

# Accurate and Diverse Sampling of Sequences based on a “Best of Many” Sample Objective

Apratim Bhattacharyya, Bernt Schiele, Mario Fritz

Max Planck Institute for Informatics, Saarland Informatics Campus, Saarbrücken, Germany  
 {abhattach, schiele, mfritz}@mpi-inf.mpg.de

## Abstract

For autonomous agents to successfully operate in the real world, anticipation of future events and states of their environment is a key competence. This problem has been formalized as a sequence extrapolation problem, where a number of observations are used to predict the sequence into the future. Real-world scenarios demand a model of uncertainty of such predictions, as predictions become increasingly uncertain – in particular on long time horizons. While impressive results have been shown on point estimates, scenarios that induce multi-modal distributions over future sequences remain challenging. Our work addresses these challenges in a Gaussian Latent Variable model for sequence prediction. Our core contribution is a “Best of Many” sample objective that leads to more accurate and more diverse predictions that better capture the true variations in real-world sequence data. Beyond our analysis of improved model fit, our models also empirically outperform prior work on three diverse tasks ranging from traffic scenes to weather data.

## 1. Introduction

Predicting the future is important in many scenarios ranging from autonomous driving to precipitation forecasting. Many of these tasks can be formulated as sequence prediction problems. Given a past sequence of events, probable future outcomes are to be predicted.

Recurrent Neural Networks (RNN) especially LSTM formulations are state-of-the-art models for sequence prediction tasks [2, 23, 6, 22]. These approaches predict only point estimates. However, many sequence prediction problems are only partially observed or stochastic in nature and hence the distribution of future sequences can be highly multi-modal. Consider the task of predicting future pedestrian trajectories. In many cases, we do not have any information about the intentions of the pedestrians in the scene. A pedestrian after walking over a Zerba crossing might decide to turn either left or right. A point estimate in such a situation would be highly

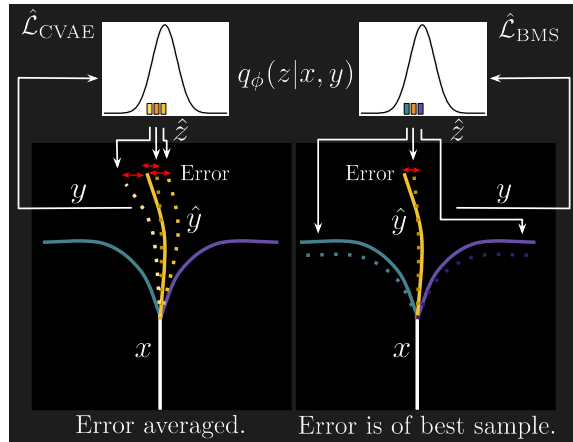


Figure 1: Comparison between our “Best of Many” sample objective and the standard CVAE objective.

unrealistic. Therefore, in order to incorporate uncertainty of future outcomes, we are interested in *structured predictions*. Structured prediction in this context implies learning a one to many mapping of a given fixed sequence to plausible future sequences [19]. This leads to more realistic predictions and enables probabilistic inference.

Recent work [14] has proposed deep conditional generative models with Gaussian latent variables for structured sequence prediction. The Conditional Variational Auto-Encoder (CVAE) framework [19] is used in [14] for learning of the Gaussian Latent Variables. We identify two key limitations of this CVAE framework. First, the currently used objectives hinder learning of diverse samples due to a marginalization over multi-modal futures. Second, a mismatch in latent variable distribution between training and testing leads to errors in model fitting. We overcome both challenges which results in more accurate and diverse samples – better capturing the true variations in data. Our main contributions are: 1. We propose a novel “best of many” sample objective; 2. We analyze the benefits of our ‘best of many’ sample objective analytically as well as show an improved fit of latent variables on models trained with this novel objec-

tive compared to prior approaches; 3. We also show for the first time that this modeling paradigm extends to full-frame images sequences with diverse multi-modal futures; 4. We demonstrate improved accuracy as well as diversity of the generated samples on three diverse tasks: MNIST stroke completion, Stanford Drone Dataset and HKO weather data. On all three datasets we consistently outperform the state of the art and baselines.

## 2. Related Work

**Structured Output Prediction.** Stochastic feed-forward neural networks (SFNN) [21] model multi-modal conditional distributions through binary stochastic hidden variables. During training multiple samples are drawn and weighted according to importance-weights. However, due to the latent variables being binary SFNNs are hard to train on large datasets. There has been several efforts to make training more efficient for binary latent variables [16, 8, 15, 13]. However, not all tasks can be efficiently modelled with binary hidden variables. In [19], Gaussian hidden variables are considered where the re-parameterization trick can be used for learning on large datasets using stochastic optimization. Inspired by this technique we model Gaussian hidden variables for structured sequence prediction tasks.

**Variational Autoencoders.** Variational learning has enabled learning of deep directed graphical models with Gaussian latent variables on large datasets [11, 10, 9]. Model training is made possible through stochastic optimization by the use of a variational lower bound of the data log-likelihood and the re-parameterization trick. In [3] a tighter lower bound on the data log-likelihood is introduced and multiple samples are used during training which are weighted according to importance weights. They show empirically that their IWAE framework can learn richer latent space representations. However, these models do not consider conditional distributions for structured output prediction. Conditional variational auto-encoders (CVAE) [19] extend the VAE framework of [11] to model conditional distributions for structured output prediction by introducing the CVAE objective which maximizes a lower bound on the conditional data log likelihood. The CVAE framework has been used for a variety of tasks. Examples include, generation of likely future frames given a single frame of a video [24], diverse images of clothed people conditioned on their silhouette [12], and trajectories of basketball players using pictorial representations [1]. However, the gap between the training and test latent variable distributions cannot be fully closed by the CVAE objective function. We consider a new multi-sample objective which relaxes the constraints on the recognition network by encouraging diverse sample generation and thus leads to a better match between the training and test latent variable distributions.

