

# Video Prediction by Modeling Videos as Continuous Multi-Dimensional Processes

Gaurav Shrivastava  
 University of Maryland, College Park  
 gauravsh@umd.edu

Abhinav Shrivastava  
 University of Maryland, College Park  
 abhinav@cs.umd.edu

## Abstract

Diffusion models have made significant strides in image generation, mastering tasks such as unconditional image synthesis, text-image translation, and image-to-image conversions. However, their capability falls short in the realm of video prediction, mainly because they treat videos as a collection of independent images, relying on external constraints such as temporal attention mechanisms to enforce temporal coherence. In our paper, we introduce a novel model class, that treats video as a continuous multi-dimensional process rather than a series of discrete frames. Through extensive experimentation, we establish state-of-the-art performance in video prediction, validated on benchmark datasets including KTH, BAIR, Human3.6M, and UCF101.<sup>1</sup>

## 1. Introduction

In the evolving landscape of machine learning and generative models, particularly in the domain of video representation [5, 32, 34–37], there exists a pivotal challenge in adequately capturing the dynamic transitions between consecutive frames. In this paper, we introduce a novel approach to video representation that treats the video as a continuous process in multi-dimensions. This methodology is anchored in the observation that transitions between consecutive frames in a video do not uniformly contain the same amount of motion. Modeling these transitions with a single-step process often leads to suboptimal quality in sampling. Our method, therefore, involves multiple predefined steps between two consecutive frames, drawing inspiration from recent advancements in diffusion models for image data. This multi-step diffusion process has been instrumental in better modeling image data, and we aim to extend this success to video data.

Previous efforts in video modeling with diffusion models have tended to approach videos as a series of images, generating separate volumes of video frame sequences and applying external constraints such as applying temporal attention to maintain the temporal coherence. We argue that

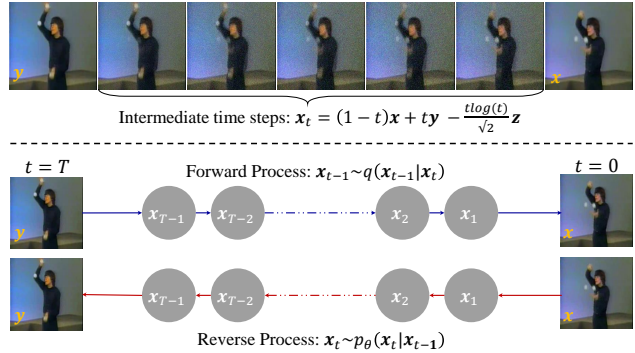


Figure 1. The figure is divided into two parts. The top portion of the figure illustrates the intermediate frames  $\mathbf{x}_t$  between two consecutive frames.  $\mathbf{x}, \mathbf{y}$  represents consecutive frames from a video sequence where  $\mathbf{y} = \mathbf{x}^{j+1}$  and  $\mathbf{x} = \mathbf{x}^j$ .  $\mathbf{x}^j$  denotes some frame at timestep  $j$  in the video sequence  $\mathcal{V} = \{\mathbf{x}^i\}_{i=1}^N$ .  $\mathbf{z}$  denotes the white noise. The lower portion of the figure represents the directed graphical model considered in this work to represent the continuous video process.

this approach overlooks the inherent continuity in video data, which can be more naturally conceptualized as a continuous multi-dimensional process. Our proposed method defines this continuous process, beginning with two consecutive frames from a video sequence as endpoints this can be observed in Fig. 1. We delineate the forward process through interpolation between these endpoints, with a predefined number of steps guiding the transition from one point to another. To ensure the existence of  $p(\mathbf{x}_t)$  at all points, we introduce a novel noise schedule that applies zero noise at both endpoints.

We approximate each step between these endpoints using a Gaussian distribution, following the assumptions made in diffusion models for images by the paper [15, 21, 38, 39]. In defining this forward process, we also lay the groundwork for estimating a reverse process. This paper presents a novel lower variational bound for estimating this reverse process.

To summarize, our contribution in this work is as follows:

- We introduce a novel approach for representing videos as multi-dimensional continuous processes.
- We derive a novel variational bound that efficiently estimates the reverse process in our proposed ‘Continuous

<sup>1</sup>Navigate to the [webpage](#) for video results.

Video Process (CVP)’ model.

- Our method employs a unique noise schedule for the continuous video process, characterized by zero noise at both endpoints, ensuring the existence of  $p(\mathbf{x}_t)$  at all intermediate timesteps.
- We demonstrate the efficacy of our approach through state-of-the-art results in video prediction tasks across four different datasets namely, KTH action recognition, BAIR robot push, Human3.6M, and UCF101 datasets. Additionally, our model requires 75% fewer sampling steps when sampling a frame compared to a diffusion-based baseline.

## 2. Related Works

Understanding and predicting future states based on observed past events is a cornerstone challenge in the domain of video understanding, crucial for applications where capturing the inherent multi-modality of future states is vital, such as in autonomous vehicles. Early methods in this field, as noted by Yuen et al.[52] and Walker et al.[47], primarily focused on matching past frames within datasets to extrapolate future states, although these predictions were constrained to either symbolic trajectories or directly retrieved future frames. The advent of deep learning has significantly propelled advancements in this area. One of the seminal works by Srivastava et al.[41] leveraged a multi-layer LSTM network for deterministic representation learning of video sequences. Subsequent studies [8, 11, 17, 30, 43, 45, 49], have expanded the scope of this research by constructing models that account for the stochastic nature of future states, marking a notable shift from earlier deterministic approaches.

Recent research in this domain has explored both implicit and explicit probabilistic modeling approaches. Implicit probabilistic modeling, typified by GAN [20]-based models, has a substantial history. Nonetheless, these models [10, 26, 28] often grapple with training stability issues and mode collapse(where model only focuses on a few modes in the dataset) issues. On the other hand, explicit probabilistic modeling for video prediction encompasses a range of methodologies, including Variational Autoencoders (VAEs) [24], Gaussian processes, and Diffusion models. VAE-based video prediction methods [7, 14, 26] tend to average results to align with all potential future scenarios, which undermines the fidelity of predictions. Gaussian process-based models [4, 35] exhibit proficiency with smaller datasets but encounter scalability issues owing to matrix inversion limitations when calculating training likelihood. While workarounds exist, they tend to compromise result fidelity.

Recent advancements in diffusion models [12, 22, 23, 46] have positioned them as the preferred choice for video prediction tasks. These multi-step models offer superior sample quality and are resilient to mode collapse. However, even with such lucrative advantages, modeling videos with

these models tends to have downsides. Majorly methods falling under this category enforce temporal consistency using artificial external constraints such as the introduction of temporal attention blocks. This might be effective but comes at a cost of significant computing power.

