian pdf. We then check how accurately the model (1) reproduces the data distribution $p(x,y|C_{train})$ it was trained on and (2) predict new distributions $p_{model}(x,y|C_{unseen})$, $C_{unseen} := \{C_6,...,C_9\}$ that the model did not see before.

The results of the generative task with CVAE-H are presented in Figure 7 and in Table 2. Note that the cross-entropy $H(p,p_{model})$ is lower-bounded by $H(p)$ as

Generative Task with CVAE-H



(a) Training data ($C_{1:5}$)  (b) Pdfs generated for **seen** $C_{1:5}$  (c) Pdfs generated for **unseen** $C_{6:9}$
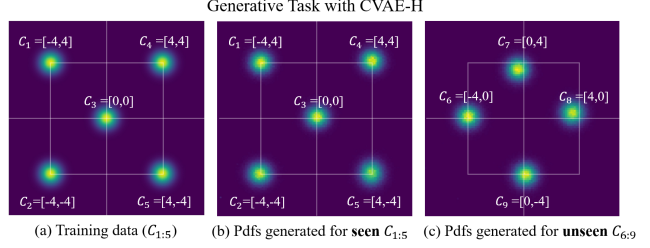
Figure 7: Qualitative results of the Gaussian 2 experiment for the generation task with CVAE-H. (a) data distribution $p(x,y|C_{train})$, b) the data distribution reproduced by CVAE-H $p_{model}(x,y|Z,C_{train})$, and c) CVAE-H's sample generation in *unseen* conditions $p_{model}(x,y|Z,C_{unseen})$.

$H(p,p_{model}) = H(p) + D_{KL}(p||p_{model})$. The differential entropy $H(p)$ of an isotropic bi-variate Gaussian distribution $p(x,y)$ is analytically obtained as $H(p) = 0.5 \cdot ln(2\pi e(\sigma)^2)^2$. The generative cross-entropy for seen conditions (i.e., $H(p,p_{model}(X_{target}|Z,C_{train}))$) are very close to the target differential entropy. On the other side, the generative cross-entropy for unseen conditions is relatively higher. This quantification agrees with the generated pdfs illustrated in Figure 7 that shows CVAE-H's capability to generalize over unseen conditions.

Table 2: Differences between the target and *generative* distributions by CVAE-H in terms of cross-entropy and KL divergence for Figure 7.

|                     | $p(x,y)$ | $p_{model}(x,y|Z,C_i)$ | |
|---------------------|----------|------------------------|------------------------|
| $C$                 | -        | $C_i \in C_{train}$    | $C_i \in C_{unseen}$   |
| $H(p)$              | 1.452    | -                      | -                      |
| $H(p,p_{model})$    | -        | 1.480                  | 2.256                  |
| $D_{KL}(p||p_{model})$ | -     | 0.028                  | 0.804                  |

In summary, the Gaussian experiments empirically verified that CVAE-H satisfies all four requirements of the ideal prediction model; probabilistic, multi-modal, context-driven, and generalizable.

## 6. Experiments: Multi-agent Forecasting

In this section, we demonstrate the performance of CVAE-H in urban driving scenarios using PRECOG-Carla dataset [22], which a publicly available dataset created using an open-source simulator called Carla [44].

### 6.1. The Dataset

PRECOG-Carla dataset includes roughly 76,000 urban driving examples each consists of a 6-second scenario, with

up to five road-agents (i.e., simulated human-driven vehicles) driving in urban areas. The first two seconds (i.e., $t = -2 : 0s$) are used as inputs to predict the future 4-second (i.e., $t = 0 : 4s$) trajectories of the road agents. Lidar point clouds of the environment is provided with three overhead lidar channels (i.e., two above ground and one ground level inputs). The goal is to produce accurate predictions of all road-agents using CVAE-H.

The 76,000 clips in the PRECOG-Carla Town01 dataset are split into a training set (80%), validation set (10%), and test set (10%). CVAE-H is trained on the PRECOG-Carla Town01-train dataset and the progress was validated over Town01-val dataset. Town01-test dataset is used to evaluate the performance of the models. It is important to note that *all qualitative and quantitative results* presented in this paper are based on *the test set*. In other words, the models were tested with *unseen* data.

We apply coordinate transformations to obtain *target-centric* coordinates as we found this helps to improve the accuracy of the multi-agent forecasting. Specifically, the coordinates of all give vehicles in the scene are transformed so that $A_k$, the target vehicle of the prediction, is positioned at $(0,0)$ at zero heading angle at $t = 0$. That is, $X_{t=0}^{A_k} := [0,0]$ and $V_{t=0}^{A_k} \cdot [0,1] = 0$. The other road-agents' positions are transformed using the same transformation matrix. As there exists five road-agents in each example, we create 5 instances with each and every one of the road-agents to be located at the center $(0,0)$ at $t = 0$. That is, $[X_{-2:4s}^{A_{\forall k}}]^{A_m}, m \in 1, 2, 3, 4, 5$, where $[X_{-2:4s}^{A_{\forall k}}]^{A_m}$ indicates that $X_{-2:4s}^{A_{\forall k}}$ is transformed in the $A_m$-centric coordinates.

