

SLAMP: Stochastic Latent Appearance and Motion Prediction

Adil Kaan Akan¹ Erkut Erdem² Aykut Erdem¹ Fatma Güney¹

¹ Koç University Is Bank AI Center, Istanbul, Turkey

² Hacettepe University Computer Vision Lab, Ankara, Turkey

{kakan20, aerdem, fguney}@ku.edu.tr erkut@cs.hacettepe.edu.tr

<https://kuis-ai.github.io/slamp>

Abstract

Motion is an important cue for video prediction and often utilized by separating video content into static and dynamic components. Most of the previous work utilizing motion is deterministic but there are stochastic methods that can model the inherent uncertainty of the future. Existing stochastic models either do not reason about motion explicitly or make limiting assumptions about the static part. In this paper, we reason about appearance and motion in the video stochastically by predicting the future based on the motion history. Explicit reasoning about motion without history already reaches the performance of current stochastic models. The motion history further improves the results by allowing to predict consistent dynamics several frames into the future. Our model performs comparably to the state-of-the-art models on the generic video prediction datasets, however, significantly outperforms them on two challenging real-world autonomous driving datasets with complex motion and dynamic background.

1. Introduction

Videos contain visual information enriched by motion. Motion is a useful cue for reasoning about human activities or interactions between objects in a video. Given a few initial frames of a video, our goal is to predict several frames into the future, as realistically as possible. By looking at a few frames, humans can predict what will happen next. Surprisingly, they can even attribute semantic meanings to random dots and recognize motion patterns [17]. This shows the importance of motion to infer the dynamics of the video and to predict the future frames.

Motion cues have been heavily utilized for future frame prediction in computer vision. A common approach is to factorize the video into static and dynamic components [35, 23, 25, 7, 10, 24, 16, 33]. First, most of the previous methods are deterministic and fail to model the uncertainty of the future. Second, motion is typically interpreted as local



Figure 1: Comparison of the first prediction frames (11th) SLAMP (left) vs. state-of-the-art method, SRVP [9] (right) on KITTI [11] (top) and Cityscapes [4] (bottom) datasets. Our method can predict both foreground and background objects better than SRVP. Full sequence predictions can be seen in Fig. 32 and 33.

changes from one frame to the next. However, changes in motion follow certain patterns when observed over some time interval. Consider scenarios where objects move with near-constant velocity, or humans repeating atomic actions in videos. Regularities in motion can be very informative for future frame prediction. In this work, we propose to explicitly model the change in motion, or *the motion history*, for predicting future frames.

Stochastic methods have been proposed to model the inherent uncertainty of the future in videos. Earlier methods encode the dynamics of the video in stochastic latent variables which are decoded to future frames in a deterministic way [5]. We first assume that both appearance and motion are encoded in the stochastic latent variables and decode them separately into appearance and motion predictions in a deterministic way. Inspired by the previous deterministic methods [8, 23, 10], we also estimate a mask relating the two. Both appearance and motion decoders are expected to predict the full frame but they might fail due to occlusions around motion boundaries. Intuitively, we predict a probabilistic mask from the results of the appearance and motion decoders to combine them into a more accurate final prediction. Our model learns to use motion cues in the dynamic parts and relies on appearance in the occluded regions.

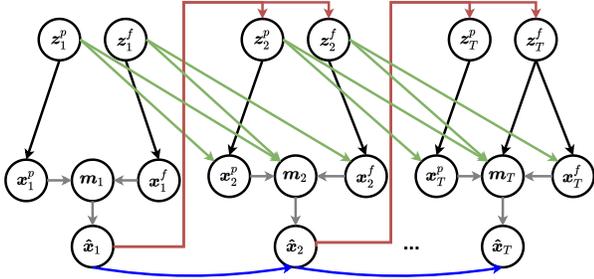


Figure 2: **Generative Model of SLAMP.** The graphical model shows the generation process of SLAMP with motion history. There are two separate latent variables for appearance z_t^p and motion z_t^f generating frames x_t^p and x_t^f (black). Information is propagated between time-steps through the recurrence between frame predictions (blue), corresponding latent variables (green), and from frame predictions to latent variables (red). The final prediction \hat{x}_t is a weighted combination of the x_t^p and x_t^f according to the mask $m(x_t^p, x_t^f)$. Note that predictions at a time-step depend on all of the previous time-steps recurrently, but only the connections between consecutive ones are shown for clarity.

The proposed stochastic model with deterministic decoders cannot fully utilize the motion history, even when motion is explicitly decoded. In this work, we propose a model to recognize regularities in motion and remember them in the motion history to improve future frame predictions. We factorize stochastic latent variables as static and dynamic components to model the motion history in addition to the appearance history. We learn two separate distributions representing appearance and motion and then decode static and dynamic parts from the respective ones.

Our model outperforms all the previous work and performs comparably to the state-of-the-art method, SRVP, [9] without any limiting assumptions on the changes in the static component on the generic video prediction datasets, MNIST, KTH and BAIR. However, our model outperforms all the previous work, including SRVP, on two challenging real-world autonomous driving datasets with dynamic background and complex object motion.

2. Related Work

Appearance-Motion Decomposition: The previous work explored motion cues for video generation either explicitly with optical flow [35, 34, 22, 23, 25, 7, 10] or implicitly with temporal differences [24] or pixel-level transformations [16, 33]. There are some common factors among these methods such as using recurrent models [29, 24, 7], specific processing of dynamic parts [16, 22, 7, 10], utilizing a mask [8, 23, 10], and adversarial training [33, 25]. We also use recurrent models, predict a mask, and separately process motion, but in a stochastic way.

The previous work which explored motion for video gen-

eration are mostly deterministic, therefore failing to capture uncertainty of the future. There are a couple of attempts to learn multiple future trajectories from a single image with a conditional variational autoencoder [34] or to capture motion uncertainty with a probabilistic motion encoder [22]. The latter work uses separate decoders for flow and frame similar to our approach, however, predicts them only from the latent vector. We incorporate information from previous frames with additional modelling of the motion history.

Stochastic Video Generation: SV2P [1] and SVG [5] are the first to model the stochasticity in video sequences using latent variables. The input from past frames are encoded in a posterior distribution to generate the future frames. In a stochastic framework, learning is performed by maximizing the likelihood of the observed data and minimizing the distance of the posterior distribution to a prior distribution, either fixed [1] or learned from previous frames [5]. Since time-variance in the model is proven crucial by the previous work, we sample a latent variable at every time step [5]. Sampled random variables are fed to a frame predictor, modelled recurrently using an LSTM. We model appearance and motion distributions separately and train two frame predictors for static and dynamic parts.

Typically, each distribution, including the prior and the posterior, is modeled with a recurrent model such as an LSTM. Villegas et al. [32] replace the linear LSTMs with convolutional ones at the cost of increasing the number of parameters. Castrejon et al. [3] introduce a hierarchical representation to model latent variables at different scales, by introducing additional complexity. Lee et al. [20] incorporate an adversarial loss into the stochastic framework to generate sharper images, at the cost of less diverse results. Our model with linear LSTMs can generate diverse and sharp-looking results without any adversarial losses, by incorporating motion information successfully into the stochastic framework. Recent methods model dynamics of the keypoints to avoid errors in pixel space and achieve stable learning [26]. This offers an interesting solution for videos with static background and moving foreground objects that can be represented with keypoints. Our model can generalize to videos with changing background without needing keypoints to represent objects.

Optical flow has been used before in future prediction [21, 25]. Li et al. [21] generate future frames from a still image by using optical flow generated by an off-the-shelf model, whereas we compute flow as part of prediction. Lu et al. [25] use optical flow for video extrapolation and interpolation without modeling stochasticity. Long-term video extrapolation results show the limitation of this work in terms of predicting future due to relatively small motion magnitudes considered in extrapolation. Differently from flow, Xue et al. [36] model the motion as image differences using cross convolutions.

State-Space Models: Stochastic models are typically auto-regressive, i.e. the next frame is predicted based on the frames generated by the model. As opposed to interleaving process of auto-regressive models, state-space models separate the frame generation from the modelling of dynamics [13]. State-of-the-art method SRVP [9] proposes a state-space model for video generation with deterministic state transitions representing residual change between the frames. This way, dynamics are modelled with latent state variables which are independent of previously generated frames. Although independent latent states are computationally appealing, they cannot model the motion history of the video. In addition, content variable designed to model static background cannot handle changes in the background. We can generate long sequences with complex motion patterns by explicitly modelling the motion history without any limiting assumptions about the dynamics of the background.

3. Methodology

3.1. Stochastic Video Prediction

Given the previous frames $\mathbf{x}_{1:t-1}$ until time t , our goal is to predict the target frame \mathbf{x}_t . For that purpose, we assume that we have access to the target frame \mathbf{x}_t during training and use it to capture the dynamics of the video in stochastic latent variables \mathbf{z}_t . By learning to approximate the distribution over \mathbf{z}_t , we can decode the future frame \mathbf{x}_t from \mathbf{z}_t and the previous frames $\mathbf{x}_{1:t-1}$ at test time.

Using all the frames including the target frame, we compute a posterior distribution $q_\phi(\mathbf{z}_t|\mathbf{x}_{1:t})$ and sample a latent variable \mathbf{z}_t from this distribution at each time step. The stochastic process of the video is captured by the latent variable \mathbf{z}_t . In other words, it should contain information accumulated over the previous frames rather than only condensing the information on the current frame. This is achieved by encouraging $q_\phi(\mathbf{z}_t|\mathbf{x}_{1:t})$ to be close to a prior distribution $p(\mathbf{z})$ in terms of KL-divergence. The prior can be sampled from a fixed Gaussian

at each time step or can be learned from previous frames up to the target frame $p_\psi(\mathbf{z}_t|\mathbf{x}_{1:t-1})$. We prefer the latter one as it is shown to work better by learning a prior that varies across time [5].

The target frame \mathbf{x}_t is predicted based on the previous frames $\mathbf{x}_{1:t-1}$ and the latent vectors $\mathbf{z}_{1:t}$.

In practice, we only use the latest frame \mathbf{x}_{t-1} and the latent vector \mathbf{z}_t as input and dependencies from further previous frames are propagated with a recurrent model. The output of the frame predictor \mathbf{g}_t

contains the information required to decode \mathbf{x}_t .

Typically, \mathbf{g}_t is decoded to a fixed-variance Gaussian distribution whose mean is the predicted target frame $\hat{\mathbf{x}}_t$ [5].

3.2. SLAMP

We call the predicted target frame, appearance prediction \mathbf{x}_t^p in the pixel space. In addition to \mathbf{x}_t^p , we also estimate optical flow $\mathbf{f}_{t-1:t}$ from the previous frame $t-1$ to the target frame t . The flow $\mathbf{f}_{t-1:t}$ represents the motion of the pixels from the previous frame to the target frame. We reconstruct the target frame \mathbf{x}_t^f from the estimated optical flow via differentiable warping [15]. Finally, we estimate a mask $\mathbf{m}(\mathbf{x}_t^p, \mathbf{x}_t^f)$ from the two frame estimations to combine them into the final estimation $\hat{\mathbf{x}}_t$:

$$\hat{\mathbf{x}}_t = \mathbf{m}(\mathbf{x}_t^p, \mathbf{x}_t^f) \odot \mathbf{x}_t^p + (\mathbf{1} - \mathbf{m}(\mathbf{x}_t^p, \mathbf{x}_t^f)) \odot \mathbf{x}_t^f \quad (1)$$

where \odot denotes element-wise Hadamard product and \mathbf{x}_t^f is the result of warping the source frame to the target frame according to the estimated flow field $\mathbf{f}_{t-1:t}$. Especially in the dynamic parts with moving objects, the target frame can be reconstructed accurately using motion information. In the occluded regions where motion is unreliable, the model learns to rely on the appearance prediction. The mask prediction learns a weighting between the appearance and the motion predictions for combining them.

We call this model SLAMP-Baseline because it is limited in the sense that it only considers the motion with respect to the previous frame while decoding the output. In SLAMP, we extend the stochasticity in the appearance space to the motion space as well. This way, we can model appearance changes and motion patterns in the video explicitly and make better predictions of future. Fig. 3 shows an illustration of SLAMP (see Appendix Section A for SLAMP-Baseline).

In order to represent appearance and motion, we compute two separate posterior distributions $q_{\phi_p}(\mathbf{z}_t^p|\mathbf{x}_{1:t})$ and $q_{\phi_f}(\mathbf{z}_t^f|\mathbf{x}_{1:t})$, respectively. We sample two latent variables \mathbf{z}_t^p and \mathbf{z}_t^f from these distributions in the pixel space and the flow space. This allows a decomposition of the video into static and dynamic components. Intuitively, we expect the dynamic component to focus on changes and the static to what remains constant from the previous frames to the target frame. If the background is moving according to a camera motion, the static component can model the change in the background assuming that it remains constant throughout the video, e.g. ego-motion of a car.

The Motion History: The latent variable \mathbf{z}_t^f should contain motion information accumulated over the previous frames rather than local temporal changes between the last frame and the target frame. We achieve this by encouraging $q_{\phi_f}(\mathbf{z}_t^f|\mathbf{x}_{1:t})$ to be close to a prior distribution in terms of KL-divergence. Similar to [5], we learn the motion prior conditioned on previous frames up to the target frame: $p_{\psi_f}(\mathbf{z}_t^f|\mathbf{x}_{1:t-1})$. We repeat the same for the static part represented by \mathbf{z}_t^p with posterior $q_{\phi_p}(\mathbf{z}_t^p|\mathbf{x}_{1:t})$ and the learned prior $p_{\psi_p}(\mathbf{z}_t^p|\mathbf{x}_{1:t-1})$.

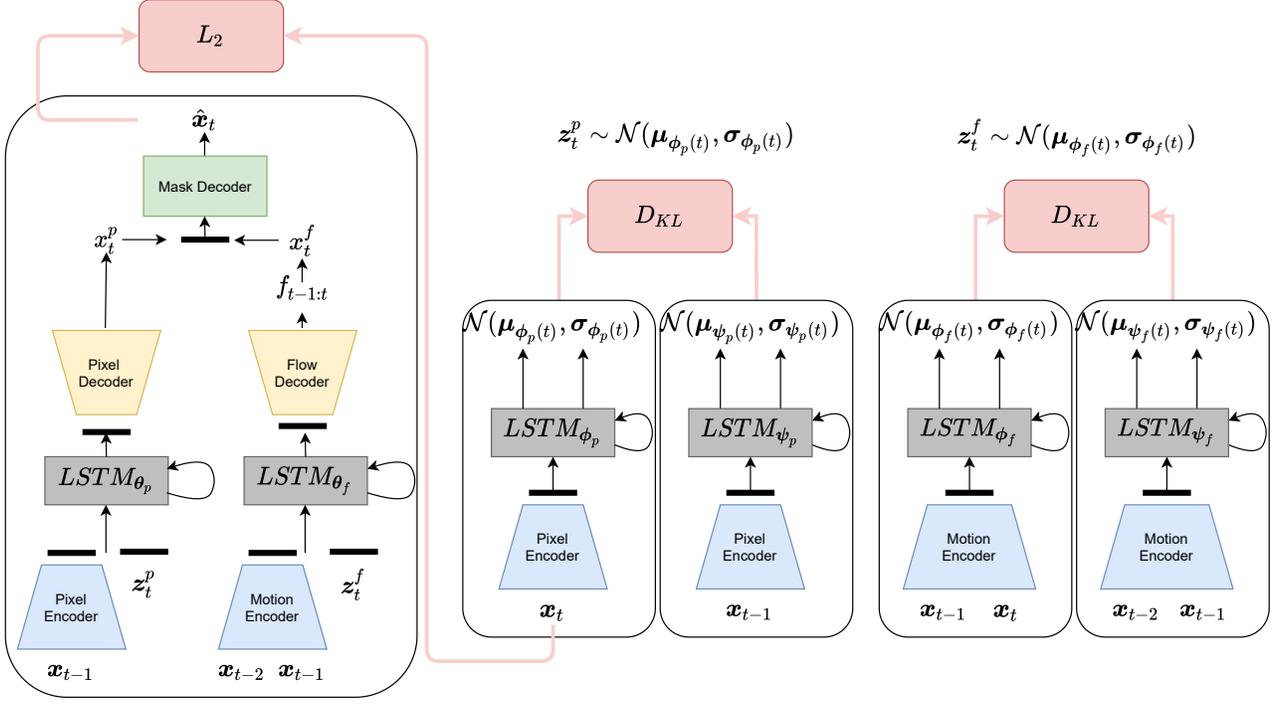


Figure 3: **SLAMP**. This figure shows the components of our SLAMP model including the prediction model, inference and learned prior models for pixel and then flow from left to right. Observations \mathbf{x}_t are mapped to the latent space by using a pixel encoder for appearance on each frame and a motion encoder for motion between consecutive frames. The blue boxes show encoders, yellow and green ones decoders, gray ones recurrent posterior, prior, and predictor models, and lastly red ones show loss functions during training. Note that L_2 loss is applied three times for appearance prediction \mathbf{x}_t^p , motion prediction \mathbf{x}_t^f , and the combination of the two $\hat{\mathbf{x}}_t$ according to the mask prediction $\mathbf{m}(\mathbf{x}_t^p, \mathbf{x}_t^f)$. We only show L_2 loss between the actual frame \mathbf{x}_t and the final predicted frame $\hat{\mathbf{x}}_t$ in the figure. For inference, only the prediction model and learned prior models are used.

