# A unified model for continuous conditional video prediction

Xi Ye
Polytechnique Montreal
xi.ye@polymtl.ca

Guillaume-Alexandre Bilodeau
Polytechnique Montreal
gabilodeau@polymtl.ca

## Abstract

*Different conditional video prediction tasks, like video future frame prediction and video frame interpolation, are normally solved by task-related models even though they share many common underlying characteristics. Furthermore, almost all conditional video prediction models can only achieve discrete prediction. In this paper, we propose a unified model that addresses these two issues at the same time. We show that conditional video prediction can be formulated as a neural process, which maps input spatio-temporal coordinates to target pixel values given context spatio-temporal coordinates and context pixel values. Specifically, we feed the implicit neural representation of coordinates and context pixel features into a Transformer-based non-autoregressive conditional video prediction model. Our task-specific models outperform previous work for video future frame prediction and video interpolation on multiple datasets. Importantly, the model is able to interpolate or predict with an arbitrary high frame rate, i.e., continuous prediction. Our source code is available at https://npvp.github.io.*

## 1. Introduction

The human ability to anticipate dynamic changes in a scene is remarkable, for instance, we can effortlessly anticipate the possible position of a car in the next few seconds or envision the most recent movement of a jogger. Conditional video prediction models are essential for creating human-like intelligent agents, which have numerous applications, such as autonomous driving, robotics, and more. In this paper, we focus on two closely related conditional video prediction tasks, video future frame prediction (VFP), which consists in predicting future frames given some past frames, and video frame interpolation (VFI) with long temporal gap. Even though VFP and VFPI share some similarities, they have been tackled by totally different methods for a long time. For example, most VFP models depend on Convolutional-LSTMs (ConvLSTMs) to predict the future frames autoregressively [6, 11, 25, 26], and they are inca-

pable of performing video interpolation. Meanwhile, most VFI methods capture the motion between input frames by estimating optical flow [20, 32] or local convolution kernels [33, 34]. None of these methods are able to solve the VFP problem.

Therefore, we address the problem of unifying multiple conditional video prediction tasks to solve them with a single model. Our motivation to propose a unified model is that multi-task learning is a good regularization for a better representation learning [16]. Thus, we believe that a unified model is beneficial for each individual task. Additionally, one common problem for almost all conditional video prediction models is that they can only generate video with a fixed frame rate, i.e., discrete prediction. However, the real world is continuous over the spatio-temporal space. Therefore, we also aim to develop a conditional video prediction model that is able to recover the underlying continuous signal of the real world given a discrete dataset, and thus enable many useful applications, e.g., generating videos with an arbitrary high frame rate, or generating a climate video with irregular time interval [35].

To address these two problems, we propose a novel unsupervised continuous conditional video prediction method based on neural processes (NPs) [14] and implicit neural representations (INRs) [40, 45]. NPs have been successfully applied for image completion [14, 40], but to the best of our knowledge, this is the first work that successfully achieves conditional video prediction based on neural processes. In addition to VFP and VFI, the flexibility of NPs also enables our model to achieve video past frame extrapolation (VPE) and video random missing frames completion (VRC). By formulating conditional video prediction as a neural process, we build a supervised mapping from any target spatio-temporal coordinate of frames to target pixel values, given observed context coordinates and pixel values. The spatio-temporal coordinates are encoded by an implicit neural representation model to achieve continuous generation. More specifically, we firstly extract the features of each video frame by training a convolutional neural network (CNN) autoencoder. Then, a Transformer-based prediction model parameterizes an attentive neural process,

which takes the target coordinates as inputs, conditions on context coordinates and context frame features, then outputs the target frame features that are finally fed into the CNN decoder to reconstruct the frame pixels. A Fourier Feature Network (FFN) learns the neural representations of coordinates, which serve as the positional information for the Transformer-based neural process model. Finally, a global latent variable is learned with variational methods to deal with prediction uncertainties. Our main contributions are:

- We propose the first neural process model for conditional video prediction (*NPVP*), which tackles VFP, VFI, VPE and VRC with one model;

- Our work is the first that successfully adapts INRs for temporal continuous VFP;

- The proposed model is able to make temporal continuous video generation, i.e., generating video with an arbitrary high frame rate;

- Our model outperforms the state-of-the-art (SOTA) models for VFP and VFI over multiple datasets.

## 2. Background

**Neural processes (NPs).** Given a set of labeled contexts $C = (X_C, Y_C) = \{(x_i, y_i)\}_{i \in \mathcal{I}(C)}$ and an unlabeled target set $T = X_T = \{x_i\}_{i \in \mathcal{I}(T)}$, Garnelo *et al.* [14] proposed (conditional) neural processes to model the predictive distribution $p(f(T)|C, T)$, where $\mathcal{I}(S)$ denotes the indices of data points in set $S$, function $f : X \rightarrow Y$ defines the mapping from domain $X$ to $Y$. Specifically, the contexts $C$ are firstly encoded and aggregated into a context embedding of fixed dimension, then $p(f(T)|C, T)$ is parameterized by a neural network with the inputs of context embedding and $T$. NPs are efficient because they preserve merits of both Gaussian processes and deep neural networks. An important property of NPs is that they are permutation invariant in $C$ and $T$ [14]. NPs can be extended to a latent variable version that accounts for the uncertainty of $f(T)$ based on VAE [24]. In order to solve the underfitting problem of NPs, Kim *et al.* [23] proposed to replace the context feature aggregation operation by an attention mechanism. NPs are required to be with scalability, flexibility and permutation invariance [23]. They have been successfully applied for image completion [14, 23, 40]. In this case, the pixel coordinates are considered as $x_i$ and pixel values are considered as $y_i$. Benefiting from the permutation invariance, NPs can predict missing pixel values condition on context pixels in arbitrary patterns.

**Implicit neural representations (INRs).** INRs [40, 45] are techniques which solve the spectral bias problem of neural networks and thus achieve a continuous mapping between the input coordinates and target signal values, e.g., pixel values. There are mainly two different types of INRs. The first one is a Fourier Feature Network (FFN) [45], which uses a Fourier feature mapping for the input of a normal multiple layer perceptron (MLP) to enable the learning of high-frequency signal components effectively. The second one is a SInusoidal REpresentation Network (SIREN) [40]. SIREN depends on periodic activation functions, i.e., sinusoidal activations, to continuously represent the signals with fine details. Both FFN and SIREN are efficient, and some work [3, 58] have proven that they are equivalent to each other. INRs have been adopted for many computer vision tasks, including image generation [41], unconditional video generation [42] and video interpolation [9].

## 3. Related work

Any video to video synthesis task can be considered as a conditional video prediction, including video translation between different domains [49], video super-resolution [17, 36], VFP and VFI. We particularly focus on the work related to VFI and VFP. The classical supervised VFI models take optical flow-based [20, 32] or kernel-based methods [33, 34] to learn the motion for the intermediate frames. The drawback is that those models require a high frame rate training dataset, which is relatively expensive to acquire. Some unsupervised VFI models have been proposed in recent years, for example, Reda *et al.* [37] developed a unsupervised VFI model based on cycle consistency. A more recent optical flow-based CNN model, VideoINR [9], successfully utilizes the INRs for continuous VFI.

VFP models can be categorized into many different types, such as deterministic models [8, 54], stochastic models [1, 11], pixel-direct generation models [8, 13] and transformation-based models [7, 22]. Almost all the VFP models are autoregressive models based on ConvLSTMs or Transformers [53, 56]. Recently, a few promising non-autoregressive VFP models were proposed [28, 48, 57]. VPTR [57] is a transformer-based non-autoregressive VFP model, but it only predicts future frames with a fixed frame rate. By combining ConvLSTMs with a neural ordinary differential equation (ODE) solver, Vid-ODE [35] is the first method that unifies the VFP and VFI into a single model, and it is able to generate temporally continuous video. Another work, masked conditional video diffusion (MCVD) [48] extends the 3D CNN-based diffusion models for video generation, but it is not an NP model and it is not able to do continuous video prediction. Benefiting from the flexibility of NPs, our model is able to perform video random missing frames completion (VRC) contrarily to MCVD. Furthermore, our model achieves stochastic prediction based on VAE instead of a diffusion model.

Our model is different from previous work in two aspects. Firstly, none of them are built to be a neural process. We believe that a NP is a better choice because it is per-

mutation invariant in constrast to ConvLSTMs or 3D-CNN that are not. Therefore, they do not have the flexiblilty of our model to achieve multiple conditional video prediction tasks with a single model. Secondly, implicit neural representation (together with NPs) enables our model to predict frames at any given temporal coordinate, even though they are not seen during the training. Most of the previous models can only predict video frames with a fixed frame rate, which is defined by the training dataset. Vid-ODE [35] circumvents this limitation by introducing ODE. However, our NP-based model outperforms Vid-ODE for both VFP and VFI, as it will be shown in experiments. Another exception is VideoINR [9], but VideoINR can only perform VFI and it does not satisfy the properties of NPs.

## 4. Proposed method

Figure 1 (a) depicts the proposed *NPVP* framework. Given some context frames $V_C \in \mathbb{R}^{L_C \times I_h \times I_w \times I_c}$, *NPVP* is trained to generate some target frames $V_T \in \mathbb{R}^{L_T \times I_h \times I_w \times I_c}$, where $L_C$ and $L_T$ denote the number of context frames and the number of target frames, respectively. $I_h, I_w, I_c$ are the image height, width, and the number of color channels. Firstly, the visual feature $Y_C \in \mathbb{R}^{L_C \times H \times W \times D}$ of context frames are extracted by a frame encoder, then a predictor predicts the target visual feature $\hat{Y}_T \in \mathbb{R}^{L_T \times H \times W \times D}$ given $Y_C, X_C$ and $X_T$. $X_C$ and $X_T$ denotes the context and target spatio-temporal coordinate encodings learned by a Fourier Feature Network (Figure 1 (b)). $X_C \in \mathbb{R}^{L_C \times H \times W \times D}$ and $X_T \in \mathbb{R}^{L_T \times H \times W \times D}$, $H, W, D$ denote the visual feature height, width and channels. Finally, the target frames are reconstructed by a frame decoder given $\hat{Y}_T$.