**Recurrent Neural Networks.** Recurrent Neural Networks

(RNNs) are state of the art methods for variety of sequence learning tasks [7, 20]. In this work, we focus on sequence to sequence regression tasks, in particular, trajectory prediction and image sequence prediction. RNNs have been used for pedestrian trajectory prediction. In [2], trajectories of multiple people in a scene are jointly modelled in a social context. However, even though the distribution of pedestrian trajectories are highly multimodal (with diverse futures), only one mean estimate is modelled. [14] jointly models multiple future pedestrian trajectories using a recurrent CVAE sampling module. Samples generated are refined and ranked using image and social context features. While our trajectory prediction model is similar to the sampling module of [14], we focus on improving the sampling module by our novel multi-sample objective function. Convolutional RNNs [22] have been used for image sequence prediction. Examples include, robotic arm movement prediction [6] and precipitation now-casting [22, 18]. In this work, we extend the model of [22] for structured sequence prediction by conditioning predictions on Gaussian latent variables. Furthermore, we show that optimization using our novel multi-sample objective leads to improved results over the standard CVAE objective.

## 3. Structured Sequence Prediction with Gaussian Latent Variables

We begin with an overview of deep conditional generative models with gaussian latent variables and the CVAE framework with the corresponding objective [19] used for training. Then, we introduce our novel “best-of-many” samples objective function. Thereafter, we introduce the conditional generative models which serve as the test bed for our novel objective. We first describe our model for structured trajectory prediction which is similar to the sampling module of [14] and consider extensions which additionally conditions on visual input and generates full image sequences.

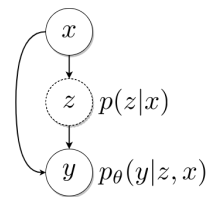


Figure 2: Conditional generative models.

We consider deep conditional generative models of the form shown in Figure 2. Given an input sequence  $x$ , a latent variable  $\hat{z}$  is drawn from the conditional distribution  $p(z|x)$  (assumed Gaussian). The output sequence  $\hat{y}$  is then sampled from the distribution  $p_\theta(y|x, z)$  of our conditional generative model with parameterized by  $\theta$ . The latent variables  $z$  enables one-to-many mapping and the learning of multiple

modes of the true posterior distribution  $p(y|x)$ . In practice, the simplifying assumption is made that  $z$  is independent of  $x$  and  $p(z|x)$  is  $\mathcal{N}(0, I)$ . Next, we discuss the training of such models.

### 3.1. Conditional Variational Auto-encoder Based Training Objective

We would like to maximize the data log-likelihood  $p_\theta(y | x)$ . To estimate the data log-likelihood of our model  $p_\theta$ , one possibility is to perform Monte-Carlo sampling of the latent variable  $z$ . For  $T$  samples, this leads to the following estimate,

$$\hat{\mathcal{L}}_{\text{MC}} = \log \left( \frac{1}{T} \sum_{i=1}^T p_\theta(y|\hat{z}_i, x) \right), \quad \hat{z}_i \sim \mathcal{N}(0, I). \quad (1)$$

This estimate is unbiased but has high variance [15]. We would underestimate the log-likelihood for some samples and overestimate for others, especially if  $T$  is small. This would in turn lead to high variance weight updates.

We can reduce the variance of updates by estimating the log-likelihood through importance sampling during training. As described in [19], we can sample the latent variables  $z$  from a recognition network  $q_\phi$  using the re-parameterization trick [11]. The data log-likelihood is,

$$\begin{aligned} \log(p_\theta(y | x)) = \\ \log \left( \int p_\theta(y|z, x) \frac{p(z|x)}{q_\phi(z|x, y)} q_\phi(z|x, y) dz \right). \end{aligned} \quad (2)$$

The integral in (2) is computationally intractable. In [19], a variational lower bound of the data log-likelihood (2) is derived, which can be estimated empirically using Monte-Carlo integration (also used in [14]),

$$\begin{aligned} \hat{\mathcal{L}}_{\text{CVAE}} = \frac{1}{T} \sum_{i=1}^T \log p_\theta(y|\hat{z}_i, x) \\ - D_{\text{KL}}(q_\phi(z|x, y) \| p(z|x)), \quad \hat{z}_i \sim q_\phi(z|x, y). \end{aligned} \quad (3)$$

The lower bound in (3) weights all samples ( $\hat{z}_i$ ) equally and so they must all ascribe high probability to the data point  $(x, y)$ . This introduces a strong constraint on the recognition network  $q_\phi$ . Therefore, the model is forced to trade-off between a good estimate of the data log-likelihood and the KL divergence between the training and test latent variable distributions. One possibility to close the gap introduced between the training and test pipelines, as described in [19], is to use an hybrid objective of the form  $(1 - \alpha)\hat{\mathcal{L}}_{\text{MC}} + \alpha\hat{\mathcal{L}}_{\text{CVAE}}$ . Although such an hybrid objective has shown modest improvement in performance in certain cases, we could not observe any significant improvement over the standard CVAE objective in our structured sequence prediction tasks. In the following, we derive our novel ‘‘best-of-many-samples’’ objective which on the one hand encourages sample diversity

and on the other hand aims to close the gap between the training and testing pipelines.