3.3. Variational Inference

For our basic formulation (SLAMP-Baseline), the derivation of the loss function is straightforward and provided in Appendix Section B. For SLAMP, the conditional joint probability corresponding to the graphical model in Fig. 2 is:

$$p(\mathbf{x}_{1:T}) = \prod_{t=1}^T p(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_t^p, \mathbf{z}_t^f) \quad (2)$$

$$p(\mathbf{z}_t^p | \mathbf{x}_{1:t-1}, \mathbf{z}_{t-1}^p) p(\mathbf{z}_t^f | \mathbf{x}_{1:t-1}, \mathbf{z}_{t-1}^f)$$

The true distribution over the latent variables \mathbf{z}_t^p and \mathbf{z}_t^f is intractable. We train time-dependent inference networks $q_{\phi_p}(\mathbf{z}_t^p | \mathbf{x}_{1:T})$ and $q_{\phi_f}(\mathbf{z}_t^f | \mathbf{x}_{1:T})$ to approximate the true distribution with conditional Gaussian distributions. In order to optimize the likelihood of $p(\mathbf{x}_{1:T})$, we need to infer latent variables \mathbf{z}_t^p and \mathbf{z}_t^f , which correspond to uncertainty of static and dynamic parts in future frames, respectively. We use a variational inference model to infer the latent variables.

Since \mathbf{z}_t^p and \mathbf{z}_t^f are independent across time, we can decompose Kullback-Leibler terms into individual time steps. We train the model by optimizing the variational lower bound

(see Appendix Section B for the derivation):

$$\log p_{\theta}(\mathbf{x}) \geq \mathcal{L}_{\theta, \phi_p, \phi_f, \psi_p, \psi_f}(\mathbf{x}_{1:T}) \quad (3)$$

$$= \sum_t \mathbb{E}_{\substack{\mathbf{z}_{1:t}^p \sim q_{\phi_p} \\ \mathbf{z}_{1:t}^f \sim q_{\phi_f}}} \log p_{\theta}(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{z}_{1:t}^p, \mathbf{z}_{1:t}^f)$$

$$- \beta \left[D_{\text{KL}}(q(\mathbf{z}_t^p | \mathbf{x}_{1:t}) || p(\mathbf{z}_t^p | \mathbf{x}_{1:t-1})) \right.$$

$$\left. + D_{\text{KL}}(q(\mathbf{z}_t^f | \mathbf{x}_{1:t}) || p(\mathbf{z}_t^f | \mathbf{x}_{1:t-1})) \right]$$

The likelihood p_{θ} , can be interpreted as an L_2 penalty between the actual frame \mathbf{x}_t and the estimation $\hat{\mathbf{x}}_t$ as defined in (1). We apply the L_2 loss to the predictions of appearance and motion components as well.

The posterior terms for uncertainty are estimated as an expectation over $q_{\phi_p}(\mathbf{z}_t^p | \mathbf{x}_{1:t})$, $q_{\phi_f}(\mathbf{z}_t^f | \mathbf{x}_{1:t})$. As in [5], we also learn the prior distributions from the previous frames up to the target frame as $p_{\psi_p}(\mathbf{z}_t^p | \mathbf{x}_{1:t-1})$, $p_{\psi_f}(\mathbf{z}_t^f | \mathbf{x}_{1:t-1})$. We train the model using the re-parameterization trick [19]. We classically choose the posteriors to be factorized Gaussian so that all the KL divergences can be computed analytically.

3.4. Architecture

We encode the frames with a feed-forward convolutional architecture to obtain appearance features at each time-step. In SLAMP, we also encode consecutive frame pairs into a feature vector representing the motion between them. We then train linear LSTMs to infer posterior and prior distributions at each time-step from encoded appearance and motion features.

Stochastic video prediction model with a learned prior [5] is a special case of our baseline model with a single pixel decoder, we also add motion and mask decoders. Next, we describe the steps of the generation process for the dynamic part.

At each time step, we encode \mathbf{x}_{t-1} and \mathbf{x}_t into \mathbf{h}_t^f , representing the motion from the previous frame to the target frame. The posterior LSTM is updated based on the \mathbf{h}_t^f :

$$\begin{aligned} \mathbf{h}_t^f &= \text{MotionEnc}(\mathbf{x}_{t-1}, \mathbf{x}_t) \quad (4) \\ \boldsymbol{\mu}_{\phi_f(t)}, \boldsymbol{\sigma}_{\phi_f(t)} &= \text{LSTM}_{\phi_f}(\mathbf{h}_t^f) \end{aligned}$$

For the prior, we use the motion representation \mathbf{h}_{t-1}^f from the previous time step, i.e. the motion from the frame $t-2$ to the frame $t-1$, to update the prior LSTM:

$$\begin{aligned} \mathbf{h}_{t-1}^f &= \text{MotionEnc}(\mathbf{x}_{t-2}, \mathbf{x}_{t-1}) \quad (5) \\ \boldsymbol{\mu}_{\psi_f(t)}, \boldsymbol{\sigma}_{\psi_f(t)} &= \text{LSTM}_{\psi_f}(\mathbf{h}_{t-1}^f) \end{aligned}$$

At the first time-step where there is no previous motion, we assume zero-motion by estimating the motion from the previous frame to itself.

The predictor LSTMs are updated according to encoded features and sampled latent variables:

$$\begin{aligned} \mathbf{g}_t^f &= \text{LSTM}_{\theta_f}(\mathbf{h}_{t-1}^f, \mathbf{z}_t^f) \quad (6) \\ \boldsymbol{\mu}_{\theta_f} &= \text{FlowDec}(\mathbf{g}_t^f) \end{aligned}$$

There is a difference between the train time and inference time in terms of the distribution the latent variables are sampled from. At train time, latent variables are sampled from the posterior distribution. At test time, they are sampled from the posterior for the conditioning frames and from the prior for the following frames. The output of the predictor LSTMs are decoded into appearance and motion predictions separately and combined into the final prediction using the mask prediction (Eq. (1)).

4. Experiments

We evaluate the performance of the proposed approach and compare it to the previous methods on three standard video prediction datasets including Stochastic Moving MNIST, KTH Actions [28] and BAIR Robot Hand [6]. We

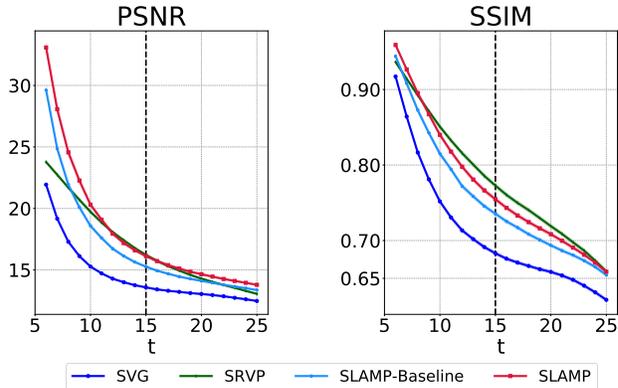


Figure 4: **Quantitative Results on MNIST.** This figure compares SLAMP to SLAMP-Baseline, SVG [5], and SRVP [9] on MNIST in terms of PSNR (left) and SSIM (right). SLAMP clearly outperforms our baseline model and SVG, and performs comparably to SRVP. Vertical bars mark the length of the training sequences.

specifically compare our baseline model (SLAMP-Baseline) and our model (SLAMP) to SVG [5] which is a special case of our baseline with a single pixel decoder, SAVP [20], SV2P [1], and lastly to SRVP [9]. We also compare our model to SVG [5] and SRVP [9] on two different challenging real world datasets, KITTI [12, 11] and Cityscapes [4], with moving background and complex object motion. We follow the evaluation setting introduced in [5] by generating 100 samples for each test sequence and report the results according to the best one in terms of average performance over the frames. Our experimental setup including training details and parameter settings can be found in Appendix Section C. We also share the code for reproducibility.

Table 1: **FVD Scores on KTH and BAIR.** This table compares all the methods in terms of FVD scores with their 95%-confidence intervals over five different samples from the models. Our model is the second best on KTH and among top three methods on BAIR.

Dataset	KTH	BAIR
SV2P	636 ± 1	965 ± 17
SAVP	374 ± 3	152 ± 9
SVG	377 ± 6	255 ± 4
SRVP	222 ± 3	163 ± 4
SLAMP-Baseline	236 ± 2	245 ± 5
SLAMP	<u>228 ± 5</u>	—

Evaluation Metrics: We compare the performance using three frame-wise metrics and a video-level one. Peak Signal-to-Noise Ratio (PSNR), *higher better*, based on L_2 distance between the frames penalizes differences in dynamics but also favors blur predictions. Structured Similarity (SSIM),

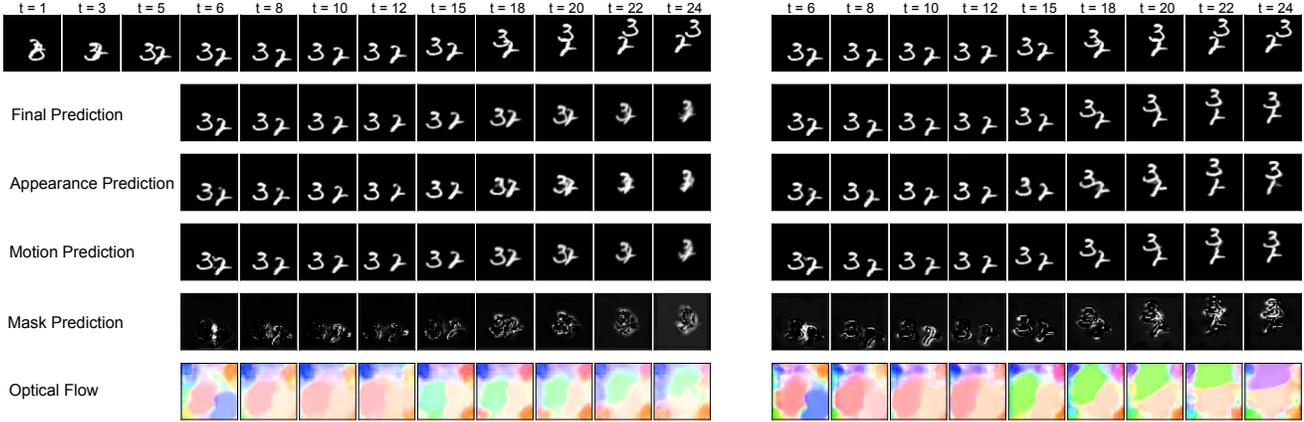


Figure 5: **SLAMP-Baseline (left) vs. SLAMP (right) on MNIST.** The top row shows the ground truth, followed by the frame predictions by the final, the appearance, the motion, and the last two rows show the mask and the optical flow predictions with false coloring. In this challenging case with bouncing and collisions, the baseline confuses the digits and cannot predict last frames correctly whereas SLAMP can generate predictions very close to the ground truth by learning smooth transitions in the motion history, as can be seen from optical flow predictions. See Fig. 13 for the color wheel showing the direction of flow.

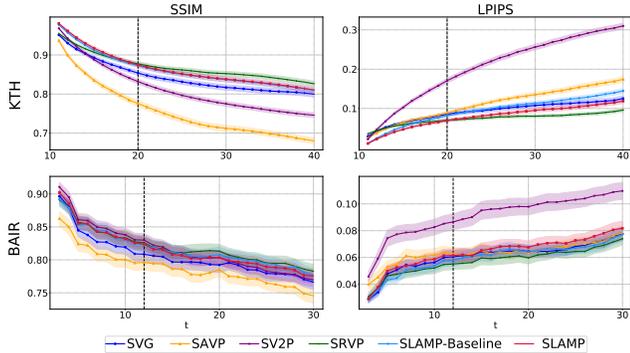


Figure 6: **Quantitative Results on KTH and BAIR.** We compare our results to previous work in terms of PSNR, SSIM, and LPIPS metrics with respect to the time steps on KTH (**top**), and BAIR (**bottom**) datasets, with 95%-confidence intervals. Vertical bars mark the length of training sequences. SLAMP outperforms previous work including SVG [5], SAVP [20], SV2P [1] and performs comparably to the state of the art method SRVP [9] on both datasets.

higher better, compares local patches to measure similarity in structure spatially. Learned Perceptual Image Patch Similarity (LPIPS) [37], *lower better*, measures the distance between learned features extracted by a CNN trained for image classification. Frchet Video Distance (FVD) [31], *lower better*, compares temporal dynamics of generated videos to the ground truth in terms of representations computed for action recognition.

Stochastic Moving MNIST: This dataset contains up to two MNIST digits moving linearly and bouncing from walls with a random velocity as introduced in [5]. Following the same training and evaluation settings as in the previous work,

we condition on the first 5 frames during training and learn to predict the next 10 frames. During testing, we again condition on the first 5 frames but predict the next 20 frames.

Fig. 4 shows quantitative results on MNIST in comparison to SVG [5] and SRVP [9] in terms of PSNR and SSIM, omitting LPIPS as in SRVP. Our baseline model with a motion decoder (SLAMP-Baseline) already outperforms SVG on both metrics. SLAMP further improves the results by utilizing the motion history and reaches a comparable performance to the state of the art model SRVP. This shows the benefit of separating the video into static and dynamic parts in both state-space models (SRVP) and auto-regressive models (ours, SLAMP). This way, models can better handle challenging cases such as crossing digits as shown next.

We qualitatively compare SLAMP to SLAMP-Baseline on MNIST in Fig. 5. The figure shows predictions of static and dynamic parts as appearance and motion predictions, as well the final prediction as the combination of the two. According to the mask prediction, the final prediction mostly relies on the dynamic part shown as black on the mask and uses the static component only near the motion boundaries. Moreover, optical flow prediction does not fit the shape of the digits but expands as a region until touching the motion region of the other digit. This is due to the uniform black background. Moving a black pixel in the background randomly is very likely to result in another black pixel in the background, which means zero-loss for the warping result. Both models can predict optical flow correctly for the most part and resort to the appearance result in the occluded regions. However, continuity in motion is better captured by SLAMP with the colliding digits whereas the baseline model cannot recover from it, leading to blur results, far from the ground truth. Note that we pick the best sample for both

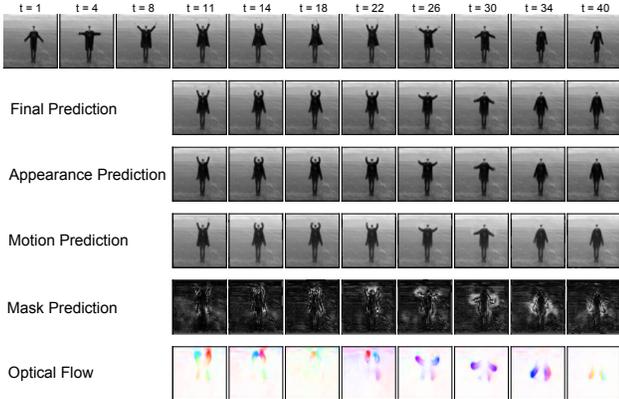


Figure 7: **Qualitative Results on KTH** We visualize the results of SLAMP on KTH dataset. The top row shows the ground truth, followed by the frame predictions by the final, the appearance, the motion, and the last two rows show the mask and the optical flow predictions. The mask prediction combines the appearance prediction (white) and the motion prediction (black) into the final prediction.

models among 100 samples according to LPIPS.

