

# Fast Fourier Inception Networks for Occluded Video Prediction

Ping Li, *Member, IEEE*, Chenhan Zhang, Xianghua Xu

**Abstract**—Video prediction is a pixel-level task that generates future frames by employing the historical frames. There often exist continuous complex motions, such as object overlapping and scene occlusion in video, which poses great challenges to this task. Previous works either fail to well capture the long-term temporal dynamics or do not handle the occlusion masks. To address these issues, we develop the fully convolutional Fast Fourier Inception Networks for video prediction, termed *FFINet*, which includes two primary components, i.e., the occlusion inpainter and the spatiotemporal translator. The former adopts the fast Fourier convolutions to enlarge the receptive field, such that the missing areas (occlusion) with complex geometric structures are filled by the inpainter. The latter employs the stacked Fourier transform inception module to learn the temporal evolution by group convolutions and the spatial movement by channel-wise Fourier convolutions, which captures both the local and the global spatiotemporal features. This encourages generating more realistic and high-quality future frames. To optimize the model, the recovery loss is imposed to the objective, i.e., minimizing the mean square error between the ground-truth frame and the recovery frame. Both quantitative and qualitative experimental results on five benchmarks, including Moving MNIST, TaxiBJ, Human3.6M, Caltech Pedestrian, and KTH, have demonstrated the superiority of the proposed approach. Our code is available at [GitHub](#).

**Index Terms**—Video prediction, occlusion, temporal dynamics, inpainting, Fourier transform.

## I. INTRODUCTION

VIDEO prediction is the pixel-level task of predicting future frames given past video frames. It has great potential in real-world applications, e.g., climate forecast [1] [2], autonomous driving [3] [4], and robot control [5] [6]. From the randomness perspective, the video prediction methods are divided into two categories, i.e., *deterministic* and *stochastic*. Given a video, the former assumes that the future is deterministic and yields the single prediction, while the latter assumes that there are multiple predictions. However, the latter requires large computations, which hinders its wide applications in highly-demanding scenarios. Hence, this work mainly focuses on the deterministic video prediction.

Most previous works either adopt Recurrent Neural Networks (RNNs) [7] or Convolutional Neural Networks (CNNs) [8] to predict future frames. Typically, Predictive RNN (Pre-RNN) [9], Eidetic 3D Long-Short Term Memory (E3D-LSTM) [10], and Motion-Aware Unit (MAU) [11] adopt RNNs to

capture the temporal dynamics; Deep Voxel Flow (DVF) [12], Disentangling Propagation & Generation (DPG) [13], and Simple Video Prediction (SimVP) [14] use CNNs to capture the inter-frame dependency. In addition, Vision Transformer (ViT) [15] has been proved ineffective in video prediction [14], which is because ViT requires large-scale data for training but the training videos in video prediction task are usually insufficient.

When there exist heavy overlaps among multiple moving objects in video, it is usually difficult to completely recover the object appearances and requires more past frames without overlap. To address this issue, E3D-LSTM [10] and CrevNet (Conditionally Reversible Network) [16] use 3D convolution layers to enlarge the temporal receptive field with larger kernel sizes, which are computationally expensive; 3D convolutions are also used in video saliency prediction [17]. Besides, MAU [11] computes the attention map between the current and the past features for capturing long-term dependency. However, these methods adopt the costly RNNs which desire long training time, e.g., it requires 10 cycles for predicting 10 frames and can not use parallelization. This inspires us to employ the Fast Fourier Transform (FFT) [18] to capture the global spatiotemporal receptive field by designing the **FFT Inception** module. It employs the group convolutions to model the local temporal dynamics and uses the channel-wise fast Fourier transform [19] to capture the global motion trend.

Moreover, previous video prediction methods assume the frames are clean, failing to consider the frames with occlusions, possibly caused by the camera pollution or the object occlusion. To tackle this challenge, one should first inpaint the occluded area and then predict the future frames. Particularly, we adopt the Fast Fourier Convolution (FFC) [20] module to recover the missing area by an **inpainter**, and impose the recovery loss (i.e., MSE-Mean Square Error) onto the inpainter to minimize the error between the recovery frame and the source clean frame.

Hence, we propose the fully-convolutional **Fast Fourier Inception Networks (FFINet)** for video prediction with occlusion, as illustrated in Fig. 1. It is composed of an encoder, an inpainter, a translator, and a decoder. Here, the **translator** is stacked with a series of FFT Inception modules, which are used for capturing both the local and global spatiotemporal features.

Our main contributions are highlighted in the following:

- To our best, we are the first to explore the occluded video prediction, and propose the fully-convolutional fast Fourier Inception Networks (FFINet) for this task.

P. Li, C. Zhang, and X. Xu are with the School of Computer Science and Technology, Hangzhou Dianzi University, Hangzhou 310018, China, and P. Li is also with the Guangdong Laboratory of Artificial Intelligence and Digital Economy (SZ), Shenzhen 518132, China (e-mail: patriclouis.lee@gmail.com, zch2020@hdu.edu.cn, xhxu@hdu.edu.cn).