### 3.2. Best of Many Samples Objective

Here, we propose our objective which unlike (3) does not weight each sample equally. Consider the functions  $f_1(z) = p(z|x)/q_\phi(z|x, y)$  and  $f_2(z) = p_\theta(y|z, x) \times q_\phi(z|x, y)$  in (2). We cannot evaluate  $f_2(z)$  directly for Monte-Carlo samples. Notice, however that both  $f_1(z)$  and  $f_2(z)$  are continuous and positive. As  $q_\phi(z|x, y)$  is normally distributed, the integral above can be very well approximated on a large enough bounded interval  $[a, b]$ . Therefore, we can use the First Mean Value Theorem of Integration [4], to separate the functions  $f_1(z)$  and  $f_2(z)$  in (2),

$$\begin{aligned} \log(p_\theta(y|x)) = \log \left( \int_a^b p_\theta(y|z, x) q_\phi(z|x, y) dz \right) \\ + \log \left( \frac{p(z'|x)}{q_\phi(z'|x, y)} \right), \quad z' \in (a, b). \end{aligned} \quad (4)$$

We can lower bound (4) with the minimum of the term on the right,

$$\begin{aligned} \log(p_\theta(y|x)) \geq \log \left( \int_a^b p_\theta(y|z, x) q_\phi(z|x, y) dz \right) \\ + \min_{z' \in (a, b)} \left( \log \left( \frac{p(z'|x)}{q_\phi(z'|x, y)} \right) \right) \end{aligned} \quad (5)$$

We can estimate the first term on the right of (5) using Monte-Carlo integration. The minimum in the second term on the right of (5) is difficult to estimate, therefore we approximate it by the KL divergence over the full distribution. The KL divergence heavily penalizes  $q_\phi(z|x, y)$  when it is high for low values  $p(z|x)$  (which leads to low value of the ratio of the distributions). This leads to the following ‘‘many-sample’’ objective, (more details in the supplementary section),

$$\begin{aligned} \hat{\mathcal{L}}_{\text{MS}} = \log \left( \frac{1}{T} \sum_{i=1}^T p_\theta(y|\hat{z}_i, x) \right) \\ - D_{\text{KL}}(q_\phi(z|x, y) \| p(z|x)), \quad \hat{z}_i \sim q_\phi(z|x, y). \end{aligned} \quad (6)$$

Compared to the CVAE objective (2), the recognition network  $q_\phi$  now has multiple chances to draw samples with high posterior probability ( $p_\theta(y | z, x)$ ). This encourages diversity in the generated samples. Furthermore, the data log-likelihood (2) estimate in this objective is tighter as  $\hat{\mathcal{L}}_{\text{MS}} \geq \hat{\mathcal{L}}_{\text{CVAE}}$  follows from the Jensen’s inequality. Therefore, this bound loosens the constraints on the recognition network  $q_\phi$  and allows it more closely match the latent variable distribution  $p(z|x)$ . However, as we focus on regression tasks, probabilities are of the form  $e^{-\text{MSE}(\hat{y}, y)}$ . Therefore, in practice the Log-Average term can cause numerical instabilities due to limited machine precision in representing

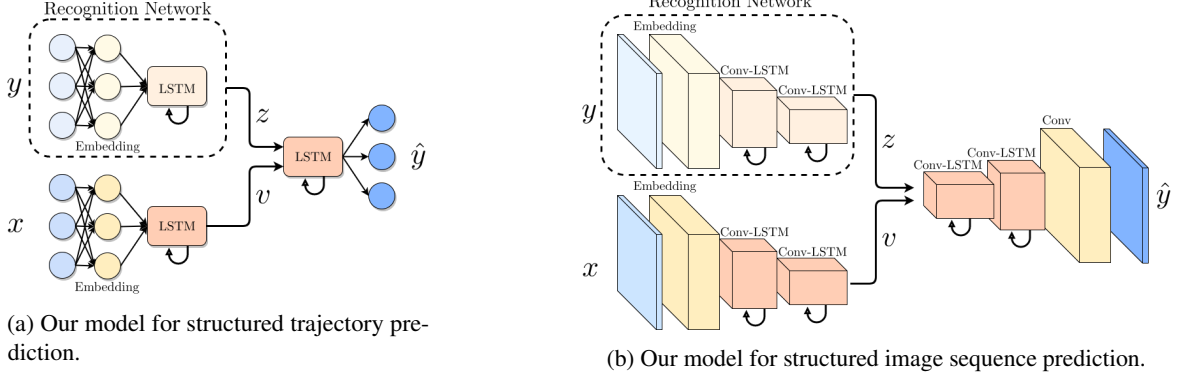


Figure 3: Our model architectures. The recognition networks are only available during training.

the probability  $e^{-\text{MSE}(\hat{y}, y)}$ . Therefore, we use a ‘‘Best of Many Samples’’ approximation  $\hat{\mathcal{L}}_{\text{BMS}}$  of (6). We can pull the constant  $1/T$  term outside the average in (6) and approximate the sum with the maximum,

$$\hat{\mathcal{L}}_{\text{MS}} = \log \left( \sum_{i=1}^T p_{\theta}(y|\hat{z}_i, x) \right) - \log(T) \quad (7)$$

$$- D_{\text{KL}}(q_{\phi}(z|x, y) \parallel p(z|x)), \hat{z}_i \sim q_{\phi}(z|x, y)$$

$$\hat{\mathcal{L}}_{\text{MS}} \geq \hat{\mathcal{L}}_{\text{BMS}} = \max_i (\log(p_{\theta}(y|\hat{z}_i, x))) - \log(T) \quad (8)$$

$$- D_{\text{KL}}(q_{\phi}(z|x, y) \parallel p(z|x)), \hat{z}_i \sim q_{\phi}(z|x, y).$$

Similar to (6), this objective encourages diversity and loosens the constraints on the recognition network  $q_{\phi}$  as only the best sample is considered. During training, initially  $p_{\theta}$  assigns low probability to the data for all samples  $\hat{z}_i$ . The  $\log(T)$  difference between (6) and (8) would be dominated by the low data log-likelihood. Later on, as both objectives promote diversity, the Log-Average term in (6) would be dominated by one term in the average. Therefore, (6) would be well approximated by the maximum of the terms in the average. Furthermore, (8) avoids numerical stability issues.