KTH Action Dataset: KTH dataset contains real videos where people perform a single action such as walking, running, boxing, etc. in front of a static camera [28]. We expect our model with motion history to perform very well by exploiting regularity in human actions on KTH. Following the same training and evaluation settings used in the previous work, we condition on the first 10 frames and learn to predict the next 10 frames. During testing, we again condition on the first 10 frames but predict the next 30 frames.

Fig. 6 and Table 1 show quantitative results on KTH in comparison to previous approaches. Both our baseline and SLAMP models outperform previous approaches and perform comparably to SRVP, in all metrics including FVD. A detailed visualization of all three frame predictions as well as flow and mask are shown in Fig. 7. Flow predictions are much more fine-grained than MNIST by capturing fast motion of small objects such as hands or thin objects such as legs (see Appendix Section E). The mask decoder learns to identify regions around the motion boundaries which cannot be matched with flow due to occlusions and assigns more weight to the appearance prediction in these regions.

On KTH, the subject might appear after the conditioning frames. These challenging cases can be problematic for some previous work as shown in SRVP [9]. Our model can generate samples close to the ground truth despite very little information on the conditioning frames as shown in Fig. 8. The figure shows the best sample in terms of LPIPS, please see Appendix Section E for a diverse set of samples with subjects of various poses appearing at different time steps.

BAIR Robot Hand: This dataset contains videos of a

Table 2: **Results with a Moving Background.** We evaluate our model SLAMP in comparison to SVG and SRVP on KITTI [11] and Cityscapes [4] datasets by conditioning on 10 frames and predicting 20 frames into the future.

Models	PSNR (\uparrow)	SSIM (\uparrow)	LPIPS (\downarrow)
SVG [5]	12.70 ± 0.70	0.329 ± 0.030	0.594 ± 0.034
SRVP [9]	13.41 ± 0.42	0.336 ± 0.034	0.635 ± 0.021
SLAMP	13.46 ± 0.74	0.337 ± 0.034	0.537 ± 0.042
KITTI [12, 11]			
Models	PSNR (\uparrow)	SSIM (\uparrow)	LPIPS (\downarrow)
SVG [5]	20.42 ± 0.63	0.606 ± 0.023	0.340 ± 0.022
SRVP [9]	20.97 ± 0.43	0.603 ± 0.016	0.447 ± 0.014
SLAMP	21.73 ± 0.76	0.649 ± 0.025	0.2941 ± 0.022
Cityscapes [4]			

robot hand moving and pushing objects on a table [6]. Due to uncertainty in the movements of the robot arm, BAIR is a standard dataset for evaluating stochastic video prediction models. Following the training and evaluation settings used in the previous work, we condition on the first 2 frames and learn to predict the next 10 frames. During testing, we again condition on the first 2 frames but predict the next 28 frames.

We show quantitative results on BAIR in Fig. 6 and Table 1. Our baseline model achieves comparable results to SRVP, outperforming other methods in all metrics except SV2P [1] in PSNR and SAVP [20] in FVD. With 2 conditioning frames only, SLAMP cannot utilize the motion history and performs similarly to the baseline model on BAIR (see Appendix Section D). This is simply due to the fact that there is only one flow field to condition on, in other words, no motion history. Therefore, we only show the results of the baseline model on this dataset.

Real-World Driving Datasets: We perform experiments on two challenging autonomous driving datasets: KITTI [12, 11] and Cityscapes [4] with various challenges. Both datasets contain everyday real-world scenes with complex dynamics due to both background and foreground motion. KITTI is recorded in one town in Germany while Cityscapes is recorded in 50 European cities, leading to higher diversity.

Cityscapes primarily focuses on semantic understanding of urban street scenes, therefore contains a larger number of dynamic foreground objects compared to KITTI. However, motion lengths are larger on KITTI due to lower frame-rate. On both datasets, we condition on 10 frames and predict 10 frames into the future to train our models. Then at test time, we predict 20 frames conditioned on 10 frames.

As shown in Table 2, SLAMP outperforms both methods on all of the metrics on both datasets, which shows its ability to generalize to the sequences with moving background. Even SVG [5] performs better than the state of

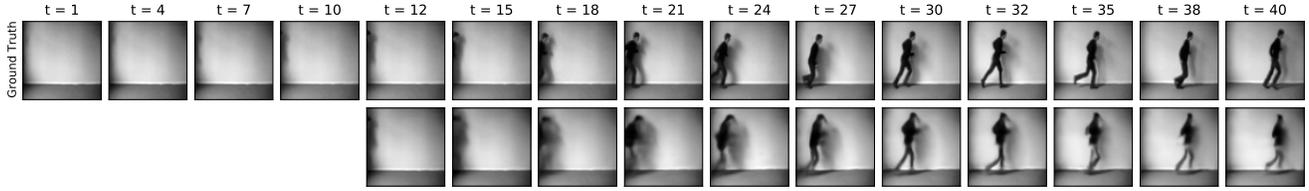


Figure 8: **Subject Appearing after the Conditioning Frames.** This figure shows a case where the subject appears after conditioning frames on KTH with ground truth (**top**) and a generated sample by our model (**bottom**). This shows our model’s ability to capture dynamics of the dataset by generating samples close to the ground truth, even conditioned on empty frames.



Figure 9: **Qualitative Comparison.** We compare SLAMP to SVG [5] and SRVP [9] on KITTI (**top**) and Cityscapes (**bottom**). Our model can better capture the changes due to ego-motion thanks to explicit modeling of motion history.

the art SRVP [9] in LPIPS metric for KITTI and on both SSIM and LPIPS for Cityscapes, which shows the limitations of SRVP on scenes with dynamic backgrounds. We also perform a qualitative comparison to these methods in Fig. 1 and Fig. 9. SLAMP can better preserve the scene structure thanks to explicit modeling of ego-motion history in the background.

Visualization of Latent Space: We visualize stochastic latent variables of the dynamic component on KTH compared to the static and SVG. (see Fig. 14 and Fig. 15)

5. Conclusion

We presented a stochastic video prediction framework to decompose video content into appearance and dynamic components. Our baseline model with deterministic motion and mask decoders outperforms SVG, which is a special case of our baseline model. Our model with motion history, SLAMP, further improves the results and reaches the performance of the state of the art method SRVP on the previously used datasets. Moreover, it outperforms both SVG and SRVP on two real-world autonomous driving datasets with dynamic background and complex motion. We show that motion his-

tory enriches model’s capacity to predict future, leading to better predictions in challenging cases.

Our model with motion history cannot realize its full potential in standard settings of stochastic video prediction datasets. A fair comparison is not possible on BAIR due to the little number of conditioning frames. BAIR holds a great promise with changing background but infrequent, small changes are not reflected in current evaluation metrics.

An interesting direction is stochastic motion decomposition, maybe with hierarchical latent variables, for modelling camera motion and motion of each object in the scene separately.

Acknowledgements. We would like to thank Jean-Yves Franceschi and Edouard Delasalles for providing technical and numerical details for the baseline performances, and Deniz Yuret for helpful discussions and comments. K. Akan was supported by KUIS AI Center fellowship, F. Güney by TUBITAK 2232 International Fellowship for Outstanding Researchers Programme, E. Erdem in part by GEBIP 2018 Award of the Turkish Academy of Sciences, A. Erdem by BAGEP 2021 Award of the Science Academy.

References

- [1] Mohammad Babaeizadeh, Chelsea Finn, Dumitru Erhan, Roy H. Campbell, and Sergey Levine. Stochastic variational video prediction. In *Proc. of the International Conf. on Learning Representations (ICLR)*, 2018. 2, 5, 6, 7, 17
- [2] Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. In *Advances in Neural Information Processing Systems*, 2015. 15
- [3] Lluís Castrejon, Nicolas Ballas, and Aaron Courville. Improved conditional vrns for video prediction. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 2
- [4] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 5, 7
- [5] Emily Denton and Rob Fergus. Stochastic video generation with a learned prior. In *Proc. of the International Conf. on Machine Learning (ICML)*, 2018. 1, 2, 3, 4, 5, 6, 7, 8, 11, 15, 16, 17, 19, 20, 28, 29
- [6] Frederik Ebert, Chelsea Finn, Alex X. Lee, and Sergey Levine. Self-supervised visual planning with temporal skip connections. In *1st Annual Conference on Robot Learning, CoRL 2017, Mountain View, California, USA, November 13-15, 2017, Proceedings*, 2017. 5, 7
- [7] Hehe Fan, Linchao Zhu, and Yi Yang. Cubic lstms for video prediction. In *Proc. of the Conf. on Artificial Intelligence (AAAI)*, 2019. 1, 2
- [8] Chelsea Finn, Ian Goodfellow, and Sergey Levine. Unsupervised learning for physical interaction through video prediction. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016. 1, 2
- [9] Jean-Yves Franceschi, Edouard Delasalles, Mickaël Chen, Sylvain Lamprier, and Patrick Gallinari. Stochastic latent residual video prediction. In *Proc. of the International Conf. on Machine Learning (ICML)*, 2020. 1, 2, 3, 5, 6, 7, 8, 15, 17, 28, 29
- [10] Hang Gao, Huazhe Xu, Qi-Zhi Cai, Ruth Wang, Fisher Yu, and Trevor Darrell. Disentangling propagation and generation for video prediction. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 1, 2
- [11] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The KITTI dataset. *International Journal of Robotics Research (IJRR)*, 2013. 1, 5, 7
- [12] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2012. 5, 7
- [13] Karol Gregor and Frederic Besse. Temporal difference variational auto-encoder. In *Proc. of the International Conf. on Learning Representations (ICLR)*, 2018. 3
- [14] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. 2018. 15
- [15] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and koray kavukcuoglu. Spatial transformer networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2015. 3
- [16] Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc V Gool. Dynamic filter networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016. 1, 2
- [17] Gunnar Johansson. Visual perception of biological motion and a model for its analysis. *Perception & Psychophysics*, 14(2):201–211, jun 1973. 1
- [18] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations*, 2015. 16
- [19] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *Proc. of the International Conf. on Learning Representations (ICLR)*, 2014. 4
- [20] Alex X. Lee, Richard Zhang, Frederik Ebert, Pieter Abbeel, Chelsea Finn, and Sergey Levine. Stochastic adversarial video prediction. *arXiv.org*, 2018. 2, 5, 6, 7, 17, 20
- [21] Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. Flow-grounded spatial-temporal video prediction from still images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 600–615, 2018. 2
- [22] Xiaodan Liang, Lisa Lee, Wei Dai, and Eric P. Xing. Dual motion gan for future-flow embedded video prediction. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2017. 2
- [23] Ziwei Liu, Raymond A. Yeh, Xiaoou Tang, Yiming Liu, and Aseem Agarwala. Video frame synthesis using deep voxel flow. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2017. 1, 2
- [24] William Lotter, Gabriel Kreiman, and David Cox. Deep predictive coding networks for video prediction and unsupervised learning. In *Proc. of the International Conf. on Learning Representations (ICLR)*, 2017. 1, 2
- [25] Chaochao Lu, Michael Hirsch, and Bernhard Scholkopf. Flexible spatio-temporal networks for video

- prediction. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, July 2017. 1, 2
- [26] Matthias Minderer, Chen Sun, Ruben Villegas, Forrester Cole, Kevin P Murphy, and Honglak Lee. Un-supervised learning of object structure and dynamics from videos. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 2
- [27] Alec Radford, Luke Metz, and Soumith Chintala. Un-supervised representation learning with deep convolutional generative adversarial networks. In *4th International Conference on Learning Representations*, 2016. 15
- [28] Christian Schüldt, Ivan Laptev, and Barbara Caputo. Recognizing human actions: A local svm approach. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2004. 5, 7
- [29] Xingjian Shi, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai kin Wong, and Wang chun Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2015. 2
- [30] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015. 15
- [31] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges. *arXiv.org*, 2019. 6
- [32] Ruben Villegas, Arkanath Pathak, Harini Kannan, Dumitru Erhan, Quoc V Le, and Honglak Lee. High fidelity video prediction with large stochastic recurrent neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 2
- [33] C. Vondrick and A. Torralba. Generating the future with adversarial transformers. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1, 2
- [34] Jacob Walker, Carl Doersch, Abhinav Gupta, and Martial Hebert. An uncertain future: Forecasting from static images using variational autoencoders. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2016. 2
- [35] Jacob Walker, Abhinav Gupta, and Martial Hebert. Dense optical flow prediction from a static image. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2015. 1, 2
- [36] Tianfan Xue, Jiajun Wu, Katherine L Bouman, and William T Freeman. Visual dynamics: Probabilistic future frame synthesis via cross convolutional networks. In *Advances In Neural Information Processing Systems*, 2016. 2
- [37] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 6

In this part, we provide additional illustrations, derivations, and results for our paper “SLAMP: Stochastic Latent Appearance and Motion Prediction”. We first show the model illustrations of our proposed model (SLAMP) and our baseline model (SLAMP-Baseline) in comparison to the previous work by Denton et al. [5] (SVG) in Section A. In Section B, we provide the full derivations of the variational inference, evidence lower bounds of our baseline model (Section B.1) and our proposed model (Section B.2). We explain the architectural choices and training details in Section C. In Section D, we present detailed versions of the quantitative results in the main paper. In addition, we present the ablation experiments for our model’s mask component. We evaluate and compare the predictions of static, dynamic heads of the model, simple averaging of the two without a mask, and our full model with learned mask. In Section E, we first provide the color wheel to interpret optical flow predictions, and a comparison of static and dynamic latent variables. We then present several qualitative results both with details as in the main paper and with random samples, showing the diversity of the generated samples on all datasets.

For video examples, please visit <https://kuis-ai.github.io/slamp/>.

A. Model Illustrations

In Fig. 10, we provide the inference procedure of our model SLAMP, in addition to the training procedure provided in the main paper. Moreover, we present graphical illustrations of the training (Fig. 11) and the inference procedures (Fig. 12) of our baseline model, SLAMP-Baseline, in comparison to SVG [5].

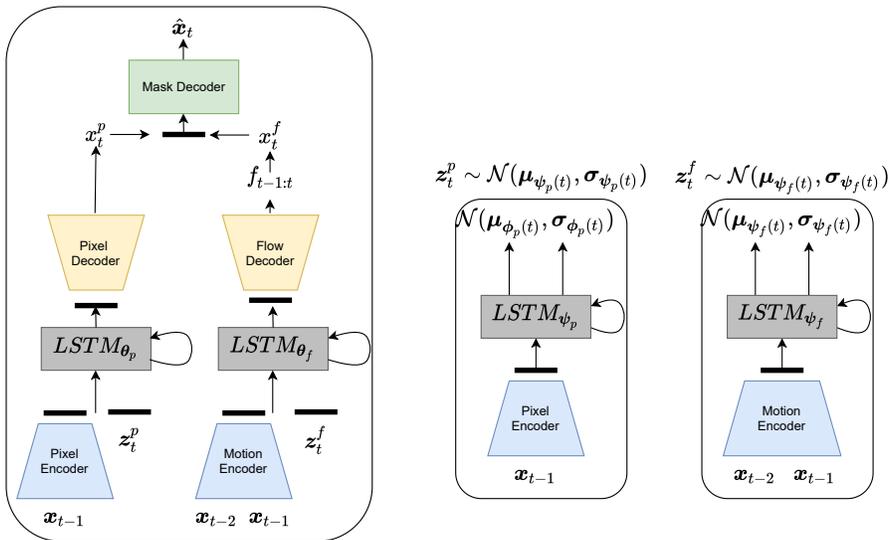


Figure 10: **Illustration of the Inference Procedure for SLAMP.** This figure illustrates the difference between the inference time and the train time in terms of the distributions the latent variables are sampled from. While at train time, latent variables are sampled from the posterior distribution, at test time, they are sampled from the posterior for the conditioning frames and from the prior for the following frames.