Another class of popular video prediction models is hierarchical prediction [5, 6, 44, 48, 50] models. These models are multistage models that decompose the problems into two stages. They first predict a high-level structure of a video, like a human pose, and then leverage that structure to make predictions at the pixel level. These models generally require additional annotation for the high-level structure for training, unlike ours that predicts future frames utilizing only the pixel-level information of context frames.

We also want to highlight some very recent works like InDI [13], and Cold diffusion [3] that provide an alternate approach to denoising diffusion models that is similar to our approach. However, their works only explored such formulation for image-based computational photography and image generation tasks.

## 3. Method

Instead of introducing noise iteratively to the frames until they conform to a Gaussian distribution, and adopting a reverse process such as denoising diffusion, a commonly employed technique for video prediction, we introduce a novel model category designed to depict videos as continuous processes. This section delves into the modeling of this continuous video process.

Suppose we have a video sequence denoted by  $\mathcal{V} = \{\mathbf{x}^t\}_1^N$  where  $\mathbf{x}^j \in \mathbb{R}^{c \times h \times w}$  is the frame at the timestep  $j$ . We represent this video sequence as a continuous process. The intermediate frames between  $\mathbf{x} = \mathbf{x}^j$  and  $\mathbf{y} = \mathbf{x}^{j+1}$  are given by the following equation.

$$\mathbf{x}_t = (1 - t)\mathbf{x} + t\mathbf{y} - \frac{t \log(t)}{\sqrt{2}}\mathbf{z} \quad (1)$$

Here,  $\mathbf{z} \sim \mathcal{N}(0, I)$  denotes the white noise. From the above Eqn, it can be seen that at  $t = 0$ , we get the frame  $\mathbf{x}^j$  and at  $t = 1$ , we get the frame  $\mathbf{x}^{j+1}$ . We utilize this continuous process of evolving  $\mathbf{x}^j \rightarrow \mathbf{x}^{j+1}$  given by Eqn. 1 and derive both the forward and reverse processes. For defining the forward process, we take steps in the direction  $t : T \rightarrow 0$  instead of the other way, which happens in denoising diffusion process [21]. The reason for this is we want the reverse process to start from past frame  $\mathbf{x}$  and according to the Eqn. 1  $\mathbf{x}_t = \mathbf{x}$  at  $t = 0$ .

We can write the forward process, i.e., going from the start point  $\mathbf{y}$  at  $t = T$  to endpoint  $\mathbf{x}$  at  $t = 0$ ,

$$\mathbf{x}_{t+\Delta t} = \mathbf{x}_t + (\mathbf{y} - \mathbf{x})\Delta t - t \log(t)\mathbf{z} \quad (2)$$

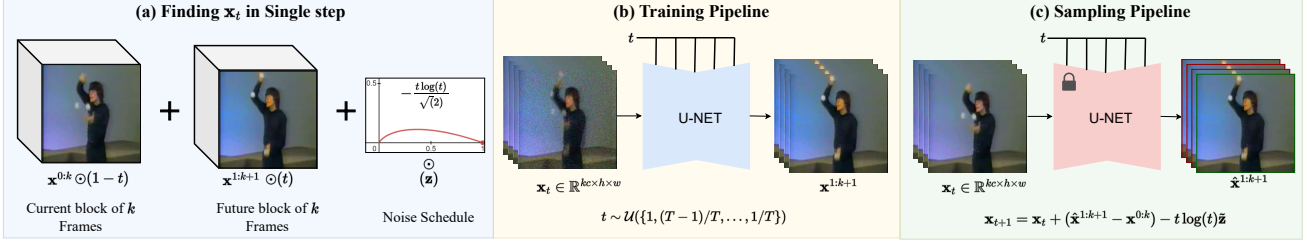


Figure 2. Fig. (a) demonstrates the methodology for estimating  $\mathbf{x}_t$  in a single step, showcasing the specific computational process involved. Fig. (b) details the training pipeline of our Continuous Video Process (CVP) model, where  $\mathbf{x}_t$  and  $t$  are fed as inputs to the U-Net architecture, and the anticipated output is  $\hat{\mathbf{y}}$ , with  $\hat{\mathbf{y}} = \mathbf{x}^{1:k+1}$  in this scenario. Fig. (c) provides an overview of the sampling pipeline utilized in our CVP method, illustrating the sequential steps to predict the next frame of the video sequence given the context frames.

From the above equation, we can write the posterior for the forward process as  $q(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{x}, \mathbf{y}) = \mathcal{N}(\mathbf{x}_{t+1} : \tilde{\mu}(\mathbf{x}_t, \mathbf{x}, \mathbf{y}), g^2(t)I)$ . Where  $g(t) = -t \log t$ . The whole derivation is provided in the appendix.

For modeling our video diffusion process, we like to model the likelihood function  $p_\theta(\mathbf{x}_T) := \int p_\theta(\mathbf{x}_{0:T}) d\mathbf{x}_{0:T-1}$  and minimize the negative log-likelihood to obtain the best fit for our model. Here,  $p_\theta(\mathbf{x}_{0:T})$  is the probability of the reverse process, and it is defined as a Markov chain with learned Gaussian transitions starting at  $p(\mathbf{x}_0) = p_{data}(\mathbf{x})$ . Important note about the notations  $\mathbf{x}_0, \mathbf{x}_T$ , unless specified consider  $\mathbf{x}_0 = \mathbf{x}$  and  $\mathbf{x}_T = \mathbf{y}$  where  $\mathbf{x}$  is the frame in the video sequence at  $j^{th}$  position and  $\mathbf{y}$  is the frame at  $(j+1)^{th}$  position. One important assumption about the continuous video process is we assume the transition between the frames  $\mathbf{x}$  and  $\mathbf{y}$  to follow Markov chain, i.e., the current state at timestep  $t$  only depends on the previous state at timestep  $t-1$ . Leveraging this assumption we can define the reverse process as follows,

$$p_\theta(\mathbf{x}_{0:T}) := p(\mathbf{x}_0) \prod_{t=1}^T p_\theta(\mathbf{x}_t|\mathbf{x}_{t-1}) \quad (3)$$

where,  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_\theta(\mathbf{x}_{t-1}, t-1), \boldsymbol{\Sigma}_\theta(\mathbf{x}_{t-1}, t-1))$ . We are interested in learning the reverse process to perform our video prediction task.