### 3.3. Model Architectures for Structured Sequence Prediction

We base our model architectures on RNN Encoder-Decoders. We use LSTM formulations as RNNs for structured trajectory prediction tasks (Figure 3a) and Convolutional LSTM formulations (Figure 3b) for structured image sequence prediction tasks. During training, we consider LSTM recognition networks in case of trajectory prediction (Figure 3a) and for image sequence prediction, we consider Conv-LSTM recognition networks (Figure 3b). Note that, as we make the simplifying assumption that  $z$  is independent of  $x$ , the recognition networks are conditioned only on  $y$ .

**Model for Structured Trajectory Prediction.** Our model for structured trajectory prediction (see Figure 3a) is simi-

lar to the sampling module of [14]. The input sequence  $x$  is processed using an embedding layer to extract features and the embedded sequence is read by the encoder LSTM. The encoder LSTM produces a summary vector  $v$ , which is its internal state after reading the input sequence  $x$ . The decoder LSTM is conditioned on the summary vector  $v$  and additionally a sample of the latent variable  $z$ . The decoder LSTM is unrolled in time and a prediction is generated by a linear transformation of its output. Therefore, the predicted sequence at a certain time-step  $\hat{y}^t$  is conditioned on the output at the previous time-step, the summary vector  $v$  and the latent variable  $z$ . As the summary  $v$  is deterministic given  $x$ , we have,

$$\begin{aligned} \log(p_{\theta}(y|x)) &= \sum_t \log(p_{\theta}(y^{t+1}|y^t, v) p(v|x)) \\ &= \sum_t \log(p_{\theta}(y^{t+1}|y^t, x)) \\ &= \int \sum_t \log(p_{\theta}(y^{t+1}|y^t, z, x) p_{\theta}(z|x)) dz. \end{aligned}$$

Conditioning the predicted sequence at all time-steps upon a single sample of  $z$  enables  $z$  to capture global characteristics (e.g. speed and direction of motion) of the future sequence and generation of temporally consistent sample sequences  $\hat{y}$ .

**Extension with Visual Input.** In case of dynamic agents e.g. pedestrians in traffic scenes, the future trajectory is highly dependent upon the environment e.g. layout of the streets. Therefore, additionally conditioning samples on sensory input (e.g. visuals of the environment) would enable more accurate sample generation. We use a CNN to extract a summary of a visual observation of a scene. This visual summary is given as input to the decoder LSTM, ensuring that the generated samples are additionally conditioned on the visual input.

**Model for Structured Image Sequence Prediction.** If the

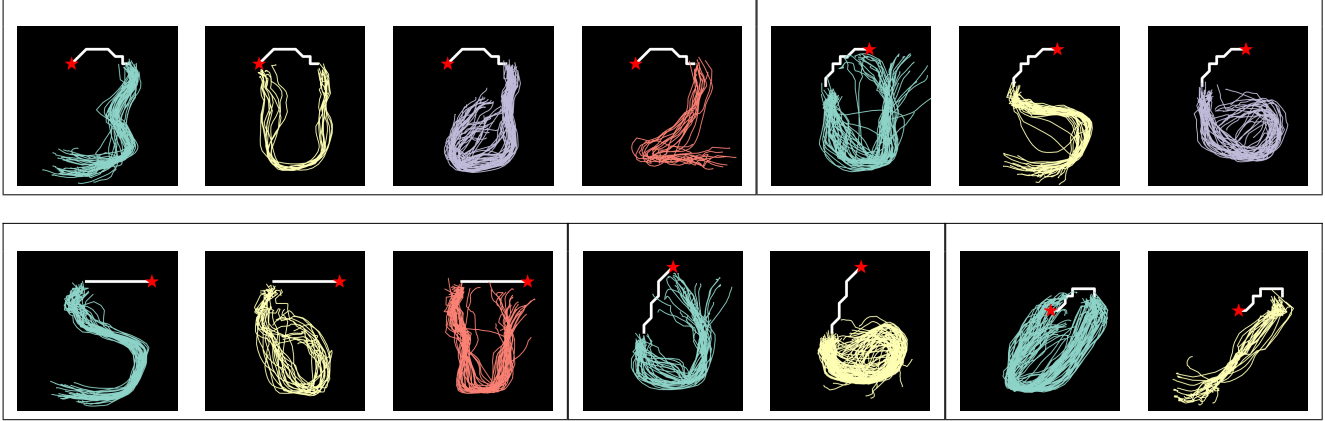


Figure 4: Diverse samples drawn from our LSTM-BMS model trained using the  $\hat{\mathcal{L}}_{\text{BMS}}$  objective, clustered using k-means. The number of clusters is set manually to the number of expected digits based on the initial stroke.