Benefiting from the flexibility of NPs, given a video clip, we can randomly select any number of frames at any time step $i$ to be the context frames $V_C$, and the remaining ones are the frames $V_T$ to be predicted. In this way, our model can be trained as a general model for multiple conditional video prediction tasks. For example, if we make all the context frames have a smaller time coordinates than all the target frames, the model is trained for VFP specifically, but if a model is trained with random context, it is able to solve multiple conditional prediction tasks at the same time.

We divided our model into two parts, a frame autoencoder and a NP-based predictor that operates over the feature space. The reason for our choice is that directly learning a NP-based model over the pixel space is expensive. Operating over the feature space allows us to train the model in two stages. We firstly train the frame encoder and frame decoder by ignoring the NP-based predictor, and then fix the parameters of the frame autoencoder to learn the predictor. The detail architectures of the autoencoder, Fourier feature network (Figure 1 (b)), and the NP-based Predictor (Figure 2) are described in the following.
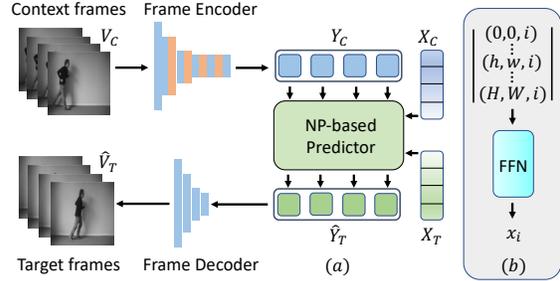


Figure 1. (a) *NPVP* framework. $Y_C$: context features; $\hat{Y}_T$: predicted target features; $X_C/X_T$: context/target spatio-temporal coordinate encodings. (b) Implicit Neural Represenations (INRs).

### 4.1. Autoencoder

We use a custom autoencoder that is adapted from Pix2Pix [19]. Specifically, we integrate non-local 2d attention layers (orange layers of the frame encoder in Figure 1 (a)) from SAGAN [59] into the CNN encoder to improve its performance. There is no modification for the frame decoder of Pix2Pix. The autoencoder is trained with a simple $L_1$ loss between input frame $I$ and reconstructed frame $\hat{I}$, i.e., $\mathcal{L}_1(I, \hat{I}) = |I - \hat{I}|$. Recall that the predictor is ignored during the learning of the autoencoder, and the autoencoder is fixed during the learning of the predictor.

### 4.2. Fourier feature network for INRs

We selected a FFN [45] instead of SIREN [40] because a FFN is easier to train. For a visual feature $y_i \in \mathbb{R}^{H \times W \times D}$ of one frame, where $i$ is the temporal coordinate, the FFN takes the coordinate $(h, w, i)$ of each feature vector at different spatio-temporal location as input, and outputs a $D$-dimensional coordinate encoding for $(h, w, i)$. The INRs is shown Figure 1 (b). $x_i \in \mathbb{R}^{H \times W \times D}$ denotes all the spatial-temporal coordinate encodings of a frame feature $y_i$. Then, $X_C$ and $X_T$ contains all the $x_i$ of the context and target coordinates respectively. Specifically, for an input 3D coordinate vector $(h, w, i)$, the FFN firstly projects it to a higher dimensional space by a Gaussian random noise matrix, then the projections are fed into a MLP with ReLU activation functions to get the output coordinate encoding. The spatio-temporal coordinates are normalized to the range $[0, 1]$. The FFN is jointly learned with the NP-based predictor.

The implicit neural representations $X_C$ and $X_T$ generated by the FFN encode the spatio-temporal location information of context features $Y_C$ and target features $Y_T$. They are critical for the learning of our Transformer-based predictor (section 4.3) because a Transformer is permutation invariant. After training, INRs are able to generalize to unseen input coordinates, which means that we can get the coordinate encoding $x_i$ at any real-number temporal coordinate $i$. Because the model predicts different $y_i$ given con-

texts and different target $x_i$ (section 4.3), we can achieve a continuous generation. For the VFP task, if we need to predict target frames beyond the maximum temporal coordinates used during training, we can perform a "block-wise" autoregressive prediction. Specifically, take the predicted future (target) frames as the past (context) frames for the next block of future (target) frames.

### 4.3. NP-based Predictor

Our NP-based predictor is designed as an attentive neural process [23] based on a video representation learning Transformer, VidHRFormer [57]. This is motivated by the fact that an attentive neural process preserves the permutation invariance and solves the underfitting problem of a vanilla NP, and VidHRFormer satisfies all the requirements of attentive neural processes.

**Loss function.** Given contexts $(X_C, Y_C)$ and $X_T$, a NP learns to maximize conditional log-likelihood $\log p(Y_T|X_C, Y_C, X_T)$. In other words, a NP makes probabilistic predictions for $Y_T$, which is normally assumed to follow a factorized Gaussian distribution $p$ [14, 23]. However, we argue that a simpler point prediction is better for video prediction. Firstly, $Y_T$ has a much higher dimensionality than the simple regression datasets or images in [14, 23], therefore it is expensive to predict the covariance even if it is diagonal (factorized). Preliminary experiments show that the predicted diagonal covariance only captures the variation of high-frequency noise instead of the desired temporal uncertainty of motion. Furthermore, even if the predicted diagonal covariance successfully estimates the motion uncertainty at each time step, we cannot use it to sample coherent video sequences because each time step is independent of each other [14]. Considering that learning full covariance is infeasible, one solution is to enforce causality among different time steps by making autoregressive predictions, which however suffers from accumulated errors and low inference speed.

Therefore, we propose a better solution which is to model $p$ as a Laplacian distribution with a constant scale parameter. This corresponds to an efficient point prediction. In order to achieve coherent sequence sampling, we introduce a global latent variable $z_e$ into the vanilla NP [14], where $z_e$ explains the uncertainty of the whole sequence $Y_T$ and thus is named "event variable". Thus, we can describe the generative process of $Y_T$ as:

$$p(Y_T|X_C, Y_C, X_T) =$$
$$\int p(Y_T|X_T, X_C, Y_C, z_e)q(z_e|X_C, Y_C)dz_e, \quad (1)$$

where $q(z_e|X_C, Y_C)$ defines a conditional prior distribution for $z_e$. By adapting a VAE [24], we can approximate Eq. 1

by maximizing the evidence lower bound (ELBO):

$$ELBO = \mathbb{E}_{q_\phi(z_e|X_T, Y_T)}[\log p(Y_T|X_T, X_C, Y_C, z_e)]$$
$$- \beta D_{KL}(q_\phi(z_e|X_T, Y_T)||q_\psi(z_e|X_C, Y_C)), \quad (2)$$

where $\beta$ is a hyperparameter. $\phi$ and $\psi$ denote the parameters of two factorized Gaussian distribution $\mathcal{N}(\mu_\phi(X_T, Y_T), \sigma_\phi(X_T, Y_T))$ and $\mathcal{N}(\mu_\psi(X_C, Y_C), \sigma_\psi(X_C, Y_C))$ respectively. Specifically, the first term in the RHS of Eq. 2 forces the predictor to reconstruct $Y_T$. The KL divergence is a regularization term that prevents the $z_e$ sampled from targets to deviate too far from the $z_e$ sampled from the context, with the assumption that both context frames and target frames are generated from the same latent event space.

As $p$ follows a Laplacian distribution with a constant scale parameter, maximizing the log-likelihood (first term in the RHS of Eq. 2) is equivalent to minimizing an $L_1$ loss, that is $|Y_T - \hat{Y}_T|$. We can derive the loss function as

$$\mathcal{L} = |Y_T - \hat{Y}_T| + \beta D_{KL}(q_\phi(z_e|X_T, Y_T)||q_\psi(z_e|X_C, Y_C)). \quad (3)$$

However, in practice, we find that learning with Eq. 3 cannot generate predictions with good visual quality, because the $L1$ loss in the RHS of Eq. 3 does not consider the curvature of the latent feature manifold learned by the frame autoencoder [39]. Therefore, we feed $\hat{Y}_T$ to the fixed frame decoder to reconstruct target frames $\hat{V}_T$ and minimize another pixel reconstruction $L_1$ loss, i.e., $|V_T - \hat{V}_T|$, at the same time. In this way, the supervisory signal from the pixel $L_1$ loss minimizes the geodesic distance between $Y_T$ and $\hat{Y}_T$ [4]. Then, the final loss function is

$$\mathcal{L} = \gamma|V_T - \hat{V}_T| + |Y_T - \hat{Y}_T|$$
$$+ \beta D_{KL}(q_\phi(z_e|X_T, Y_T)||q_\psi(z_e|X_C, Y_C)), \quad (4)$$

where $\gamma$ is a hyperparameter.

**Network architecture.** The architecture of the proposed NP-based predictor is shown in Figure 2. It is composed of a context event encoding module, a target event encoding module and a Transformer Decoder $\mathcal{T}_D$. The context event encoding shares the same architecture as the target event encoding, they are responsible for learning $q_\psi$ and $q_\phi$ respectively. $\mathcal{T}_D$ is responsible for predicting $\hat{Y}_T$.

In detail, the architecture of the context event encoding path can be formalized by these operations,

$$M_C = \mathcal{T}_E(X_C, Y_C) \quad (5)$$
$$\mu_\psi, \sigma_\psi = E_C(mean(M_C)), \quad (6)$$

where $\mathcal{T}_E : X_C \times Y_C \rightarrow M_C \in \mathbb{R}^{L_C \times H \times W \times D}$ denotes a Transformer encoder, which is composed by multiple VidHRFormer blocks [57]. The spatio-temporal separated attention mechanism of VidHRFormer block ensures
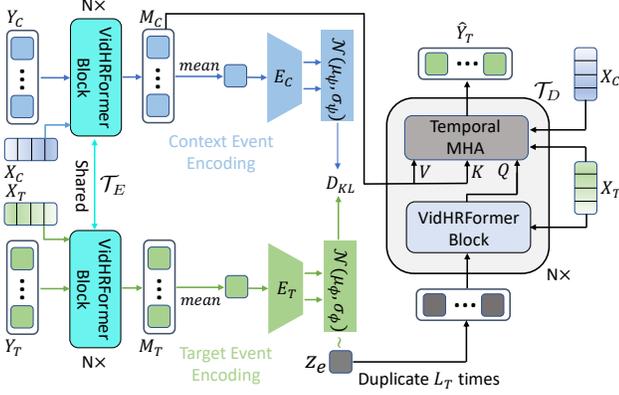
Figure 2. Architecture of the NP-based predictor. Target event encoding is only used during training. $z_e$ is sampled from $\mathcal{N}(\mu_\psi, \sigma_\psi)$ during test.

the permutation invariance along the temporal dimension. $X_C$ is fused into $\mathcal{T}_E$ as the positional encodings. In order to generate $(\mu_\psi, \sigma_\psi)$, $M_C$ is firstly averaged along the temporal dimension, then passed through an event encoder $E_C : \mathbb{R}^{H \times W \times D} \to \mathbb{R}^{H \times W \times D}$, where $E_C$ is a small CNN with two output heads for $\mu_\psi$ and $\sigma_\psi$, respectively. We use the *mean* operation because it is an efficient aggregation method and it is permutation invariant. Please see the supplementary material for the detailed architectures of VidHRFormer block [57] and $E_C$.