## 6.2. Hyper-network Design

The hypernetwork of CVAE-H takes the large condition $C$, processes it, and encodes it into the network parameters for the VAE. Unlike the Gaussian experiments that work with low-dimensional $C$, the condition for the self-driving prediction task is much more high-dimensional. $C$ is defined as $C := [\Omega, X_{-2:0s}^{A_{\forall k}}]$, where $\Omega \in \mathbb{R}^{3 \times 200 \times 200}$ indicates the 200 by 200 3-channel lidar scans of the environment and $X_{-2:0s}^{A_{\forall k}} \in \mathbb{R}^{5 \times 10 \times 2}$ represents the labelled trajectories of 5 road agents over 10 time-steps. In this regard, a dedicated hypernetwork architecture must be used to effectively extract the spatial and social features of the inputs. Through changes in the hyper-network, CVAE-H can scale up to tackle the high-dimensional multi-agent forecasting problems for autonomous driving.

Figure 8 depicts the customized hypernetwork design for the multi-agent forecasting. The hypernetwork takes perception inputs as the aforementioned condition $C$ and outputs a set of network parameters $\theta = [\theta_{enc}, \theta_{dec}]$ that are network parameters for the encoder $Z = f_{enc}(X; \theta_{enc})$ and decoder $X = f_{dec}(Z; \theta_{dec})$.
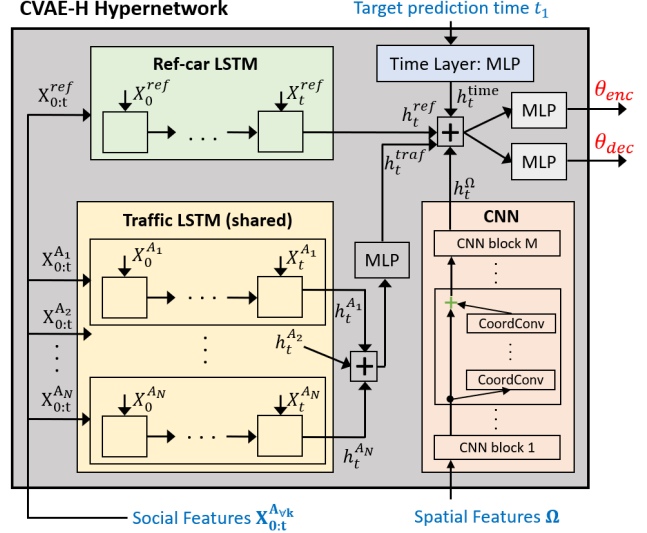


Figure 8: The hypernetwork design for the multi-agent forecasting for autonomous driving. The inputs to the hypernetwork are highlighted in blue and the outputs are colored red.

Our hypernetwork design was inspired by HCNAF [5]. Precisely, we leverage the social, spatial, and time modules of [5]. The difference is that our hypernetwork works with *target-centric* coordinates $[X_{-2:4s}^{A_{\forall k}}]^{A_m}$ whereas HCNAF works with *AV-centric* coordinate as HCNAF performs an occupancy forecasting that is independently obtained from the road agents. Another difference is that our hypernetwork models network parameters for VAE $\theta_{enc}$ and $\theta_{dec}$, as opposed to the parameters of normalizing flow of HCNAF.

The details of the hypernetwork design is as follows. As aforementioned, $C$ is formed with spatial features $\Omega$ (i.e., lidar scan) and social features $X_{-2:0s}^{A_{\forall k}}$ (i.e., labelled trajectories of the road agents). Technically, the spatial inputs can come from any sensor outputs (e.g., lidar, camera, or map), but we use lidar data only since that is what the PRECOG-Carla dataset provides. The hyper-network consists of three components: (1) a social module, (2) a spatial module, and 3) a time module. The outputs of the three modules $h^S, h^\Omega, h^T$ are concatenated and fed into two separate 3-layer MLPs with Relu activations, each dedicated to produce network parameters for encoder and decoder, as shown in Figure 8.