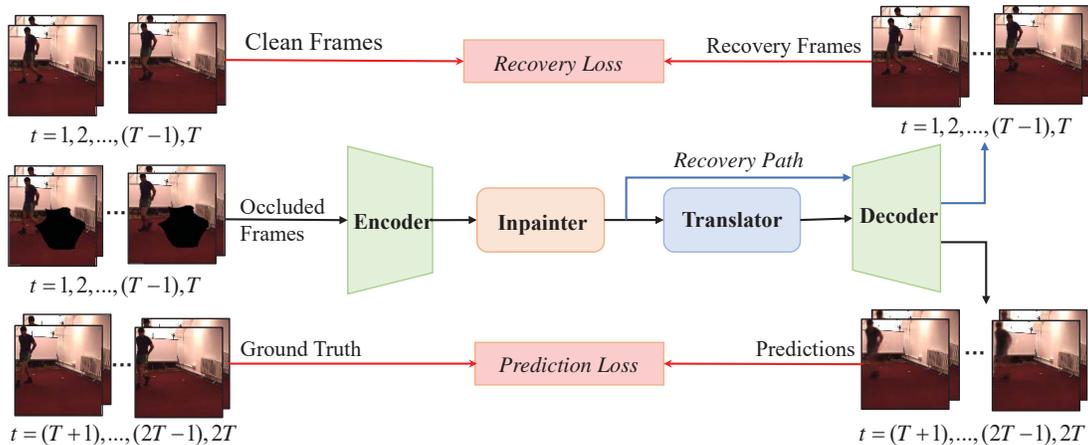


Fig. 1. Illustration of the occluded video prediction task. Given  $T$  input frames, it predicts the future  $T$  frames. Both the recovery loss and the prediction loss adopt the Mean Squared Error (MSE).

- Our framework includes an inpainter that consists of the Fast Fourier Convolution modules to recover the missing areas in video, and a translator that consists of the Fast Fourier Transform inception modules to learn the spatiotemporal features for predicting future frames.
- Empirical studies on several benchmarks including Moving MNIST [21], TaxiBJ [22], Human3.6M [23], Caltech Pedestrian [24], and KTH [25], have demonstrated the superiority of the proposed FFINet framework for video prediction with occlusions.

The rest of this paper is organized as follows. Section II reviews some closely related works and Section III describes the proposed FFINet framework. Then, experimental results are reported in Section IV with rigorous analysis. Finally, we conclude this work in Section V.

## II. RELATED WORK

This section mainly discusses several encoder-translator-decoder architectures of video prediction, including full RNN, RNN or ViT (as translator) with CNN, and full CNN.

### A. Full RNN

The full-RNN methods employ the stacked RNNs for encoding, spatiotemporal feature translation, and decoding. For example, Convolutional LSTM [1] extends the fully-connected LSTMs to the architecture with convolutions, and the memory updates only along the temporal dimension; PredRNN [9] uses the spatiotemporal memory flow to update the memory states along both the spatial and the temporal dimensions, but it suffers from the gradient propagation difficulty when capturing the long-term dependency, which is addressed by PredRNN++ [26]. Moreover, MIM (Memory In Memory) [27] models both the non-stationary and stationary properties by self-updated memory; to reduce the computational and memory costs, fRNN (folded RNN) [28] shares the state between the encoder and decoder layers, while the representation is stratified during learning; Conv-TT-LSTM [29] (Convolutional Tensor-Train LSTM) extends ConvLSTM from the first-order to the higher-order scenario by combining convolutional features across

time, e.g., update the current features by employing the past continuous features. In addition, Motion RNN [2] divides the physical motion into the transient variation and motion trend, which are updated simultaneously; PredRNNv2 [30] develops a decoupling loss and a reverse scheduled sampling strategy to extend the original PredRNN.

### B. RNN or ViT with CNN

This type uses RNN or ViT as the translator, and CNN as the encoder and decoder. CNN is used to learn the spatial feature and RNN or ViT is used to model temporal dynamics. For example, E3D-LSTM [10] uses 3D convolution to extract spatiotemporal features, while CrevNet [16] includes the flow-based encoder and decoder to capture information-preserving features. But they fail to predict the motion trend during more complex scenarios. So, PhyDNet (Physical Dynamics Network) [31] and Vid-ODE (Video generation with Ordinary Differential Equation) [32] use the partial-differential equation or ODE to enhance the physical dynamics modeling; DDME (Dynamic Motion Estimation and Evolution) [33] employs the dynamic convolution to generate convolution kernel at each moment to estimate the temporal dynamics; some works use the dynamic convolution kernel in the fully convolutional network to generate the talking face video [34]. The above methods are good at capturing the short-term dependency, but the predicted frames become obscure when the time gets longer. To capture the long-term dependency, LMC (Long-term Motion Context) [35] uses a memory to save the long-term motion context, which is used for predicting the future; MAU [11] captures the inter-frame motion by broadening the temporal receptive field of the predictive units.

Furthermore, some attempts have been devoted to the sub-tasks in video prediction. For example, Chang *et al.* [36] have developed the spatiotemporal residual predictive model for high-resolution video prediction by focusing more on frame details; CPL (Continual Predictive Learning) [37] presents the mixture world model and the predictive experience replay strategy to alleviate the continual learning problem; MAC (Modular Action Concept network) [38] considers the seman-

tic action-conditional video prediction, which predicts future frames according to the semantic labels that describe the action interactions.

In addition, Vision Transformer (ViT) [15] has been used for modeling the latent dynamics, e.g., Weissenborn *et al.* [39] design a simple auto-regressive video generation model with a 3D self-attention mechanism to yield continuous frames; Rakhimov *et al.* [40] directly uses ViT to model the latent dynamics; TCTN (Temporal Convolutional Transformer Network) [41] employs the transformer-based encoder with temporal convolution layers to capture both the short-term and long-term dependencies. However, ViT model requires large-scale training data to achieve satisfying performance, and it does not bring about much gains on video prediction since the training data in this task is usually insufficient.