The target event encoding path has the same architecture as the context event encoding path, which is formalized as,

$$M_T = \mathcal{T}_E(X_T, Y_T) \tag{7}$$

$$\mu_\phi, \sigma_\phi = E_T(mean(M_T)), \tag{8}$$

where $M_T \in \mathbb{R}^{L_T \times H \times W \times D}$. Note that the target event encoding shares the same $\mathcal{T}_E$ with the context event encoding, but it has a different target event encoder $E_T : \mathbb{R}^{H \times W \times D} \to \mathbb{R}^{H \times W \times D}$.

We hypothesize that all target visual features are generated by event variable $z_e \in \mathbb{R}^{H \times W \times D}$, and $z_e \sim \mathcal{N}(\mu_\phi, \sigma_\phi)$ during training. During test, the ground truth $Y_T$ is not accessible, then $z_e$ is sampled from the learned context prior event space $\mathcal{N}(\mu_\psi, \sigma_\psi)$, which is generated by the context event encoding path.

Finally, $\hat{Y}_T$ is generated by conditioning on $(X_T, X_C, M_C, z_e)$ through another Transformer $\mathcal{T}_D$,

$$\hat{Y}_T = \mathcal{T}_D(X_T, X_C, M_C, z_e). \tag{9}$$

The architecture of a $\mathcal{T}_D$ block is the same as the Transformer decoder block of used in VPTR [57] (see supplementary material for the detailed architecture). $M_C$ and $X_C$ are fed into $\mathcal{T}_D$ as the key/value and positional encodings of the encoder-decoder temporal multi-head attention layer respectively. They provide context information for $\hat{Y}_T$. $X_T$

is also injected into $\mathcal{T}_D$ as positional encodings. Note that event variable $z_e$ is duplicated $L_T$ times and fed into $\mathcal{T}_D$ as the initial query for each target frame feature. In this way, we can generate $\hat{Y}_T$ with arbitrary frame rate, i.e., continuous prediction, as long as we input the desired $X_T$, which is produced by the trained FFN for free.

## 5. Experiments

We evaluated the proposed predictor on multiple realistic video datasets, KTH [38], BAIR [12], KITTI [15], Cityscapes [10], and a synthetic video dataset, Stochastic Moving MNIST (SM-MNIST) [11]. KITTI and Cityscapes are resized to $128 \times 128$, all others are resized to the resolution of $64 \times 64$. Following the experimental configuration of previous work, we present the quantitative results of Peak Signal-to-Noise Ratio (PSNR), Fréchet Video Distance (FVD) [46], Learned Perceptual Image Patch Similarity (LPIPS) [60] and Structural Similarity Index Measure (SSIM). The LPIPS is reported in $10^{-3}$ scale. Same as previous stochastic methods, 100 different predictions are sampled for each test example, then the best SSIM, LPIPS, PSNR, and the average FVD of the generated samples are reported.

For a fair comparison with previous task-specific models, we first train different models for VFP and VFI respectively, i.e., we train task-specific models following the same training procedures that they used. Then, we present results with a unified model that is not task-specific, but since the training is different, results are not fully comparable with task-specific models. The supplementary material includes more qualitative examples and implementation details.

### 5.1. Task-specific models

**VFP.** The VFP experimental results are summarized in Table 1. For the KTH dataset, our *NPVP* model is trained to predict 10 future frames given 10 past frames. During test, the performance is evaluated on predicting 20 future frames conditioned on 10 past frames, which is achieved by a block-wise autoregressive inference. Our *NPVP* achieves the best SSIM and outperforms previous methods by a large margin in terms of LPIPS. For KITTI, *NPVP* is trained to predict 5 future frames given 4 past frames. Compared with previous methods, our *NPVP* reaches the best performance in terms of both SSIM and LPIPS. Qualitative results (Figure 3) show that *NPVP* predicts future frames with good visual quality despite the large motion of KITTI dataset, which has a low frame rate of 10 fps. The results on KITTI dataset demonstrate that *NPVP* is capable of challenging real-world traffic video prediction.

For the Cityscapes dataset, *NPVP* is trained to predict 10 future frames given 2 past frames, but 28 future frames are predicted by block-wise autoregressive inference during test. *NPVP* achieves the best SSIM and the second-best

| Models | KTH, *10 → 20* | | | | Models | KITTI, *4 → 5* | | | Models | Cityscapes, *2 → 28* | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS ↓ | | | SSIM↑ | LPIPS↓ | | | FVD↓ | SSIM↑ | LPIPS↓ |
| PredRNN++ [50] | 28.47 | 0.865 | 228.9 | | PredNet [31] | 47.56 | 629.5 | | SVG-LP [11] | 1300.26 | 0.574 | 549.0 |
| STMFANet [21] | **29.85** | 0.893 | 118.1 | | MCNet [47] | 55.48 | 373.9 | | VRNN 1L [5] | 682.08 | 0.609 | 304.0 |
| E3D-LSTM [51] | *29.31* | 0.879 | - | | Voxel Flow [29] | 42.62 | 415.9 | | Hier-VRNN [5] | 567.51 | 0.628 | 264.0 |
| Conv-TT-LSTM [44] | 28.36 | *0.907* | 133.4 | | FVS [54] | 60.77 | *304.9* | | GHVAEs [52] | *418.00* | *0.740* | 194.0 |
| Vid-ODE [35] | 28.19 | 0.878 | *80.0* | | SADM [2] | *64.72* | 311.6 | | MCVD-concat [48] | **141.31** | 0.690 | **112.0** |
| VPTR-NAR [57] | 26.96 | 0.879 | 86.1 | | | | | | | | | |
| *NPVP* (ours) | 27.66 | **0.909** | **66.0** | | *NPVP* (ours) | **66.12** | **279.0** | | *NPVP* (ours) | 768.04 | **0.744** | *183.2* |

Table 1. VFP results on KTH, KITTI and Cityscapes. **Boldface**: best results. *Blue*: second best results.

LPIPS. We suspect that the gap between *NPVP* and MCVD-concat in terms of FVD is due to the fact that a vanilla VAE is not expressive enough [5]. All methods for Cityscapes in Table 1 are VAE-based, except for the MCVD-concat, which uses a denoising diffusion model (DDM) [18, 43]. It has a brightness-changing problem [48], but outperforms all VAE-based methods for the FVD score. Castrejón *et al*. [5] has shown that increasing the levels of VAE latent variables is beneficial for the expressiveness of VFP models. Meanwhile, DDM can be considered as a special form of hierarchical VAEs with a large levels of latent variable, and all the latent variables have the same dimensionality as the data, i.e., a video clip. Therefore, MCVD-concat achieves the best FVD probably because of the high expressiveness of its latent variables. Our proposed method may suffer from weaker temporal coherence caused by its non-autoregressive prediction. All VAE-based models in Table 1 are auto-regressive models, but the loss function of *NPVP* assumes that frames at different time steps are independent of each other to preserve the unified model flexibility, even though the temporal attention exchange information between different frames. Nevertheless, *NPVP* achieves a comparable or better performance than the SOTAs. Visual examples of VFP on Cityscapes (Figure 3) shows that our method predicts better than MCVD [48] that suffers from brightness change problems.

**VFI.** We follow the experimental protocol in [48]. For the VFI task, given $p$ past frames and $f$ future frames as the context, the model is trained to generate $k$ intermediate frames, i.e., target frames. For VFI, the context spans across the past and future, which limits the event stochasticity.

The VFI results are summarized in Table 2. For the KTH dataset, our task-specific *NPVP* (*15→10*) model outperforms the SOTA MCVD (*15→10*) in terms of SSIM. It also outperforms Vid-ODE in terms of all metrics by a large margin. Note that a smaller $p + f$ and a larger $k$ means a harder VFI task. Therefore, we can draw the conclusion that our model is better than SVG-LP and SDVI-*full*, because (*15→10*) is harder than (*18→7*). Our *NPVP* (*10→5*) also outperforms MCVD (*10→5*) model in terms of both SSIM and PSNR. The *NPVP* (*10→5*) outperforms
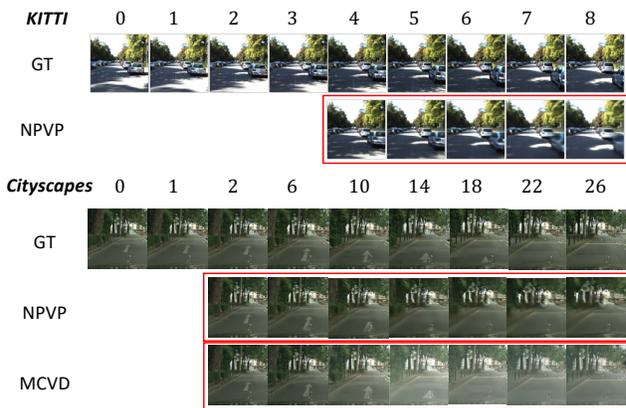


Figure 3. VFP examples on KITTI and Cityscapes. Frames inside the red boxes are future frames predicted by the model.

*NPVP* (*15→10*). We believe that it is because 10 context frames are enough to provide good context information for the KTH dataset, and interpolating 5 frames is much easier than interpolating 10 frames.