The social module is designed using Long Short Term Memory (LSTM) [45]. Each LSTM module takes $X_{-2:0s}^{A_k} := [x_t^{A_k}, y_t^{A_k}]$, the historical states of a road agent in the scene and encodes temporal dependencies and trends among the state parameters. We use two 2-layer LSTM modules; one for encoding the hidden states of road-agents except the target vehicle ($h^{traffic} \in \mathbb{R}^{40}$) and the other for encoding the hidden state of the target vehicle ($h^{ref} \in \mathbb{R}^{24}$). This is be-

cause the target is positioned at (0,0) and others somewhere else. The resulting output latent vector is the concatenation of the two; $h^{social} = [h^{ref} \oplus h^{traffic}]$.

The spatial module takes in the processed spatial data, which is denoted as $\Omega$, from external perception modules. Our spatial module is based on convolutional neural networks that use residual connections. Recall that the backbone network works with Cartesian (x,y) coordinates space and pixel (image) space. To strengthen the association between the 1D numerical coordinates and 2D image data, we use the coordinate convolution layers [46] as opposed to regular convolution layers. Overall, the spatial module consists of 4 blocks where each block consists of 5 coordconv layers with residual connections, max-pooling layers, and batch-normalization layers. The output latent variable for spatial feature is $h^\Omega \in \mathbb{R}^{96}$.

Lastly, the time module encodes the temporal information into a latent vector $h^{time} \in \mathbb{R}^{10}$ and consists of a 3-layer MLP. The input to the module is the target forecasting time $t_1 \in \mathbb{R}^1$.

### 6.3. Qualitative Results

The following two subsection share the results of the trajectory forecasting experiment. We divide the evaluation into qualitative (visualizations) and quantitative evaluations (accuracy metrics).

In this subsection, we present the qualitative results in Figure 9 and 10. In Figure 9, we show how CVAE-H outputs are contextual to spatial features of environments. Figure 10 illustrates how CVAE-H outputs are contextual to social features.

In each figure, two examples from the PRECOG-Carla Town01 test set are depicted. Each example includes the inputs to the model (i.e., lidar scans of the environment and the 2-second historical states of the road agents), the 4-second label future trajectories, and the outputs of the model (i.e., CVAE-H predictions). We produce 100 sample predictions per car so that the diversity and multi-modality of the outputs are clearly displayed. The estimated road boundaries are overlaid using black lines.

Figure 9 visualizes that CVAE-H predictions are contextual to the spatial features of the environments. The top figure portrays a 4-way intersection environment where the orange car is entering the intersection, and the green and yellow cars are cruising. A vehicle that approaches a 4-way intersection typically has 3 distinct options (or modes): leftturn, right-turn, or pass-through. The figure suggests that CVAE-H captured all three modes. For green and yellow cars, CVAE-H predicted them to cruise, as reflected in the predictions that progressed towards the direction of travels.

The bottom plot of Figure 9 depicts a curved-road environment with five cars in the scene. The orange car is driving in the corner of the road, and the green, yellow, and pur-

◊ past ($X^{A_k}_{-2:0s}$), □ future ($X^{A_k}_{0:4s}$), ○ pred ($\hat{X}^{A_k}_{0:4s}$)

3 modes predicted: left, right turn, or straight
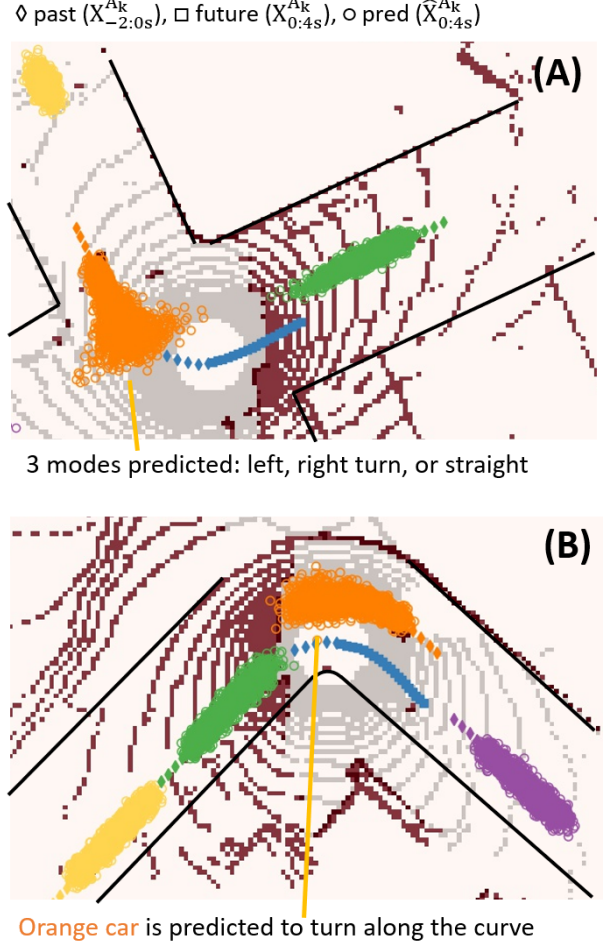
Orange car is predicted to turn along the curve

Figure 9: The first set of examples qualitatively shows that CVAE-H predictions are contextual to *spatial* features of the environments (i.e., road topologies) and *multi-modal*.

ple cars are either approaching or leaving the curved road. The output prediction for the orange car shows that CVAE-H accounts for the spatial constraint of the curve and makes accurate predictions along the curvature of the road.