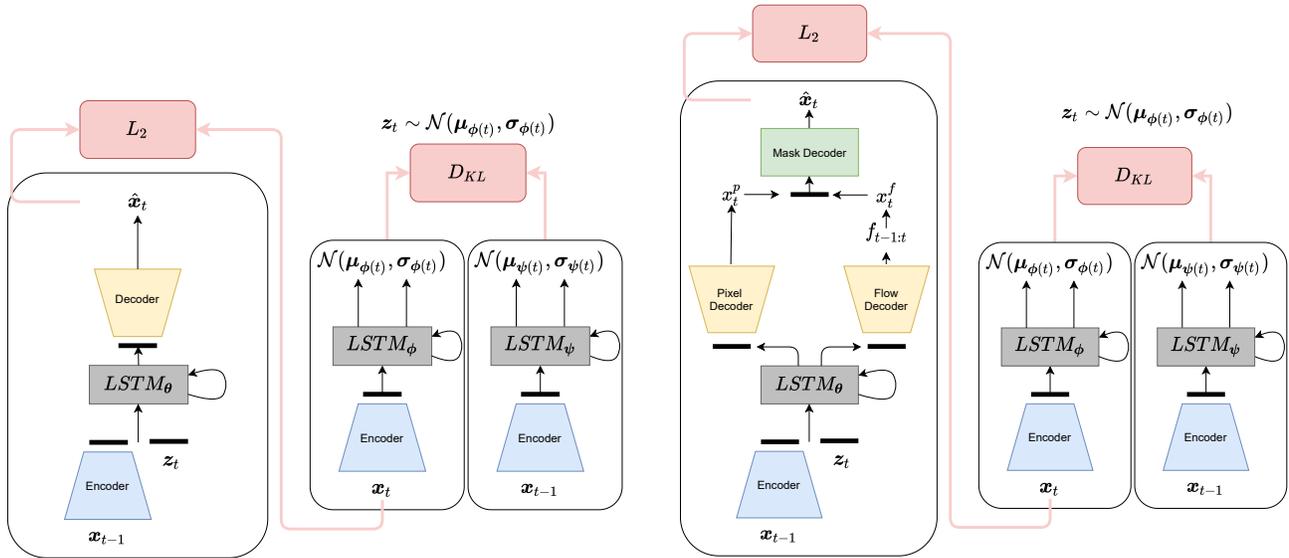


Figure 11: **Illustration of the Training Procedure for SVG (left) and SLAMP-Baseline (right).** The main difference between SVG and SLAMP-Baseline is that SLAMP-Baseline has three decoders instead of one pixel decoder. In SLAMP-Baseline, in addition to the appearance prediction \mathbf{x}_t^p , we also estimate flow $f_{t-1:t}$ and warp the previous frame according to the estimated flow to obtain motion prediction \mathbf{x}_t^f . Mask decoder takes appearance and motion predictions as input and generates a weighted combination of the two, $\hat{\mathbf{x}}_t$ as the final prediction. Note that SVG corresponds to only appearance prediction case of our baseline model.

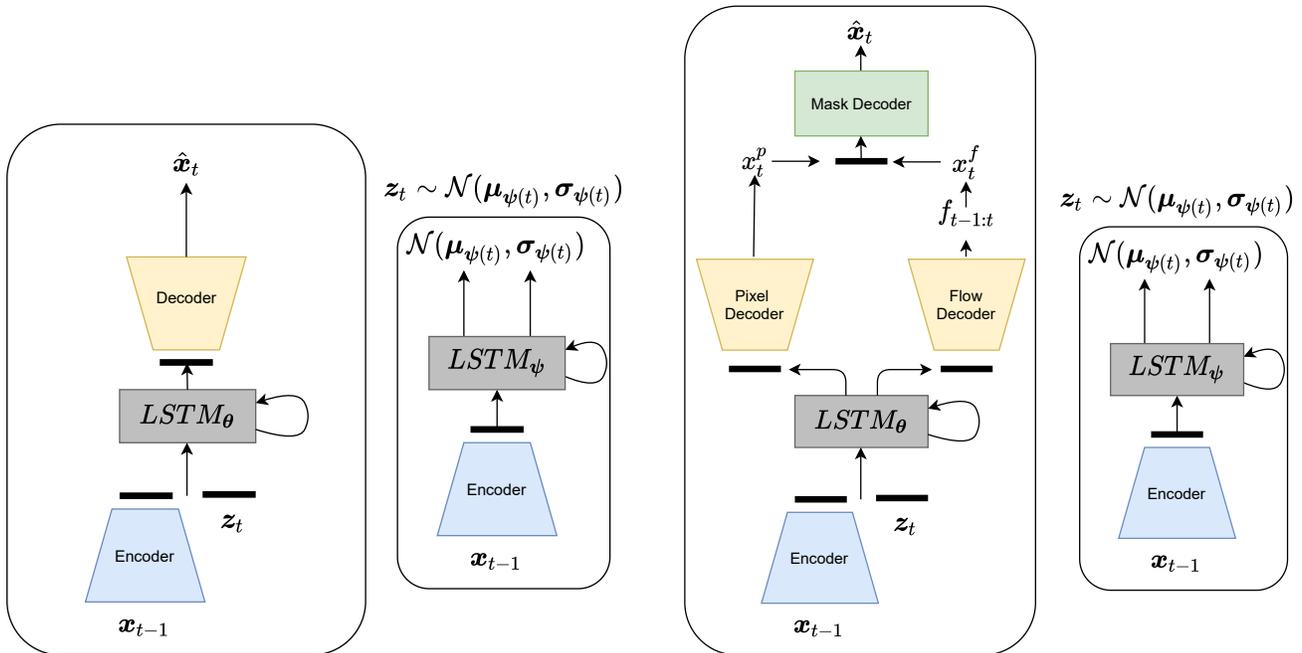


Figure 12: **Illustration of the Inference Procedure for SVG (left) and SLAMP-Baseline (right).** This figure illustrates the inference time in comparison to the train time in terms of the distribution the latent variables are sampled from. While at train time, latent variables are sampled from the posterior distribution, at test time, they are sampled from the posterior for the conditioning frames and from the prior for the following frames.

B. Derivations

Here, we provide derivations of inference steps and variational lower bounds of the baseline method, SLAMP-Baseline (Section B.1), and our method SLAMP (Section B.2).

B.1. Derivation of the ELBO for SLAMP-Baseline

We first derive the variational lower bound for the baseline model with one posterior and one learned prior distribution.

$$\begin{aligned}
\log p_{\theta}(\mathbf{x}) &= \log \int_{\mathbf{z}} p_{\theta}(\mathbf{x}|\mathbf{z}) p(\mathbf{z}|\mathbf{x}) \\
&= \log \int_{\mathbf{z}} p_{\theta}(\mathbf{x}|\mathbf{z}) p(\mathbf{z}|\mathbf{x}) \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \\
&= \log \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \frac{p_{\theta}(\mathbf{x}|\mathbf{z}) p(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \\
&\geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z}) p(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \\
&= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z}) - \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}|\mathbf{x})} \\
&= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z}) - D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}))
\end{aligned} \tag{7}$$

We model the posterior distribution with a recurrent network. The recurrent network outputs a different posterior distribution, $q_{\phi}(\mathbf{z}_t|\mathbf{x}_{1:t})$, at every time step. Due to independence of latent variables across time, $\mathbf{z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_T]$, we can derive the estimation of posterior distribution across time steps as follows:

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \prod_t q_{\phi}(\mathbf{z}_t|\mathbf{x}_{1:t}) \tag{8}$$

Since the latent variables, $\mathbf{z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_T]$, are independent across time, we can further decompose Kullback-Leibler term in the evidence lower bound into individual time steps:

$$\begin{aligned}
&D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}_{1:t-1})) \\
&= \int_{\mathbf{z}} q_{\phi}(\mathbf{z}|\mathbf{x}) \log \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}|\mathbf{x}_{1:t-1})} \\
&= \int_{\mathbf{z}_1} \dots \int_{\mathbf{z}_T} q_{\phi}(\mathbf{z}_1|\mathbf{x}_1) \dots q_{\phi}(\mathbf{z}_T|\mathbf{x}_{1:T}) \log \frac{q_{\phi}(\mathbf{z}_1|\mathbf{x}_1) \dots q_{\phi}(\mathbf{z}_T|\mathbf{x}_{1:T})}{p(\mathbf{z}_1|\mathbf{x}_1) \dots p(\mathbf{z}_T|\mathbf{x}_{1:T-1})} \\
&= \int_{\mathbf{z}_1} \dots \int_{\mathbf{z}_T} q_{\phi}(\mathbf{z}_1|\mathbf{x}_1) \dots q_{\phi}(\mathbf{z}_T|\mathbf{x}_{1:T}) \sum_t \log \frac{q_{\phi}(\mathbf{z}_t|\mathbf{x}_{1:t})}{p(\mathbf{z}_t|\mathbf{x}_{1:t-1})} \\
&= \sum_t \int_{\mathbf{z}_1} \dots \int_{\mathbf{z}_T} q_{\phi}(\mathbf{z}_1|\mathbf{x}_1) \dots q_{\phi}(\mathbf{z}_T|\mathbf{x}_{1:T}) \log \frac{q_{\phi}(\mathbf{z}_t|\mathbf{x}_{1:t})}{p(\mathbf{z}_t|\mathbf{x}_{1:t-1})}
\end{aligned} \tag{9}$$

And because $\int_x p(x) = 1$, this simplifies to:

$$\begin{aligned}
&= \sum_t \int_{\mathbf{z}_t} q_{\phi}(\mathbf{z}_t|\mathbf{x}_{1:t}) \log \frac{q_{\phi}(\mathbf{z}_t|\mathbf{x}_{1:t})}{p(\mathbf{z}_t|\mathbf{x}_{1:t-1})} \\
&= \sum_t D_{\text{KL}}(q_{\phi}(\mathbf{z}_t|\mathbf{x}_{1:t}) || p(\mathbf{z}_t|\mathbf{x}_{1:t-1}))
\end{aligned}$$

At each time step, our model predicts \mathbf{x}_t , conditioned on \mathbf{x}_{t-1} and \mathbf{z}_t . Since our model has recurrence connections, it considers not only \mathbf{x}_{t-1} and \mathbf{z}_t , but also $\mathbf{x}_{1:t-2}$ and $\mathbf{z}_{1:t-1}$. Therefore, we can further write our inference as:

$$\begin{aligned}\log p_{\theta}(\mathbf{x}|\mathbf{z}) &= \log \prod_t p_{\theta}(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t}) \\ &= \sum_t \log p_{\theta}(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t})\end{aligned}\quad (10)$$

Combining all of them leads to the following variational lower bound:

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &\geq \mathcal{L}_{\theta, \phi, \psi}(\mathbf{x}_{1:T}) \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z}) - D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\psi}(\mathbf{z}|\mathbf{x})) \\ &= \sum_t \left[\mathbb{E}_{q_{\phi}(\mathbf{z}_{1:t}|\mathbf{x}_{1:t})} \log p_{\theta}(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t}) \right. \\ &\quad \left. - D_{\text{KL}}(q_{\phi}(\mathbf{z}_t|\mathbf{x}_{1:t}) \parallel p_{\psi}(\mathbf{z}_t|\mathbf{x}_{1:t-1})) \right]\end{aligned}\quad (11)$$

B.2. Derivation of the ELBO for SLAMP

In this section, we derive the variational lower bound for the proposed model with two posterior and two learned prior distributions.

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &= \log \int_{\mathbf{z}^p} \int_{\mathbf{z}^f} p_{\theta}(\mathbf{x}|\mathbf{z}^p, \mathbf{z}^f) p(\mathbf{z}^p|\mathbf{x}) p(\mathbf{z}^f|\mathbf{x}) \\ &= \log \int_{\mathbf{z}^p} \int_{\mathbf{z}^f} p_{\theta}(\mathbf{x}|\mathbf{z}^p, \mathbf{z}^f) p(\mathbf{z}^p|\mathbf{x}) p(\mathbf{z}^f|\mathbf{x}) \frac{q_{\phi_p}(\mathbf{z}^p|\mathbf{x})}{q_{\phi_p}(\mathbf{z}^p|\mathbf{x})} \frac{q_{\phi_f}(\mathbf{z}^f|\mathbf{x})}{q_{\phi_f}(\mathbf{z}^f|\mathbf{x})} \\ &= \log \mathbb{E}_{\substack{\mathbf{z}^p \sim q_{\phi_p} \\ \mathbf{z}^f \sim q_{\phi_f}}} \frac{p_{\theta}(\mathbf{x}|\mathbf{z}^p, \mathbf{z}^f) p(\mathbf{z}^p|\mathbf{x}) p(\mathbf{z}^f|\mathbf{x})}{q_{\phi_p}(\mathbf{z}^p|\mathbf{x}) q_{\phi_f}(\mathbf{z}^f|\mathbf{x})} \\ &\geq \mathbb{E}_{\substack{\mathbf{z}^p \sim q_{\phi_p} \\ \mathbf{z}^f \sim q_{\phi_f}}} \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z}^p, \mathbf{z}^f) p(\mathbf{z}^p|\mathbf{x}) p(\mathbf{z}^f|\mathbf{x})}{q_{\phi_p}(\mathbf{z}^p|\mathbf{x}) q_{\phi_f}(\mathbf{z}^f|\mathbf{x})} \\ &= \mathbb{E}_{\substack{\mathbf{z}^p \sim q_{\phi_p} \\ \mathbf{z}^f \sim q_{\phi_f}}} \log p_{\theta}(\mathbf{x}|\mathbf{z}^p, \mathbf{z}^f) - \mathbb{E}_{\mathbf{z}^p \sim q_{\phi_p}} \log \frac{q_{\phi_p}(\mathbf{z}^p|\mathbf{x})}{p(\mathbf{z}^p|\mathbf{x})} - \mathbb{E}_{\mathbf{z}^f \sim q_{\phi_f}} \log \frac{q_{\phi_f}(\mathbf{z}^f|\mathbf{x})}{p(\mathbf{z}^f|\mathbf{x})} \\ &= \mathbb{E}_{\substack{\mathbf{z}^p \sim q_{\phi_p} \\ \mathbf{z}^f \sim q_{\phi_f}}} \log p_{\theta}(\mathbf{x}|\mathbf{z}^p, \mathbf{z}^f) - D_{\text{KL}}(q_{\phi_p}(\mathbf{z}^p|\mathbf{x}) \parallel p(\mathbf{z}^p|\mathbf{x})) - D_{\text{KL}}(q_{\phi_f}(\mathbf{z}^f|\mathbf{x}) \parallel p(\mathbf{z}^f|\mathbf{x}))\end{aligned}\quad (12)$$

We model the posterior distributions with two recurrent networks. The recurrent networks output two different posterior distributions, $q_{\phi_p}(\mathbf{z}_t^p|\mathbf{x}_{1:t})$ and $q_{\phi_f}(\mathbf{z}_t^f|\mathbf{x}_{1:t})$, at every time step. Due to the independence of the latent variables across time, $\mathbf{z}^p = [\mathbf{z}_1^p, \mathbf{z}_2^p, \dots, \mathbf{z}_T^p]$ and $\mathbf{z}^f = [\mathbf{z}_1^f, \mathbf{z}_2^f, \dots, \mathbf{z}_T^f]$, we can derive the estimation of posterior distributions across time steps as follows:

$$\begin{aligned}q_{\phi_p}(\mathbf{z}^p|\mathbf{x}) &= \prod_t q_{\phi_p}(\mathbf{z}_t^p|\mathbf{x}_{1:t}) \\ q_{\phi_f}(\mathbf{z}^f|\mathbf{x}) &= \prod_t q_{\phi_f}(\mathbf{z}_t^f|\mathbf{x}_{1:t})\end{aligned}\quad (13)$$

Since the latent variables, $\mathbf{z}^p = [\mathbf{z}_1^p, \mathbf{z}_2^p, \dots, \mathbf{z}_T^p]$ and $\mathbf{z}^f = [\mathbf{z}_1^f, \mathbf{z}_2^f, \dots, \mathbf{z}_T^f]$, are independent across time and independent from each other, we can further decompose Kullback-Leibler terms in the evidence lower bound into individual time steps as in Eq. (9).