For the SM-MNIST dataset, our *NPVP* (*10→10*) outperforms all the previous methods in terms of both PSNR and SSIM, even for the easier (*18→7*) and (*10→5*) tasks. It means that given the same number of context frames, our model is able to interpolate more intermediate frames with better quality. The *NPVP* (*10→5*) outperforms *NPVP* (*10→10*) as expected. In short, the results demonstrate that *NPVP* achieves a new SOTA for VFI on the SM-MNIST dataset. The randomness of SM-MNIST only occurs when the characters bounce off the boundaries. Most of the time, the character trajectories of intermediate frames can be determined based on the past and future frames, thus the model tends to ignore the event variable as little stochasticity exists. A similar phenomenon is observed for VFI on the KTH dataset because natural human motion in missing frames is mostly constrained by past and future movements.

For the BAIR dataset, our *NPVP* (*18→7*) outperforms SVG-LP and SDVI-*full* by a large margin. Compared with

| Models | KTH | | | | SM-MNIST | | | | BAIR | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | (p+f→k) | PSNR↑ | SSIM↑ | LPIPS↓ | (p+f→k) | PSNR↑ | SSIM↑ | LPIPS↓ | (p+f→k) | PSNR↑ | SSIM↑ | LPIPS↓ |
| SVG-LP [11] | (18→7) | 28.13 | 0.883 | - | (18→7) | 13.54 | 0.741 | - | (18→7) | 18.65 | 0.846 | - |
| SDVI-*full* [55] | (18→7) | 29.19 | 0.901 | - | (18→7) | 16.03 | 0.842 | - | (18→7) | 21.43 | 0.880 | - |
| Vid-ODE [35] | - | 31.77 | 0.911 | 48.0 | - | - | - | - | - | - | - | - |
| MCVD [48] | (15→10) | 34.67 | 0.943 | - | (10→10) | 20.94 | 0.854 | - | (4→5) | *25.16* | *0.932* | - |
| MCVD [48] | (10→5) | *35.61* | 0.963 | - | (10→5) | 27.69 | 0.940 | - | - | - | - | - |
| *NPVP (ours)* | (15→10)† | 33.60 | *0.969* | 22.3 | (10→10) | *28.11* | *0.958* | *17.3* | (18→7) | 22.97 | 0.909 | *21.8* |
| | (10→5) | **37.17** | **0.984** | **10.5** | (10→5) | **34.34** | **0.992** | **4.1** | (4→5) | **25.28** | **0.933** | **14.7** |

Table 2. VFI results. $p/f$: number of past/future frames; $k$: number of intermediate frames to interpolate. **Smaller $p+f$ and larger $k$ means a harder VFI task**. †: $p = 8$, $f = 7$; for our other models, $p$ equals to $f$. **Boldface**: best results. *Blue*: second best results.

MCVD ($4{\to}5$), *NPVP* ($4{\to}5$) achieves a slightly better performance in terms of both PSNR and SSIM.

## 5.2. A unified model for VFP, VFI, VPE and VRC

We trained a unified *NPVP* on the KTH dataset to demonstrate that our NP-based conditional video prediction model is flexible enough to perform VFP, VFI, video past frame extrapolation (VPE), and video random missing frames completion (VRC) with one single learned model. More importantly, all the aforementioned tasks can be solved with an arbitrary high frame rate, i.e., a continuous prediction. In order to learn one model for all tasks, *NPVP* is trained with random contexts, i.e., given a video clip with $L$ frames, we draw $L_C$ frames (the probability of each frame follows a uniform distribution) as contexts and the remaining $L_T = L - L_C$ frames are target frames, together with their corresponding coordinates. $L = 20$, and the value of $L_C$ varies in the range of $[4, 16]$.

In Figure 4, we present examples of one model for all four different conditional video prediction tasks. The first row is the ground-truth (GT) frames. The frames inside a red box are target frames generated by the model given the other context frames. We can observe that the target frames visual quality of VRC and VFI is better than the VPE and VFP, because VPE and VFP only observe the future or past frames, thus there are more uncertainties and it is harder to predict target frames. On the contrary, the context frames of VRC and VFI are scattered across the temporal dimension, which provides more accurate context motion information about the ground-truth event. Because the model tends to minimize the loss quickly by solving the easier tasks, the model trained with random contexts needs more epochs to reach a comparable performance on VFP and VPE tasks compared to a model trained specifically for that task.

Besides, we investigated the quantitative change of target frames visual quality w.r.t. the varying number of context frames for VFP and VFI with our unified model (Figure 5). The results demonstrate that all three metrics of quality monotonically improves as more context frames are fed into the model. That is, the model generates more accurate tar-
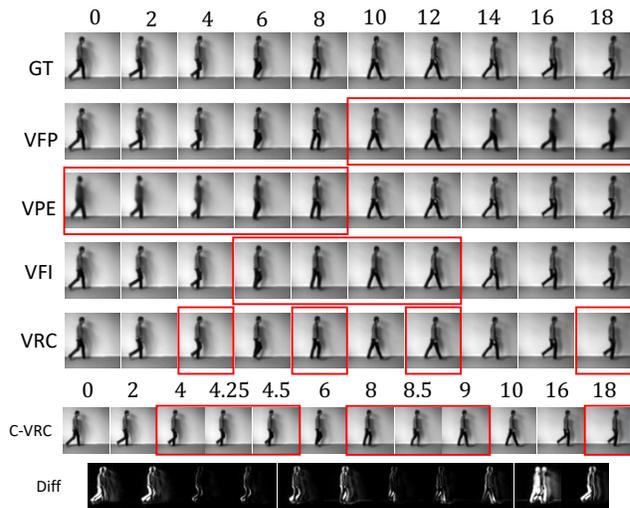


Figure 4. One model for all tasks. Frames inside the red boxes are target frames generated by the model. C-VRC denotes continuous VRC. Diff are the difference images between neighboring frames of C-VRC to show that they are all different and that the temporal coordinates are taken into account.
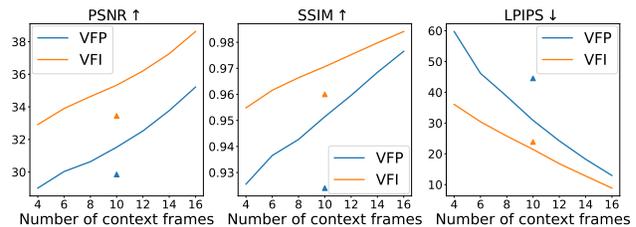


Figure 5. Metric curves of VFP and VFI on KTH for an increasing number of context frames. ▲ and ▲ denote results of task-specific *NPVP* ($10{\to}10$) for VFI and VFP respectively.

get frames given more context frames, which is in alignment with the property of NPs [14]. Figure 5 also shows the performance gap between VFP and VFI, which indicates that VFI is an easier task. Results of the task-specific *NPVP* ($10{\to}10$) for VFI and VFP are plotted in Figure 5. The uni-

| Models | | | | | | | VFI, *5+5 →10* | | | VFP, *10 → 10* | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| INR | $\mathcal{T}_D$-6L | $\mathcal{T}_D$-8L | $\mathcal{T}_E$ | $FL_1$ | $PL_1$ | $D_{KL}$ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| 1 | | ✓ | | | ✓ | | | 30.58 | 0.937 | 50.07 | 28.05 | 0.904 | 71.75 |
| 2 ✓ | ✓ | | | ✓ | | | 31.34 | 0.953 | 37.03 | 28.83 | 0.920 | 71.04 |
| 3 ✓ | ✓ | | | ✓ | ✓ | | 33.17 | 0.958 | 31.11 | 29.60 | 0.920 | 62.96 |
| 4 ✓ | | ✓ | | ✓ | ✓ | | 33.20 | 0.959 | 29.64 | 29.85 | 0.922 | 57.49 |
| 5 ✓ | | ✓ | ✓ | ✓ | ✓ | | *33.77* | *0.962* | *26.83* | *30.11* | *0.927* | *53.89* |
| *NPVP* ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | **34.07** | **0.972** | **25.59** | **30.37** | **0.941** | **52.18** |

Table 3. Ablation Study on KTH dataset, trained with random contexts. $FL_1$ denotes feature space $L_1$ loss. $PL_1$ denotes pixel space $L_1$ loss. *NPVP* is the stochastic counterpart of model 5. **Boldface**: best results. *Blue*: second best results.

fied model outperforms the task-specific models in terms of all metrics for both VFI and VFP. Therefore, the results validate our motivation that multi-task learning is beneficial.

Finally, a continuous video random missing frames completion (C-VRC) (Figure 4, bottom) experiment was conducted to show the continuous generation ability of our model. We observe that it is capable of generating frames at unseen temporal coordinates, such as $4.25, 4.5, 8.5$. Therefore, we can do conditional video prediction at an arbitrary high frame rate. We observe from the generated videos that the continuous predictions by the unified model have a much better temporal consistency than task-specific models. This is because the random contexts help the model to learn more complex conditional distributions [23], thus improving the generalization ability of INRs.

### 5.3. Ablation study

We summarize the ablation study results in Table 3. The base model (model 1) is a deterministic model (no target event encoding) with a $\mathcal{T}_D$ with 6 layers ($\mathcal{T}_D$-6L), which is trained only by a $L_1$ loss over feature space, and spatio-temporal coordinates are directly fed into the model, i.e., without INR. We gradually modify the architecture and loss function to investigate the influence of the various components. All models are trained with random contexts and thus we evaluate the performance of both VFI and VFP.

**INR.** We observe a significant performance boost on all metrics in model 2, which validates the effectiveness of INR. Qualitative results also show that predictions of model 1 lack temporal consistency, and have unnatural, choppy motion. *We believe INR improves performance in three ways*: 1) learned Fourier features represent high-frequency details better, contributing to better visual quality. 2) INR learns continuous mapping from coordinates to video pixels, improving continuous prediction. 3) Learnable Fourier features are better multi-dimensional positional encoding, capturing more complex relationships [27].

**Pixel $L_1$ loss.** By introducing the pixel $L_1$ loss, we observe performance improvement in terms of almost all metrics. Particularly for LPIPS, model 3 outperforms model 2 by a large margin. It proves that pixel $L_1$ loss is beneficial

for the target frames visual quality.