CVAE-H considers the social contexts of the environment as displayed in Figure 10. First, let us pay attention to the orange and green cars in the figure. In the top plot, the orange car is the first in the queue, whereas the orange car is not in the bottom plot. CVAE-H understands this social context and accordingly predicts that the orange car in the top environment may move in the x-direction, as shown in the elongated uncertainty in the x-direction, whereas the orange car in the bottom is most likely to stay idle. Indeed, this is how a queue is cleared. Secondly, the green car in both top and bottom environments are located last in the queue; CVAE-H understands that there are two other vehicles in front and accordingly predicts that the green car would not move. Lastly, CVAE-H captures the multi-modality of the

◊ past ($X^{A_k}_{-2:0s}$), □ future ($X^{A_k}_{0:4s}$), ○ pred ($\hat{X}^{A_k}_{0:4s}$)



(A)

Green car is predicted to hold position
Orange car (first in the queue) to possibly move

(B)

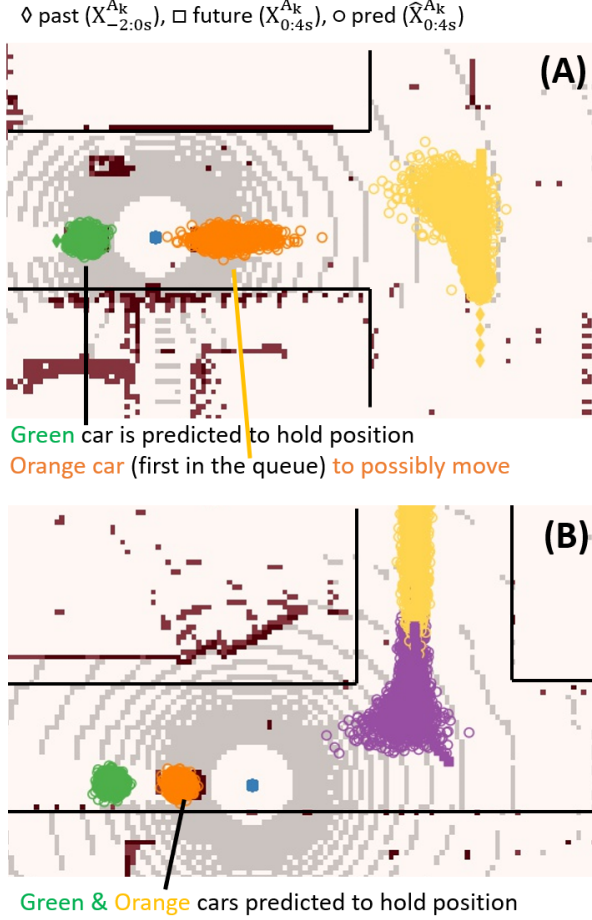Green & Orange cars predicted to hold position

Figure 10: Qualitative demonstrations of *socially contextual* and *multi-modal* CVAE-H outputs.

yellow car in the top environment and the purple car in the bottom environment.

We emphasize that CVAE-H not only outputs predictions, but also the probability associated with them. Although the certainty of the prediction is not used in this work as it is out of scope, the access to the probability is essential for planning tasks. The prediction probability is obtained through the final output layer of CVAE-H, which is a Gaussian mixture. In this regard, CVAE-H is *probabilitic*. As aforementioned, all qualitative examples presented in this section were produced using a *single* CVAE-H model. This empirically shows that CVAE-H is *general*. In this sense, CVAE-H possesses all four attributes of the ideal prediction model.

### 6.4. Quantitative Results

For quantitative evaluations, we use the minimum mean squared deviation of sampled trajectories from the ground-truth (*minMSD*). MinMSD is a measure of diversity or multi-modality of the model outputs. It is often used for

evaluating self-driving prediction models, in particular, generative prediction models that construct probabilistic models with sampling capability [4, 16, 22, 23]. The definition of MinMSD is as follows.

$$\hat{m}_K = min_{k \in 1..K}[||X^k - X^{gt}||^2 / T], \quad (4)$$

where $K$, $X^k$, $X^{gt}$, and $T$ corresponds to the number of trajectory samples, k-th sample trajectory, ground-truth trajectory, and the prediction horizon. While there exist similar metrics like minimum average displacement error (minADE) and minimum final displace error (minFDE) [16, 23, 47], we report minMSD to compare our models against the available benchmarks.