At each time step, our model predicts \mathbf{x}_t , conditioned on \mathbf{x}_{t-1} , \mathbf{z}_t^p , \mathbf{z}_t^f . Since our model has recurrence connections, it considers not only \mathbf{x}_{t-1} , \mathbf{z}_t^p and \mathbf{z}_t^f , but also $\mathbf{x}_{1:t-2}$, $\mathbf{z}_{1:t-1}^p$ and $\mathbf{z}_{1:t-1}^f$. Therefore, we can further write our inference as:

$$\begin{aligned} \log p_{\theta}(\mathbf{x}|\mathbf{z}^p, \mathbf{z}^f) &= \log \prod_t p_{\theta}(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t}^p, \mathbf{z}_{1:t}^f) \\ &= \sum_t \log p_{\theta}(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t}^p, \mathbf{z}_{1:t}^f) \end{aligned} \quad (14)$$

Combining all of them leads to the following variational lower bound:

$$\begin{aligned} \log p_{\theta}(\mathbf{x}) &\geq \mathcal{L}_{\theta, \phi_p, \psi_f}(\mathbf{x}_{1:T}) \\ &= \mathbb{E}_{\mathbf{z}^p \sim q_{\phi_p}} \log p_{\theta}(\mathbf{x}|\mathbf{z}^p, \mathbf{z}^f) - D_{\text{KL}}(q_{\phi_p}(\mathbf{z}^p|\mathbf{x}) \parallel p_{\psi_p}(\mathbf{z}^p|\mathbf{x})) - D_{\text{KL}}(q_{\phi_f}(\mathbf{z}^f|\mathbf{x}) \parallel p_{\psi_f}(\mathbf{z}^f|\mathbf{x})) \\ &\quad \mathbf{z}^f \sim q_{\phi_f} \\ &= \sum_t \mathbb{E}_{\mathbf{z}^p \sim q_{\phi_p}} \log p_{\theta}(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{z}_{1:t}^p, \mathbf{z}_{1:t}^f) \\ &\quad \mathbf{z}^f \sim q_{\phi_f} \\ &\quad - D_{\text{KL}}(q_{\phi_p}(\mathbf{z}_t^p|\mathbf{x}_{1:t}) \parallel p_{\psi_p}(\mathbf{z}_t^p|\mathbf{x}_{1:t-1})) \\ &\quad - D_{\text{KL}}(q_{\phi_f}(\mathbf{z}_t^f|\mathbf{x}_{1:t}) \parallel p_{\psi_f}(\mathbf{z}_t^f|\mathbf{x}_{1:t-1})) \end{aligned} \quad (15)$$

C. Training Details

We provide training details including scheduled sampling (Section C.1), architecture details (Section C.2), and the hyper-parameters used in the optimization (Section C.3).

C.1. Scheduled Sampling

Scheduled sampling proposed for sequence prediction [2] has been proven useful for several tasks where predictions need to be made based on the generated results from the previous time steps. We also experiment with scheduled sampling as part of our training procedure. Scheduled sampling prevents the model from conditioning on ground-truth perfect samples which are not available at test time. This is achieved by allowing the model to slowly encounter generated samples instead of ground-truth perfect samples. The ratio of ground-truth perfect samples over generated samples is decreased throughout the training. Specifically, we apply inverse sigmoid decay. We report the scores with and without scheduled sampling for the proposed models, both SLAMP-Baseline and SLAMP, on all datasets. As can be seen from Table 3, scheduled sampling is not crucial but it improves the results on KTH, especially for SLAMP.

C.2. Architecture Details

Encoders and Decoders: For all encoders and decoders, we use the same architectures as the previous work [5, 9]: a DCGAN [27] generator and discriminator for MNIST, and a VGG16 architecture [30] for KTH and BAIR datasets. In all datasets, we encode the image into an appearance feature vector of size $h_{\text{appearance}} = 128$ and the two consecutive images into a motion feature vector of size $h_{\text{motion}} = 128$. Compared to SVG, there are two more decoders for predicting flow and mask in our models. For flow decoder, we use the same decoder with two output channels representing motion in horizontal and vertical direction. See below for the details of the mask decoder. For SLAMP, motion encoder takes concatenated frame pair as input and outputs a feature vector encoding motion from one frame to the next.

Similar to previous work [5, 9], we also use skip connections but with a minor modification. In the previous work, the skip connection from either the last conditioning frame or last generated frame is used. Instead, we take the running average of all the skip connections from seen or generated frames. For example, at time step 15, we use the average of previous 14 skip connections that are generated.

Mask Predictor: For mask predictor, we use a 5-layer CNN with 2 Squeeze and Excitation layers (SE-Layer) [14] after each two convolutional layers. In the CNN, we use 64-channel filters at each layer and do not reduce the resolution by using 3×3 kernels with padding. We simply concatenate the output of pixel decoder and warped prediction along their channel axis and feed it into mask predictor which outputs a one-channel image. We apply sigmoid at the end to map the output to the range between 0 and 1, representing the weight to combine the appearance and the motion predictions.

Table 3: **Ablation Study.** This table shows the quantitative results comparing SLAMP-Baseline and SLAMP with scheduled sampling (+ SS) and without during training. Following the previous work, we report the results as the mean and the 95%-confidence interval in terms of PSNR, SSIM, and LPIPS on all the datasets except LPIPS on MNIST.

	Models	PSNR	SSIM	LPIPS
MNIST	SLAMP-Baseline	16.83 ± 0.06	0.7537 ± 0.0017	—
	SLAMP-Baseline + SS	16.32 ± 0.06	0.7343 ± 0.0016	—
	SLAMP	18.07 ± 0.07	0.7736 ± 0.0019	—
	SLAMP + SS	17.54 ± 0.08	0.7567 ± 0.0018	—
KTH	SLAMP-Baseline	28.47 ± 0.27	0.8527 ± 0.0053	0.0896 ± 0.0038
	SLAMP-Baseline + SS	29.20 ± 0.28	0.8633 ± 0.0048	0.0951 ± 0.0036
	SLAMP	28.91 ± 0.28	0.8604 ± 0.0049	0.0860 ± 0.0037
	SLAMP + SS	29.39 ± 0.30	0.8646 ± 0.0049	0.0795 ± 0.0033
BAIR	SLAMP-Baseline	19.60 ± 0.26	0.8175 ± 0.0083	0.0596 ± 0.0031
	SLAMP-Baseline + SS	19.55 ± 0.26	0.8171 ± 0.0083	0.0634 ± 0.0034
	SLAMP	19.67 ± 0.26	0.8161 ± 0.0086	0.0639 ± 0.0037
	SLAMP + SS	19.75 ± 0.26	0.8160 ± 0.0084	0.0661 ± 0.0035

LSTMs and Latent Variables: For prior, posterior, and frame predictor LSTMs, we use the settings proposed in SVG [5]. All LSTMs have 256 neurons and all prior and posterior LSTMs have one layer whereas the frame predictor LSTMs have two layers. For the size of the latent variables, we use 20, 50, 64 for MNIST, KTH, and BAIR, respectively.

C.3. Optimization Hyper-Parameters

All the models are trained with Adam optimizer [18], with decay rates $\beta_1 = 0.9$ and $\beta_2 = 0.999$. We train each model for 300 epochs where each epoch consists of 1000 updates, unless otherwise is specified. We take the model which performs the best in the validation set. We will share the trained models upon publication for replicating the results. Dataset-specific parameters for each dataset are as follows:

MNIST: The batch size is chosen to be 32, learning rate is 3×10^{-4} and $\beta = 1 \times 10^{-4}$. We continue training the models on MNIST with a lower learning rate, 1×10^{-5} , a lower $\beta = 5 \times 10^{-5}$, and a lower $\beta_1 = 0.7$.

KTH: The batch size is chosen to be 20, learning rate is 1×10^{-4} and $\beta = 1 \times 10^{-6}$. We apply scheduled sampling with inverse sigmoid decay.

BAIR: The batch size is chosen to be 20, learning rate is 1×10^{-4} and $\beta = 1 \times 10^{-4}$.

Training details for KITTI and Cityscapes:

We use 92×310 image resolution for KITTI and 128×256 for Cityscapes. We replaced LSTMs with ConvLSTMs and used 3×10 intermediate feature size for KITTI, 4×8 for Cityscapes. We used a shared encoder to downsample the image first and then, use two separate encoders for pixel and motion encoders to make model less powerful. We increased the number of layers in the shared encoder to downsample the higher resolution image, and preserve the VGG-based structure.

For SVG, we use the same settings as SLAMP. For SRVP, we use the same shared encoder and use a channel pooling at the end to make the convolutional feature vector compatible with the rest of the architecture.

We train all the models until the models see $2.4M$ video samples. We use the largest batch size that we could use and choose the learning rate 1×10^{-4} for all the models.

D. Detailed Quantitative Results

In this section, we provide a detailed version of the quantitative results presented in the main paper in Figure 3 and 5.

We compare the performance of SLAMP-Baseline and SLAMP to the previous work in terms of PSNR, SSIM, and LPIPS averaged over all time steps on MNIST (Table 4), KTH (Table 5), and BAIR (Table 6) datasets. Confirming the results in the main paper, the proposed model SLAMP with motion history outperforms both the baseline model, SLAMP-Baseline, and the previous work [5, 1, 20] and performs comparably to the state of the art model SRVP [9]. See the main paper for a detailed analysis.

Table 4: **Results on MNIST.** This table compares the results of SLAMP and SLAMP-Baseline to the previous work on MNIST dataset. Following the previous work, we report the results as the mean and the 95%-confidence interval in terms of PSNR and SSIM. Bold and underlined scores indicate the best and the second best performing method, respectively.

Models	PSNR	SSIM
SVG [5]	14.50 ± 0.04	0.7090 ± 0.0015
SRVP [9]	<u>16.93 ± 0.07</u>	0.7799 ± 0.0020
SLAMP-Baseline	16.83 ± 0.06	0.7537 ± 0.0018
SLAMP	18.07 ± 0.08	<u>0.7736 ± 0.0019</u>

Table 5: **Results on KTH.** This table compares the results of SLAMP and SLAMP-Baseline to the previous work on KTH dataset. Following the previous work, we report the results as the mean and the 95%-confidence interval in terms of PSNR, SSIM, and LPIPS. Bold and underlined scores indicate the best and the second best performing method, respectively.

Models	PSNR	SSIM	LPIPS
SV2P [1]	28.19 ± 0.31	0.8141 ± 0.0050	0.2049 ± 0.0053
SAVP [20]	26.51 ± 0.29	0.7564 ± 0.0062	0.1120 ± 0.0039
SVG [5]	28.06 ± 0.29	0.8438 ± 0.0054	0.0923 ± 0.0038
SRVP [9]	29.69 ± 0.32	0.8697 ± 0.0046	0.0736 ± 0.0029
SLAMP-Baseline	29.20 ± 0.28	0.8633 ± 0.0048	0.0951 ± 0.0036
SLAMP	<u>29.39 ± 0.30</u>	<u>0.8646 ± 0.0050</u>	<u>0.0795 ± 0.0034</u>

Table 6: **Results on BAIR.** This table compares the results of SLAMP and SLAMP-Baseline to the previous work on BAIR dataset. Following the previous work, we report the results as the mean and the 95%-confidence interval in terms of PSNR, SSIM, and LPIPS. Bold and underlined scores indicate the best and the second best performing method, respectively.

Models	PSNR	SSIM	LPIPS
SV2P [1]	20.39 ± 0.27	0.8169 ± 0.0086	0.0912 ± 0.0053
SAVP [20]	18.44 ± 0.25	0.7887 ± 0.0092	0.0634 ± 0.0026
SVG [5]	18.95 ± 0.26	0.8058 ± 0.0088	0.0609 ± 0.0034
SRVP [9]	19.59 ± 0.27	0.8196 ± 0.0084	0.0574 ± 0.0032
SLAMP-Baseline	19.60 ± 0.26	<u>0.8175 ± 0.0084</u>	<u>0.0596 ± 0.0032</u>
SLAMP	<u>19.67 ± 0.26</u>	0.8161 ± 0.0086	0.0639 ± 0.0037

In addition, we provide detailed results corresponding to the components of our model. We evaluate the result of the static head, the dynamic head, and simply their average without the learned mask and compare them to our full model with the learned mask in Table 7. Our full model, SLAMP, performs the best in all datasets according to all three evaluation metrics by combining the two predictions according to the mask prediction.

Table 7: **Mask Ablation Study.** This table shows the quantitative results comparing the components of SLAMP model. **Static** refers to the evaluation of the prediction of the static head directly, **Dynamic** refers to the evaluation of the prediction of the dynamic head which uses optical flow to predict the next frame, **Average** refers to average of static and dynamic heads without using the mask prediction. The last row for each dataset show the results of our model which uses the mask prediction to fuse the static and dynamic predictions. Following the previous work, we report the results as the mean and the 95%-confidence interval in terms of PSNR, SSIM, and LPIPS on all the datasets except LPIPS on MNIST.

	Models	PSNR	SSIM	LPIPS
MNIST	SLAMP - Static	15.03 ± 0.04	0.7273 ± 0.0014	—
	SLAMP - Dynamic	17.64 ± 0.09	0.7639 ± 0.0019	—
	SLAMP - Average	15.96 ± 0.05	0.7377 ± 0.0017	—
	SLAMP	18.07 ± 0.07	0.7736 ± 0.0019	—
KTH	SLAMP - Static	28.12 ± 0.28	0.8410 ± 0.0056	0.0844 ± 0.0037
	SLAMP - Dynamic	16.14 ± 0.14	0.7614 ± 0.0064	0.3689 ± 0.0080
	SLAMP - Average	21.61 ± 0.11	0.8359 ± 0.0048	0.2039 ± 0.0054
	SLAMP	29.39 ± 0.30	0.8646 ± 0.0049	0.0795 ± 0.0033
BAIR	SLAMP - Static	19.38 ± 0.26	0.8119 ± 0.00837	0.0606 ± 0.0033
	SLAMP - Dynamic	16.65 ± 0.20	0.7643 ± 0.0010	0.1176 ± 0.0048
	SLAMP - Average	18.51 ± 0.20	0.8073 ± 0.0088	0.0897 ± 0.0044
	SLAMP	19.75 ± 0.26	0.8160 ± 0.0084	0.0661 ± 0.0035

E. Additional Visualizations and Qualitative Results

E.1. Optical Flow Visualization

Fig. 13 shows the color wheel used to visualize the optical flow with false coloring. Colors show the direction of motion and the intensity of color in the visualizations show the magnitude of motion, i.e. intense colors for large motions. By following the usual practice in optical flow, we predict flow from the target frame to the current frame and apply inverse warping to obtain the target frame. Therefore, the direction of motion is also inverse, i.e. the opposite direction on the wheel shows the motion from the current frame to the target frame.

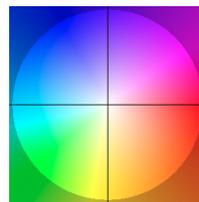


Figure 13: **Optical Flow False Coloring.** Colors on the wheel indicate the direction of motion in 2D.

E.2. Comparison of Static and Dynamic Latent Variables

In Fig. 8 of the main paper, we provide a visualization of stochastic latent variables of the dynamic component on KTH using t-SNE. Here, we provide both the static and the dynamic components for a comparison. The same colors from the main paper show the semantic classes of video frames plotted. As can be seen from Fig. 14, static variables on the right are more scattered and do not form clusters according to semantic classes as in the dynamic variables on the left (and in the main paper). This shows that our model can learn video dynamics according to semantic classes with separate modelling of the dynamic component.

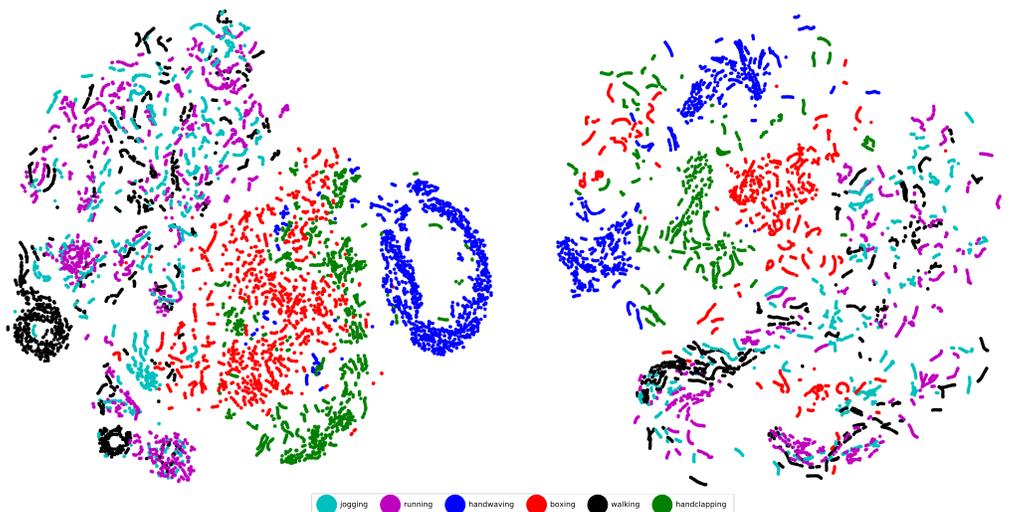


Figure 14: **Dynamic (left) vs. Static (right) Latent Variables.** This figure shows the T-SNE visualization of dynamic and static latent variables on 300 test videos from KTH dataset. In dynamic latent variables, different classes with similar repetitive movements such as walking, running, and jogging are clustered together. However, in static latent variables, points are more scattered and do not form clusters according to semantic actions.

