

Quantifying Uncertainty in Motion Prediction with Variational Bayesian Mixture

Supplementary Material

A. Implementation Details

We implement our model using PyTorch, trained for 20 epochs on the INTERACTION dataset with a batch size of 64 and 25 epochs on the Argoverse 2 dataset with a batch size of 64. With only 1.3M parameters, the model balances scalability and performance. We set $\alpha = 1$ and use the Adam optimization solver with a learning rate of 0.0001 and the learning rate decay schedule with a step size of 5 epochs and a rate of 0.3 to ensure efficient convergence. We train and evaluate our model using only a single NVIDIA GeForce RTX 3090 Ti.

B. Evaluation Metrics

B.1. Uncertainty Quantification

In our formulation, the random variables s_f and v are both Gaussian variables, and z follows a categorical distribution. Therefore, we can compute the total uncertainty for a predicted distribution by its entropy. Given the generative model in equation 4, the total entropy can be estimated by the summation of three individual expected entropy:

$$\begin{aligned} & \sum_z \int_v \int_{s_f} p(s_f, v, z | \mathbf{x}) \log p(s_f, v, z | \mathbf{x}) ds_f dv \\ &= \mathbb{E}_{v, z \sim p(v, z | \mathbf{x})} \text{Entropy}(p(s_f | v, x)) \\ & \quad + \mathbb{E}_{s_f \sim p(z \sim p(z))} \text{Entropy}(p(v | \mathbf{x}, z)) \\ & \quad + \text{Entropy}(p(z)). \end{aligned} \quad (17)$$

Since we have a fixed prior $p(z)$, the comparison of the total entropy reduces to comparing the sum of the first two terms. In our experiment, we use Monte-Carlo sampling to generate N_{mc} samples of v for entropy calculation.

B.2. Motion Prediction

For motion prediction, we use the standard Minimum Average Displacement Error (minADE), Minimum Final Displacement Error (minFDE), and Miss Rate (MR) to assess the accuracy and effectiveness of our approach. minADE and minFDE are distance-based metrics commonly used in multi-modal trajectory prediction (i.e., trajectory prediction with multiple possible outcomes) tasks. The minADE calculates the average Euclidean distance between predicted and ground truth trajectories at each time step, taking the minimum across all trajectories in the prediction set:

$$\text{minADE}(\hat{x}_n^k, x_n) = \frac{1}{NT} \sum_{n=1}^N \min_{k=1, \dots, K} \sum_{t=1}^T \|\hat{x}_{n,t}^k - x_{n,t}\|_2. \quad (18)$$

On the other hand, the minFDE measures the Euclidean distance between predicted and ground truth final positions, effectively assessing the long-term prediction performance of the model:

$$\text{minFDE}(\hat{x}_n^k, x_n) = \frac{1}{N} \sum_{n=1}^N \min_{k=1, \dots, K} \|\hat{x}_{n,T}^k - x_{n,T}\|_2. \quad (19)$$

MR represents the ratio of 'miss' cases over all cases. The definitions of MR are significantly different for the INTERACTION dataset and the Argoverse 2 dataset.

In the INTERACTION dataset, if its prediction at the final timestamp ($T=30$) is out of a given lateral or longitudinal threshold of the ground truth, it will be assumed as a 'miss.' In the INTERACTION dataset, we need to align both the ground truth and the prediction by rotating them based on the yaw angle of the ground truth at the final timestamp, ensuring that the x-axis represents the longitudinal direction and the y-axis corresponds to the lateral direction. The lateral threshold is established as 1 meter, while the longitudinal threshold is a piecewise function set as:

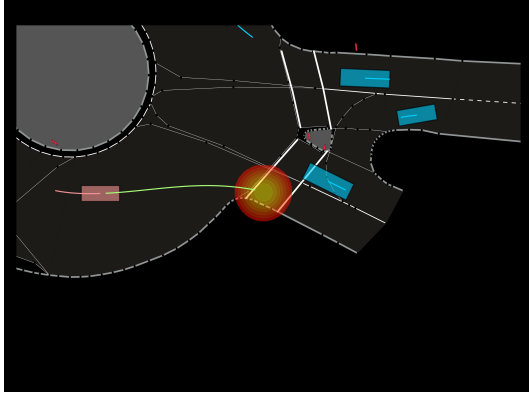
$$\text{Threshold}_{lon} = \begin{cases} 1 & v < 1.4m/s \\ 1 + \frac{v-1.4}{11-1.4} & 1.4m/s \leq v \leq 11m/s \\ 2 & v \geq 11m/s \end{cases} \quad (20)$$

For the Argoverse 2 dataset, the MR indicates the proportion of test samples where none of the predicted trajectories fall within a 2-meter range of the ground truth, as measured through the endpoint error measurement.