Table 3 presents the evaluation results against available benchmarks. Our model outperforms most of these prediction models [4, 6, 10, 22, 23], except MFP [16]. Although MFP achieved the best performance, MFP is a deterministic model that outputs a fixed number of deterministic modes (i.e., trajectories), compared to other methods including ours that outputs a non-deterministic trajectories.

Table 3: Quantitative evaluation of the proposed prediction models using MinMSD with K=12 on PRECOG-Carla town01 testset, averaged over all 5 cars. Lower is better.

| Method | minMSD($m^2$), K=12 | |
|---|---|---|
| GAN-based | SocialGAN [10] | 1.464 |
| Flow-based | R2P2 [4] | 0.843 |
| | PRECOG-ESP [22] | 0.716 |
| GMM-based | MultiPath [23] | 0.680 |
| AE-based (deterministic) | MFP [16] | **0.279** |
| VAE-based | DESIRE [6] | 2.599 |
| | CVAE-H (ours) | **0.312** |

While Table 3 presented minMSD averaged over all 5 cars, Table 4 shares the detailed minMSD numbers per road agent in the scene against available benchmarks.

Table 4: The detailed minMSD results for individual agents.

| Method | Car 1 | Car 2 | Car 3 | Car 4 | Car 5 |
|---|---|---|---|---|---|
| DESIRE [6] | 2.621 | 2.422 | 2.710 | 2.969 | 2.391 |
| PRECOG-ESP [22] | 0.340 | 0.759 | 0.809 | 0.851 | 0.828 |
| CVAE-H (ours) | **0.151** | **0.346** | **0.304** | **0.355** | **0.404** |

# 7. Conclusion

We present CVAE-H, a novel variational autoencoder conditionalized using a hypernetwork, and demonstrate that CVAE-H is an ideal model for the high-dimensional multi-agent forecasting problem of self-driving. Unlike the traditional embedding approach to make a VAE conditional, CVAE-H offers an unique way of conditioning a VAE using the hypernetwork architecture. This facilitates the informational flow of the conditions to the VAE network across all layers and thus allows the VAE to scale up to large conditional inputs of the prediction problem. We first evaluated CVAE-H in simple generative tasks where we demonstrated the data generation and generalization capabilities of the proposed model. Then, we scaled up the hypernetwork to work with a high-dimensional trajectory forecasting task for autonomous driving. We showed that CVAE-H is probabilistic, multi-modal, context-driven, general, and that it achieves a comparable performance to the state-of-the-art prediction model.

# References

[1] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2446–2454, 2020. 1

[2] Christopher M Bishop. Mixture density networks. 1994. 2, 5

[3] Geunseob Oh and Huei Peng. Impact of traffic lights on trajectory forecasting of human-driven vehicles near signalized intersections. *arXiv preprint arXiv:1906.00486*, 2019. 2, 3

[4] Nicholas Rhinehart, Kris M Kitani, and Paul Vernaza. R2p2: A reparameterized pushforward policy for diverse, precise generative path forecasting. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 772–788, 2018. 2, 10

[5] Geunseob Oh and Jean-Sebastien Valois. Hcnaf: Hyper-conditioned neural autoregressive flow and its application for probabilistic occupancy map forecasting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14550–14559, 2020. 2, 4, 5, 6, 7, 8

[6] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 336–345, 2017. 2, 3, 10

[7] Sergio Casas, Wenjie Luo, and Raquel Urtasun. Intentnet: Learning to predict intention from raw sensor data. In Aude Billard, Anca Dragan, Jan Peters, and Jun Morimoto, editors, *Proceedings of The 2nd Conference on Robot Learning*, volume 87 of *Proceedings of Machine Learning Research*, pages 947–956. PMLR, 29–31 Oct 2018. 2

[8] Wenjie Luo, Bin Yang, and Raquel Urtasun. Fast and furious: Real time end-to-end 3d detection, tracking and motion forecasting with a single convolutional net. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2

[9] Alexandre Alahi, Kratarth Goel, Vignesh Ramanathan, Alexandre Robicquet, Li Fei-Fei, and Silvio Savarese. Social lstm: Human trajectory prediction in crowded spaces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 961–971, 2016. 2

[10] Agrim Gupta, Justin Johnson, Li Fei-Fei, Silvio Savarese, and Alexandre Alahi. Social gan: Socially acceptable trajectories with generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2255–2264, 2018. 2, 10

[11] Nauman Sohani, Geunseob GS Oh, and Xinpeng Wang. A data-driven, falsification-based model of human driver behavior. In *2020 American Control Conference (ACC)*, pages 3892–3899. IEEE, 2020. 2