### C. Full CNN

This type fully uses CNN to learn and update the spatiotemporal features. Some works will additionally use optical flow to assist the video prediction, such as DVF (Deep Voxel Flow) [12] employs the CNN auto-encoder to learn the voxel flow to reconstruct the frame by using nearby frame voxel flow; Wu *et al.* [42] treat video prediction as the video frame interpolation optimization, which is affected by the optical flow quality; DPG (Disentangling Propagation and Generation) [13] and LCVG (Layered Controllable Video Generation) [43] both separate the foreground from the background, where the former uses optical flow and the latter uses CNN for discrimination. Moreover, Generative Adversarial Network (GAN) [44] is applied to video prediction for increasing the authenticity and the continuity of predicted frames, e.g., rCycleGAN (Retrospective Cycle GAN) [45] enhances the temporal consistency by learning cycle GAN; Xu *et al.* [46] develop a progressive multiple granularity analysis framework to match the prototype motion dynamics with the input sequence. However, the above methods desire complex modules and training skills to improve the performance. Recently, SimVP [14] is a simple video prediction model using common convolution modules, and achieves the State-Of-The-Art (SOTA) performance. But its temporal receptive field is still narrow, limiting its performance upgrade. This inspires us to design an efficient module to enlarge the temporal receptive field to capture the long-term motion dynamics.

## III. OUR FFINET METHOD

This section mainly describes the proposed FFINet framework as depicted in Fig. 2, which includes encoder, inpainter, translator, and decoder.

### A. Problem Formulation

Given a video sequence with  $T$  frames, i.e.,  $\mathcal{X}_T = \{\mathbf{X}_t \in \mathbb{R}^{C \times H \times W} | t = 1, 2, \dots, T\}$ , where  $\mathbf{X}_t$  denotes the  $t$ -th RGB frame with width  $W$ , height  $H$ , and  $C$  channels, the video prediction task aims to generate the future  $T'$  frames, i.e.,  $\mathcal{Y}_{T'} = \{\mathbf{X}_t\}_{t=T+1}^{T+T'}$ , by learning a mapping function  $\mathcal{F}_\theta : \mathcal{X}_T \mapsto \mathcal{Y}_{T'}$ , where  $\theta$  is the model parameter. The vanilla

goal is to minimize the loss between the ground-truth frames  $\mathcal{Y}_{T'}$  and the predicted frames  $\mathcal{F}_\theta(\mathcal{X}_T)$ . For occluded video prediction, the model should first inpaint the corrupted frames and then predict the future frames. Given the occluded video sequence  $\hat{\mathcal{X}}_T = \{\hat{\mathbf{X}}_t \in \mathbb{R}^{C \times H \times W} | t = 1, 2, \dots, T\}$ , the goal is to minimize the loss between the ground-truth frames  $\mathcal{Y}_{T'}$  and the predicted frames  $\mathcal{F}'_\theta(\hat{\mathcal{X}}_T)$  with inpainting.

### B. Encoder

The encoder is used to learn spatial features of video frames, and consists of  $\tilde{N}$  stacked  $3 \times 3$  Conv2D blocks with group normalization  $\text{GN}(\cdot)$  to speed up the convergence and leaky ReLU function  $\sigma(\cdot)$  to enhance the feature nonlinearity. Mathematically, the hidden feature is computed by

$$\mathbf{Z}_i = \sigma(\text{GN}(\text{Conv2D}(\mathbf{Z}_{i-1}))), 1 \leq i \leq \tilde{N}, \quad (1)$$

where  $\mathbf{Z}_0 \in \mathbb{R}^{(B \cdot T) \times C \times H \times W}$  is the input video sequence. Here  $B$  denotes the batch size, and we do the downsampling every two convolution layers by halving the height and the width of the feature map. At the last convolution layer, we concatenate the spatial features along the temporal dimension and reshape them to the tensor  $\mathbf{Z}_{\tilde{N}} \in \mathbb{R}^{B \times (T \cdot \tilde{C}) \times H' \times W'}$ , where  $\tilde{C}$ ,  $H'$ ,  $W'$  are the feature channel, height, and width. Generally, the feature dimension is reduced, and thus the computational cost can be saved a lot. The obtained features will be fed into the inpainter for recovering the missing areas.

### C. Inpainter

The inpainter is used to recover the occluded frames and it adopts two Fast Fourier Convolution (FFC) [20] modules. Each FFC module consists of two inter-connected branches, i.e., common convolution layer on the half of the feature channels, and channel-wise fast Fourier transform on the rest channels. The former captures the local spatial features and the latter captures the global context.

The input of the inpainter is the encoded feature  $\mathbf{Z}_{\tilde{N}}$ , and we learn the local feature  $\mathbf{Z}^l \in \mathbb{R}^{B \times \frac{T \cdot \tilde{C}}{2} \times H' \times W'}$  according to

$$\mathbf{Z}^l = \sigma(\text{GN}(\text{Conv2D}(\mathbf{Z}^l) + \text{Conv2D}(\mathbf{Z}^g))), \quad (2)$$

where Conv2D adopts the kernel size of  $3 \times 3$ . The first convolution layer is used to obtain the local feature, and the second convolution layer is used to fuse both the local and the global features.

Similarly, we learn the global feature  $\mathbf{Z}^g \in \mathbb{R}^{B \times \frac{T \cdot \tilde{C}}{2} \times H' \times W'}$  according to

$$\mathbf{Z}^g = \sigma(\text{GN}(\text{Conv2D}(\mathbf{Z}^l) + \text{FU}(\mathbf{Z}^g))), \quad (3)$$

where the  $3 \times 3$  Conv2D fuses the local and the global features, and  $\text{FU}(\cdot)$  denotes the Fourier Unit that includes the channel-wise fast Fourier transform and  $1 \times 1$  Conv2D, which has the global receptive field. The fused features from the two branches are concatenated along the channel, leading to the recovery feature tensor  $\bar{\mathbf{Z}}_{\tilde{N}} \in \mathbb{R}^{B \times (T \cdot \tilde{C}) \times H' \times W'}$ .