In Fig. 15, we visualize the latent variables of SVG [5] to show the difference between our architecture’s latent variables and SVG’s latent variables. In our architecture, dynamic branch learns the similar repetitive movements whereas static branch learns the general image information. Therefore, dynamic latent variables form cluster around similar movements. However, in SVG, there is only one branch to predict the future frames, which only encodes the general image information rather than motion cues. Therefore, samples do not form clusters according to semantic classes as in our case for dynamic latent variables shown on the left in Fig. 14.

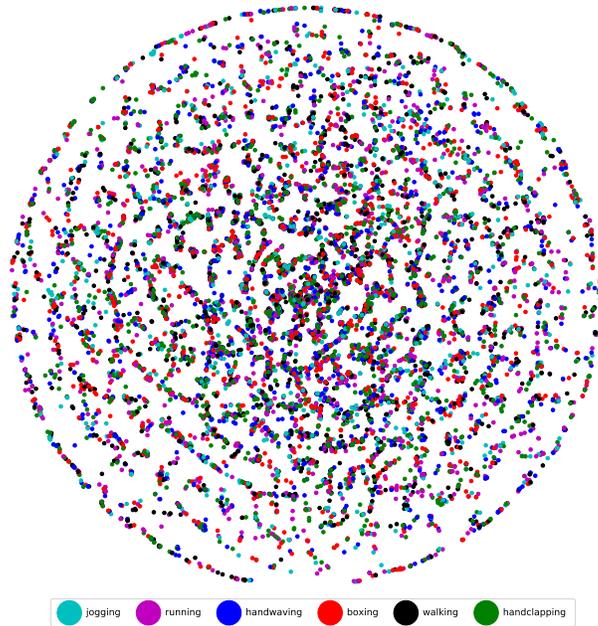


Figure 15: **Latent Variables of SVG.** This figure shows the T-SNE visualization of latent variables of SVG [5] on 300 test videos from KTH dataset. The latent variables form a normal distribution when visualized with T-SNE because SVG method uses a standard normal distribution as a fixed prior in KTH dataset. The model learns the general image information instead of motions groupings.

E.3. Diversity of Generated Samples

As proposed in SAVP [20], as a measure of diversity, we visualize the average over 100 generated samples. According to this measure, if a model is able to generate diverse results, generated samples should differ where there is motion, e.g. a moving object appearing in different positions and moving in different directions, leading to blurring out of the moving object. Therefore, we expect to see the background without moving objects in the average of the generated samples. The average samples confirm this for our model as shown for MNIST Fig. 16, KTH Fig. 17, and BAIR Fig. 18.

There is a special case on KTH which further supports our diversity claim as shown in Fig. 7 of the main paper. When subject appears after conditioning frames, our model can handle stochasticity of this challenging case and can generate **diverse** sequences. We show the best prediction and three random predictions in Fig. 19. Generated samples differ in terms of pose and speed of the subject as well as the time step that the subject appears.

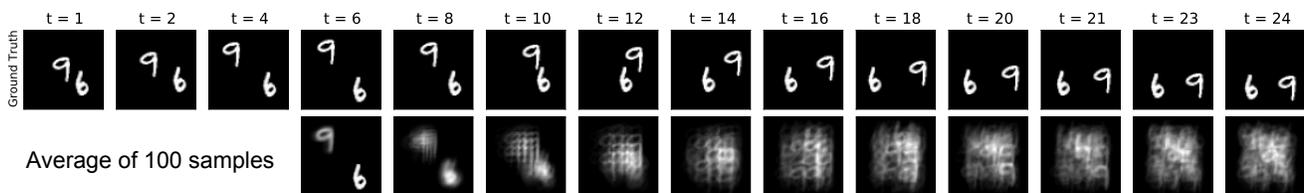


Figure 16: **Diversity on MNIST.** After a digit hits the wall, it can move in any direction. Our model successfully models the stochasticity of this case and generates diverse results, resulting in blurry average images.

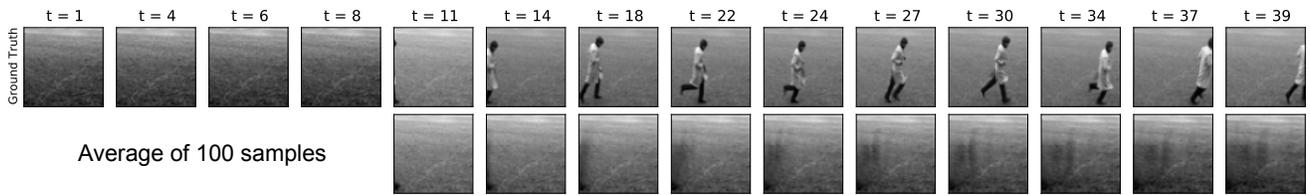


Figure 17: **Diversity on KTH.** Since the running person appears after the conditioning frames, the model should generate different results for each sample. The average of the generated samples does not contain any human because our model can generate diverse results, e.g. person in various poses appearing at different time steps with different speed.

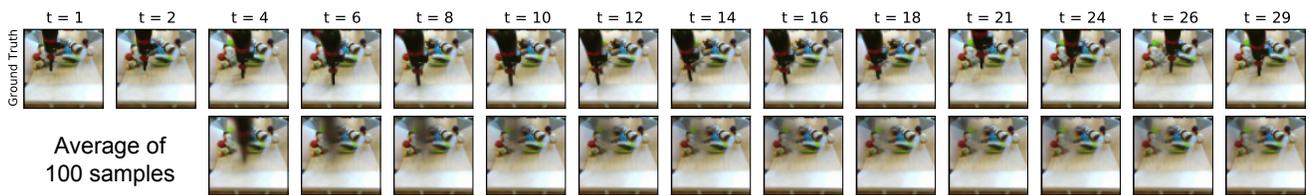


Figure 18: **Diversity on BAIR.** The robot hand can move in any direction at each time step, therefore the generated samples should differ from each other in terms of the position of the robot hand. The moving robot hand becomes invisible in the average images after the first few frames, which is an indication of the diversity of the generated samples.

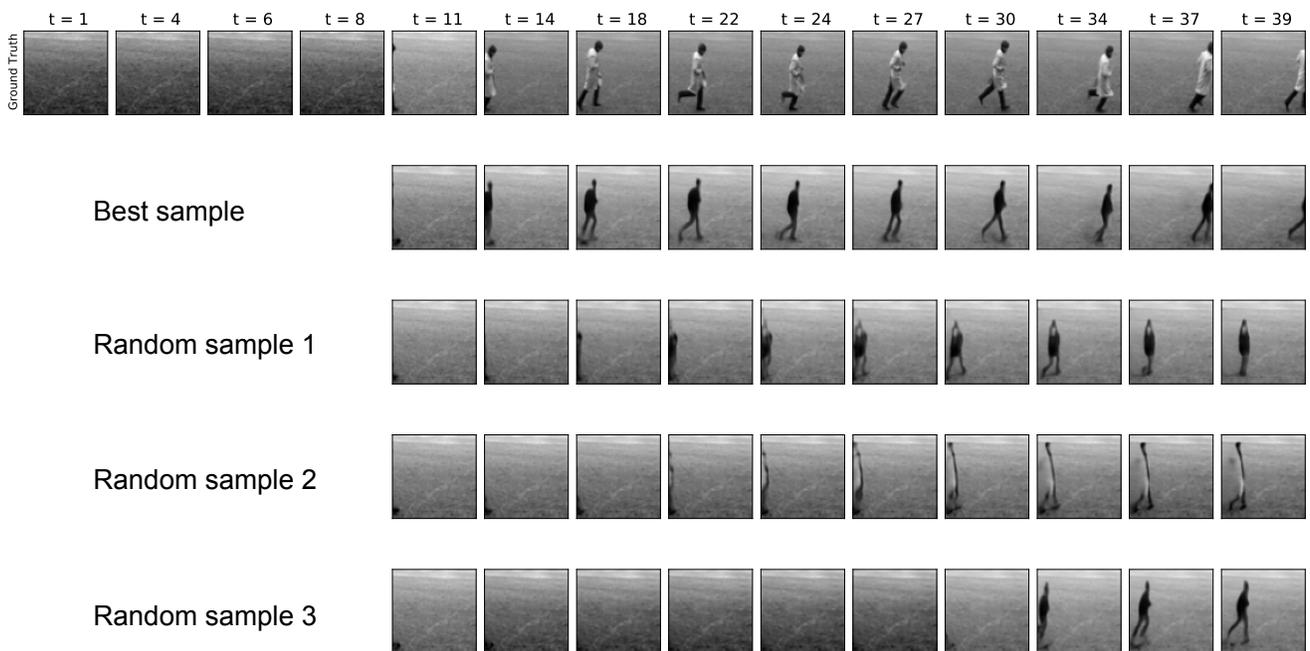


Figure 19: **Subject Appearing after Conditioning Frames on KTH.** We show the best sample and three random samples for the case where the subject appears after conditioning frames. Our model generates different results at each random sample by learning dataset dynamics.

E.4. Additional Qualitative Results

For each dataset, we show random examples with details of the best sample as well as three random samples generated. The detailed visualizations show appearance and motion prediction separately as well as the mask prediction and optical flow with false coloring. Random sample visualizations show the best sample and three random samples. We also provide full sequences of the samples in Fig. 1 in Fig. 32 and 33.

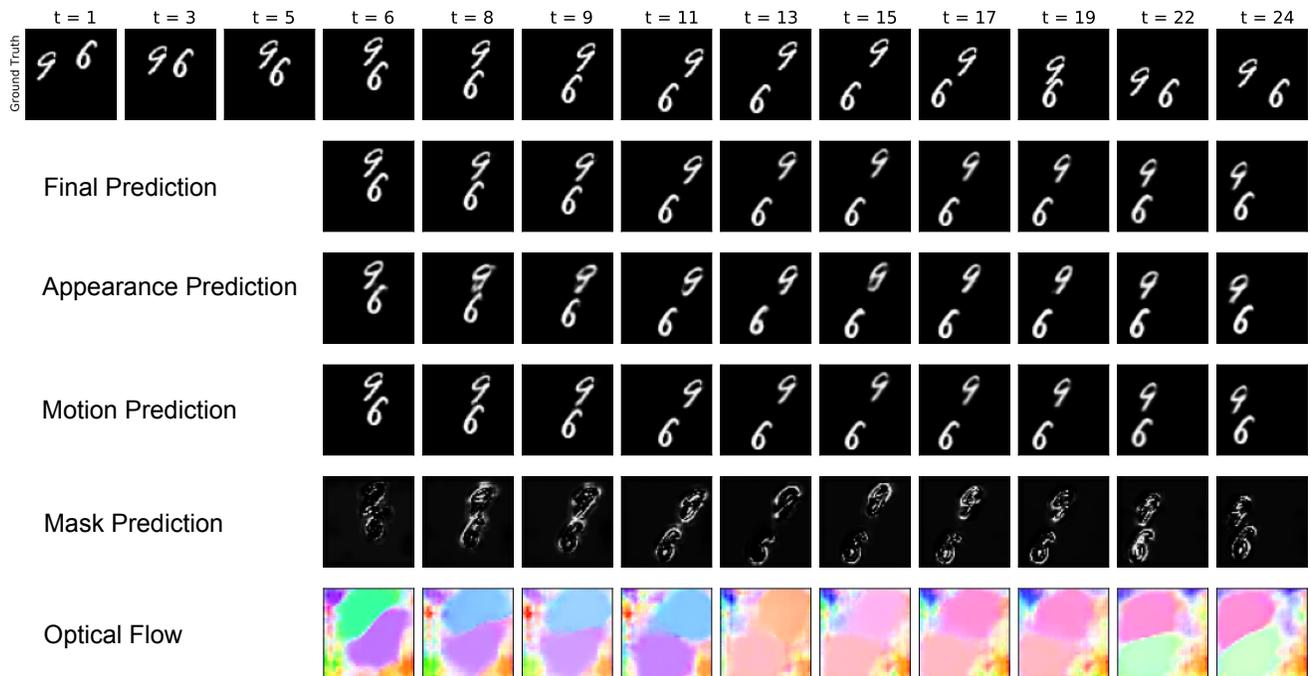


Figure 20: **No Overlapping Digits.** This figure show a regular case with two non-overlapping digits. Note that predicted flow is from the target frame to the current frame since we apply inverse warping. The correctness of the optical flow estimation can be verified by inspecting Fig. 13.

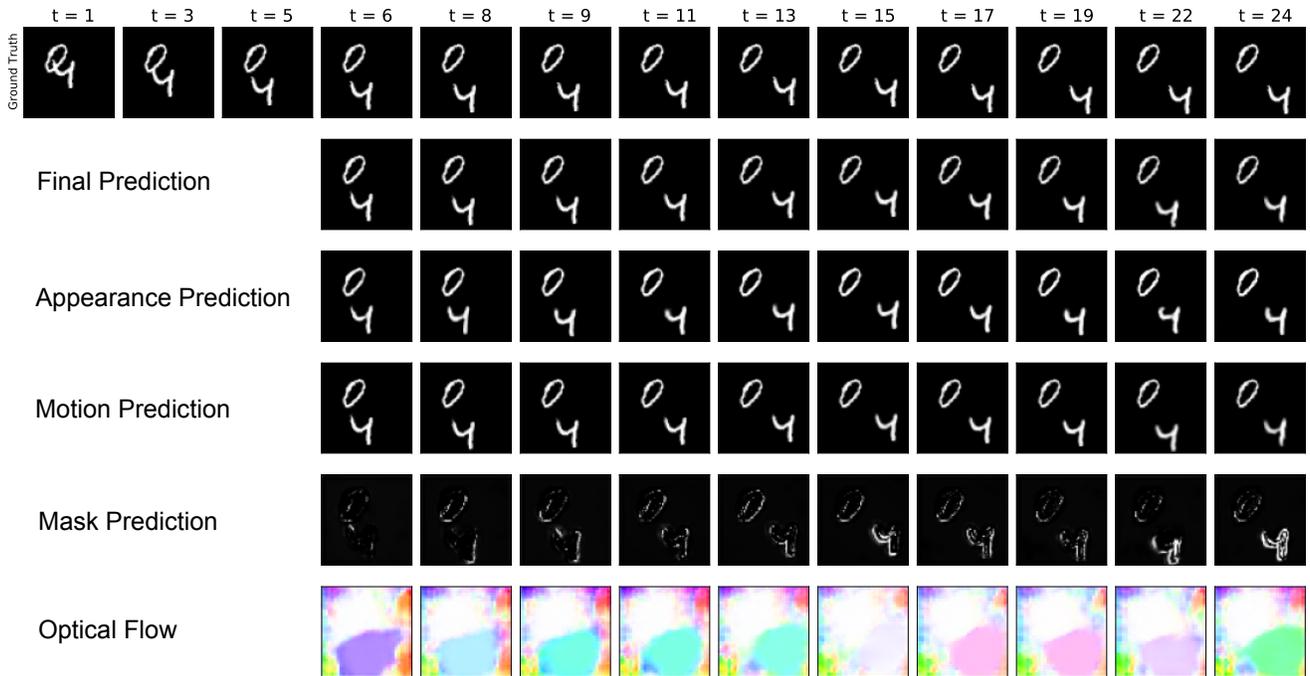


Figure 21: **Stationary Digit.** This figure shows a case where a digit, i.e. 0, is not moving. As can be seen from the last row, optical flow is correctly estimated as zero for that digit.