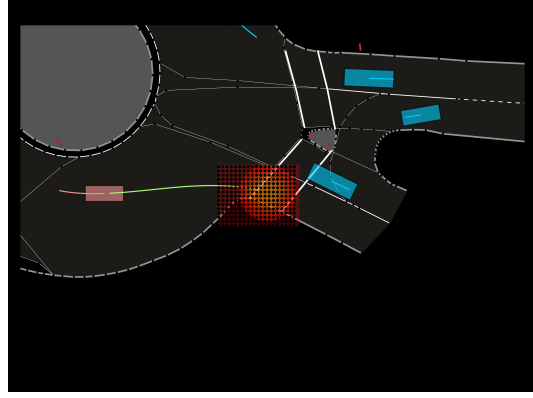
C. Methodology Details

C.1. Derivation of Evidence Lower Bound (ELBO)

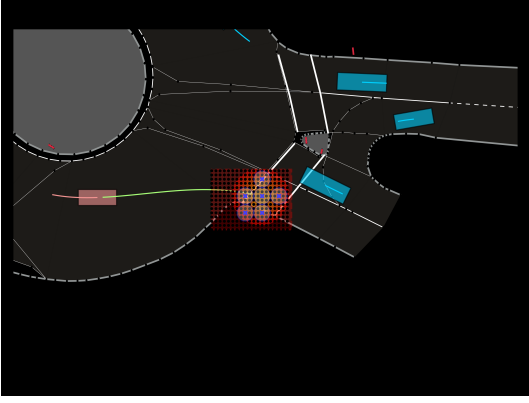
The standard and intuitive objective for training the probabilistic model in SeNeVA is to let the modeled conditional distribution $p(s_f | \mathbf{x})$ to match ground-truth data distribution through maximizing the likelihood. However, direct computation on the likelihood function is intractable since it involves calculating the integration given as $p(s_f | \mathbf{x}) = \sum_z \int_v p(s_f, v, z | \mathbf{x}) dv dz$, which is hard to estimate and optimize. To address this issue, we follow the popular *variational inference* method and introduce a tractable, closed-form, and easy-sampling proxy posterior $q(v, z | s_f, \mathbf{x})$ of the latent variables conditioned on the observed variables, and the lower bound of the log-likelihood can be derived with



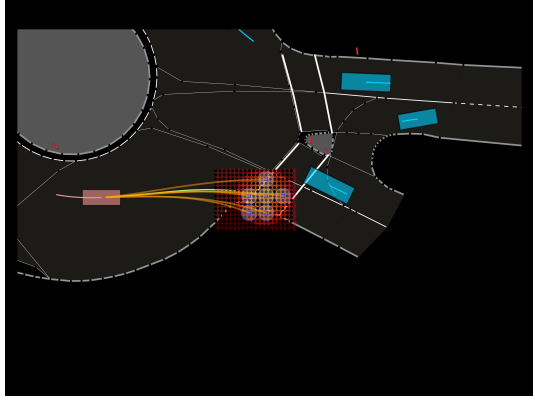
(a) **Distribution evaluation output.** The heatmap illustrate the distribution of y_{H+T} in this case quantified by the SeNeVA model.



(b) **Dense grid generation.** For sampling, we generate a grid of candidates that covers 2 standard deviation area of the y_{H+T} distribution.



(c) **NMS sampling results.** Following Algorithm 1, we can sample a set of top- M candidates (blue) regarding circular buffers defined by radius r and their IoU threshold γ .



(d) **Backward completion.** Starting from the last location, we assume homogeneous uncertainty over time and compute intermediate waypoints to obtain the final trajectories (orange).

Figure 6. **Example visualization of the backward sampling process.** The multi-modal trajectory prediction is generated through (a) evaluating distribution, (b) generating dense candidates, (c) applying NMS sampling, and finally (d) completing intermediate trajectories. To better illustrate the effectiveness of our method, we plot the history (red) and ground-truth future trajectory (green) of the target agent.