**Size of $\mathcal{T}_D$.** In order to investigate the influence of the size of $\mathcal{T}_D$, we increase the number of layers of $\mathcal{T}_D$ from 6 to 8 ($\mathcal{T}_D$-8L). Comparing the results of model 4 with model 3, we observe a performance boost for all metrics of VFI and VFP, which indicates that a larger $\mathcal{T}_D$ is useful.

**Context Transformer encoder $\mathcal{T}_E$.** For Model 1-4, there is no explicit temporal relationship learning among context frame features. Because a good aggregation of the context information is critical to improve the performance of NPs, we now include the context Transformer $\mathcal{T}_E$, which models the temporal relationship of $Y_C$ and generates $M_C$ for the $\mathcal{T}_D$ and $z_e$. Comparing model 5 with the previous models, we observe further improvement over all metrics for both VPF and VFI, especially for the LPIPS metric.

**Stochastic vs Deterministic.** Finally, we modified model 5 to be stochastic by introducing the VAE architecture to give our proposed method *NPVP*. *NPVP* outperforms model 4 for both VFI and VFP in terms of all metrics as we expected, because the event variable takes into account the randomness of prediction instead of only predicting the average of all possible outcomes as its deterministic counterpart [1]. We propose a two-stage training strategy to stabilize the training of *NPVP*. For the initial stage, we ignore target event encoding path and $D_{KL}$, in other words, we trained a model 5 first, then the target event encoding path and $D_{KL}$ are included for training.

## 6. Conclusion

We proposed a novel continuous conditional video prediction model based on a neural process and implicit neural representations. By training with random contexts, we can address multiple conditional video prediction tasks with only one model, including future frame prediction, frame interpolation, past frame extrapolation and random missing frame completion. Importantly, all the tasks can be tackled with an arbitrary high frame rate. Results show that our model achieves the SOTA for video future frame prediction and video frame interpolation over multiple datasets.

# References

[1] Mohammad Babaeizadeh, Chelsea Finn, Dumitru Erhan, Roy Campbell, and Sergey Levine. Stochastic variational video prediction. In *ICLR*, 2018. 2, 8

[2] Xinzhu Bei, Yanchao Yang, and Stefano Soatto. Learning Semantic-Aware Dynamics for Video Prediction. 2021. 6, 1

[3] Nuri Benbarka, Timon Höfer, Hamd ul-Moqeet Riaz, and Andreas Zell. Seeing Implicit Neural Representations As Fourier Series. In *WACV*, 2022. 2

[4] Sarthak Bhagat, Shagun Uppal, Zhuyun Yin, and Nengli Lim. Disentangling Multiple Features in Video Sequences Using Gaussian Processes in Variational Autoencoders. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *ECCV*, pages 102–117, Cham, 2020. Springer International Publishing. 4

[5] L. Castrejon, N. Ballas, and A. Courville. Improved Conditional VRNNs for Video Prediction. In *ICCV*, Oct. 2019. 6

[6] Zheng Chang, Xinfeng Zhang, Shanshe Wang, Siwei Ma, Yan Ye, Xinguang Xiang, and Wen Gao. MAU: A Motion-Aware Unit for Video Prediction and Beyond. In *NeurIPS*, May 2021. 1

[7] Xiongtao Chen, Wenmin Wang, Jinzhuo Wang, and Weimian Li. Learning object-centric transformation for video prediction. In *MM 2017 - Proceedings of the 2017 ACM Multimedia Conference*, pages 1503–1512, 2017. 2

[8] Xinyuan Chen, Chang Xu, Xiaokang Yang, and Dacheng Tao. Long-term video prediction via criticization and retrospection. *IEEE Transactions on Image Processing*, 29:7090–7103, 2020. 2

[9] Zeyuan Chen, Yinbo Chen, Jingwen Liu, Xingqian Xu, Vidit Goel, Zhangyang Wang, Humphrey Shi, and Xiaolong Wang. VideoINR: Learning Video Implicit Neural Representation for Continuous Space-Time Super-Resolution. In *CVPR*, 2022. 2, 3

[10] 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. pages 3213–3223, 2016. 5

[11] Emily Denton and Rob Fergus. Stochastic Video Generation with a Learned Prior. In *ICML*, 2018. 1, 2, 5, 6, 7

[12] Frederik Ebert, Chelsea Finn, Alex X. Lee, and Sergey Levine. Self-supervised visual planning with temporal skip connections. In *CoRL*, 2017. 5

[13] Jean-Yves Franceschi, Edouard Delasalles, Mickael Chen, Sylvain Lamprier, and Patrick Gallinari. Stochastic Latent Residual Video Prediction. In *ICLR*, 2020. 2

[14] Marta Garnelo, Dan Rosenbaum, Christopher Maddison, Tiago Ramalho, David Saxton, Murray Shanahan, Yee Whye Teh, Danilo Rezende, and S. M. Ali Eslami. Conditional Neural Processes. In *ICML*, 2018. 1, 2, 4, 7

[15] A Geiger, P Lenz, C Stiller, and R Urtasun. Vision meets robotics: The KITTI dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, Sept. 2013. 5

[16] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, Nov. 2016. 1

[17] Muhammad Haris, Gregory Shakhnarovich, and Norimichi Ukita. Recurrent Back-Projection Network for Video Super-Resolution. In *CVPR*, 2019. 2

[18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models. In *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates, Inc., 2020. 6

[19] Phillip Isola, Jun Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017. 3

[20] Huaizu Jiang, Deqing Sun, Varun Jampani, Ming-Hsuan Yang, Erik Learned-Miller, and Jan Kautz. Super SloMo: High Quality Estimation of Multiple Intermediate Frames for Video Interpolation. In *CVPR*, 2018. 1, 2

[21] Beibei Jin, Yu Hu, Qiankun Tang, Jingyu Niu, Zhiping Shi, Yinhe Han, and Xiaowei Li. Exploring Spatial-Temporal Multi-Frequency Analysis for High-Fidelity and Temporal-Consistency Video Prediction. In *CVPR*, 2020. 6

[22] B. Jin, Y. Hu, Y. Zeng, Q. Tang, S. Liu, and J. Ye. VarNet: Exploring Variations for Unsupervised Video Prediction. In *IROS*, 2018. 2

[23] Hyunjik Kim, Andriy Mnih, Jonathan Schwarz, Marta Garnelo, Ali Eslami, Dan Rosenbaum, Oriol Vinyals, and Yee Whye Teh. Attentive Neural Processes. In *ICLR*, 2019. 2, 4, 8

[24] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *ICLR*, 2014. 2, 4

[25] Y. Kwon and M. Park. Predicting Future Frames Using Retrospective Cycle GAN. In *CVPR*, 2019. 1

[26] Sangmin Lee, Hak Gu Kim, Dae Hwi Choi, Hyung-Il Kim, and Yong Man Ro. Video Prediction Recalling Long-Term Motion Context via Memory Alignment Learning. In *CVPR*, 2021. 1

[27] Yang Li, Si Si, Gang Li, Cho-Jui Hsieh, and Samy Bengio. Learnable fourier features for multi-dimensional spatial positional encoding. In *Advances in Neural Information Processing Systems*, volume 34, 2021. 8

[28] Zhouyong Liu, Shun Luo, Wubin Li, Jingben Lu, Yufan Wu, Chunguo Li, and Luxi Yang. ConvTransformer: A Convolutional Transformer Network for Video Frame Synthesis. In *arXiv:2011.10185 [cs]*, Nov. 2020. 2

[29] Z. Liu, R. A. Yeh, X. Tang, Y. Liu, and A. Agarwala. Video Frame Synthesis Using Deep Voxel Flow. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 4473–4481, Oct. 2017. ISSN: 2380-7504. 6

[30] Ilya Loshchilov and Frank Hutter. SGDR: Stochastic Gradient Descent with Warm Restarts. In *ICLR*, 2017. 1

[31] William Lotter, Gabriel Kreiman, and David Cox. Deep predictive coding networks for video prediction and unsupervised learning. In *5th international conference on learning representations, ICLR 2017 - conference track proceedings*, pages 1–18, 2017. arXiv: 1605.08104 tex.arxivid: 1605.08104. 6

[32] S. Niklaus and F. Liu. Softmax Splatting for Video Frame Interpolation. In *CVPR*, pages 5436–5445, 2020. 1, 2

[33] Simon Niklaus, Long Mai, and Feng Liu. Video Frame Interpolation via Adaptive Convolution. pages 670–679, 2017. 1, 2

[34] Simon Niklaus, Long Mai, and Oliver Wang. Revisiting Adaptive Convolutions for Video Frame Interpolation. pages 1099–1109, 2021. 1, 2

[35] Sunghyun Park, Kangyeol Kim, Junsoo Lee, Jaegul Choo, Joonseok Lee, Sookyung Kim, and Edward Choi. Vid-ODE: Continuous-Time Video Generation with Neural Ordinary Differential Equation. In *AAAI*, May 2021. 1, 2, 3, 6, 7

[36] E. Pérez-Pellitero, M. S. M. Sajjadi, M. Hirsch, and B. Schölkopf. Photorealistic video super resolution. In *Workshop and challenge on perceptual image restoration and manipulation (PIRM) at the 15th european conference on computer vision (ECCV)*, 2018. 2

[37] Fitsum A. Reda, Deqing Sun, Aysegul Dundar, Mohammad Shoeybi, Guilin Liu, Kevin J. Shih, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Unsupervised Video Interpolation Using Cycle Consistency. pages 892–900, 2019. 2

[38] C. Schuldt, I. Laptev, and B. Caputo. Recognizing human actions: a local SVM approach. In *ICPR*, 2004. 5

[39] Hang Shao, Abhishek Kumar, and P. Thomas Fletcher. The Riemannian Geometry of Deep Generative Models. pages 315–323, 2018. 4

[40] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit Neural Representations with Periodic Activation Functions. In *NeurIPS*, 2020. 1, 2, 3

[41] Ivan Skorokhodov, Savva Ignatyev, and Mohamed Elhoseiny. Adversarial Generation of Continuous Images. In *CVPR*, 2021. 2

[42] Ivan Skorokhodov, Sergey Tulyakov, and Mohamed Elhoseiny. StyleGAN-V: A Continuous Video Generator With the Price, Image Quality and Perks of StyleGAN2. In *CVPR*, 2022. 2