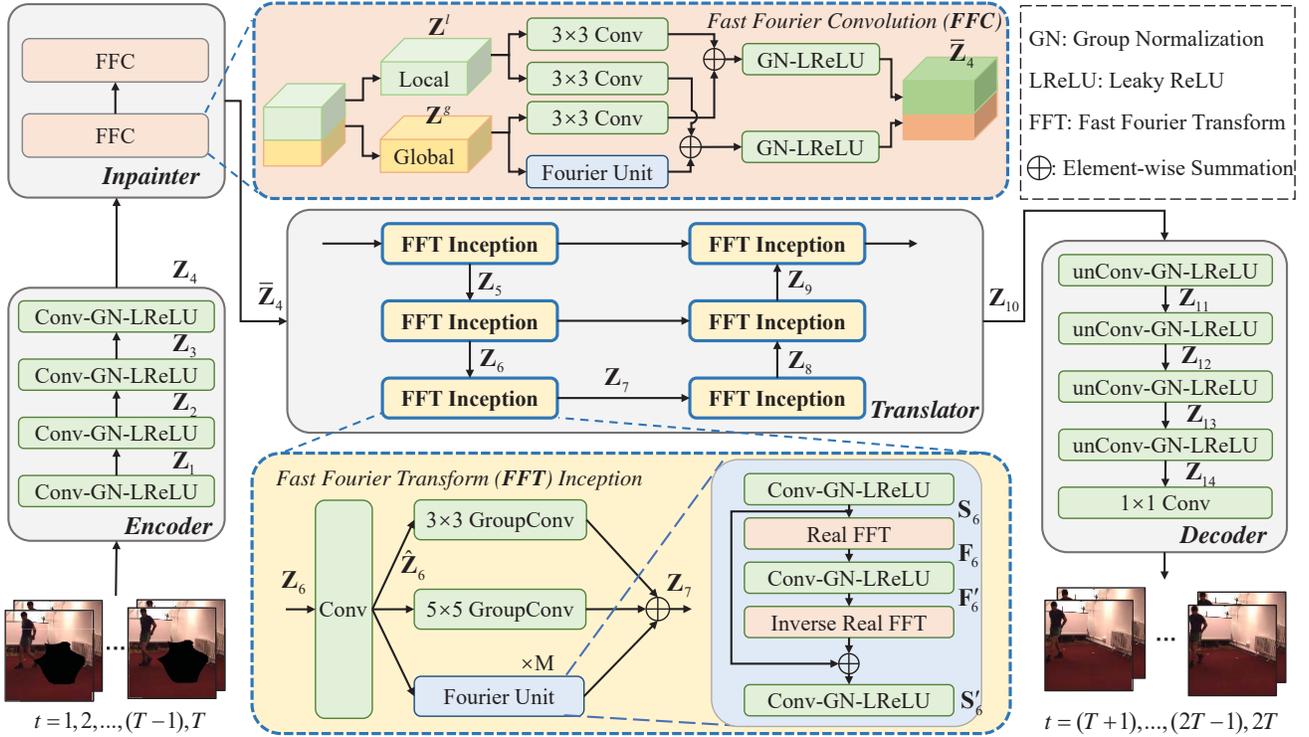


Fig. 2. Framework of the Fast Fourier Inception Networks (FFINet) for occluded video prediction. Here,  $(\tilde{N}, \hat{N})=(4,6)$ ,  $M$  is the number of Fourier Units.

#### D. Translator

The translator learns the temporal evolution by capturing and updating the spatiotemporal features. Inspired by [14], it is composed of  $\hat{N}$  stacked Fast Fourier Transform (FFT) Inception blocks  $\text{FFTI}(\cdot)$  (See the middle bottom in Fig. 2), and the hidden feature  $\mathbf{Z}_j \in \mathbb{R}^{B \times \hat{C} \times H' \times W'}$  is obtained by

$$\mathbf{Z}_j = \text{FFTI}(\mathbf{Z}_{j-1}), \tilde{N} < j \leq \tilde{N} + \hat{N}, \quad (4)$$

where  $\hat{C}$  is the channel number, which keeps still across succeeding FFT Inception blocks in the translator. FFT Inception block mainly involves the convolution layers with the kernel sizes of  $1 \times 1$ ,  $3 \times 3$ , and  $5 \times 5$ , as well as the Fourier Units. Here,  $1 \times 1$  Conv2D reduces the feature dimension from  $T \cdot \hat{C}$  (the first block) or  $\hat{C}$  (the remaining blocks) to  $\frac{\hat{C}}{2}$ , while the  $3 \times 3$  and  $5 \times 5$  Conv2D adopt the group convolutions, i.e., equally dividing the feature channels into 8 groups and each group captures distinct local pattern of the features. For the last block, the output channel is recovered to  $T' \cdot \hat{C}$ .

Note that our FFT Inception block abandons the large convolution kernels to save the computational costs, and introduces the Fourier Unit  $\text{FU}(\cdot)$  to capture the global context in early layers. The Fourier Unit (FU) has a receptive field covering the entire frame by including the channel-wise fast Fourier transform and  $1 \times 1$  Conv2D at lower cost. In particular, FU is composed of three convolution layers with group normalization and leaky ReLU function, real FFT, and inverse real FFT (IFFT( $\cdot$ )) layers. The convolution layers capture the local context and the real FFT models the global context. The first and the third convolution layers are used to align the feature dimension, while the second convolution layer is used

to update the spatiotemporal features in the frequency domain. Here, the real FFT (using half of the FFT spectrum) is applied to real-valued spatial features and the inverse real FFT makes the recovered spatial feature be real valued [18]. The details can be expressed by