Jensen's Inequality:

$$\begin{aligned}
 \log p(\mathbf{s}_f | \mathbf{x}) &= \log \int_z \int_{\mathbf{v}} p(\mathbf{s}_f, \mathbf{v}, z | \mathbf{x}) d\mathbf{v} dz \\
 &= \log \mathbb{E}_{q(\mathbf{v}, z | \mathbf{s}_f, \mathbf{x})} \left[\frac{p(\mathbf{s}_f, \mathbf{v}, z | \mathbf{x})}{q(\mathbf{v}, z | \mathbf{s}_f, \mathbf{x})} \right] \quad (21) \\
 &\geq \mathbb{E}_{q(\mathbf{v}, z | \mathbf{s}_f, \mathbf{x})} \log \left[\frac{p(\mathbf{s}_f, \mathbf{v}, z | \mathbf{x})}{q(\mathbf{v}, z | \mathbf{s}_f, \mathbf{x})} \right].
 \end{aligned}$$

Since we factorize the joint distribution $p(\mathbf{s}_f, \mathbf{v}, z | \mathbf{x})$ and the posterior $q(\mathbf{v}, z | \mathbf{s}_f, \mathbf{x})$ in equation 4 and equation 6, respectively, we can leverage the factorization and expand the expectation term to compute the analytical solution of the lower bound. To simplify the following notation, we denote $p(\mathbf{s}_f, \mathbf{v}, z | \mathbf{x}) = p_{\mathbf{s}_f, \mathbf{v}, z}$, $q(\mathbf{v}, z | \mathbf{s}_f, \mathbf{x}) = q_{\mathbf{v}, z}$, $p(\mathbf{s}_f | \mathbf{v}, \mathbf{x}) = p_{\mathbf{s}_f}$, $p(\mathbf{v} | \mathbf{x}, z) = p_{\mathbf{v}}$, $p(z) = p_z$, $q(\mathbf{v} | \mathbf{s}_f, \mathbf{x}) = q_{\mathbf{v}}$, and $q(z | \mathbf{v}, \mathbf{x}) = q_z$. The expansion of the expectation term

writes:

$$\begin{aligned}
 \mathbb{E}_{q_{\mathbf{v}, z}} \log \left[\frac{p_{\mathbf{s}_f, \mathbf{v}, z}}{q_{\mathbf{v}, z}} \right] &= \int_z \int_{\mathbf{v}} \log \left[\frac{p_{\mathbf{s}_f} p_{\mathbf{v}} p_z}{q_{\mathbf{v}} q_z} \right] \cdot q_{\mathbf{v}, z} d\mathbf{v} dz \\
 &= \mathbb{E}_{q_{\mathbf{v}}} \log p_{\mathbf{s}_f} + \int_z \int_{\mathbf{v}} q_{\mathbf{v}} \log \left[\frac{p_{\mathbf{v}}}{q_{\mathbf{v}}} \right] d\mathbf{v} dz \\
 &\quad + \int_{\mathbf{v}} q_{\mathbf{v}} \int_z q_z \log \left[\frac{p_z}{q_z} \right] dz d\mathbf{v} \\
 &= \mathbb{E}_{q_{\mathbf{v}}} (\log p_{\mathbf{s}_f} - D_{\text{KL}}(q_z \| p_z)) - \mathbb{E}_{q_z} D_{\text{KL}}(q_{\mathbf{v}} \| p_{\mathbf{v}}). \quad (22)
 \end{aligned}$$

The expansion above is the ELBO objective we maximize during training equivalent to the formula given in equation 10. Therefore, maximizing the lower bound is equal to minimizing the KL divergence, driving the variational posterior $q(\mathbf{v}, z | \mathbf{s}_f, \mathbf{x})$ towards the ground-truth posterior. As a result, maximizing the ELBO objective can effectively maximize the likelihood.