The forward process or the diffusion process is a fixed Markov chain that gradually transforms the frame  $\mathbf{y}$  to frame  $\mathbf{x}$ .

$$q(\mathbf{x}_{0:T-1}|\mathbf{x}_T) := \prod_{t=1}^T q(\mathbf{x}_{t-1}|\mathbf{x}_t), \quad (4)$$

Training is performed by minimizing the variational

bound on the negative log-likelihood.

$$\begin{aligned} \mathbb{E}[-\log p_\theta(\mathbf{x}_T)] &\leq \mathbb{E}_q \left[ -\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{0:T-1}|\mathbf{x}_T)} \right] \quad (5) \\ &\leq \mathbb{E}_q \left[ -\log p(\mathbf{x}_0) - \sum_{t \geq 1} \log \frac{p_\theta(\mathbf{x}_t|\mathbf{x}_{t-1})}{q(\mathbf{x}_{t-1}|\mathbf{x}_t)} \right] \quad (6) \\ &=: L(\theta) \quad (7) \end{aligned}$$

This variational bound can be simplified to the following (we refer the readers to the appendix to follow the simplification of from Eqn. 7 to the following equation),

$$L(\theta) =: \sum_{t \geq 1} D_{\text{KL}}(q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}, \mathbf{y}) \| p_\theta(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x})) \quad (8)$$

In the above Eqn, the KL divergence term utilizes the comparison of  $p_\theta(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x})$  with forward process posterior term, which is tractable under the process given by Eqn. 2. The forward process posterior term is given by

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}, \mathbf{y}) = \mathcal{N}(\mathbf{x}_t : \tilde{\mu}(\mathbf{x}_{t-1}, \mathbf{x}, \mathbf{y}), g^2(t)I) \quad (9)$$

where,  $\tilde{\mu}(\mathbf{x}_t, \mathbf{x}, \mathbf{y}) = \mathbf{x}_t + (\mathbf{y} - \mathbf{x})$  and  $g(t) = -t \log(t)$ . Consequently, all KL divergences in Eqn. 8 are comparisons between Gaussians, so they can be calculated in a Rao-Blackwellized fashion with closed-form expressions instead of high-variance Monte Carlo estimates. It is important to note while deriving the Eqn. 8, we ignore some terms that purely involve the forward process posteriors as  $q$  has no learnable parameters, so such terms are constants during training.

Now we discuss our choices in  $p_\theta(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}) = \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_\theta(\mathbf{x}_{t-1}, t-1, \mathbf{x}), \boldsymbol{\Sigma}_\theta(\mathbf{x}_{t-1}, t-1, \mathbf{x}))$  for  $1 < t \leq T$ . First, we set  $\boldsymbol{\Sigma}_\theta(\mathbf{x}_{t-1}, t-1) = g^2(t)I$  to untrained time dependent constants. Experimentally, the choice of  $g(t) = -t \log(t)$  works the best. This noise function has an interesting property that noise is absent both at the start and end points, i.e.,  $g(t) = 0 \quad \forall t = \{0, 1\}$ .

Second, to represent the mean  $\mu_\theta(\mathbf{x}_t, t, \mathbf{x})$ , we propose a specific parameterization motivated by the forward process posterior given by Eqn. 9. With  $p_\theta(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}) = \mathcal{N}(\mathbf{x}_t; \mu_\theta(\mathbf{x}_{t-1}, t-1, \mathbf{x}), g^2(t)\mathbf{I})$ , we can write:

$$L(\theta) := \mathbb{E}_q \left[ \frac{1}{2g^2(t)} \|\tilde{\mu}(\mathbf{x}_t, \mathbf{x}, \mathbf{y}) - \mu_\theta(\mathbf{x}_t, t, \mathbf{x})\|^2 \right] + C \quad (10)$$

where  $C$  is a constant that does not depend on  $\theta$ . So, we see that the most straightforward parameterization of  $\mu_\theta$  is a model that predicts  $\tilde{\mu}_t$ , the forward process posterior mean.

However, we can simplify Eqn. 10 further and obtain a very simple training loss objective by delving in the term  $\tilde{\mu}$ . We further parameterize the term  $\mu_\theta$  as follows,

$$\mu_\theta(\mathbf{x}_t, t, \mathbf{x}) = \mathbf{x}_t + (\mathbf{y}_\theta(\mathbf{x}_t) - \mathbf{x}) \quad (11)$$

When we substitute this  $\mu_\theta(\mathbf{x}_t, t, \mathbf{x})$  parameterization in the Eqn. 10 we get the simplified version of the loss  $L(\theta)$  as follows,

$$L_{\text{simple}}(\theta) := \mathbb{E}_{t, \mathbf{x}_t} \left[ \frac{1}{2g^2(t)} \|\mathbf{y} - \mathbf{y}_\theta((\mathbf{x}_t, t))\|^2 \right] \quad (12)$$

For training the video prediction model utilizing the above Eqn. 12 we obtain the  $\mathbf{x}_t$  as a function of  $t$  by leveraging the Eqn. 1. The following equation gives a more generic form of the final loss function utilized to train the video prediction model,

$$\arg \min_{\theta} \mathbb{E}_{t, \mathbf{x}, \mathbf{y}} \left[ \frac{1}{2g^2(t)} \left\| \mathbf{y} - \mathbf{y}_\theta((1-t)\mathbf{x} + t\mathbf{y} + \frac{g(t)}{\sqrt{2}}\mathbf{z}, t) \right\|^2 \right] \quad (13)$$

The whole training and sampling pipeline is described in the training Alg. 1, sampling Alg. 2 and depicted in Fig. 2.

## 4. Experiments

Video prediction task can be defined as given a few context frames, the model has to predict the subsequent future frames. In this section, we empirically demonstrate that our approach yields superior results in modeling the video prediction task.

### 4.1. Datasets

We chose 4 different types of datasets to demonstrate the efficacy of our approach. These are standard benchmarks for video prediction tasks. Dataset lists include KTH action recognition dataset [33], BAIR robot pushing dataset [16], Human3.6M [9] and UCF101 [40] datasets. Training and architecture-specific details about the approach are included in the appendix.

**KTH Action Recognition Dataset.** The KTH action dataset [33] consists of video sequences of 25 people performing six different actions: walking, jogging, running,

---

### Algorithm 1 Training of CVP model

---

- 1: **repeat**
  - 2:  $\mathbf{x}, \mathbf{y} \sim q_{\text{data}}(\mathbf{x}, \mathbf{y})$
  - 3:  $t \sim \text{Uniform}(\{1, \dots, T\})$
  - 4:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - 5: Take gradient descent step on  $\nabla_{\theta} \frac{1}{2g^2(t)} \|\mathbf{y} - \mathbf{y}_\theta((1-t)\mathbf{x} + t\mathbf{y} - (t \log(t)/\sqrt{2})\mathbf{z}, t)\|^2$
  - 6: **until** converged
- 

---

### Algorithm 2 Sampling Algorithm

---

- 1:  $\mathbf{x} \sim q_{\text{data}}(\mathbf{x})$
  - 2:  $\mathbf{x}_0 = \mathbf{x}$
  - 3:  $d = \frac{1}{N}$ , Here  $N$  denotes number of steps.
  - 4: **for**  $t = 1, \dots, N$  **do**
  - 5:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \text{Id})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$
  - 6:  $\mathbf{x}_{t+1} = \mathbf{x}_t + (\hat{y}(\mathbf{x}_t, t) - \mathbf{x})d - t \log(t)\mathbf{z}$
  - 7: **end for**
  - 8: **return**  $\mathbf{x}_T$
- 

boxing, hand-waving, and hand-clapping. The background is uniform, and a single person is performing actions in the foreground. The foreground motion of the person in the frame is fairly regular. The frames in the video for this dataset consist of a single channel. The spatial resolution of the frames in the video is downsampled to the size of  $64 \times 64$ .