[43] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-Based Generative Modeling through Stochastic Differential Equations. 2021. 6

[44] Jiahao Su, Wonmin Byeon, Jean Kossaifi, Furong Huang, Jan Kautz, and Animashree Anandkumar. Convolutional Tensor-Train LSTM for Spatio-temporal Learning. In *NeurIPS*, 2020. 6

[45] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier Features Let Networks Learn High Frequency Functions in Low Dimensional Domains. In *NeurIPS*, 2020. 1, 2, 3

[46] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphaël Marinier, Marcin Michalski, and Sylvain Gelly. FVD: A new Metric for Video Generation. In *ICLR Workshop*, 2019. 5

[47] Ruben Villegas, Jimei Yang, Seunghoon Hong, Xunyu Lin, and Honglak Lee. Decomposing motion and content for natural video sequence prediction. In *ICLR*, 2017. 6

[48] Vikram Voleti, Alexia Jolicoeur-Martineau, and Christopher Pal. Masked Conditional Video Diffusion for Prediction, Generation, and Interpolation. In *Advances in Neural Information Processing Systems*, 2022. 2, 6, 7, 4, 5

[49] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-Video Synthesis. In *Advances in Neural Information Processing Systems*, volume 31, 2018. 2

[50] Yunbo Wang, Zhifeng Gao, Mingsheng Long, Jianmin Wang, and Philip S. Yu. PredRNN++: Towards A Resolution of the Deep-in-Time Dilemma in Spatiotemporal Predictive Learning. In *ICML*, 2018. 6

[51] Yunbo Wang, Lu Jiang, Ming-Hsuan Yang, Li-Jia Li, Mingsheng Long, and Li Fei-Fei. Eidetic 3D LSTM: A Model for Video Prediction and Beyond. In *ICLR*, 2018. 6

[52] Bohan Wu, Suraj Nair, Roberto Martin-Martin, Li Fei-Fei, and Chelsea Finn. Greedy Hierarchical Variational Autoencoders for Large-Scale Video Prediction. In *CVPR*, pages 2318–2328, 2021. 6

[53] Chenfei Wu, Jian Liang, Lei Ji, Fan Yang, Yuejian Fang, Daxin Jiang, and Nan Duan. N\"UWA: Visual Synthesis Pre-training for Neural visUal World creation. *arXiv:2111.12417 [cs]*, Nov. 2021. 2

[54] Y. Wu, R. Gao, J. Park, and Q. Chen. Future Video Synthesis With Object Motion Prediction. In *CVPR*, 2020. 2, 6

[55] Qiangeng Xu, Hanwang Zhang, Weiyue Wang, Peter Belhumeur, and Ulrich Neumann. Stochastic Dynamics for Video Infilling. pages 2714–2723, 2020. 7

[56] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. VideoGPT: Video Generation using VQ-VAE and Transformers. In *arXiv:2104.10157 [cs]*, Sept. 2021. arXiv: 2104.10157. 2

[57] Xi Ye and Guillaume-Alexandre Bilodeau. VPTR: Efficient Transformers for Video Prediction. In *26th International Conference on Pattern Recognition (ICPR)*, 2022. 2, 4, 5, 6, 1

[58] Gizem Yüce, Guillermo Ortiz-Jiménez, Beril Besbinar, and Pascal Frossard. A Structured Dictionary Perspective on Implicit Neural Representations. In *CVPR*, 2022. 2

[59] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-Attention Generative Adversarial Networks. In *ICML*, 2019. 3

[60] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, pages 586–595, 2018. 5

# Supplementary Material: A unified model for continuous conditional video prediction

## A. Table of important acronyms and notations

| | |
|---|---|
| NPVP: | Neural process for conditional video prediction |
| VFI: | Video frame interpolation |
| VFP: | Video future frame prediction |
| VPE: | Video past frame extrapolation |
| VRC: | Video random missing frames completion |
| NPs: | Neural processes |
| INRs: | Implicit neural representations |
| FFN: | Fourier feature network |
| SIREN: | Sinusoidal representation networks |
| MLP: | Multiple layer perceptron |
| CNN: | Convolutional neural network |
| ConvLSTMs: | Convolutional-LSTMs |
| $V_C$: | Context video frames |
| $V_T$: | Target video frames |
| $X_C$: | Context coordinate representations |
| $Y_C$: | Context video frame features |
| $X_T$: | Target coordinate representations |
| $Y_T$: | Target video frame features |
| $M_C$: | Output feature of $\mathcal{T}_E$ given $X_C$ and $Y_C$ |
| $M_T$: | Output feature of $\mathcal{T}_E$ given $X_T$ and $Y_T$ |
| $z_e$: | event variable |
| $\mathcal{T}_E$: | Transformer encoder |
| $\mathcal{T}_D$: | Transformer decoder |
| $E_C$: | Context event CNN encoder |
| $E_T$: | Target event CNN encoder |

Table S1. Table of important acronyms and notations

## B. Implementation details

### B.1 Datasets

**KTH.** KTH dataset includes grayscale videos of 6 different human actions. Following the experimental setup of previous work, we take persons 1-16 as training set, and persons 17-25 as test set. Random horizontal flips and vertical flips are applied to each video clip as data augmentation.

**BAIR.** BAIR dataset includes RGB video clips of a robot arm randomly moving over a table with small objects. The training and test sets are defined by the creators of BAIR. Random horizontal flips and vertical flips are applied to each video clip as data augmentation.

**SM-MNIST.** Stochastic Moving MNIST (SM-MNIST) is a synthetic dataset includes videos of two randomly moving MNIST characters within a square region. There is no data augmentation for SM-MNIST during training.

**Cityscapes.** Cityscapes dataset includes high-resolution urban traffic videos of many cities. Note that we do not use any annotation provided by Cityscapes, for example, object classes or segmentation masks. Same as previous work, we use the raw video clips from the "$leftImg8bit\_sequence\_trainvaltest.zip$" of Cityscapes. The frames are firstly center-cropped to be square, then we resize the frames to be the resolution of $128 \times 128$. There is no data augmentation for the Cityscapes dataset during training.

**KITTI.** KITTI dataset includes traffic videos across multiple scenarios, including city, residential, road etc. We follow the experimental setup of previous works [2], i.e., randomly select 4 sequences from the raw data of KITTI for testing and use the remaining videos for training. The frames are firstly center-cropped and then resize to be the resolution of $128 \times 128$. Random horizontal flips and vertical flips are applied to each video clip as data augmentation.

### B.2 Training details

**Training of the autoencoder.** For all datasets, the dimension of visual features is set to be $H = 8, W = 8, D = 512$. For input with a resolution of $64 \times 64$, the frame encoder includes 3 downsampling blocks and 2 residual blocks. For input with a resolution of $128 \times 128$, the frame encoder includes 4 downsampling blocks and 3 residual blocks. The number of upsampling blocks for the frame decoder equals to the number of downsampling blocks in the corresponding frame encoder. An Adam optimizer with a learning rate of $1e^{-4}$ is used for the training.

**Training of the NPs-based predictor.** For all datasets, $\gamma = 0.01$. For BAIR and SM-MNIST, $\beta = 1e^{-6}$. For KTH, $\beta = 1e^{-8}$. The predictors are trained by AdamW, we take a cosine annealing learning rate scheduler with warm restarts [30] at every 150 epochs, the maximum learning rate is $1e^{-4}$ and the minimum learning rate is $1e^{-7}$. Gradient clipping is applied to $\mathcal{T}_E$ and $\mathcal{T}_D$ during training. Please visit https://npvp.github.io for the code.

### B.3 Architecture of VidHRFormer block

For the convenience of the readers, we have redrawn the detail architecture of VidHRFormer block [57] and the VPTR decoder block in Figure S1.

### B.4 Architecture of Event encoder $E_C$ and $E_T$

$E_C$ and $E_T$ share the same architecture, see Figure S2. They are implemented by a small neural network with three $Conv - BN - ReLU$ layers and two $Conv$ heads to output $\mu$ and $\sigma$ respectively.
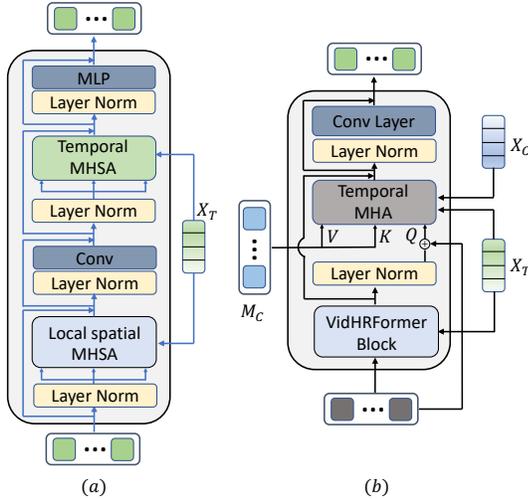
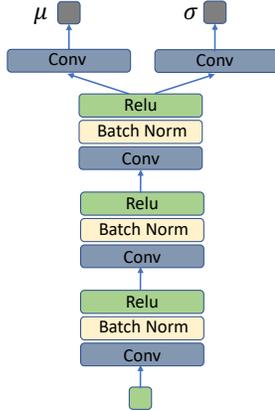Figure S1. (a) VidHRFormer block [57]. (b) Decoder block of VPTR [57].



Figure S2. Architecture of the Event encoders.

amples. We also present the VFI results of MCVD [48] on KTH and SM-MNIST datasets for qualitative comparison. Please visit `https://npvp.github.io` for video examples.

### C.3 Task-specific VFP

We present VFP examples by task-specific *NPVP* models, see Figure S7, Figure S8 and Figure S9. For Cityscapes dataset, we show the results of MCVD for comparison. Please visit `https://npvp.github.io` for video examples.

## C. Qualitative examples

### C.1 Unified model

Here we show another example (see Figure S3) of the unified model on Cityscapes dataset for all four different conditional video prediction tasks. In order to demonstrate the continuous prediction ability of NPVP, we take the trained unified model to solve different tasks with different rates, please visit `https://npvp.github.io` for video examples of a unified model for KTH dataset.