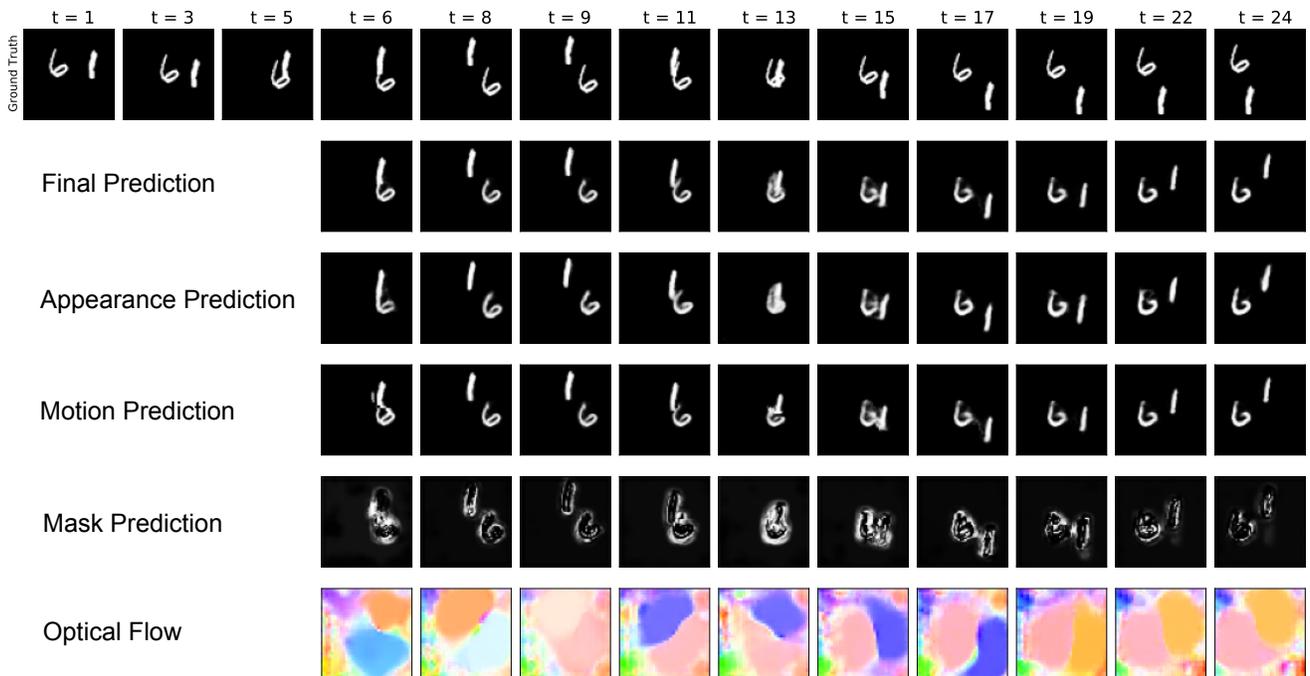


Figure 22: **Overlapping Digits.** This figure shows a challenging case where two digits cross each other and continue moving. The digits start overlapping at around $t = 13$. Our model can successfully handle this challenging case by preserving the appearance of digits.

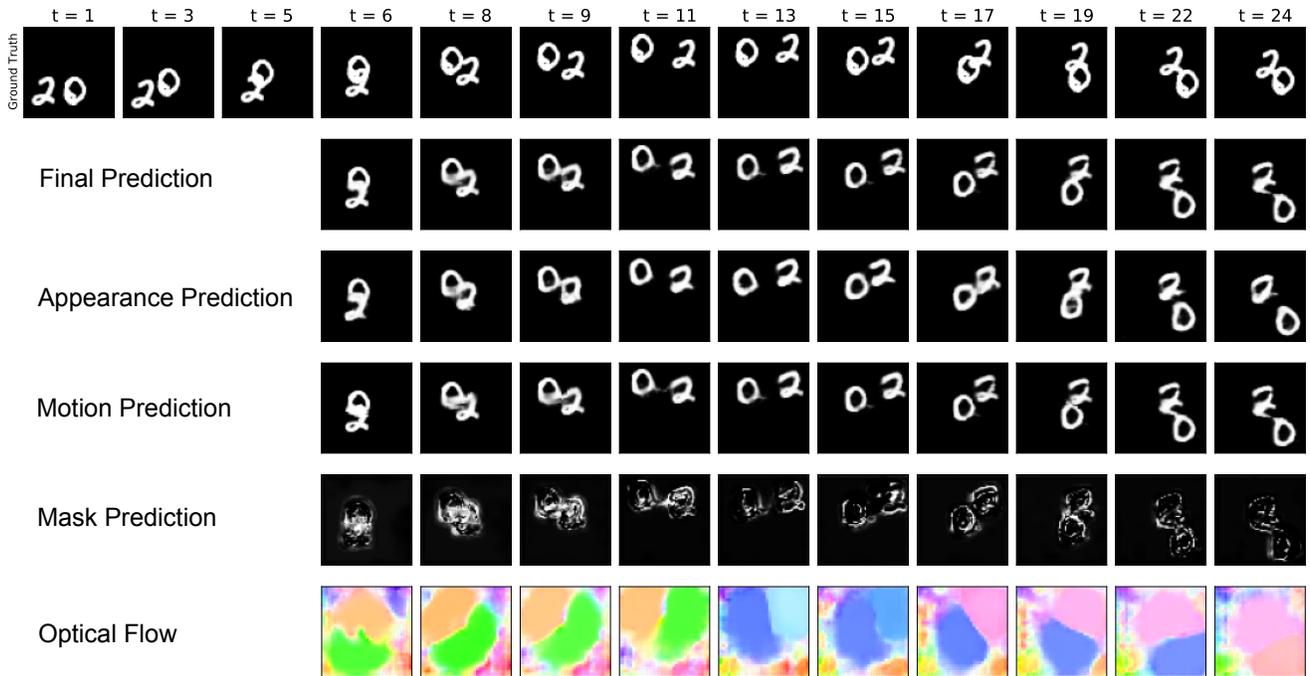


Figure 23: **Overlapping Digits.** This figure shows a challenging case where two digits cross each other and continue moving. The digits start overlapping at around $t = 6$. Our model can successfully handle this challenging case by preserving the appearance of digits.

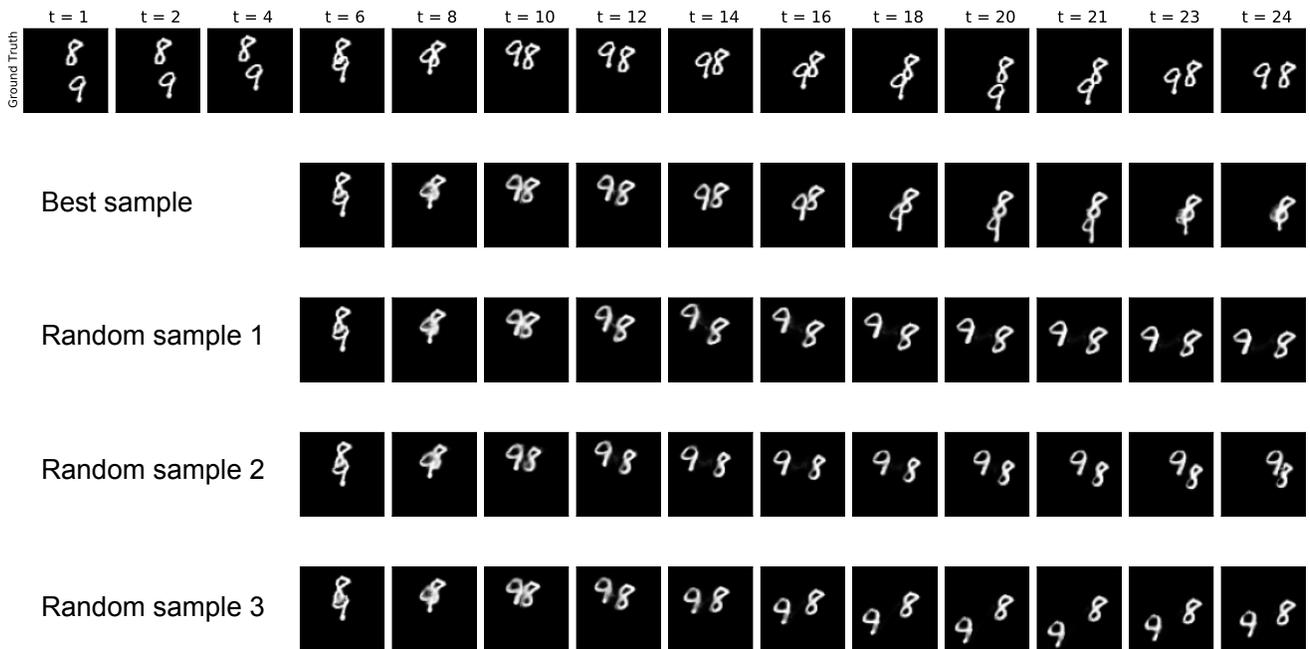


Figure 24: **Random Samples.** We show the best sample and three random samples generated. All of the predictions are sharp-looking and different than each other, which proves that our model can generate diverse results.

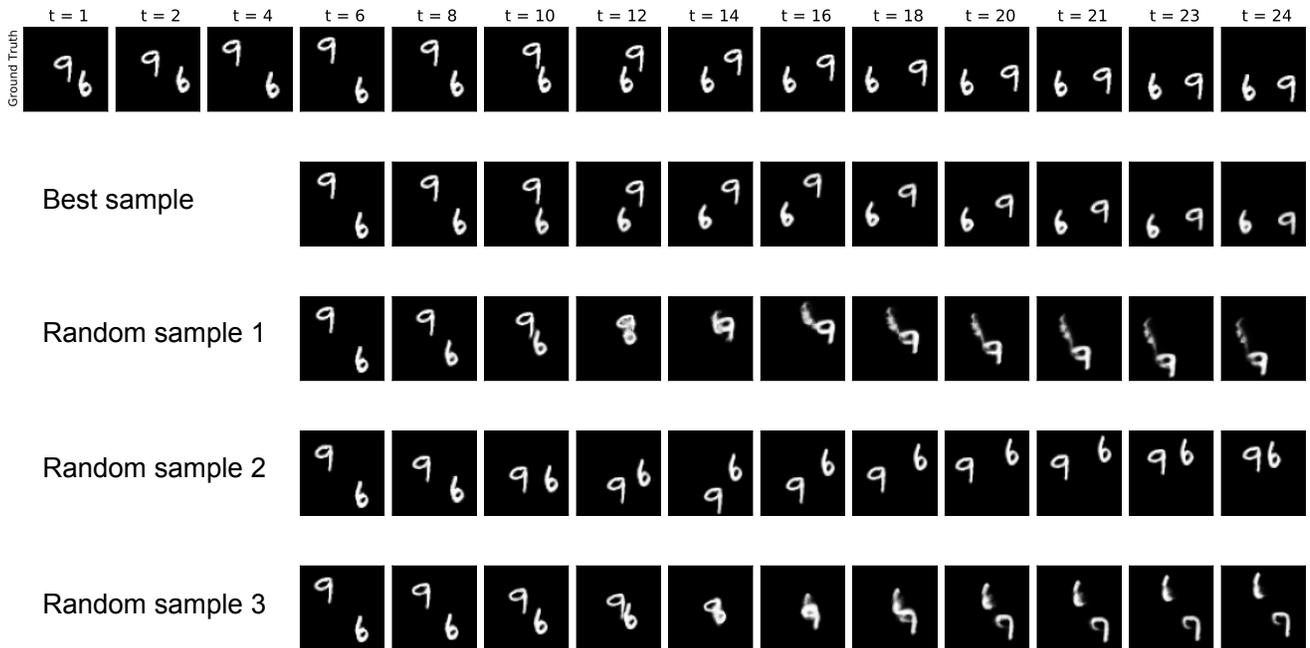


Figure 25: **Random Samples.** We show the best sample and three random samples generated. The first and the third samples cannot preserve the shape of the digits, however, the best sample and the second sample are still sharp-looking.

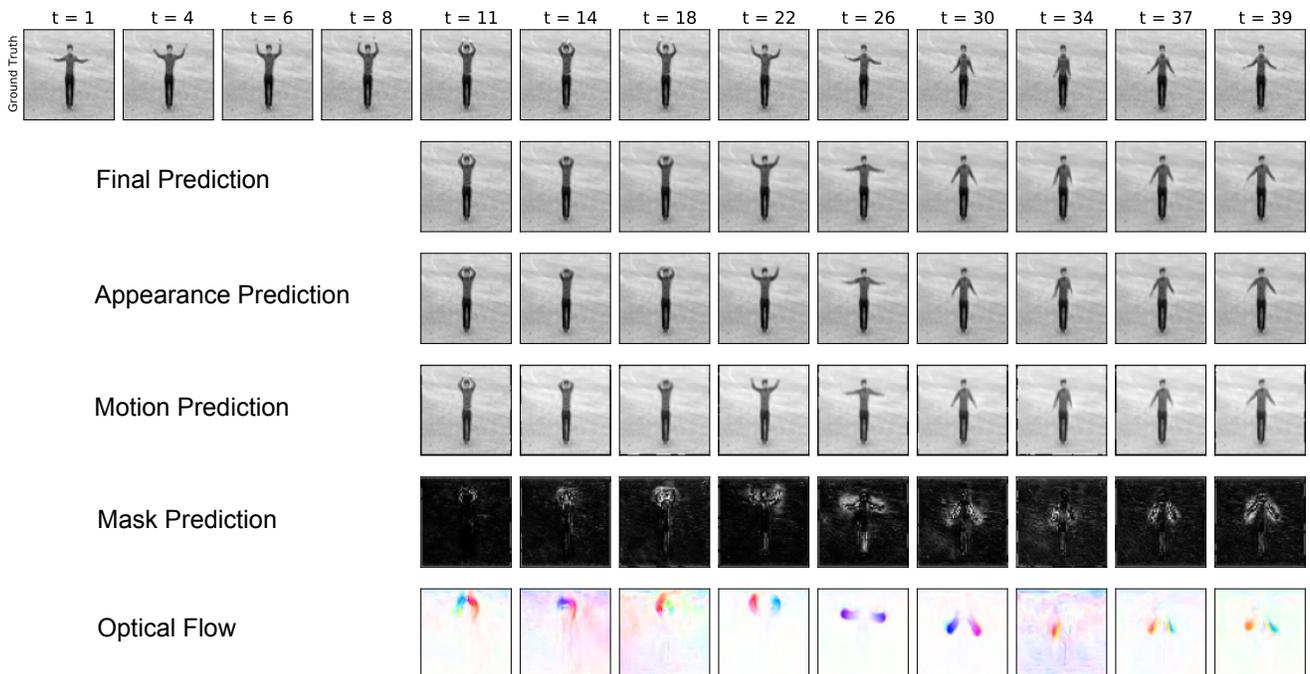


Figure 26: **Person Waving Hands.** This figure shows the challenging case of hand waving with fast motion in a small region. Note that predicted flow is from the target frame to the current frame since we apply inverse warping. The correctness of the optical flow estimation can be verified by checking Fig. 13. Our model focuses on the motion prediction for moving hands but it recovers the occluded motion boundaries from the appearance prediction.

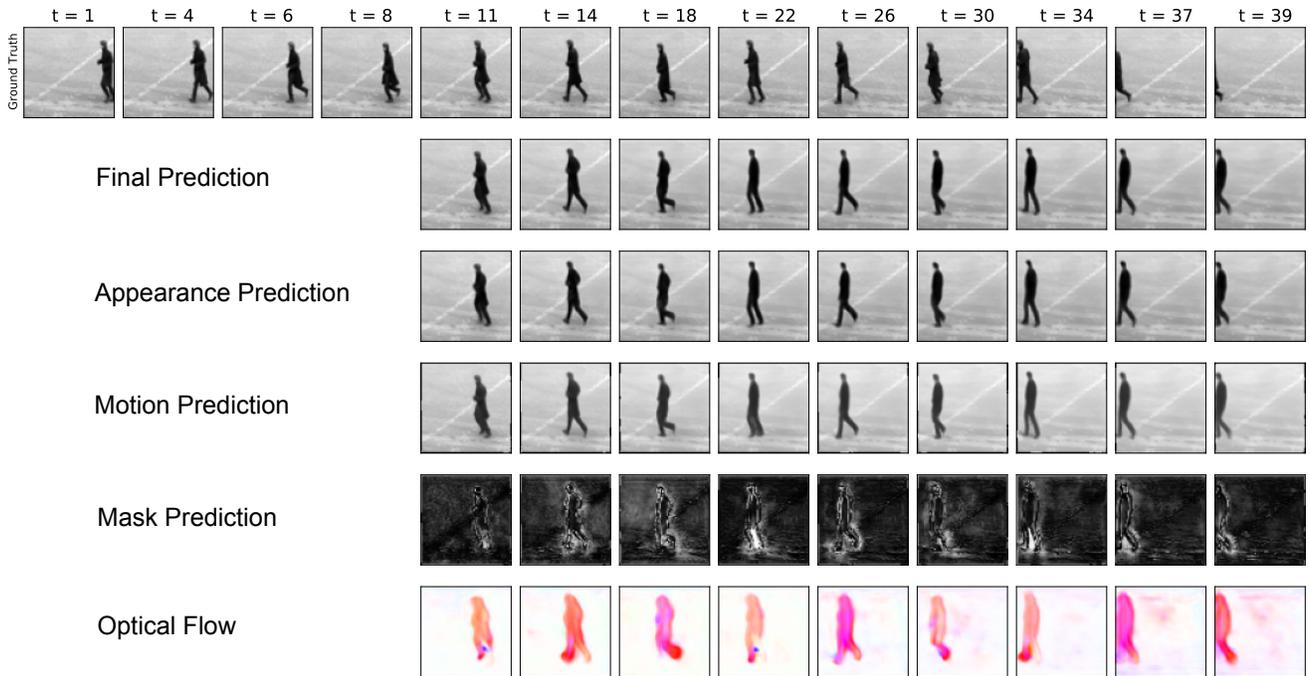


Figure 27: **Person Running**. Our model correctly estimates the optical flow showing the motion of a person running and it can recover the occluded pixels around the legs from the appearance prediction.