**BAIR pushing Dataset.** The BAIR robot pushing dataset [16] contains the videos of table mounted sawyer robotic arm pushing various objects around. The BAIR dataset consists of different actions given to the robotic arm to perform. The spatial resolution of the frames in the video is kept to be  $64 \times 64$ .

**Human3.6M Dataset.** Human3.6M [9] dataset consists of 10 subjects performing 15 different actions. The pose information from the dataset was not used in predicting next frame. The background is uniform, and a single person is performing actions in the foreground. The foreground motion of the person in the frame is fairly regular. The frames in the video for this dataset consist of ‘RGB’ channels. The spatial resolution of the frames in the video is downsampled to the size of  $64 \times 64$ .

**UCF101 Dataset.** This dataset [40] consists of 13,320 videos belonging to 101 different action classes. The video seems to have a variety of backgrounds and the frames of the video have three channels, namely ‘RGB’. We reshape the resolution of frames from the original size of  $320 \times 240$  down to  $128 \times 128$  for our video prediction tasks. The downsampling is done utilizing the bicubic downsampling.

### 4.2. Metrics

We primarily use the FVD [42] metric to determine the best-performing baseline when evaluating a video prediction

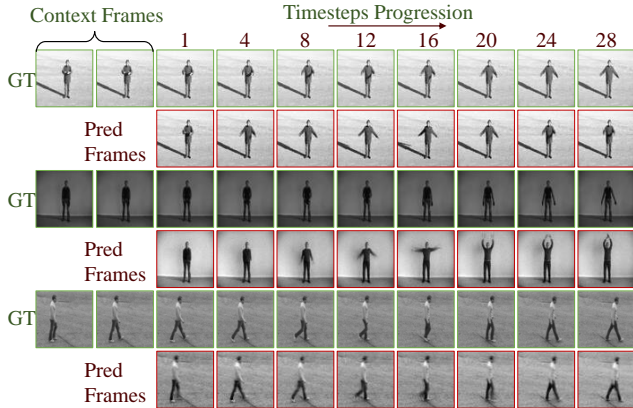


Figure 3. Figure represents qualitative results of our CVP model on the KTH dataset. The number of context frames used in the above setting is 4 for all three sequences. Every 4<sup>th</sup> predicted future frame is shown in the figure.

Table 1. Video prediction results on KTH (64 × 64), predicting 30 and 40 frames using models trained to predict  $k$  frames at a time. All models condition on 10 past frames on 256 test videos.

KTH [10 → #pred; trained on $k$ ]	$k$	#pred	FVD↓	PSNR↑	SSIM↑
SVG-LP [14]	10	30	377	28.1	0.844
SAVP [26]	10	30	374	26.5	0.756
MCVD [46]	5	30	323	27.5	0.835
SLAMP [1]	10	30	228	29.4	0.865
SRVP [18]	10	30	222	29.7	0.870
RIVER [12]	10	30	180	<b>30.4</b>	0.86
<b>CVP (Ours)</b>	<b>1</b>	30	<b>140.6</b>	29.8	<b>0.872</b>
Struct-vRNN [29]	10	40	395.0	24.29	0.766
SVG-LP [14]	10	40	157.9	23.91	0.800
MCVD [46]	5	40	276.7	26.40	0.812
SAVP-VAE [26]	10	40	145.7	26.00	0.806
Grid-keypoints [19]	10	40	144.2	27.11	0.837
RIVER [12]	10	40	170.5	29.0	0.82
<b>CVP (Ours)</b>	<b>1</b>	40	<b>120.1</b>	<b>29.2</b>	<b>0.841</b>

task. FVD metric evaluates a baseline on both terms, the reconstruction quality and diversity of the generated samples. FVD is calculated as the frechet distance between the I3D embeddings of generated video samples and real samples. The I3D network used for obtaining the embeddings for real and generated video is trained on the Kinetics-400 dataset.

## 5. Setup and Results

Below, we describe in detail how the setup for our experiment looks compared to baselines. We also showcase our findings about the performance of our method and comparison to baselines in this section.

**KTH action recognition dataset:** For this dataset, we adhered to the baseline setup [46], which utilizes the first 10 frames as context frames. In baseline setup, these 10 frames are utilized to predict the subsequent 30 and 40 frames. A

notable aspect of our experiment is we only used the last 4 frames from this sequence of 10 frames as context frames in our CVP model, while disregarding the information in the remaining 6 frames. This decision was taken to maintain consistency with the experimental setups used in prior baseline methodologies. The outcomes of this evaluation are summarized in Table 1.

It can be observed from the Table 1, our model’s unique approach requires a significantly reduced number of frames for training. Contrary to other methods that train on an additional set of  $k$  frames (10[context frames]+ $k$ [future frames]), our model uses just one frame (effectively 4[context frames]+1[future frames]). We employ the 4 context frames to predict the immediate next frame and then autoregressively generate either 30 or 40 frames, depending on the evaluation requirement. This methodology is supported by our model’s efficient handling of video sequences as continuous processes, which eliminates the need for external artificial constraints, such as temporal attention mechanisms.

The results, as shown in Table 1, clearly indicate that our method delivers state-of-the-art performance when compared to other baseline models. Additionally, the qualitative results for our CVP model on the KTH dataset can be observed in Fig. 3.

**BAIR Robot Push dataset:** The BAIR Robot Push dataset is characterized by highly stochastic video sequences. In our study, we adhered to a baseline setup [46] with three main experimental settings: 1) using only one context frame to predict the next 15 frames, 2) employing two context frames to predict 14 future frames, and 3) utilizing two context frames to forecast the next 28 frames. The outcomes of these approaches are summarized in Table 2.