[12] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 2, 3

[13] David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016. 2, 3

[14] Yuexin Ma, Xinge Zhu, Sibo Zhang, Ruigang Yang, Wenping Wang, and Dinesh Manocha. Trafficpredict: Trajectory prediction for heterogeneous traffic-agents. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6120–6127, 2019. 2, 3

[15] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1349–1358, 2019. 2

[16] Yichuan Charlie Tang and Ruslan Salakhutdinov. Multiple futures prediction. *arXiv preprint arXiv:1911.00997*, 2019. 2, 10

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

[18] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. In *Advances in neural information processing systems*, pages 3483–3491, 2015. 2

[19] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 2

[20] Ian Goodfellow. Nips 2016 tutorial: Generative adversarial networks. *arXiv preprint arXiv:1701.00160*, 2016. 2

[21] Boris Ivanovic, Karen Leung, Edward Schmerling, and Marco Pavone. Multimodal deep generative models for trajectory prediction: A conditional variational autoencoder approach. *IEEE Robotics and Automation Letters*, 6(2):295–302, 2020. 2, 3

[22] Nicholas Rhinehart, Rowan McAllister, Kris Kitani, and Sergey Levine. Precog: Prediction conditioned on goals in visual multi-agent settings. *arXiv preprint arXiv:1905.01296*, 2019. 2, 7, 10

[23] Yuning Chai, Benjamin Sapp, Mayank Bansal, and Dragomir Anguelov. Multipath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction. *arXiv preprint arXiv:1910.05449*, 2019. 3, 10

[24] Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural networks. *science*, 313(5786):504–507, 2006. 3

[25] Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An introduction to variational methods for graphical models. *Machine learning*, 37(2):183–233, 1999. 3, 5

[26] Junyoung Chung, Kyle Kastner, Laurent Dinh, Kratarth Goel, Aaron C Courville, and Yoshua Bengio. A recurrent latent variable model for sequential data. *Advances in neural information processing systems*, 28:2980–2988, 2015. 3

[27] Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improved variational inference with inverse autoregressive flow. In *Advances in neural information processing systems*, pages 4743–4751, 2016. 3, 5

[28] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. 2016. 3

[29] Tomer Galanti and Lior Wolf. On the modularity of hypernetworks. *Advances in Neural Information Processing Systems*, 33, 2020. 4

[30] Luca Bertinetto, João F Henriques, Jack Valmadre, Philip Torr, and Andrea Vedaldi. Learning feed-forward one-shot learners. *Advances in neural information processing systems*, 29:523–531, 2016. 4

[31] Joseph Suarez. Language modeling with recurrent highway hypernetworks. *Advances in neural information processing systems*, 30:3267–3276, 2017. 4

[32] Andrew Brock, Theodore Lim, James Millar Ritchie, and Nicholas J Weston. Smash: One-shot model architecture search through hypernetworks. In *6th International Conference on Learning Representations*, 2018. 4

[33] Chris Zhang, Mengye Ren, and Raquel Urtasun. Graph hypernetworks for neural architecture search. In *International Conference on Learning Representations*, 2018. 4

[34] Johannes von Oswald, Christian Henning, João Sacramento, and Benjamin F Grewe. Continual learning with hypernetworks. In *International Conference on Learning Representations*, 2019. 4

[35] Jonathan Lorraine and David Duvenaud. Stochastic hyperparameter optimization through hypernetworks. *arXiv preprint arXiv:1802.09419*, 2018. 4

[36] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *European conference on computer vision*, pages 630–645. Springer, 2016. 4

[37] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 4, 5

[38] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4

[39] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 4

[40] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. 4

[41] Aaron van den Oord, Oriol Vinyals, et al. Neural discrete representation learning. In *Advances in Neural Information Processing Systems*, pages 6306–6315, 2017. 5

[42] George Papamakarios, Theo Pavlakou, and Iain Murray. Masked autoregressive flow for density estimation. In *Advances in Neural Information Processing Systems*, pages 2338–2347, 2017. 5

[43] Chin-Wei Huang, David Krueger, Alexandre Lacoste, and Aaron Courville. Neural autoregressive flows. In *International Conference on Machine Learning*, pages 2078–2087, 2018. 5, 6

[44] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. Carla: An open urban driving simulator. *arXiv preprint arXiv:1711.03938*, 2017. 7

[45] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997. 8

[46] Rosanne Liu, Joel Lehman, Piero Molino, Felipe Petroski Such, Eric Frank, Alex Sergeev, and Jason Yosinski. An intriguing failing of convolutional neural networks and the coordconv solution. In *Advances in Neural Information Processing Systems*, pages 9605–9616, 2018. 9

[47] Tung Phan-Minh, Elena Corina Grigore, Freddy A Boulton, Oscar Beijbom, and Eric M Wolff. Covernet: Multimodal behavior prediction using trajectory sets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14074–14083, 2020. 10

# A. More Qualitative Results of the Multi-agent Forecasting

In this section, we present more qualitative results of the multi-agent forecasting. All results are based on the Town01 test set of the PRECOG-Carla dataset.