C.2. Derivation of Assignment Network Loss

The assignment network directly approximates $p(z|\mathbf{x})$ to avoid tedious sampling at inference time from the latent \mathbf{v} space to estimate the posterior $q(z|\mathbf{x}) = \int_{\mathbf{v}} q(z|\mathbf{v}, \mathbf{x}) d\mathbf{v}$. We can obtain the distribution over z given in equation 13 by applying Bayes’ rule. Herein, we estimate the conditional distribution $p(\mathbf{s}_t|\mathbf{x}, z)$ by applying Monte-Carlo sampling over the latent \mathbf{v} space at training:

$$\begin{aligned} p(\mathbf{s}_t|\mathbf{x}, z) &= \int_{\mathbf{v}} p(\mathbf{s}_t, \mathbf{v}|\mathbf{x}, z) d\mathbf{v} \\ &\approx \frac{1}{N_{\text{mc}}} \sum_{n=1}^{N_{\text{mc}}} p(\mathbf{s}_t, \mathbf{v}^{(n)}|\mathbf{x}) p(\mathbf{v}^{(n)}, z|\mathbf{x}). \end{aligned} \quad (23)$$

C.3. Backward Sampling

As mentioned in section 4.5, we propose the backward sampling procedure to generate a collection of trajectories leveraging the distribution information learned by the model. The idea is first to sample the final location y_{H+T} that accounts for most uncertainty in the trajectory. The backward sampling procedure consists of three steps: Evaluation, Sampling, and Completion.

Evaluation At this stage, we leverage the output $\hat{\pi}$ from the assignment network to determine how we evaluate the distribution of y_{H+T} . One can use the component corresponding to $\hat{\pi}_{\text{max}}$. In our case, we promote multi-modality by computing the distribution as a mixture of top-6 components. For handling the latent space \mathbf{v} , one can apply Monte-Carlo sampling to approximate the integral. In our case, we choose to use the maximum likelihood samples (i.e., $\mathbf{v}^{\text{ml}} = \underset{\mathbf{v}}{\text{argmax}} p(\mathbf{v}, z|\mathbf{x})$) in equation 16 to evaluate the distribution, as shown in Figure 6a.

Sampling To allow full exploitation of the distribution information, we first generate a dense grid of candidates that covers the area within 2 standard deviations of the distribution mean. We adopt a rectangular grid with a resolution of 0.5 meters for simplicity, as illustrated in Figure 6b. One can quickly improve precision by choosing a smaller resolution or clipping the grid area. We then apply the NMS sampling given in Algorithm 1 to sample M candidates from the dense grid considering their circular buffers determined by hyperparameter r and the IoU threshold γ (see Figure 6c). Together, the two hyperparameters determine the density of selected candidates. In our practice, we choose $r = 1.4$ meters and $\gamma = 0\%$.

Completion The last step is to complete the intermediate trajectory from the target agent’s current position to the

sampled final locations. One can easily apply random sampling on each timestep to get the waypoints. Nevertheless, we find trajectories generated by this approach lack auto-consistency and can be non-smooth. To address the problem, we propose a strong assumption that displacement uncertainty is uniform over time. Hence, we can first parameterize an uncertainty distance parameter $u^{(m)}$ for each selected candidate and then use it for computing waypoints for all previous timesteps. Specifically, for sampled candidate $y_{H+T}^{(m)}$ from the distribution $\mathcal{N}(\mu_{H+T}, \Sigma_{H+T})$, we have

$$u^{(m)} = L_{H+T}^{-1} \left(y_{H+T}^{(m)} - \mu_{H+T} \right) : \Sigma_{H+T} = LL^T, \quad (24)$$

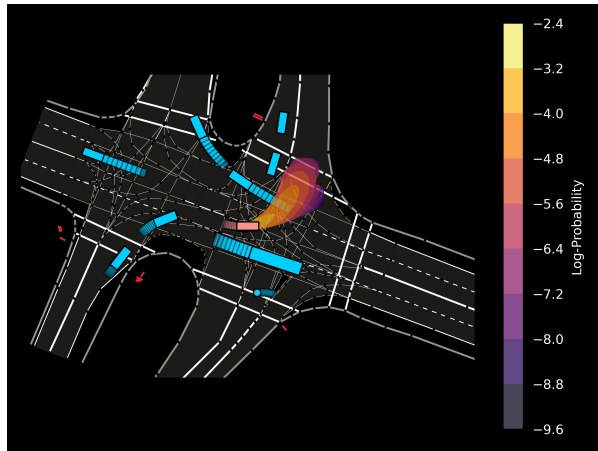
where L is the upper triangle Cholesky decomposition of the covariance. For each timestep $t = 1, \dots, T - 1$, we have