$$\mathbf{S}_j = \sigma(\text{GN}(\text{Conv2D}(\hat{\mathbf{Z}}_j))) \in \mathbb{R}^{B \times \frac{\hat{C}}{2} \times H' \times W'}, \quad (5)$$

$$\mathbf{F}_j = \text{Real FFT}(\mathbf{S}_j) \in \mathbb{C}^{B \times \hat{C} \times H' \times \frac{W'}{2}}, \quad (6)$$

$$\mathbf{F}'_j = \sigma(\text{GN}(\text{Conv2D}(\mathbf{F}_j))) \in \mathbb{C}^{B \times \hat{C} \times H' \times \frac{W'}{2}}, \quad (7)$$

$$\mathbf{S}'_j = \sigma(\text{GN}(\text{Conv2D}(\text{IFFT}(\mathbf{F}'_j) + \mathbf{S}_j))) \in \mathbb{R}^{B \times \hat{C} \times H' \times W'}, \quad (8)$$

where  $\mathbb{C}$  denotes the frequency domain, and  $\hat{\mathbf{Z}}_j = \text{Conv2D}(\mathbf{Z}_j) \in \mathbb{R}^{B \times \frac{\hat{C}}{2} \times H' \times W'}$ . Here,  $\mathbf{S}_j$  denotes the spatiotemporal feature in source domain,  $\mathbf{F}_j$  denotes the spatiotemporal feature in frequency domain,  $\mathbf{S}'_j$  and  $\mathbf{F}'_j$  denote the updated features. In practice, we use multiple Fourier Units (e.g.,  $M=3$ ) in the model.

#### E. Decoder

The decoder is composed of  $\tilde{N}$  stacked unConv blocks, which are used to decode the updated spatiotemporal features into future frames. Following [14], we use  $\text{ConvTranspose2d}$  to serve as the  $\text{unConv}(\cdot)$  operator for upsampling the features along the spatial dimension. The hidden feature can be computed by

$$\mathbf{Z}_k = \sigma(\text{GN}(\text{unConv}(\mathbf{Z}_{k-1}))), \tilde{N} + \hat{N} < k \leq 2\tilde{N} + \hat{N}, \quad (9)$$

where  $k$  indexes the unConv block.

Given the updated spatiotemporal feature  $\mathbf{Z}_{\tilde{N}+\hat{N}} \in \mathbb{R}^{B \times (T' \cdot \tilde{C}) \times H' \times W'}$ , the decoder outputs the future  $T'$  frames  $\mathcal{Y}_{T'} \in \mathbb{R}^{B \times T' \times C \times H \times W}$ .

### F. Loss Function

There are two losses in our FFINet model, including the prediction loss  $\mathcal{L}_{pre}$  and the recovery loss  $\mathcal{L}_{rec}$ , both of which adopt the MSE function, i.e.,

$$\mathcal{L} = \mathcal{L}_{pre} + \lambda \mathcal{L}_{rec}, \quad (10)$$

where the constant  $\lambda > 0$  governs the contribution of the inpainter to the objective. It becomes the traditional video prediction when the inpainter is removed.

The prediction loss minimizes the error between the ground-truth frames  $\mathcal{Y}'_T = \{\mathbf{X}_t\}_{t=T+1}^{T+T'}$  and the predicted frames  $\mathcal{F}_\theta(\mathcal{X}_T)$ , i.e.,

$$\mathcal{L}_{pre} = \min_{\theta} \sum_{t=T+1}^{T+T'} \|\mathbf{X}_t - \mathcal{F}_\theta(\mathcal{X}_T)_t\|_2^2, \quad (11)$$

where  $\mathbf{X}_t$  is the  $t$ -th ground-truth frame,  $\mathcal{F}_\theta(\mathcal{X}_T)_t$  is the  $t$ -th predicted frame,  $\mathcal{X}_T$  denotes the training frames,  $\theta$  is the model parameter, and  $\|\cdot\|$  denotes the  $\ell_2$ -norm.

The recovery loss minimizes the error between the ground-truth frames  $\mathcal{Y}_T = \{\mathbf{X}_t\}_{t=1}^T$  and the recovery frames  $\mathcal{R}_\phi(\hat{\mathcal{X}}_T)$  by the inpainter, i.e.,

$$\mathcal{L}_{rec} = \min_{\phi} \sum_{t=1}^T \|\mathbf{X}_t - \mathcal{R}_\phi(\hat{\mathbf{X}}_t)\|_2^2, \quad (12)$$

where  $\phi$  is the parameter of the model without the translator.

## IV. EXPERIMENTS

This section shows extensive experimental results on several benchmark data sets. All experiments were conducted on a machine with three NVIDIA RTX 3090 Graphics Cards, and our model was compiled using PyTorch 1.12, Python 3.10, and CUDA 11.1.

### A. Data Sets

In total, there are five publicly available video databases used in the experiments. Details are shown below.

**Moving MNIST**<sup>1</sup> [21]. It consists of paired evolving hand-written digits from the MNIST<sup>2</sup> data set. Following [9], the training set includes 10000 sequences and the test set includes 5000 sequences. Each sequence consists of 20 successive  $64 \times 64$  frames with 2 randomly appearing digits. Among them, 10 frames are the input and the rest are the output. The initial position and rate of each digit are random, but the rate keeps the same across the entire sequence.

**TaxiBJ**<sup>3</sup> [22] is collected from the real-world traffic scenario in Beijing, ranging from 2013 to 2016. The traffic flows have strong temporal dependency among nearby area, and the data pre-processing follows [27]. The data of the last four

TABLE I  
EXPERIMENTAL SETTINGS.