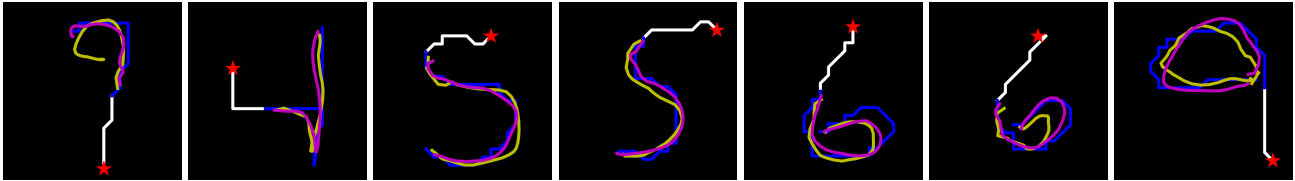


Figure 5: Top 10% of samples drawn from the LSTM-BMS model (magenta) and the LSTM-CVAE model (yellow), with the groundtruth in (blue).

sequence  $(x, y)$  in question consists of images e.g. frames of a video, the trajectory prediction model Figure 3a cannot exploit the spatial structure of the image sequence. More specifically, consider a pixel  $y_{i,j}^{t+1}$  at time-step  $t + 1$  of the image sequence  $y$ . The pixel value at time-step  $t + 1$  depends upon only the pixel  $y_{i,j}^t$  and a certain neighbourhood around it. Furthermore, spatially neighbouring pixels are correlated. This spatial structure can be exploited by using Convolutional LSTMs [22] as RNN encoder-decoders. Conv-LSTMs retain spatial information by considering the hidden states  $h$  and cell states  $c$  as 3D tensors – the cell and hidden states are composed of vectors  $c_{i,j}^t, h_{i,j}^t$  corresponding to each spatial position. New cell states, hidden states and outputs are computed using convolutional operations. Therefore, new cell states  $c_{i,j}^{t+1}$ , hidden states  $h_{i,j}^{t+1}$  depend upon only a local spatial neighbourhood of  $c_{i,j}^t, h_{i,j}^t$ , thus preserving spatial information.

We propose conditional generative models networks with Conv-LSTMs for structured image sequence prediction (Figure 3b). The encoder and decoder consists of two stacked Conv-LSTMs for feature aggregation. As before, the output is conditioned on a latent variable  $z$  to model multiple modes of the conditional distribution  $p(y | x)$ . The future states of neighboring pixels are highly correlated. However, spatially distant parts of the image sequences can evolve independently. To take into account the spatial structure of images,

we consider latent variables  $z$  which are 3D tensors. As detailed in Figure 3b, the input image sequence  $x$  is processed using a convolutional embedding layer. The Conv-LSTM reads the embedded input sequence and produces a 3D tensor  $v$  as the summary. The 3D summary  $v$  and latent variable  $z$  is given as input to the Conv-LSTM decoder at every time-step. The cell state, hidden state or output at a certain spatial position,  $c_{i,j}^t, h_{i,j}^t, y_{i,j}^t$ , it is conditioned on a sub-tensor  $z_{i,j}$  of the latent tensor  $z$ . Spatially neighbouring cell states, hidden states (and thus outputs) are therefore conditioned on spatially neighbouring sub-tensors  $z_{i,j}$ . This coupled with the spatial information preserving property of Conv-LSTMs detailed above, enables  $z$  to capture spatial location specific characteristics of the future image sequence and allows for modeling the correlation of future states of spatially neighboring pixels. This ensures spatial consistency of sampled output sequences  $\hat{y}$ . Furthermore, as in the fully connected case, conditioning the full output sequence sample  $\hat{y}$  is on a single sample of  $z$  ensures temporal consistency.

## 4. Experiments

We evaluate our models both on synthetic and real data. We choose sequence datasets which display multimodality. In particular, we evaluate on key strokes from MNIST sequence data [5] (which can be seen as trajectories in a

Method	CLL
LSTM	136.12
LSTM-MC	102.34
LSTM-CVAE	96.42
LSTM-BMS	<b>95.63</b>

Table 1: Evaluation on the MNIST Sequence dataset.

constrained space), human trajectories from Stanford Drone data [17] and radar echo image sequences from HKO [22]. All models were trained using the ADAM optimizer, with a batch size of 32 for trajectory data and 4 for the radar echo data. All experiments were conducted on a single Nvidia M40 GPU with 12GB memory. For models trained using the  $\hat{\mathcal{L}}_{\text{CVAE}}$  and  $\hat{\mathcal{L}}_{\text{BMS}}$  objectives, we use  $T = \{10, 10, 5\}$  samples during training on the MNIST Sequence, Stanford Drone, and HKO datasets respectively.

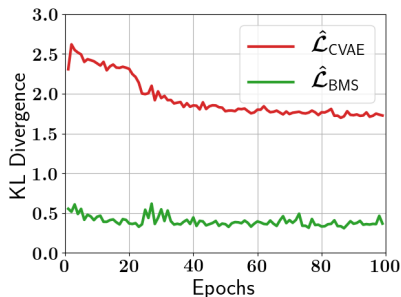


Figure 6: KL Divergence during training on the MNIST Sequence dataset.

#### 4.1. MNIST Sequence

The MNIST sequence dataset consists of pen strokes which closely approximates the skeleton of the digits in the MNIST dataset. We focus on the stroke completion task. Given an initial stroke the distribution of possible completions is highly multimodal. The digits 0, 3, 2 and 8, have the same initial stroke with multiple writing styles for each digit. Similarly for the digits 0 and 6, with multiple writing styles for each digit.

We fix the length of the initial stroke sequence at 10. We use the trajectory prediction model from Figure 3a and train it using the  $\hat{\mathcal{L}}_{\text{BMS}}$  objective (LSTM-BMS). We compare it against the following baselines: 1. A vanilla LSTM encoder-decoder regression model (LSTM) without latent variables; 2. The trajectory prediction model from Figure 3a trained using the  $\hat{\mathcal{L}}_{\text{MC}}$  objective (LSTM-MC); 3. The trajectory prediction model from Figure 3a trained using the  $\hat{\mathcal{L}}_{\text{CVAE}}$  objective (LSTM-CVAE). We use the negative conditional log-likelihood metric (CLL) and report the results in Table 1. We use  $T = 100$  samples to estimate the CLL.