As observed in Table 2, a trend emerges where increasing the number of frames predicted at a time concurrently results in a degradation of prediction quality. This phenomenon is hypothesized to stem from an augmented disparity between the blocks of context frames and predicted future frames. Specifically, consider the scenario where two context frames are designated as  $\mathbf{x}^{0:2}$ , corresponding to  $\mathbf{x}$  in the context of Eqn.1. Under the first experimental condition, where the model predicts a single frame at a time, the future frame prediction block is represented as  $\mathbf{x}^{1:3}$ , analogous to  $\mathbf{y}$  in Eqn.1. Conversely, in the second condition, where two frames are predicted simultaneously, the future frame block extends to  $\mathbf{x}^{2:4}$ , again paralleling  $\mathbf{y}$  in the equation. This setup implies that in the former setting, interpolation occurs between adjacent frames (i.e., the transition from  $\mathbf{x}^0 \rightarrow \mathbf{x}^1$  and  $\mathbf{x}^1 \rightarrow \mathbf{x}^2$ ), while in the latter, interpolation spans a two-frame interval (i.e., the transition from  $\mathbf{x}^0 \rightarrow \mathbf{x}^2$  and from  $\mathbf{x}^1 \rightarrow \mathbf{x}^3$ ). The expanded interval in the second scenario is posited as the causative factor for the observed reduction in predictive performance, particularly in configurations where  $k = 2$  and  $p = 2$ .

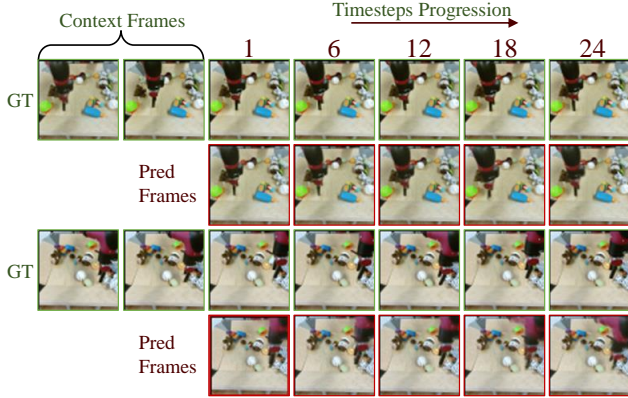


Figure 4. Figure represents qualitative results of our CVP model on the BAIR dataset. The number of context frames used in the above setting is two for both sequences. Every 6<sup>th</sup> predicted future frame is shown in the figure.

The results, as shown in Table 2, clearly indicate that our method delivers state-of-the-art performance compared to other baseline models. Additionally, the qualitative results for our CVP model on the BAIR dataset can be observed in Fig. 4.

**Human3.6M dataset:** Similar to the KTH dataset, the Human3.6M dataset features actors performing distinct actions against a static background. However, the Human3.6M dataset distinguishes itself by offering a greater variety of distinct actions within its videos and providing three-channel video frames, in contrast to the single-channel frames of the KTH dataset. For evaluating the Human3.6M dataset, we employed a similar setup to that used for the KTH dataset, where 5 frames are provided as context, and the model predicts the subsequent 30 frames based on these context frames. The results of this evaluation are summarized in Table 3.

An analysis of Table 3 reveals that our model, with its unique approach, requires a significantly lower number of frames for training, needing only a total of 6 frames per block to yield results that are considerably better than those of the baselines.

The results, as presented in Table 3, unequivocally demonstrate that our method outperforms other baseline models, establishing a new state-of-the-art on the Human3.6M dataset. Furthermore, the qualitative efficacy of our CVP model on the Human3.6M dataset is illustrated in Fig. 5, showcasing the model’s ability to effectively capture and predict the dataset’s varied actions.

**UCF101 dataset:** The UCF101 dataset presents a greater level of complexity compared to the KTH or Human3.6M datasets, owing to its substantially higher number of action categories, diverse backgrounds, and significant camera movements. Notably, we only use information from the context frames for our frame-conditional generation task. No

Table 2. BAIR dataset evaluation. Video prediction results on BAIR ( $64 \times 64$ ) conditioning on  $p$  past frames and predicting  $pred$  frames in the future, using models trained to predict  $k$  frames at a time. The common way to compute the FVD is to compare  $100 \times 256$  generated sequences to 256 randomly sampled test videos. Best results are marked in *bold*.

BAIR ( $64 \times 64$ )	$p$	$k$	#pred	FVD↓
LVT [31]	1	15	15	125.8
DVD-GAN-FP [10]	1	15	15	109.8
TrIVD-GAN-FP [28]	1	15	15	103.3
VideoGPT [51]	1	15	15	103.3
CCVS [25]	1	15	15	99.0
FitVid [2]	1	15	15	93.6
MCVD [46]	1	5	15	89.5
NÜWA [27]	1	15	15	86.9
RaMViD [23]	1	15	15	84.2
VDM [22]	1	15	15	66.9
RIVER [12]	1	15	15	73.5
<b>CVP (Ours)</b>	1	<b>1</b>	15	<b>70.1</b>
<hr/>				
DVG [35]	2	14	14	120.0
SAVP [26]	2	14	14	116.4
MCVD [46]	2	5	14	87.9
<b>CVP (Ours)</b>	2	2	14	68.2
<b>CVP (Ours)</b>	2	<b>1</b>	14	<b>65.1</b>
<hr/>				
SAVP [26]	2	10	28	143.4
Hier-vRNN [7]	2	10	28	143.4
MCVD [46]	2	5	28	118.4
<b>CVP (Ours)</b>	2	<b>2</b>	28	95.1
<b>CVP (Ours)</b>	2	<b>1</b>	28	<b>85.1</b>

Table 3. Quantitative comparisons on the Human3.6M dataset. The best results under each metric are marked in bold.

Human3.6M	$p$	$k$	#pred	FVD↓
SVG-LP [14]	5	10	30	718
Struct-VRNN [29]	5	10	30	523.4
DVG [35]	5	10	30	479.5
SRVP [18]	5	10	30	416.5
Grid keypoint [19]	8	8	30	166.1
<b>CVP (Ours)</b>	5	1	30	<b>144.5</b>

extra information, like class labels, was used for the prediction task. In evaluating the UCF101 dataset, we adopted an approach similar to that used for the Human3.6M dataset, where 5 context frames are provided, and the model is tasked with predicting the next 16 frames based on these. The outcomes of this evaluation are detailed in Table. 4.