Figure S1 illustrates 3 examples from the test set. The first example (A) depicts a three-way environment where the purple and yellow cars are entering the intersection, the green car is making a left-turn, and the orange car is waiting. For the purple car, CVAE-H predicted it to have two modes; either turning left or right. For the green car, the proposed prediction model predicted it to finish a left-turn with the correct turning curvature, which suggests that CVAE-H took account for the spatial features of the environment. CVAE-H accurately predicted the orange car to hold the position and the yellow car to cruise.

The second example (B) describes a 4-way intersection environment where the green and orange cars are entering the intersection, the purple car leaving the intersection, and the yellow car is waiting. Among the four road agents, the orange car is the most multi-modal as it could possibly turn left, right, or pass through the intersection. The predictions from CVAE-H successfully captured all three modes. CVAE-H predicted the green and purple cars to cruise. The yellow car is mostly predicted to hold the position, but a small possibility to start moving is reflected as the uncertainties in the x-axis.

The third example (C) is about another 3-way intersection environment where the green is entering the intersection, the purple car is on a left-turn, and the orange and yellow cars are cruising. For the green car which is about to enter the intersection, CVAE-H predicted that it could make either a left or right turn. For the purple car which is already on a left-turn, it was predicted to finish the turn toward the end of the prediction horizon. CVAE-H predicted the yellow and orange cars to continue to cruise.

Figure S2 illustrates another 3 examples from the test set. Example (D) describes a 4-way intersection environment where the purple car is entering the intersection, the green car is in the middle of the intersection-passing, the orange car is leaving the intersection, and the yellow is waiting for its turn. The purple car could possibly either turn left, right, or pass through the intersection. CVAE-H captured all three modes successfully. For the green car, CVAE-H predicted it to continue to pass through the intersection, as it was already on its way through the intersection. The proposed prediction model predicted the orange car to cruise and the yellow car to hold its position.

Example (E) illustrates a 4-way intersection where the orange and purple cars are approaching the intersection, and the green car is entering the intersection. As the green car just started to make a right-turn, the prediction weights more towards the right-turn and pass-through modes, but still captured a small possibility of the left-turn mode. For the orange and purple cars, CVAE-H predicted them to continue to cruise.