We observe that our LSTM-BMS model achieves the best CLL. This means that our LSTM-BMS model fits the data distribution best. Furthermore, we see that the latent variables sampled from our recognition network  $q_{\phi}(z | x, y)$  during training better matches the true distribution  $p(z | x)$  used during testing. This can be seen through the KL divergence  $D_{\text{KL}}(q_{\phi}(z | x, y) || p(z | x))$  in Figure 6 during training of the recognition network trained with the  $\hat{\mathcal{L}}_{\text{BMS}}$  objective versus that of the  $\hat{\mathcal{L}}_{\text{CVAE}}$  objective. We observe that the KL divergence of the recognition network trained with the  $\hat{\mathcal{L}}_{\text{BMS}}$  to be substantially lower, thus, reducing the mismatch in the latent variable  $z$  between the training and testing pipelines.

We show qualitative examples of generated samples in Figure 4 from the LSTM-BMS model. We show  $T = 100$  samples per test example. The initial conditioning stroke is shown in white. The samples drawn are diverse and clearly multimodal. We cluster the generated samples using k-means for better visualization. The number of clusters is set manually to the number of expected digits based on the initial stroke. In particular, our model generates corresponding to 2, 3, 0 (1st example), 0, 6 (2nd example) and so on.

We compare the accuracy of samples generated by our LSTM-BMS model versus the LSTM-CVAE model in Figure 5. We display mean of the oracle top 10% of samples (closest in euclidean distance w.r.t. groundtruth) generated by both models. Comparing the results we see that, using the  $\hat{\mathcal{L}}_{\text{BMS}}$  objective leads to the generation of more accurate samples.

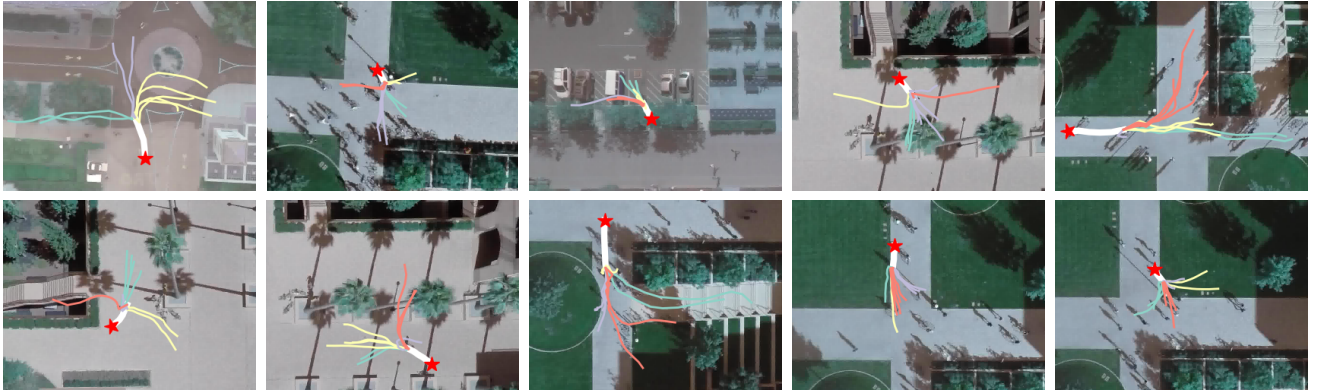
#### 4.2. Stanford Drone

The Stanford Drone dataset consists of overhead videos of traffic scenes. Trajectories of various dynamic agents including Pedestrians and Bikers are annotated. The paths of such agents are determined by various factors including the intention of the agent, paths of other agents and the layout of the scene. Thus, the trajectories are highly multimodal. As in [17, 14], we predict the trajectories of these agents 4.8 seconds into the future conditioned on the past 2.4 seconds. We use the same dataset split as in [14]. We encode trajectories as relative displacement from the initial position. The trajectory at each time-step can be seen as the velocity of the agent.

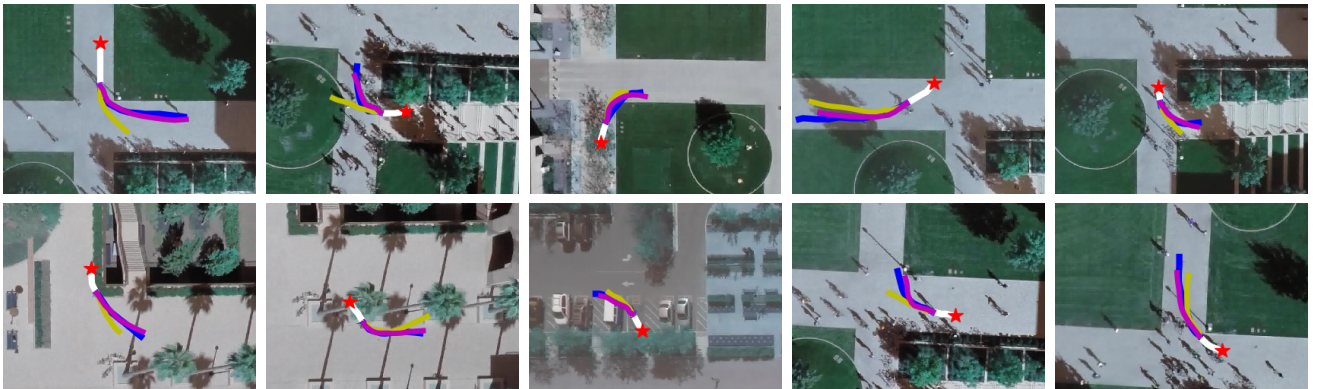
We consider the extension of our trajectory prediction model (Figure 3a) discussed in subsection 3.3 conditioned on the last visual observation from the overhead camera. We use a 6 layer CNN to extract visual features (see supplementary material). We train this model with the  $\hat{\mathcal{L}}_{\text{BMS}}$  objective and compare it to: 1. A vanilla LSTM encoder-decoder regression model with and without visual observation (LSTM); 2. The state of the art DESIRE-SI-IT4 model from [14]; 3. Our extended trajectory prediction model Figure 3a trained using the  $\hat{\mathcal{L}}_{\text{CVAE}}$  objective (LSTM-CVAE).