Dataset	$(H, W, C)$	In→Out	$\tilde{C}$	$\hat{C}$	$\tilde{N}$	$\hat{N}$	$M$	Epoch
Moving MNIST [21]	(64,64,1)	10→10	64	512	4	6	3	2000
TaxiBJ [22]	(32,32,2)	4→4	64	256	3	4	2	80
Human3.6M [23]	(128,128,3)	4→4	64	64	1	10	2	100
KITTI&Caltech [24]	(128,160,3)	10→1	64	128	1	6	2	50
KTH [25]	(128,128,1)	10→20/40	32	128	3	8	1	100

weeks are used as the test set (1334 clips) while the rest are the training set (19627 clips). Each clip has 8 frames, where 4 frames are the input and the others are the output. The size of each video frame is  $32 \times 32 \times 2$ , and the two channels indicate the in and out traffic flow.

**Human3.6M**<sup>4</sup> [23] contains the sports videos of 11 subjects in 17 scenes, involving 3.6 million human pose images from 4 distinct camera views. Following [27], we use the data in the walking scene, which includes  $128 \times 128 \times 3$  RGB frames. The subsets  $\{S1, S5-S8\}$  are for training (2624 clips) and  $\{S9, S11\}$  are for test (1135 clips). Each clip has 8 frames, and the half of them are the input.

**KITTI&Caltech Pedestrian**<sup>5</sup>. Following [27], we use 2042 clips in KITTI [24] for training and 1983 clips in Caltech Pedestrian [47] for test. Both of them are driving databases taken from a vehicle in an urban environment, and the RGB frames are resized to  $128 \times 160$  by center-cropping and downsampling. The former includes “city”, “residential”, and “road” categories, while the latter has about 10 hours of  $640 \times 480$  video. Each clip has 20 consecutive frames, where 10 frames are the input and the others are the output.

**KTH**<sup>6</sup> [25] includes six action classes, i.e., walking, jogging, running, boxing, hand waving, and hand clapping, involving 25 subjects in four different scenes. Each video clip is taken in 25 fps and is 4 seconds on average. Following [48], the gray-scale frames are resized to  $128 \times 128$ . The training set has 5200 clips (16 subjects) and the test set has 3167 clips (9 subjects). Each clip has 30 frames, where 10 frames are the input and 20 frames are the output.

### B. Evaluation Metrics

Following [31] [16] [14], we employ MSE (Mean Square Error), MAE (Mean Absolute Error), SSIM (Structure Similarity Index Measure) [49], and PSNR (Peak Signal to Noise Ratio) to evaluate the quality of the predicted frames. On Caltech Pedestrian [47], we use MSE, SSIM, and PSNR; on KTH [25], we use SSIM and PSNR; on the remaining ones, we use MAE, MSE, and SSIM. SSIM ranges from -1 to 1, and the images are more similar when it approaches 1. The larger the PSNR db value, the better quality it achieves.

### C. Experimental Setup

**Training Phase.** The FFINet model is trained on the training set of each database, and use Adam [51] optimizer with

<sup>1</sup>[https://www.cs.toronto.edu/~nitish/unsupervised\\_video/](https://www.cs.toronto.edu/~nitish/unsupervised_video/)

<sup>2</sup><http://yann.lecun.com/exdb/mnist/>

<sup>3</sup><https://github.com/TolicWang/DeepST/tree/master/data/TaxiBJ>

<sup>4</sup><http://vision.imar.ro/human3.6m/description.php>

<sup>5</sup><https://www.cvlabs.net/datasets/kitti/>

<sup>6</sup><https://www.csc.kth.se/cvap/actions/>

TABLE II  
COMPARISON RESULTS ON MOVING MNIST [21], TAXIBJ [22], AND HUMAN3.6M [23].

Method	Venue	Moving MNIST [21]			TaxiBJ [22]			Human3.6M [23]		
		MSE↓	MAE↓	SSIM↑	MSE↓	MAE↓	SSIM↑	MSE↓	MAE↓	SSIM↑
PredRNN [9]	NeurIPS'17	56.8	126.1	0.867	46.4	17.1	0.971	48.4	18.9	0.781
PredRNN++ [26]	ICML'18	46.5	106.8	0.898	44.8	16.9	0.977	45.8	17.2	0.851
MIM [27]	CVPR'19	44.2	101.1	0.910	42.9	16.6	0.971	42.9	17.8	0.790
E3D-LSTM [10]	ICLR'19	41.3	86.4	0.910	43.2	16.9	0.979	46.4	16.6	0.869
PhyDNet [31]	CVPR'20	24.4	70.3	0.947	41.9	<b>16.2</b>	<b>0.982</b>	36.9	16.2	0.901
CrevNet [16]	ICLR'20	22.3	-	<u>0.949</u>	-	-	-	-	-	-
MAU [11]	NeurIPS'21	27.6	80.3	<u>0.937</u>	42.2*	16.4*	0.982*	<u>31.2*</u>	15.0*	0.885*
STAM [50]	TMM'23	28.6	-	0.935	44.1	-	-	-	<u>13.2</u>	0.875
SimVP [14]	CVPR'22	23.8	<u>69.9</u>	0.948	<u>41.4</u>	<b>16.2</b>	<b>0.982</b>	31.6	15.1	<u>0.904</u>
PredRNNv2 [30]	TPAMI'23	<u>19.9</u>	-	0.939	45.6*	16.8*	0.980*	36.3*	17.7*	0.863*
Ours		<b>19.2</b>	<b>60.4</b>	<b>0.958</b>	<b>41.2</b>	<b>16.2</b>	<b>0.982</b>	<b>23.3</b>	<b>11.9</b>	<b>0.913</b>