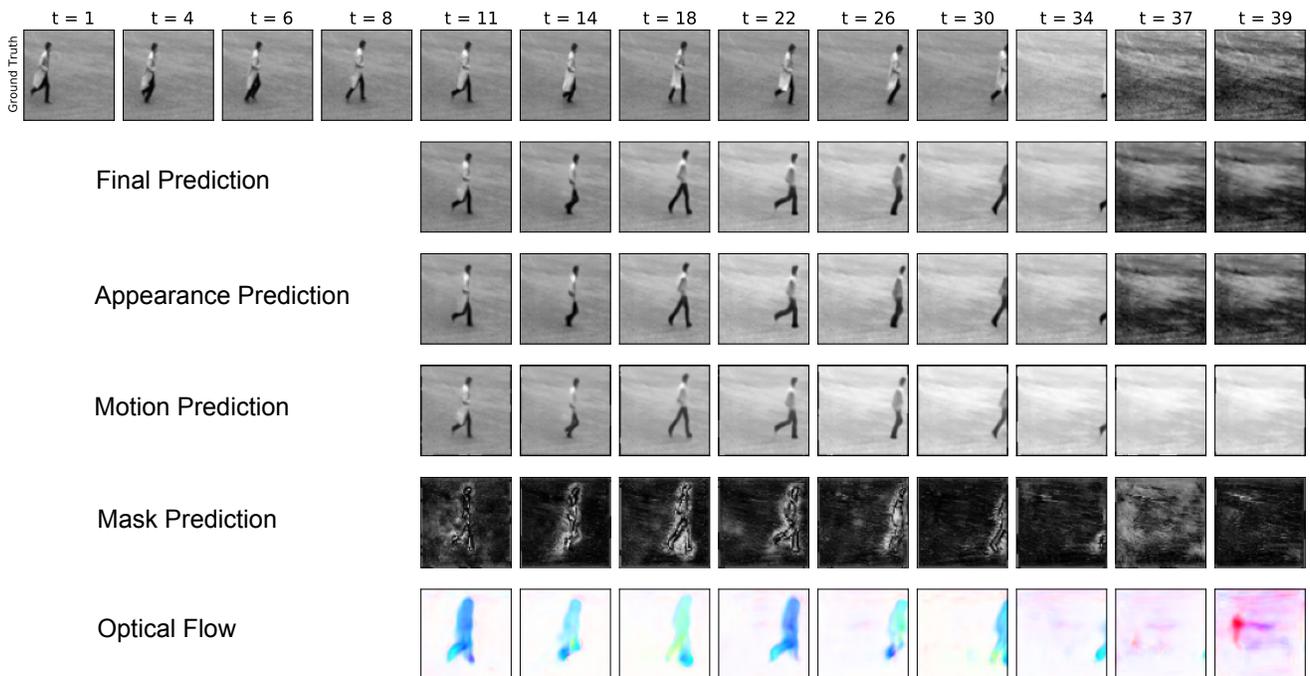


Figure 28: **Person Leaving the Frame**. Our model correctly estimates optical flow when the person leaves the scene in the middle of the sequence by predicting nearly zero flow towards the end.

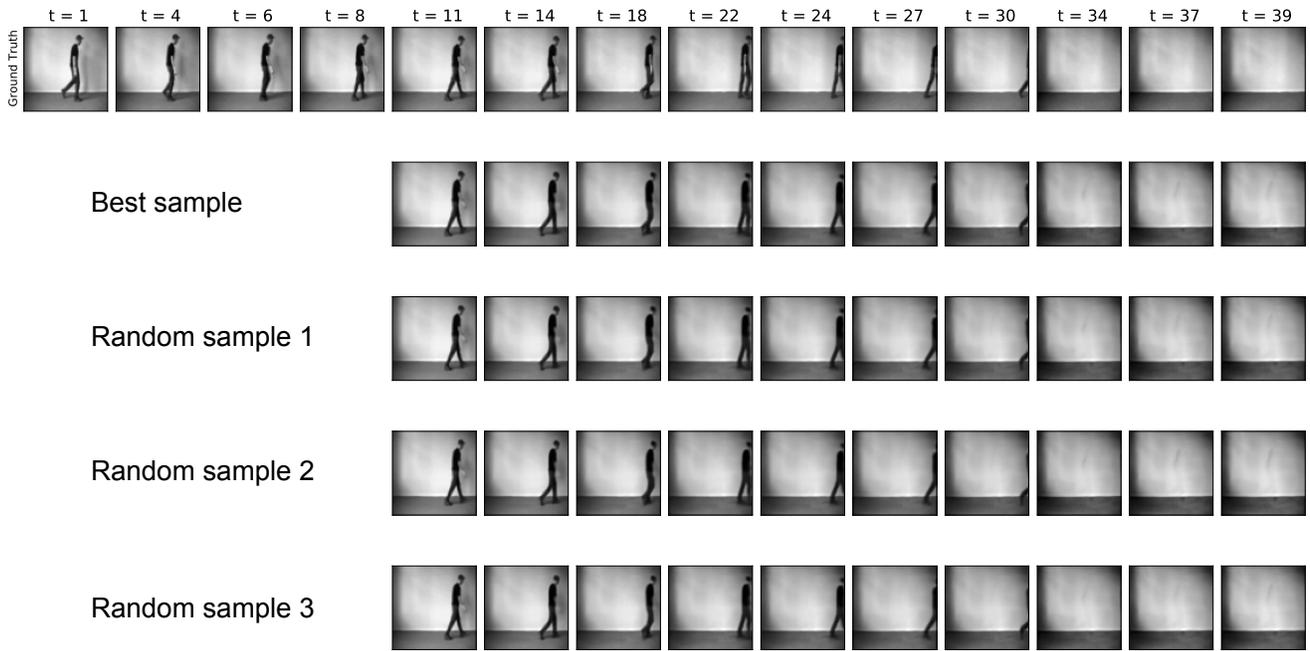


Figure 29: **Random Samples.** We show the best sample and three random samples generated. Random samples look very similar due to the regular motion in the conditioning frames. Our model can capture the motion from the conditioning frames by generating consistent samples, only with minor differences in speed.

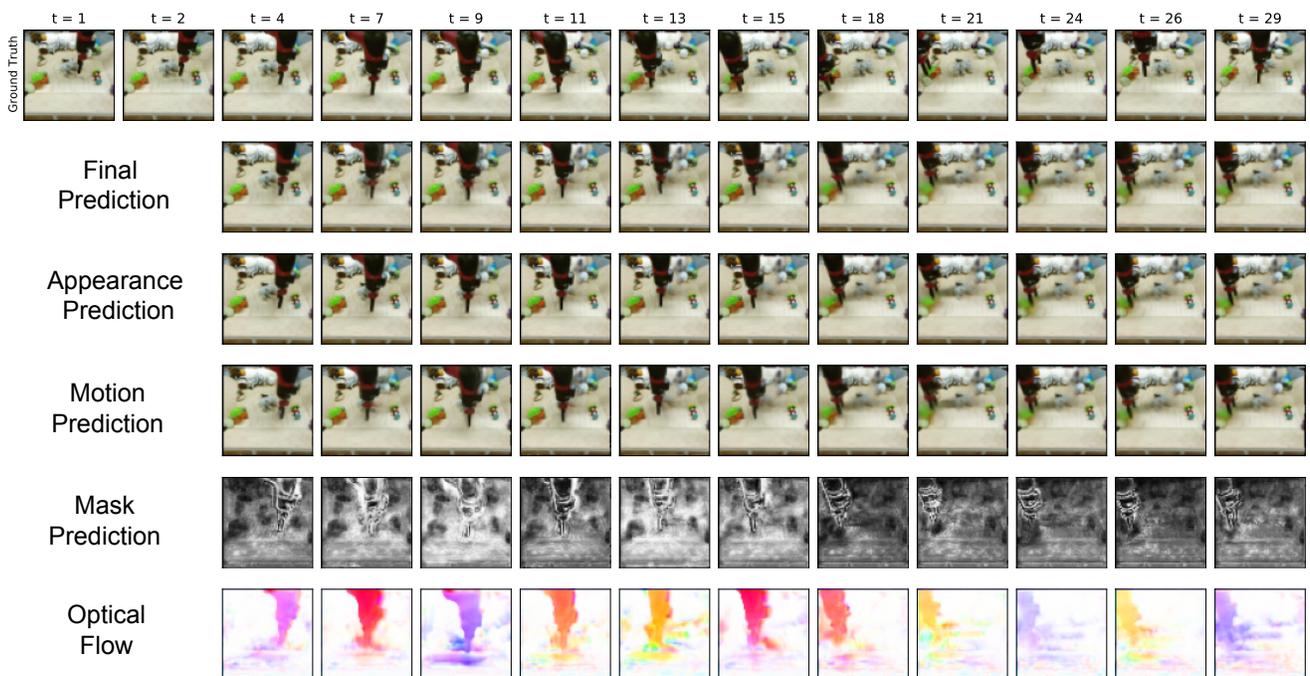


Figure 30: **Results on BAIR.** Optical flow consistently changes from one frame to the next, showing that our model can learn the dataset dynamics. Note that predicted flow is from the target frame to the current frame since we apply inverse warping.

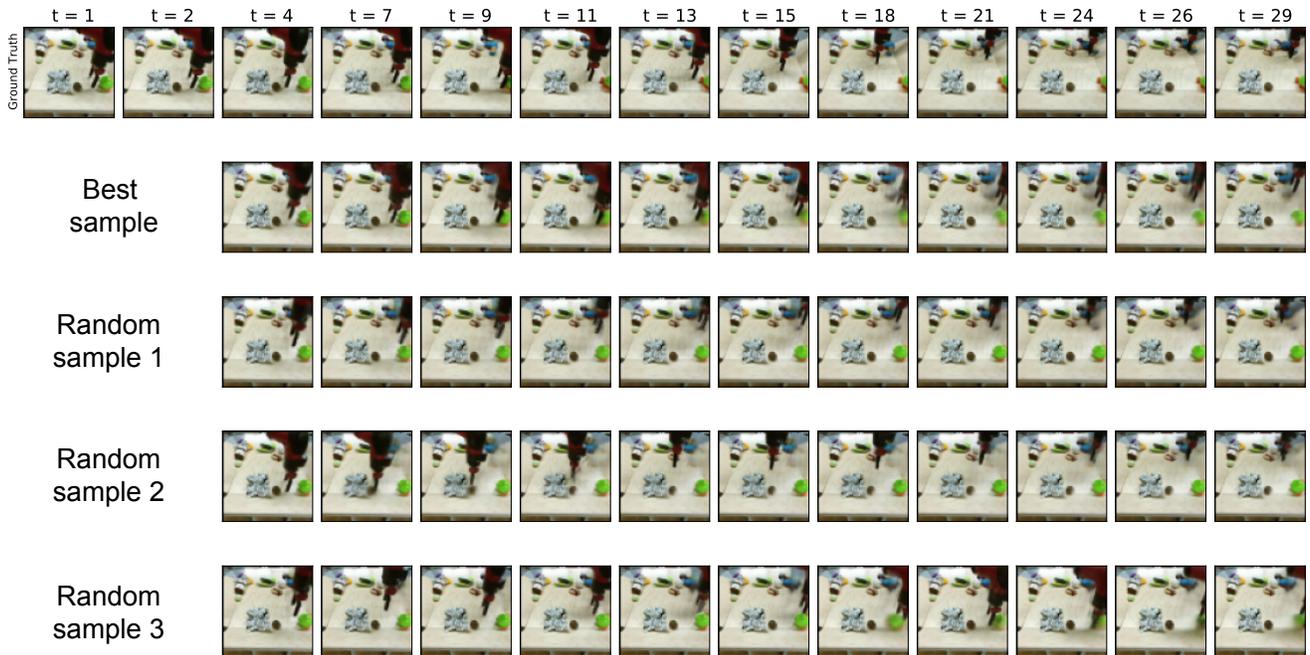


Figure 31: **Random Samples.** We show the best sample and three random samples generated. Random samples look very similar in the beginning due to regular motion in the conditioning frames. Towards the end of the sequence, samples start looking different which shows that our model can generate diverse results.

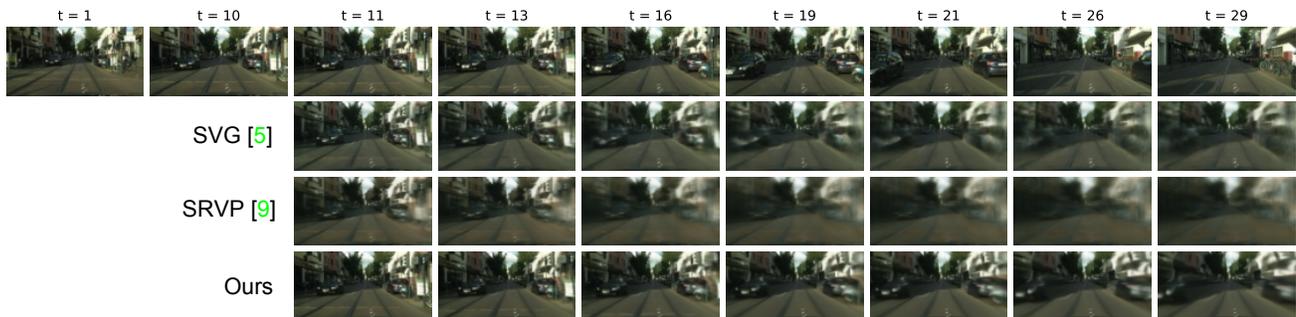


Figure 32: **Full sequence of Cityscapes in Fig. 1** We show the full sequence comparisons with baseline method, SVG [5], and state-of-the-art method, SRVP [9]. Our model can model the ego-motion while both SRVP and SVG suffers from it significantly.

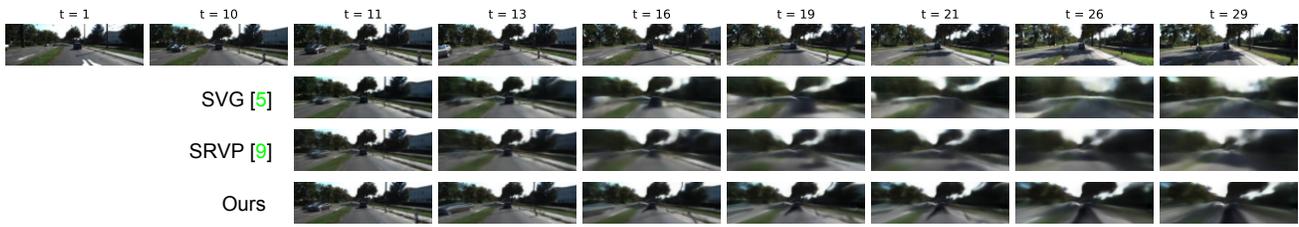


Figure 33: **Full sequence of KITTI in Fig. 1** We show the full sequence comparisons with baseline method, SVG [5], and state-of-the-art method, SRVP [9]. Our model can model both the ego-motion and independently moving objects while both SRVP and SVG cannot reconstruct the future frames successfully.