### C.2 Task-specific VFI

We present uncurated VFI examples of KTH, SM-MNIST and BAIR datasets by task-specific *NPVP* models, see Figure S4, Figure S5 and Figure S6. As there is little stochasticity for VFI on KTH and SM-MNIST, we only show the example with the best SSIM from 100 random ex-
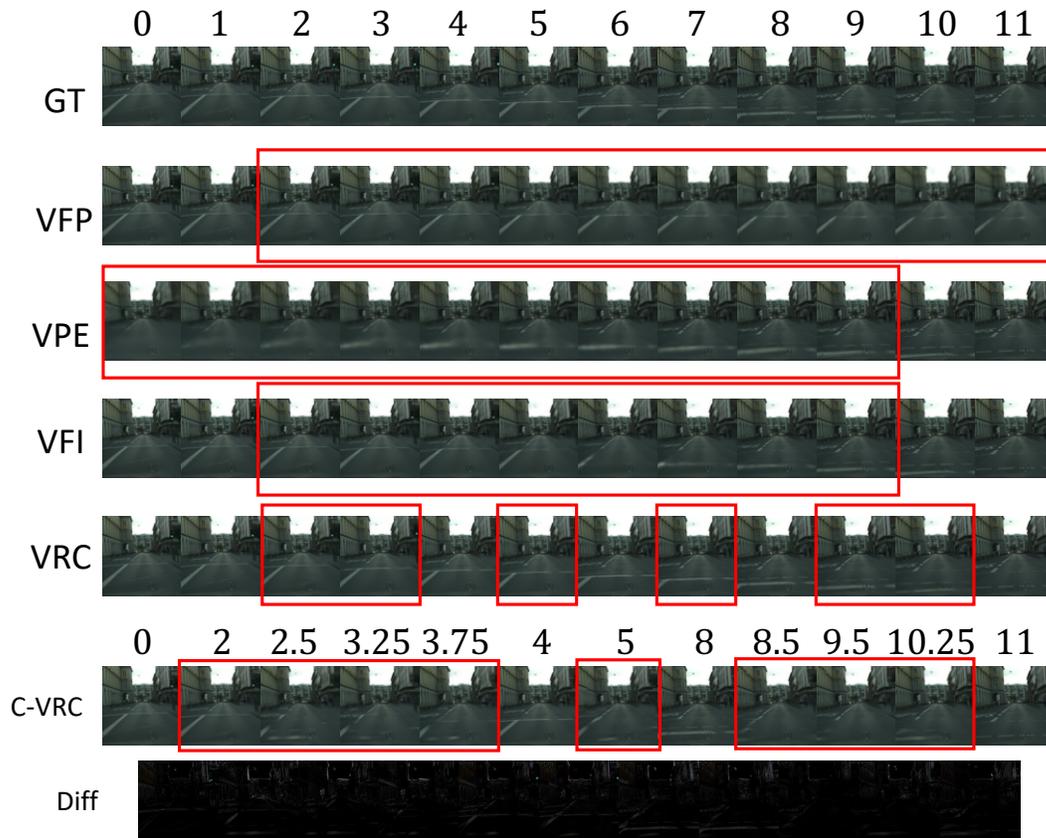
Figure S3. One model for all tasks. Frames inside the red boxes are target frames generated by the model. C-VRC denotes continuous VRC. Diff are the difference images between neighboring frames of C-VRC to show that they are all different and that the temporal coordinates are taken into account.
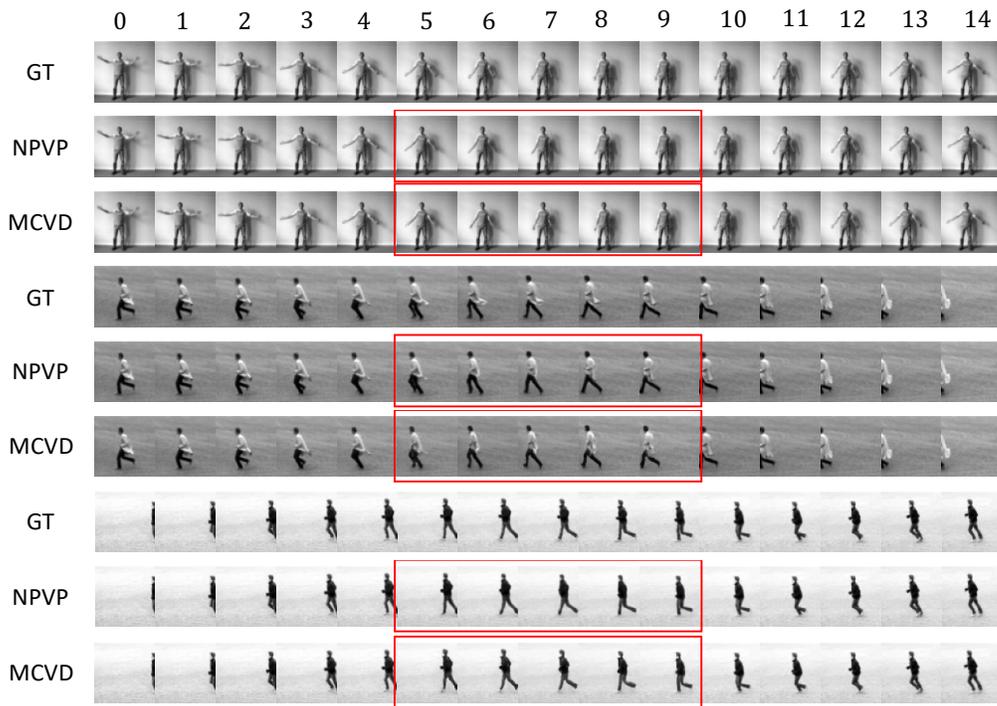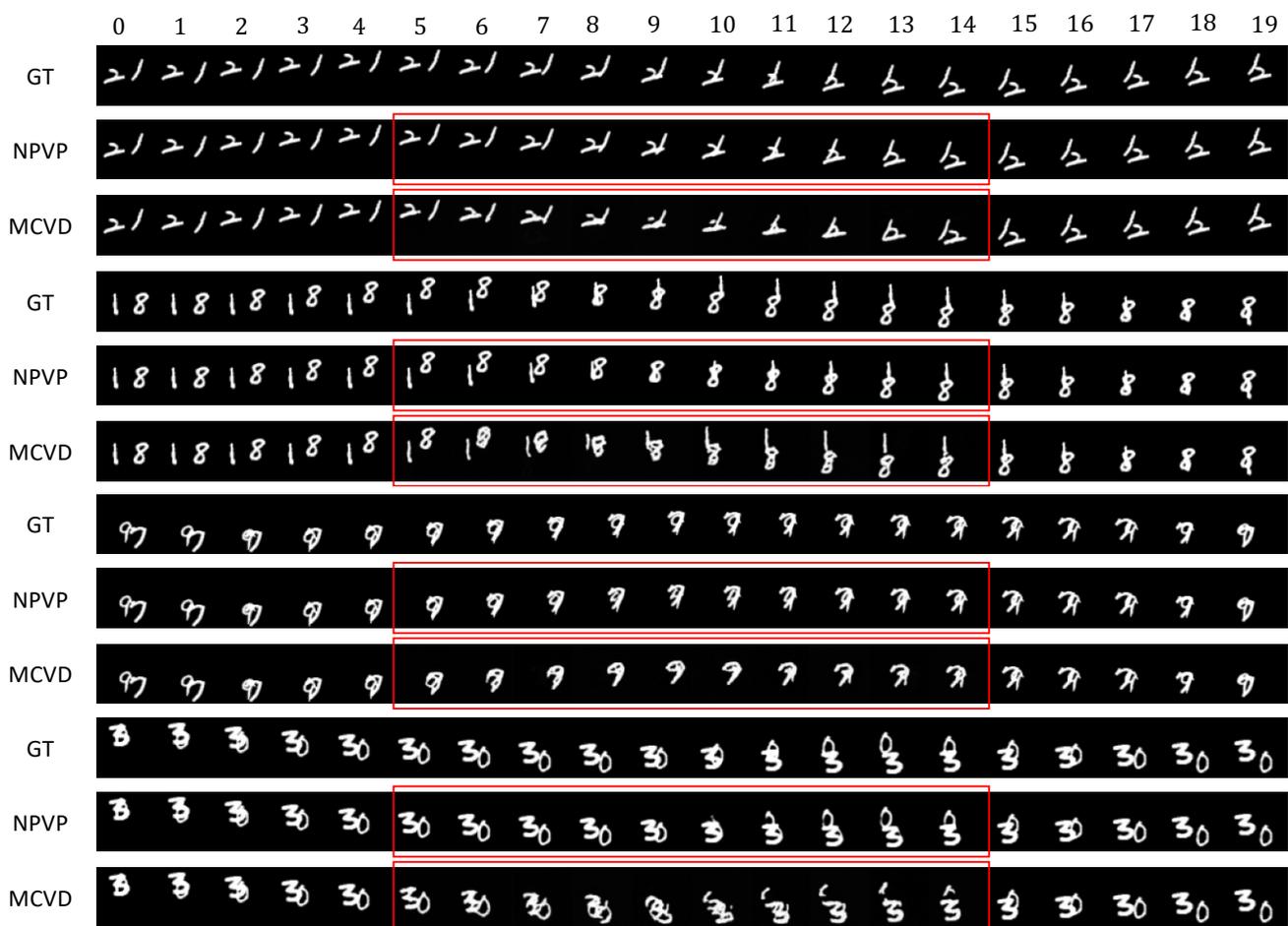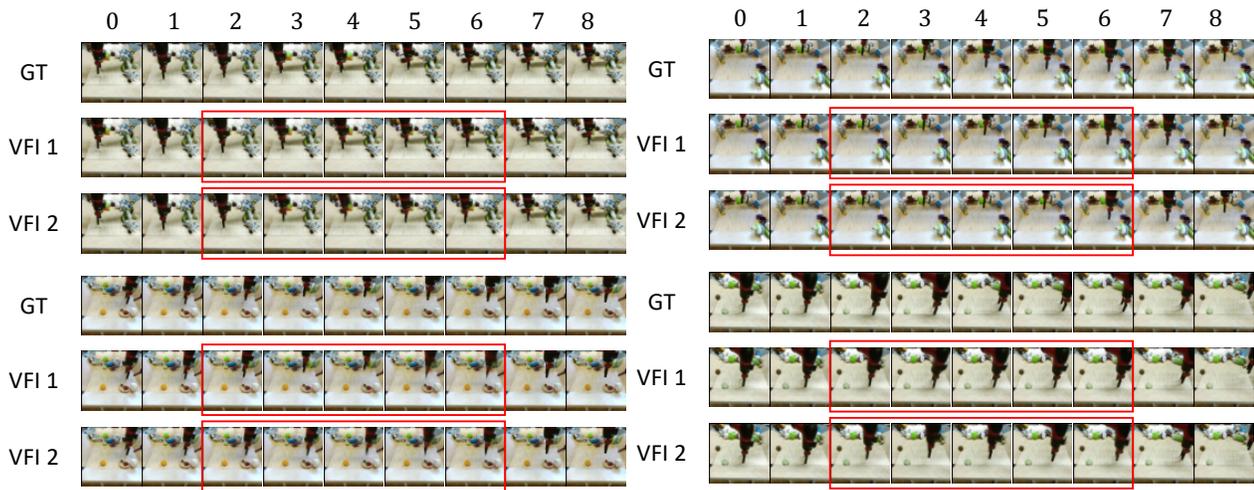
Figure S4. VFI examples on KTH by a Task-specific *NPVP* (*10→5*) model. Frames inside the red boxes are target frames generated by the models. Compared with MCVD [48], predicted moving arms or legs by *NPVP* are more realistic and more similar to the ground-truth.