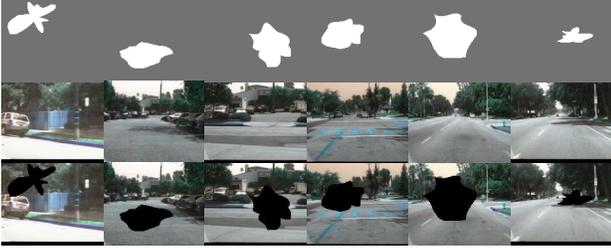


Fig. 3. Occluded examples of the Caltech Pedestrian dataset. Top row: masks, middle row: source frame, bottom row: occluded frame.

the OneCycle [52] learning rate scheduler and the momentum  $(\beta_1, \beta_2)=(0.9, 0.999)$ . The initial learning rate is 0.01, and the batch size  $B$  is 16. The constant  $\lambda$  is set to 0.5 for Moving MNIST [21] and 1.0 for the rest. The other experimental settings are listed in Table I, where  $\tilde{N}$  and  $\hat{N}$  denotes the convolution layer number and the FFT Inception block number of the encoder and the translator, respectively;  $M$  is the number of Fourier Units in the FFT Inception block. Note that in the ablation study, we use 10→20 for KTH dataset.

The feature height and width  $(H', W')=(\lfloor \frac{H}{2^{\lfloor \tilde{N}/2 \rfloor}} \rfloor, \lfloor \frac{W}{2^{\lfloor \hat{N}/2 \rfloor}} \rfloor)$ , where  $\lfloor \cdot \rfloor$  is the floor function. To generate the masks, we adopt the algorithm in [53], which randomly produces a set of control points around a unit circle and smoothly connects them to a closed cyclic contour by cubic Bezier curves. We show some occluded examples of Caltech Pedestrian in Fig. 3. The masks are randomly and instantly generated during training for enhancing the model robustness to different occlusions.

**Inference Phase.** The masks are generated in advance and keep still during frame prediction for fairness.

#### D. Quantitative Results

We show the quantitative comparison results in two situations, including video prediction with and without occlusion. On Moving MNIST [21], TaxiBJ [22], and Human3.6M [23], we compare PredRNN [9], PredRNN++ [26], PredRNNv2 [30], MIM (Memory In Memory) [27], E3D-LSTM [10], PhyDNet (Physical Dynamics Network) [31], MAU (Motion-Aware Unit) [11], CrevNet [16], and SimVP [14]; on Cal-

TABLE III  
COMPARISON RESULTS ON CALTECH PEDESTRIAN [47].

Method	Venue	Caltech Pedestrian [47](10→1)		
		MSE↓	SSIM↑	PSNR↑(db)
DVF [12]	ICCV'17	-	0.897	26.2
Dual-GAN [54]	ICCV'17	2.41	0.899	-
PredNet [55]	ICLR'17	2.42	0.905	27.6
CtrlGen [56]	CVPR'18	-	0.900	26.5
ContextVP [57]	ECCV'18	1.94	0.921	28.7
DPG [13]	ICCV'19	-	0.923	28.2
STMFANet [58]	CVPR'20	1.59	0.927	29.1
CrevNet [16]	ICLR'20	1.55	0.925	29.3
MAU [11]	NeurIPS'21	1.24	0.943	<u>30.1</u>
VPCL [59]	CVPR'22	-	0.928	-
SimVP [14]	CVPR'22	1.56	0.940	<b>33.1</b>
STAM [50]	TMM'23	<b>1.11</b>	<u>0.945</u>	29.9
Ours		<u>1.14</u>	<b>0.949</b>	<b>33.1</b>

TABLE IV  
COMPARISON RESULTS ON KTH [25].

Method	Venue	KTH [25](10→20)		KTH [25](10→40)	
		SSIM↑	PSNR↑(db)	SSIM↑	PSNR↑(db)
DFN [60]	NeurIPS'16	0.794	27.26	0.652	23.01
MCnet [48]	ICLR'17	0.804	25.95	0.730	23.89
PredRNN [9]	NeurIPS'17	0.839	27.55	0.703	24.16
fRNN [28]	ECCV'18	0.771	26.12	0.678	23.77
SV2P [5]	ICLR'18	0.838	27.79	0.789	26.12
PredRNN++ [26]	ICML'18	0.865	28.47	0.741	25.21
SAVP [61]	ICLR'19	0.852	27.77	0.811	26.18
E3D-LSTM [10]	ICLR'19	0.879	29.31	0.810	27.24
STMFANet [58]	CVPR'20	0.893	29.85	0.851	27.56
GridVP [62]	IROS'21	-	-	0.837	27.11
SimVP [14]	CVPR'22	<u>0.905</u>	<u>33.72</u>	<u>0.886</u>	<u>32.93</u>
PredRNNv2 [30]	TPAMI'23	0.838	28.37	-	-
Ours		<b>0.912</b>	<b>34.24</b>	<b>0.894</b>	<b>33.28</b>

tech Pedestrian [47], we compare PredNet [55], ContextVP [57], STMFANet (Spatial-Temporal Multi-Frequency Analysis Network) [58], CrevNet (Conditionally Reversible Network)

TABLE V  
COMPUTATION COMPARISON ON MOVING MNIST [21].