An examination of Table. 4 reveals that our CVP model surpasses the performance of other baseline models, thereby setting a new benchmark for the UCF101 dataset. Additionally, the qualitative performance of our CVP model on the UCF101 dataset is depicted in Fig. 6. This illustration

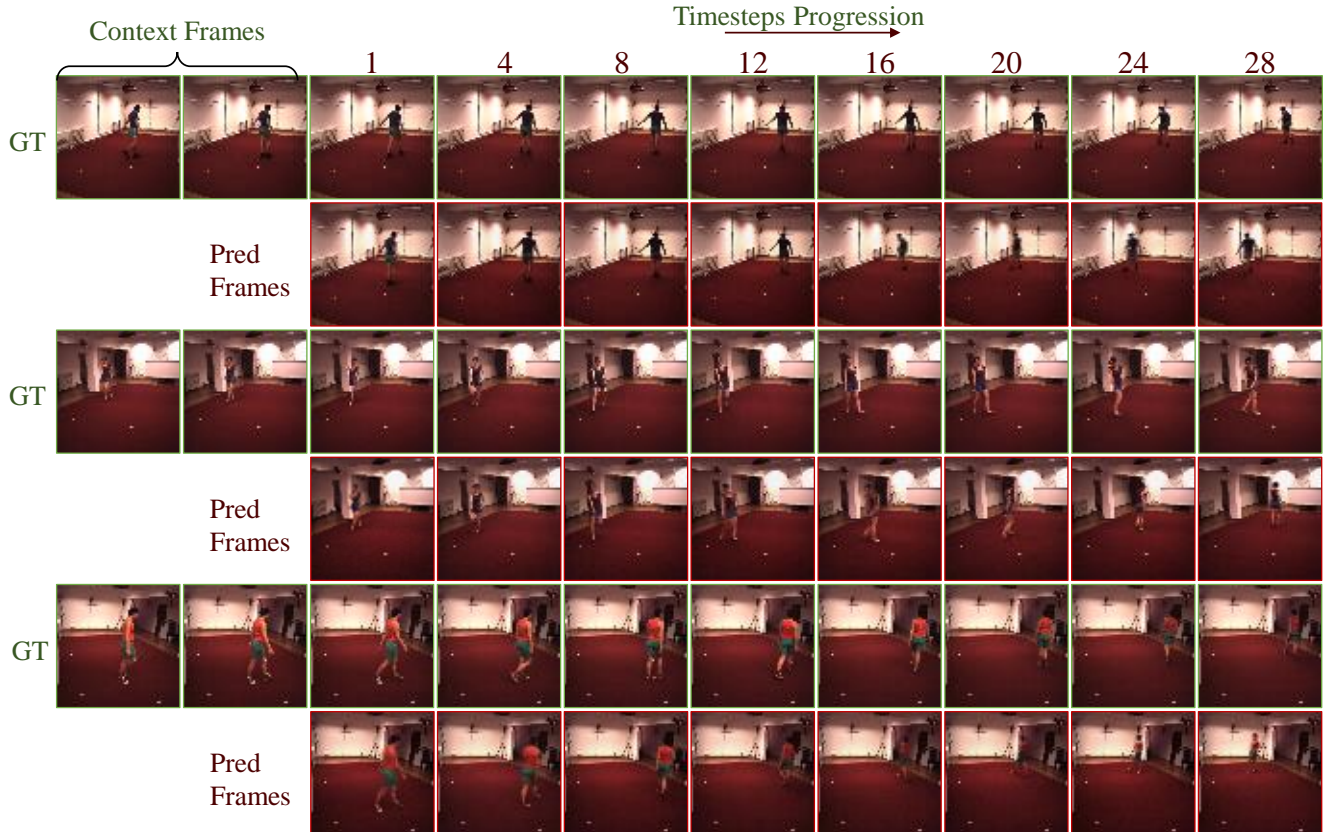


Figure 5. Figure represents qualitative results of our CVP model on the Human3.6M dataset. The number of context frames used in the above setting is 4 for all three sequences. Every 4<sup>th</sup> predicted future frame is shown in the figure.

Table 4. Video prediction results on UCF (128 × 128), predicting 16 frames. All models are conditioned on 5 past frames.

UCF101 [5 → 16]   $p$ $k$ #pred   FVD↓
SVG-LP [14]   5 10 16   1248
CCVS [25]   5 16 16   409
MCVD [46]   5 5 16   387
RaMViD [23]   5 4 16   356
<b>CVP (Ours)</b>   5 <b>1</b> 16   <b>245.2</b>

showcases the model’s proficiency in accurately capturing and predicting the diverse range of actions featured in the dataset.

## 6. Limitation

While our method demonstrates promising results in video prediction, it is important to acknowledge its limitations to guide future research and application development.

A primary limitation of our approach is its reliance on a limited context frame window for predicting the next frame. Specifically, when a context vector, denoted as  $\mathbf{x}^{0:4}$ , com-

prising 4 video frames is used, the prediction of the subsequent frame is entirely dependent on this four-frame window. This model architecture performs adequately in scenarios involving uniform video sequences. However, its efficacy diminishes in a setting that requires more context to predict the future frame. Addressing this limitation requires a more adaptive approach that can handle varying contextual information, a challenge we have earmarked for future research.

Another constraint lies in the computational efficiency of our model. Currently, it necessitates multiple steps to sample a single frame, which could become a significant bottleneck, especially when a larger number of frame predictions are required. Although our method is more efficient in terms of the number of steps needed for frame sampling compared to diffusion-based counterparts, further optimization is necessary to reduce the computational overhead associated with this process.

Additionally, our experimental setup was constrained by the computational resources available to us. The model was developed and tested using just two A6000 GPUs. This limitation raises questions about the potential improvements