Figure S5. VFI examples on SM-MNIST by a Task-specific *NPVP* (*10→10*) model. Frames inside the red boxes are target frames generated by the model. Compared with MCVD [48], the interpolation quality of *NPVP* is better as it captures the shape and motion of MNIST characters for missing frames.

Figure S6. VFI examples on BAIR by a Task-specific *NPVP* (*4→5*) model. Frames inside the red boxes are target frames generated by the model. VFI 1 and VFI 2 denote two different random interpolations given the same contexts.
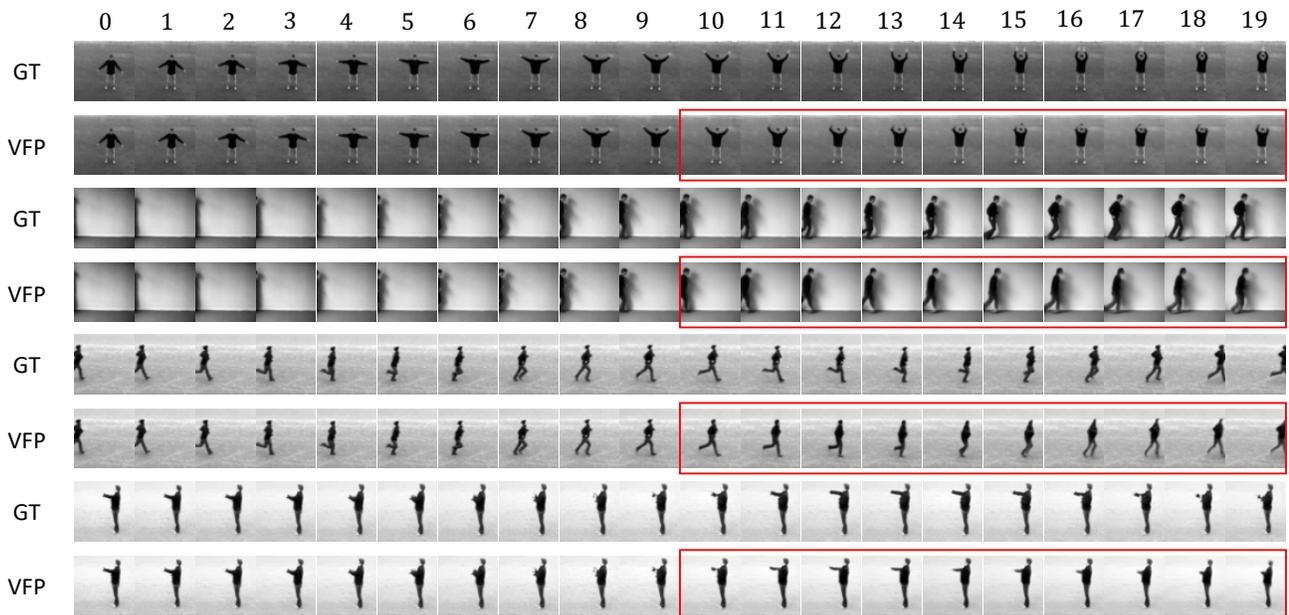


Figure S7. VFP examples on KTH by a Task-specific *NPVP* (*10→10*) model. Frames inside the red boxes are target frames generated by the model.
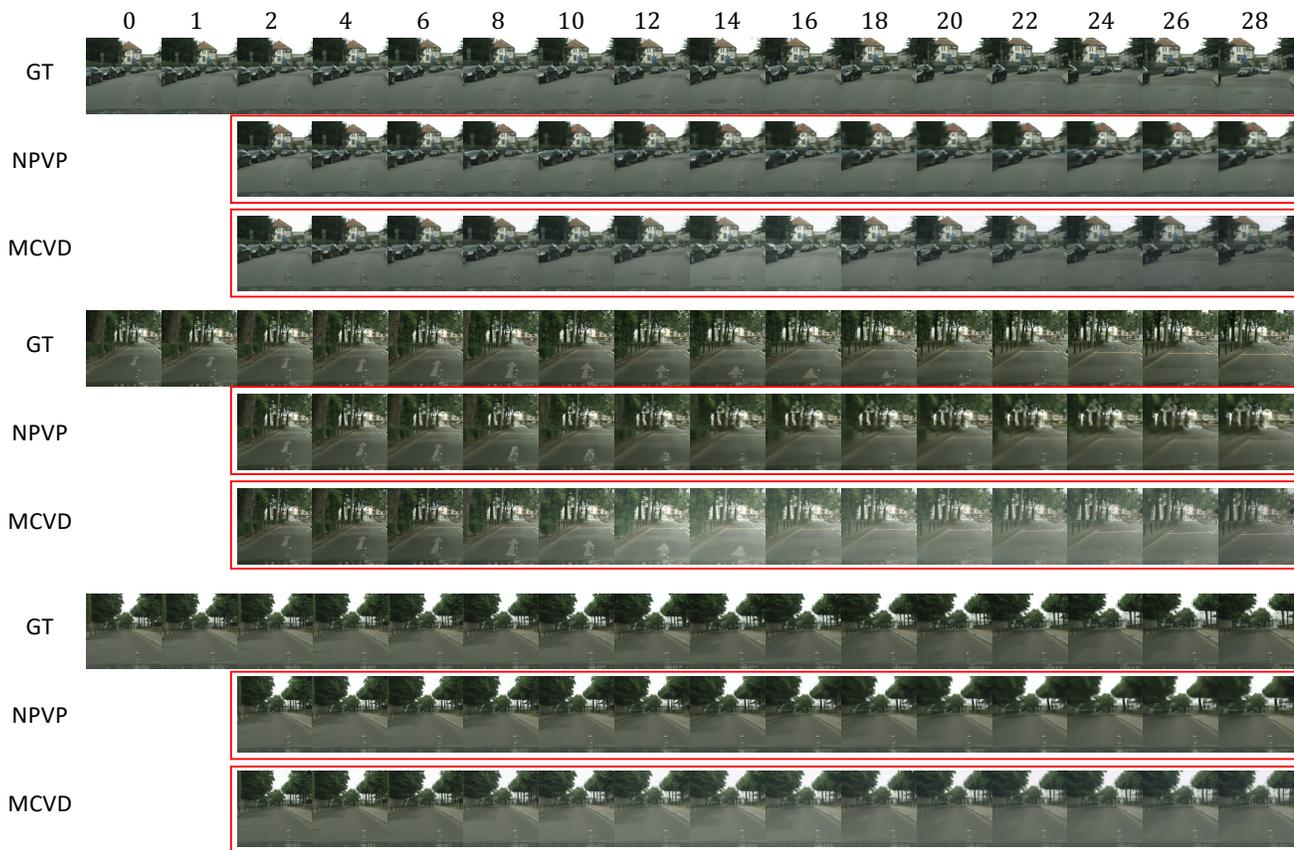
Figure S8. VFP examples on Cityscapes by a Task-specific *NPVP* (*2→28*) model. Frames inside the red boxes are target frames generated by the model. Here we also show the examples generated by MCVD [48], which suffers from a brightness-changing problem.
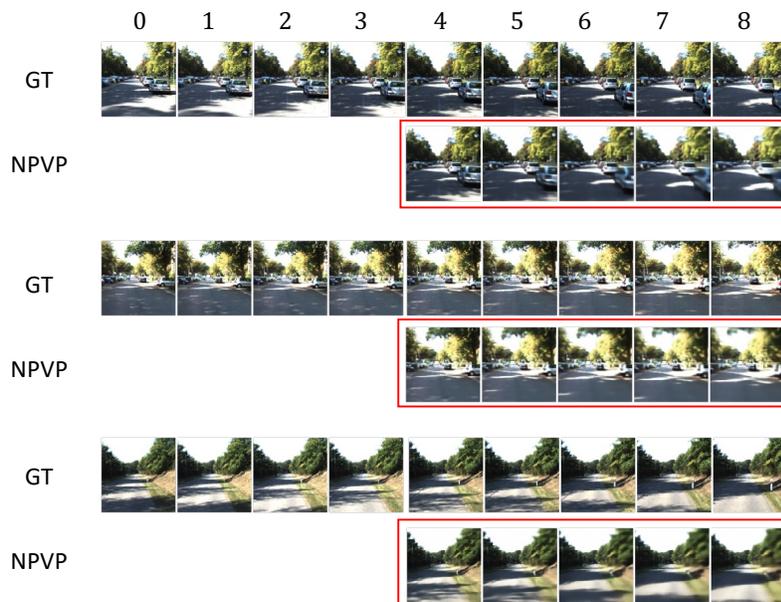


Figure S9. VFP examples on KITTI by a Task-specific *NPVP* (*4→5*) model. Frames inside the red boxes are target frames generated by the model.