$$y_{H+t}^{(m)} = \mu_{H+t} + L_{H+t} \cdot u^{(m)}. \quad (25)$$

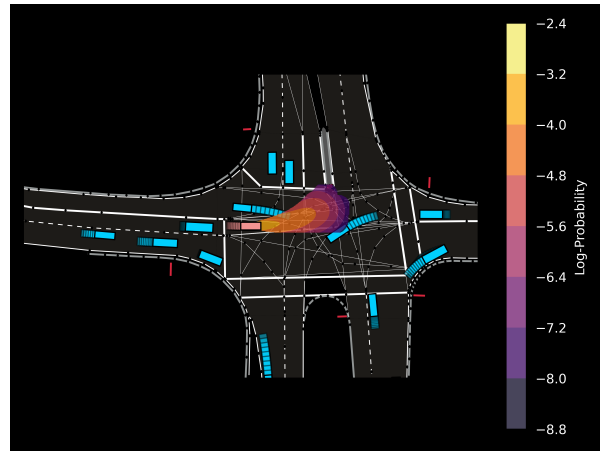
Finally, we connect intermediate waypoints with the sampled candidates to derive the required trajectory set. Figure 6d illustrates the final output from the backward sampling process.

D. Extensive Qualitative Results

We further visualize the quantified trajectory distributions on some representative cases selected from the INTERACTION dataset. Figure 7 illustrates two examples from unsignalized intersections, where SeNeVA successfully identifies the left-turn intention of the driver and quantifies the distribution of future trajectories that conform to the road geometry. In Figure 8, we visualize two cases in the expressway merging, where the SeNeVA model can anticipate the maneuver of the surrounding vehicles and predict distributions that avoid collisions.

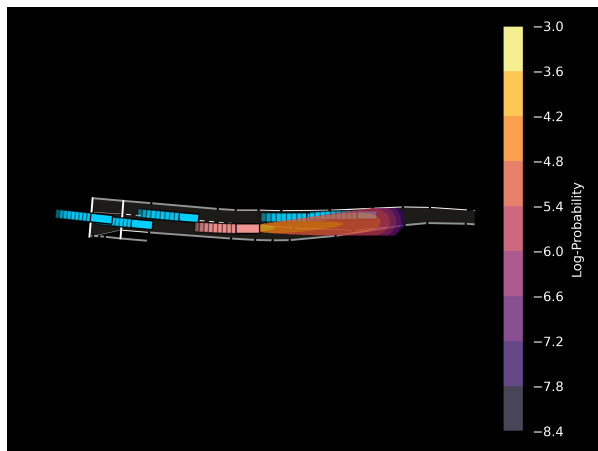


(a) Results on a case from DR_USA_Intersection_GL

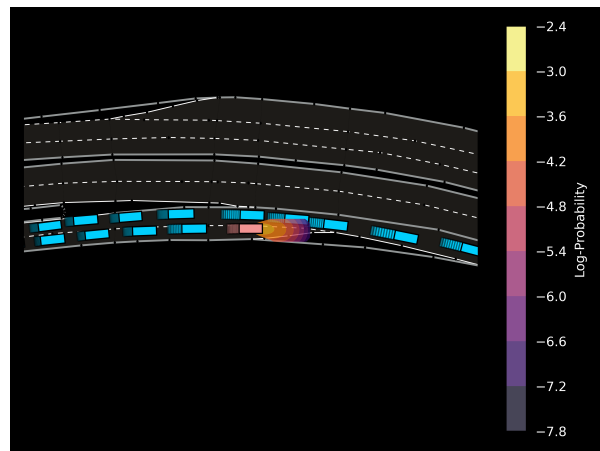


(b) Results on a case from DR_USA_Intersection_MA

Figure 7. **Representative example visualization of quantified uncertainty on intersections.** The heatmap generated by the SeNeVA model successfully identifies the left-turn intention of drivers in both cases. The predicted distributions conform to the road geometry.



(a) Results on a case from DR_DEU_Merging_MT



(b) Results on a case from DR_CHN_Merging_ZS0

Figure 8. **Representative example visualization of quantified uncertainty on intersections.** The model recognizes the existence of surrounding vehicles and predicts with higher certainty that a vehicle will stay hold to avoid collisions.