Method	Visual	Error at 1.0(sec)	Error at 2.0(sec)	Error at 3.0(sec)	Error at 4.0(sec)	CLL
LSTM	x	1.08	2.57	4.70	7.20	134.29
LSTM	RGB	0.84	1.95	3.86	6.24	133.12
DESIRE-SI-IT4 [14]	RGB	1.29	2.35	3.47	5.33	x
LSTM-CVAE	RGB	<b>0.71</b>	1.86	3.39	5.06	127.51
LSTM-BMS	RGB	0.80	<b>1.77</b>	<b>3.10</b>	<b>4.62</b>	<b>126.65</b>

Table 2: Evaluation on the Stanford Drone dataset. Euclidean distance measured at  $(1/5)$  resolution.



(a) Diverse samples drawn from our LSTM-BMS model trained using the  $\hat{\mathcal{L}}_{\text{BMS}}$  objective, color-coded after clustering using k-means with four clusters.



(b) Top 10% of samples drawn from the LSTM-BMS model (magenta) and the LSTM-CVAE model (yellow), with the groundtruth in blue.

Figure 7: Qualitative evaluation on the Stanford Drone dataset.

We report the results in Table 2. We report the CLL metric and the euclidean distance in pixels between the true trajectory and the oracle top 10% of generated samples at 1, 2, 3 and 4 seconds into the future at  $(1/5)$  resolution (as in [14]). Our LSTM-BMS model again performs best both with respect to the euclidean distance and the CLL metric. This again demonstrates that using the  $\hat{\mathcal{L}}_{\text{BMS}}$  objective enables us to better fit the groundtruth data distribution and enables the generation of more accurate samples. The performance advantage with respect to DESIRE-SI-IT4 [14] is due to 1. Conditioning the decoder LSTM in Figure 3a

directly on the visual input at higher  $(1/2)$  versus  $(1/5)$  resolution (as our LSTM-CVAE outperforms DESIRE-SI-IT4), 2. Our  $\hat{\mathcal{L}}_{\text{BMS}}$  objective (as our LSTM-BMS outperforms both DESIRE-SI-IT4 and LSTM-CVAE).

We show qualitative examples of generated samples ( $T = 10$ ) in Figure 7a. We color code the generated samples using k-means with four clusters. The qualitative examples display high plausibility and diversity. They follow the layout of the scene, the location of roads, vegetation, vehicles etc. We qualitatively compare the accuracy of samples generated by our LSTM-BMS model versus the LSTM-CVAE model

Method	Rainfall-MSE	CSI	FAR	POD	Correlation	CLL
[22]	1.420	0.577	0.195	0.660	0.908	x
Conv-LSTM-CVAE	1.259	0.651	<b>0.155</b>	0.701	0.910	132.78
Conv-LSTM-BMS	<b>1.163</b>	<b>0.670</b>	0.163	<b>0.734</b>	<b>0.918</b>	<b>132.52</b>

Table 3: Evaluation on HKO radar image sequences.

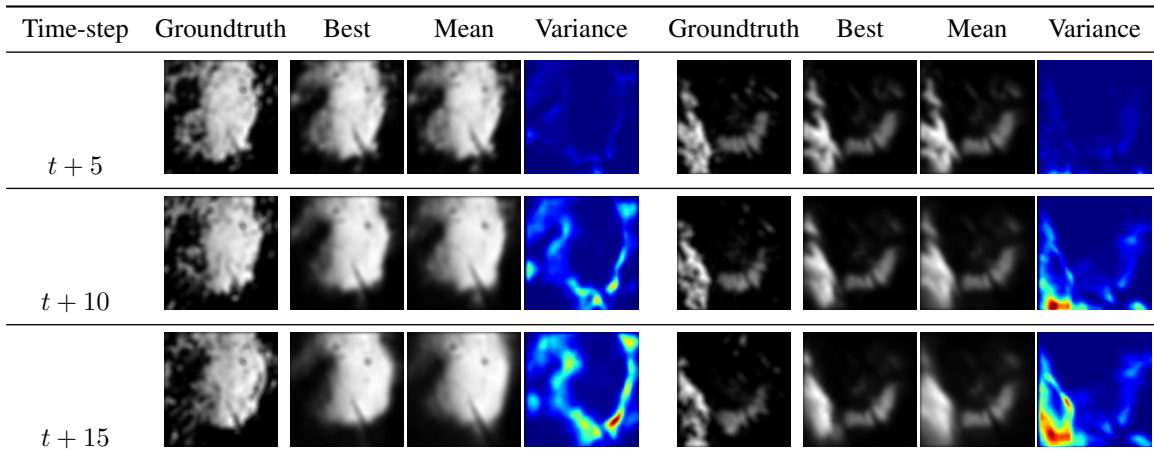


Figure 8: Statistics of samples generated by our LSTM-BMS model on the HKO dataset.

in Figure 7b. We see that the oracle top 10% of samples generated using the  $\hat{\mathcal{L}}_{\text{BMS}}$  objective are more accurate and thus more representative of the groundtruth.

### 4.3. Radar Echo

The Radar Echo dataset [22] consists of weather radar intensity images from 97 rainy days over Hong Kong from 2011 to 2013. The weather evolves due to variety of factors, which are difficult to identify using only the radar images, with varied and multimodal futures. Each sequences consists of 20 frames each of resolution  $100 \times 100$ , recorded at intervals of 6 minutes. We use the same dataset split as [22] and predict the next 15 images given the previous 5 images.