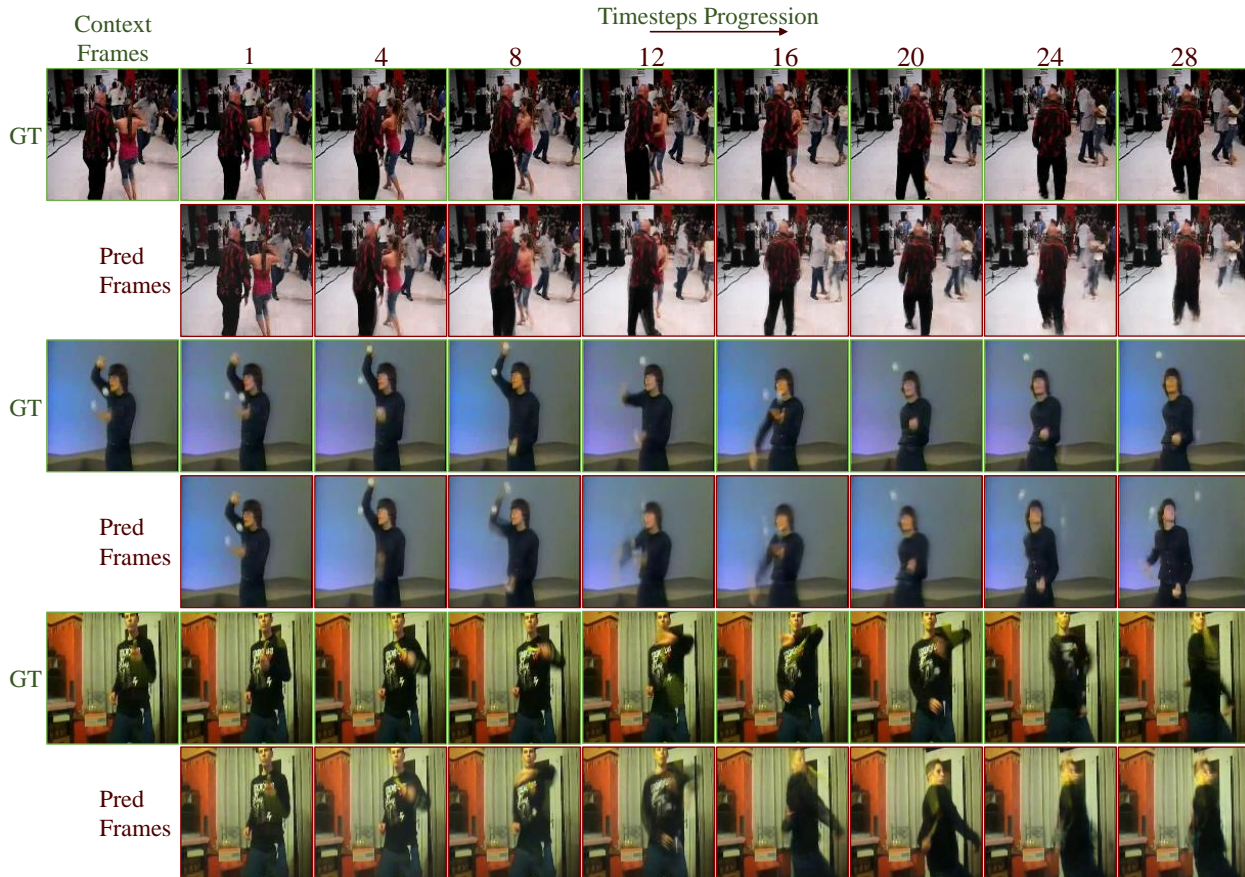


Figure 6. Figure represents qualitative results of our CVP model on the UCF dataset. The number of context frames used in the above setting is 5 for all three sequences. Every 4<sup>th</sup> predicted future frame is shown in the figure.

that could be achieved with a more powerful computational setup. A larger model with an increased number of parameters, trained on more advanced hardware, could potentially unveil further advancements in video prediction capabilities. We recognize this as an important area for investigation and encourage labs with more substantial resources to explore this avenue.

In summary, while our model represents a significant step forward in video prediction, these limitations highlight crucial areas for future research and development, paving the way for more robust and versatile video prediction models.

## 7. Conclusion

In this work, we have presented a novel model class designed specifically for video representation, marking a significant advancement in the field of video prediction tasks. Our comprehensive experimental evaluations across various datasets, including KTH, BAIR, Human3.6M, and UCF101, have not only validated the effectiveness of our model but also established new benchmarks in state-of-the-art performance for video prediction tasks.

A notable aspect of our approach is its efficiency in terms of the required number of context and future frames for training. Moreover, our model’s continuous video process capability uniquely operates without the need for additional constraints such as temporal attention, which are typically employed to ensure temporal consistency. This aspect of our model underscores its inherent ability to maintain temporal coherence, further simplifying the video prediction process while enhancing its effectiveness.

In conclusion, the innovations introduced in our model offer promising directions for future research in video representation and prediction. The achievements demonstrated in this paper not only contribute to the advancement of video prediction methodologies but also open avenues for exploring more efficient and effective ways of video representation in various real-world applications.

**Acknowledgements.** This project was partially funded by NSF CAREER Award (2238769) to AS. We also thank Shirley Huang for providing feedback on the manuscript.



## References

- [1] Adil Kaan Akan, Erkut Erdem, Aykut Erdem, and Fatma Güney. Slamp: Stochastic latent appearance and motion prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14728–14737, 2021. [5](#)
- [2] Mohammad Babaeizadeh, Mohammad Taghi Saffar, Suraj Nair, Sergey Levine, Chelsea Finn, and Dumitru Erhan. Fitvid: Overfitting in pixel-level video prediction. *arXiv preprint arXiv:2106.13195*, 2021. [6](#)
- [3] Arpit Bansal, Eitan Borgnia, Hong-Min Chu, Jie S Li, Hamid Kazemi, Furong Huang, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Cold diffusion: Inverting arbitrary image transforms without noise. *arXiv preprint arXiv:2208.09392*, 2022. [2](#)
- [4] Sarthak Bhagat, Shagun Uppal, Zhuyun Yin, and Nengli Lim. Disentangling multiple features in video sequences using gaussian processes in variational autoencoders, 2020. [2](#)
- [5] Navaneeth Bodla, Gaurav Shrivastava, Rama Chellappa, and Abhinav Shrivastava. Hierarchical video prediction using relational layouts for human-object interactions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12146–12155, 2021. [1](#), [2](#)
- [6] Haoye Cai, Chunyan Bai, Yu-Wing Tai, and Chi-Keung Tang. Deep video generation, prediction and completion of human action sequences. *Lecture Notes in Computer Science*, page 374–390, 2018. [2](#)
- [7] Lluís Castrejon, Nicolas Ballas, and Aaron Courville. Improved conditional vrns for video prediction. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7608–7617, 2019. [2](#), [6](#)
- [8] Lluís Castrejón, Nicolas Ballas, and Aaron C. Courville. Improved conditional vrns for video prediction. *CoRR*, abs/1904.12165, 2019. [2](#)
- [9] Cristian Sminchisescu Catalin Ionescu, Fuxin Li. Latent structured models for human pose estimation. In *International Conference on Computer Vision*, 2011. [4](#)
- [10] Aidan Clark, Jeff Donahue, and Karen Simonyan. Adversarial video generation on complex datasets. *arXiv preprint arXiv:1907.06571*, 2019. [2](#), [6](#)
- [11] Francesco Cricri, Xingyang Ni, Mikko Honkala, Emre Aksu, and Moncef Gabbouj. Video ladder networks. *CoRR*, abs/1612.01756, 2016. [2](#)
- [12] Aram Davtyan, Sepehr Sameni, and Paolo Favaro. Efficient video prediction via sparsely conditioned flow matching. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 23263–23274, 2023. [2](#), [5](#), [6](#)
- [13] Mauricio Delbracio and Peyman Milanfar. Inversion by direct iteration: An alternative to denoising diffusion for image restoration. *arXiv preprint arXiv:2303.11435*, 2023. [2](#)
- [14] Emily Denton and Rob Fergus. Stochastic video generation with a learned prior, 2018. [2](#), [5](#), [6](#), [7](#)
- [15] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34, 2021. [1](#)
- [16] Frederik Ebert, Chelsea Finn, Alex X. Lee, and Sergey Levine. Self-supervised visual planning with temporal skip connections, 2017. [4](#)
- [17] N. Elsayed, A. S. Maida, and M. Bayoumi. Reduced-gate convolutional lstm architecture for next-frame video prediction using predictive coding. In *2019 International Joint Conference on Neural Networks (IJCNN)*, pages 1–9, 2019. [2](#)
- [18] Jean-Yves Franceschi, Edouard Delasalles, Mickaël Chen, Sylvain Lamprier, and Patrick Gallinari. Stochastic latent residual video prediction. In *International Conference on Machine Learning*, pages 3233–3246. PMLR, 2020. [5](#), [6](#)
- [19] Xiaojie Gao, Yueming Jin, Qi Dou, Chi-Wing Fu, and Pheng-Ann Heng. Accurate grid keypoint learning for efficient video prediction. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5908–5915. IEEE, 2021. [5](#), [6](#)
- [20] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014. [2](#)
- [21] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 2020. [1](#), [2](#)
- [22] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J. Fleet. Video diffusion models. 2022. [2](#), [6](#)
- [23] Tobias Höppe, Arash Mehrjou, Stefan Bauer, Didrik Nielsen, and Andrea Dittadi. Diffusion models for video prediction and infilling, 2022. [2](#), [6](#), [7](#)
- [24] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *CoRR*, abs/1312.6114, 2013. [2](#)
- [25] Guillaume Le Moing, Jean Ponce, and Cordelia Schmid. Ccvs: context-aware controllable video synthesis. *Advances in Neural Information Processing Systems*, 34:14042–14055, 2021. [6](#), [7](#)
- [26] Alex X. Lee, Richard Zhang, Frederik Ebert, Pieter Abbeel, Chelsea Finn, and Sergey Levine. Stochastic adversarial video prediction. *arXiv preprint arXiv:1804.01523*, 2018. [2](#), [5](#), [6](#)
- [27] Jian Liang, Chenfei Wu, Xiaowei Hu, Zhe Gan, Jianfeng Wang, Lijuan Wang, Zicheng Liu, Yuejian Fang, and Nan Duan. Nuwa-infinity: Autoregressive over autoregressive generation for infinite visual synthesis. *Advances in Neural Information Processing Systems*, 35:15420–15432, 2022. [6](#)
- [28] Pauline Luc, Aidan Clark, Sander Dieleman, Diego de Las Casas, Yotam Doron, Albin Cassirer, and Karen Simonyan. Transformation-based adversarial video prediction on large-scale data. *arXiv preprint arXiv:2003.04035*, 2020. [2](#), [6](#)
- [29] Matthias Minderer, Chen Sun, Ruben Villegas, Forrester Cole, Kevin P Murphy, and Honglak Lee. Unsupervised learning of object structure and dynamics from videos. *Advances in Neural Information Processing Systems*, 32, 2019. [5](#), [6](#)
- [30] Marc Oliu, Javier Selva, and Sergio Escalera. Folded recurrent neural networks for future video prediction. *CoRR*, abs/1712.00311, 2017. [2](#)
- [31] Ruslan Rakhimov, Denis Volkhonskiy, Alexey Artemov, Denis Zorin, and Evgeny Burnaev. Latent video transformer. *arXiv preprint arXiv:2006.10704*, 2020. [6](#)