Example (F) depicts another 4-way intersection where the green and yellow cars are approaching the intersection, and the orange car is cruising behind the blue reference car. The green car is predicted to have three modes; straight, left-turn, and right-turn. The orange and yellow cars are predicted to cruise.
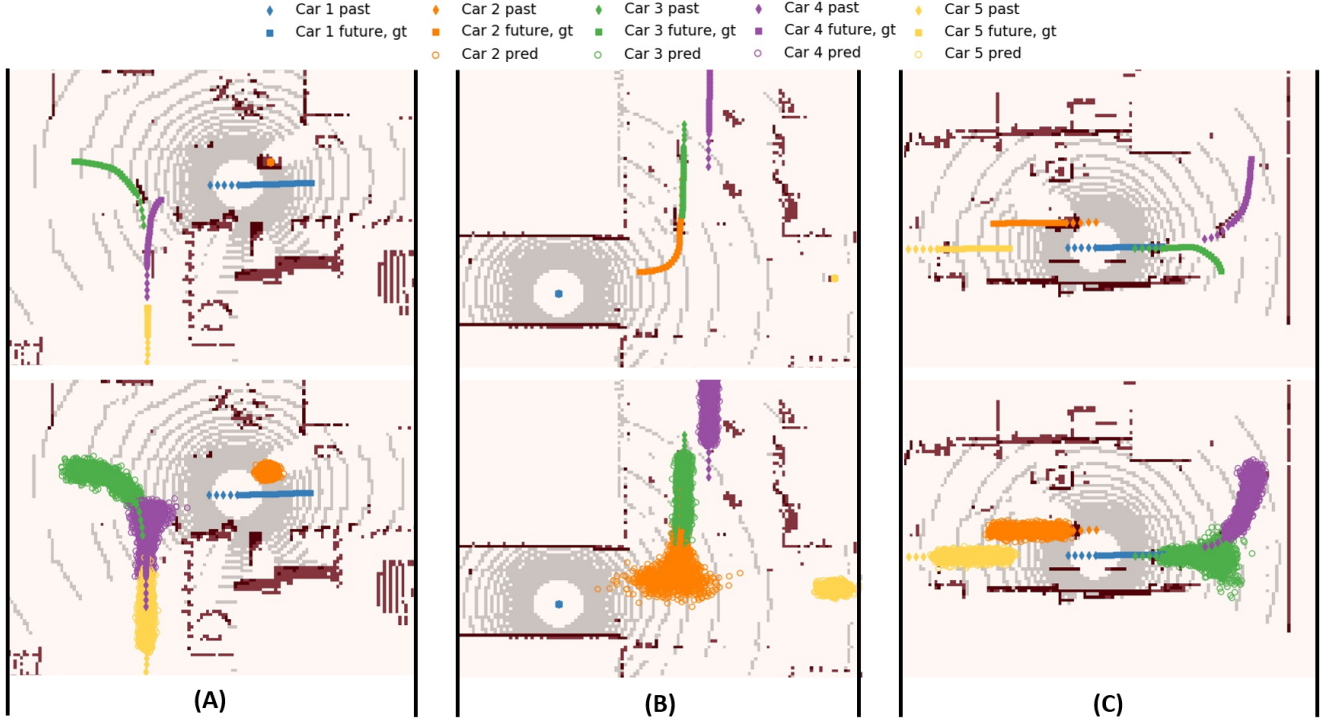
Figure S1: More qualitative examples. Each example consists of a top and bottom plot. A top plot depicts the input to the prediction model as well as the ground-truth future positions of the road agents. The input consists of a lidar scan of the environment $\Omega$ and 2-second history $X_{-2:0s}^{A_{\forall k}}$ (depicted as $\diamond$). The 4-second ground-truth future positions $X_{0:4s}^{A_{\forall k}}$ are depicted using $\square$. A bottom plot depicts $\hat{X}_{0:4s}^{A_{\forall k}}$, the outputs of the prediction model overlaid over future time-steps, $t = 0 : 4s$. In other words, a bottom plot describes how the predictions progress over time.
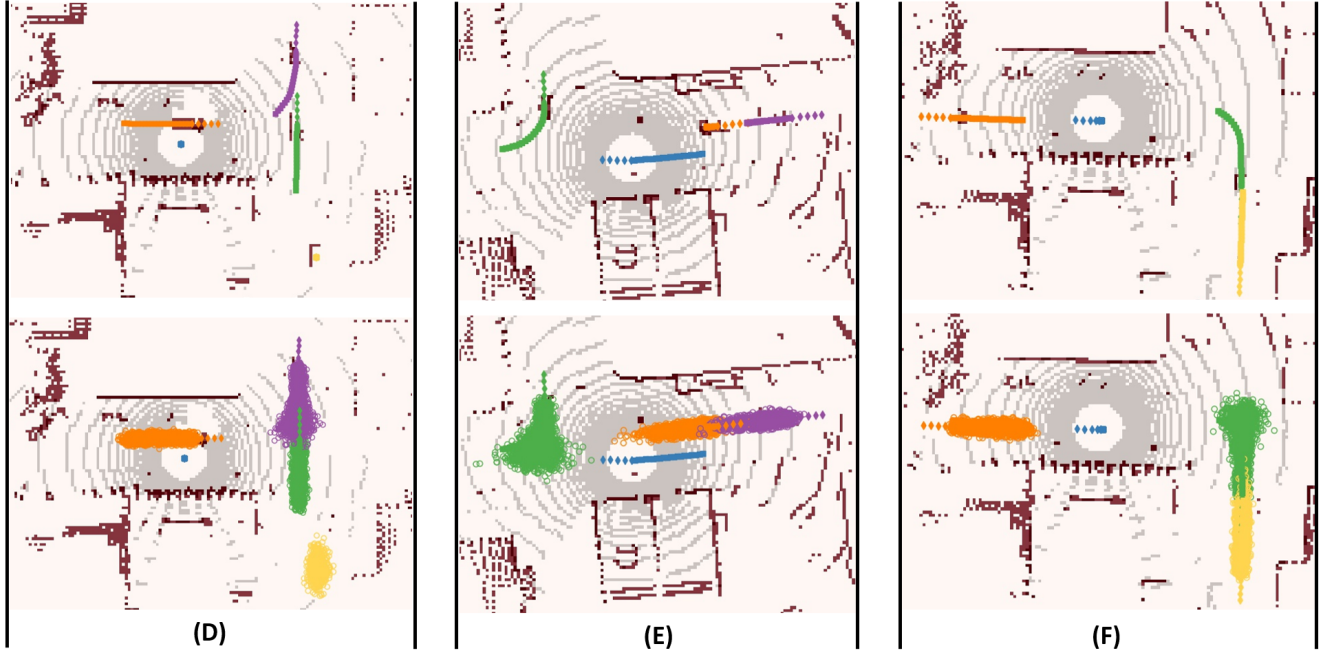


Figure S2: Another 3 examples (D-F) of the multi-agent forecasting with CVAE-H. The explanations about each example are provided in the body of the appendix.