We compare our image sequence prediction model in Figure 3b trained with the  $\hat{\mathcal{L}}_{\text{BMS}}$  (Conv-LSTM-BMS) objective to one trained with the  $\hat{\mathcal{L}}_{\text{CVAE}}$  (Conv-LSTM-CVAE) objective. We additionally compare it to the Conv-LSTM model of [22]. In addition to the CLL metric (calculated per image sequence), we use the following precipitation nowcasting metrics from [22], 1. Rainfall mean squared error (Rainfall-MSE), 2. Critical success index (CSI), 3. False alarm rate (FAR), 4. Probability of detection (POD), and 5. Correlation. For fair comparison we estimate these metrics using  $T = 1$  random samples from the Conv-LSTM-CVAE and Conv-LSTM-BMS models.

We report the results in Table 3. Both the Conv-LSTM-CVAE and Conv-LSTM-BMS models perform better compared to [22]. This is due to use of embedding layers for fea-

ture extraction and the use of  $2 \times 2$  max pooling in between two Conv-LSTM layers for feature aggregation (compared no embedding layers or pooling in [22]). Furthermore, the superior CLL of the Conv-LSTM-BMS model demonstrates it’s ability to fit the data distribution better. We show qualitative examples in Figure 8 at  $t + 5$ ,  $t + 10$  and  $t + 15$ . We generate  $T = 50$  samples and show the sample closest to the groundtruth (Best), the mean of all the samples and the per-pixel variance in the samples. The qualitative examples demonstrate that our model produces highly accurate and diverse samples.

## 5. Conclusion

We have presented a novel “best of many” sample objective for Gaussian latent variable models and show its advantages for learning conditional models on multi-modal distributions. Our analysis shows indeed the learnt latent representation is better matched between training and test time – which in turn leads to more accurate samples. We show the benefits of our model on trajectory as well as image sequence prediction using three diverse datasets: MNIST strokes, Stanford drone and HKO weather. Our proposed approach consistently outperforms baselines and state of the art in all these scenarios.

**Acknowledgments** We would like to thank Francesco Croce for his comments and feedback.



## References

- [1] D. Acuna. Unsupervised modeling of the movement of basketball players using a deep generative model. [http://www.cs.toronto.edu/~davidj/projects/unsupervised\\_modeling\\_using\\_a\\_DGM.pdf](http://www.cs.toronto.edu/~davidj/projects/unsupervised_modeling_using_a_DGM.pdf), 2017. Accessed: 2017-10-26.
- [2] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese. Social lstm: Human trajectory prediction in crowded spaces. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [3] Y. Burda, R. Grosse, and R. Salakhutdinov. Importance weighted autoencoders. *ICLR*, 2015.
- [4] M. Comenetz. *Calculus: the elements*. World Scientific Publishing Co Inc, 2002.
- [5] E. D. De Jong. The mnist sequence dataset. <https://edwin-de-jong.github.io/blog/mnist-sequence-data/>, 2016. Accessed: 2017-10-26.
- [6] C. Finn, I. Goodfellow, and S. Levine. Unsupervised learning for physical interaction through video prediction. In *Advances in Neural Information Processing Systems*, 2016.
- [7] A. Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, 2013.
- [8] S. Gu, S. Levine, I. Sutskever, and A. Mnih. Muprop: Unbiased backpropagation for stochastic neural networks. *ICLR*, 2016.
- [9] D. Jimenez Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *International Conference on Machine Learning*, 2014.
- [10] D. P. Kingma, S. Mohamed, D. J. Rezende, and M. Welling. Semi-supervised learning with deep generative models. In *Advances in Neural Information Processing Systems*, 2014.
- [11] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *ICLR*, 2013.
- [12] C. Lassner, G. Pons-Moll, and P. V. Gehler. A generative model of people in clothing. 2017.
- [13] K. Lee, J. Kim, S. Chong, and J. Shin. Simplified stochastic feedforward neural networks. *arXiv preprint arXiv:1704.03188*, 2017.
- [14] N. Lee, W. Choi, P. Vernaza, C. B. Choy, P. H. Torr, and M. Chandraker. Desire: Distant future prediction in dynamic scenes with interacting agents. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [15] A. Mnih and D. Rezende. Variational inference for monte carlo objectives. In *International Conference on Machine Learning*, 2016.
- [16] T. Raiko, M. Berglund, G. Alain, and L. Dinh. Techniques for learning binary stochastic feedforward neural networks. *Advances in Neural Information Processing Systems*, 2014.
- [17] A. Robicquet, A. Sadeghian, A. Alahi, and S. Savarese. Learning social etiquette: Human trajectory understanding in crowded scenes. In *European Conference on Computer Vision*, 2016.
- [18] X. Shi, Z. Gao, L. Lausen, H. Wang, D.-Y. Yeung, W.-k. Wong, and W.-c. Woo. Deep learning for precipitation nowcasting: A benchmark and a new model. 2017.
- [19] K. Sohn, H. Lee, and X. Yan. Learning structured output representation using deep conditional generative models. In *Advances in Neural Information Processing Systems*, 2015.
- [20] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, 2014.
- [21] Y. Tang and R. R. Salakhutdinov. Learning stochastic feed-forward neural networks. In *Advances in Neural Information Processing Systems*, 2013.
- [22] S. Xingjian, Z. Chen, H. Wang, D.-Y. Yeung, W.-K. Wong, and W.-c. Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In *Advances in Neural Information Processing Systems*, 2015.
- [23] H. Xu, Y. Gao, F. Yu, and T. Darrell. End-to-end learning of driving models from large-scale video datasets. *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [24] T. Xue, J. Wu, K. L. Bouman, and W. T. Freeman. Visual dynamics: Probabilistic future frame synthesis via cross convolutional networks. In *Advances in Neural Information Processing Systems*, 2016.