- [32] Nirat Saini, Bo He, Gaurav Shrivastava, Sai Saketh Rambhatla, and Abhinav Shrivastava. Recognizing actions using object states. In *ICLR2022 Workshop on the Elements of Reasoning: Objects, Structure and Causality*, 2022. 1
- [33] C. Schuldt, I. Laptev, and B. Caputo. Recognizing human actions: a local svm approach. In *Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004.*, pages 32–36 Vol.3, 2004. 4
- [34] Gaurav Shrivastava. *Diverse Video Generation*. PhD thesis, University of Maryland, College Park, 2021. 1
- [35] Gaurav Shrivastava and Abhinav Shrivastava. Diverse video generation using a gaussian process trigger. *arXiv preprint arXiv:2107.04619*, 2021. 2, 6
- [36] Gaurav Shrivastava, Ser-Nam Lim, and Abhinav Shrivastava. Video dynamics prior: An internal learning approach for robust video enhancements. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [37] Gaurav Shrivastava, Ser-Nam Lim, and Abhinav Shrivastava. Video decomposition prior: Editing videos layer by layer. In *The Twelfth International Conference on Learning Representations*, 2024. 1
- [38] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *International Conference on Learning Representations*, 2020. 1
- [39] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*, 2021. 1
- [40] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild, 2012. 4
- [41] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhutdinov. Unsupervised learning of video representations using lstms. In *Proceedings of the 32Nd International Conference on International Conference on Machine Learning - Volume 37*, pages 843–852. JMLR.org, 2015. 2
- [42] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric and challenges, 2018. 4
- [43] Ruben Villegas, Jimei Yang, Seunghoon Hong, Xunyu Lin, and Honglak Lee. Decomposing motion and content for natural video sequence prediction. *CoRR*, abs/1706.08033, 2017. 2
- [44] Ruben Villegas, Jimei Yang, Yuliang Zou, Sungryull Sohn, Xunyu Lin, and Honglak Lee. Learning to generate long-term future via hierarchical prediction, 2017. 2
- [45] Ruben Villegas, Arkanath Pathak, Harini Kannan, Dumitru Erhan, Quoc V. Le, and Honglak Lee. High fidelity video prediction with large stochastic recurrent neural networks, 2019. 2
- [46] Vikram Voleti, Alexia Jolicoeur-Martineau, and Chris Pal. Mcvd-masked conditional video diffusion for prediction, generation, and interpolation. *Advances in Neural Information Processing Systems*, 35:23371–23385, 2022. 2, 5, 6, 7
- [47] Jacob Walker, Abhinav Gupta, and Martial Hebert. Patch to the future: Unsupervised visual prediction. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 3302–3309, 2014. 2
- [48] Jacob Walker, Kenneth Marino, Abhinav Gupta, and Martial Hebert. The pose knows: Video forecasting by generating pose futures. In *International Conference on Computer Vision*, 2017. 2
- [49] Yunbo Wang, Lu Jiang, Ming-Hsuan Yang, Li-Jia Li, Ming-sheng Long, and Li Fei-Fei. Eidetic 3d LSTM: A model for video prediction and beyond. In *International Conference on Learning Representations*, 2019. 2
- [50] Nevan Wichers, Ruben Villegas, Dumitru Erhan, and Honglak Lee. Hierarchical long-term video prediction without supervision, 2018. 2
- [51] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. Videogpt: Video generation using vq-vae and transformers. *arXiv preprint arXiv:2104.10157*, 2021. 6
- [52] Jenny Yuen and Antonio Torralba. A data-driven approach for event prediction. In *European Conference on Computer Vision*, pages 707–720. Springer, 2010. 2