Method	Venue	FLOPs (G)↓	Train (s)↓	Test (fps)↑	#Params (M)↓	MSE ↓
PredRNN [9]	NeurIPS'17	115.6	300	107	23.8	56.8
PredRNN++ [26]	ICML'18	171.7	530	70	38.6	46.5
MIM [27]	CVPR'19	179.2	564	55	38.0	44.2
E3D-LSTM [10]	ICLR'19	298.9	1417	59	51.3	41.3
CrevNet [16]	ICLR'20	270.7	1030	10	5.0	22.3
PhyDNet [31]	CVPR'20	15.3	196	63	<b>3.1</b>	24.4
MAU [11]	NeurIPS'21	17.8	210	58	<u>4.5</u>	27.6
SimVP [14]	CVPR'22	19.4	86	190	22.3	23.8
Ours( $M=3$ )		<b>7.83</b>	<b>81</b>	<b>204</b>	19.1	<b>19.2</b>

[16], Dual-GAN (Dual Generative Adversarial Network) [54], DVF (Deep Voxel Flow) [12], CtrlGen (Controllable video Generation) [56], DPG (Disentangling Propagation and Generation) [13], VPCL (Video Prediction with Correspondence-wise Loss) [59], and SimVP [14]; on KTH [25], we compare PredRNN [9], PredRNN++ [26], fRNN [28], MCNet [48], E3D-LSTM [10], SV2P (Stochastic Variational Video Prediction) [5], STMFANet [58], GridVP [62], DFN (Dynamic Filter Network) [60], SAVP (Stochastic Adversarial Video Prediction) [61], and SimVP [14]. The best records are highlighted in bold and the second-best ones are underlined; “-” indicates the record is unavailable; “\*” indicates the record is obtained by re-implementing the code provided by the authors.

**Traditional video prediction.** The performance comparison results on Moving MNIST [21], TaxiBJ [22], and Human3.6M [23] are shown in Table II. Following [14], MSE values are enlarged by 100 times for TaxiBJ, MSE and MAE values are divided by 10 and 100 for Human3.6M. Note that the methods in the top group adopt the fixed training set, and the rest generate the training sample online, e.g., randomly select two digits and their motion path. From the table, it can be seen that our method consistently outperforms the other competitive alternatives across all evaluation metrics. Compared to the strongest baseline, our FFINet model reduces the MAE by 9.5 on Moving MNIST, while it reduces MSE by 7.9 on Human3.6M. Previous methods like PredRNN [9], PredRNN++ [26], MIM [27] use recurrent neural networks to model the temporal dynamics, failing to capture the long-term dependency and thus obtaining the poor predictions. E3D-LSTM [10] and CrevNet [16] adopt the 3D convolution to enlarge the receptive field but largely increase the computational costs. By contrast, our method employs the Fast Fourier Transform Inception blocks as the translator to better capture the spatiotemporal tendency by learning both the local and the global spatiotemporal features in video, resulting in higher-quality predictions.

Moreover, the performance comparison results on Caltech Pedestrian [47] and KTH [25] are shown in Table III and Table IV, respectively. From the tables, we observe that our approach has the most satisfying overall performance on predicting future frames in comparison with several SOTA methods. This demonstrates that the stacked Fourier transform

TABLE VI  
COMPARISON RESULTS ON MOVING MNIST [21] WITH OCCLUSION.

Method	Venue	Occ.	MSE↓	MAE↓	SSIM↑
CrevNet [16]	ICLR'20		30.2	86.3	0.935
		✓	35.0(+4.8)	94.2(+7.9)	0.915
PhyDNet [31]	CVPR'20		24.4	70.3	0.947
		✓	29.1(+4.7)	80.0(+9.7)	0.936
MAU [11]	NeurIPS'21		27.6	80.3	0.937
		✓	35.7(+8.1)	98.6(+18.3)	0.913
SimVP [14]	CVPR'22		23.8	68.9	0.948
		✓	28.7(+5.1)	80.8(+11.9)	0.936
PredRNNv2 [30]	TPAMI'23		27.4	82.2	0.937
		✓	32.8(+5.4)	92.3(+10.1)	0.925
Ours			19.2	60.4	0.958
Ours	w/o Inpainter	✓	22.7(+3.5)	68.4(+8.0)	0.950
Ours	w/ Inpainter	✓	<b>21.7(+2.5)</b>	<b>65.8(+5.4)</b>	<b>0.952</b>

inception blocks are able to learn the temporal evolution by adopting group convolutions and the channel-wise Fourier convolutions.

In addition, we show the efficiency comparison results on Moving MNIST in Table V, where the training time is computed for one epoch (per frame) using a single RTX3090 and the test time is the average fps of 10,000 samples. As seen from the table, our FFINet method enjoys the best prediction performance at the lowest training time with the fast inference speed. For example, our FLOPS is less than the half of the best candidate SimVP [14] but with much lower MSE.

**Occluded video prediction.** The results of occluded video prediction are shown in Table VI, Table VII, Table VIII, Table IX, and Table X, for Moving MNIST [21], TaxiBJ [22], Human3.6M [23], Caltech Pedestrian [47], and KTH [25], respectively. Here, “Occ.” denotes whether the video is occluded by the masks.

From these tables, we can see that the video prediction performance degenerates a lot when the video frames are occluded by random masks, which indicates the occlusions really do harm to predicting future frames. Moreover, when we use the Fast Fourier Convolution blocks to build the inpainter for recovering the missing areas in the video frames, the video prediction performance is improved, e.g., the MAE value reduces from 68.4 to 65.8 on Moving MNIST. This demonstrates that it is beneficial for the model to fill the missing areas by employing the inpainter to enlarge the receptive field. Note that the performance is slightly boosted with occlusion for some methods like MAU [11] and PredRNNv2 [30] on TaxiBJ [22], which might be the reason that there are many repeated patterns in the traffic flow and the masks are treated as noise to increase the model robustness.

### E. Ablation Study

We conduct the ablations on the FFT Inception block with different Fourier Units (without inpainter), and the recovery loss with different hyper-parameters  $\lambda$  (with inpainter) on the test data. Note that it requires two days to train for 2000 epochs until convergence on Moving MNIST, which is very

TABLE VII  
COMPARISON RESULTS ON TAXIBJ [22] WITH OCCLUSION.

Method	Venue	Occ.	MSE↓	MAE↓	SSIM↑
PhyDNet [31]	CVPR'20		41.9	16.2	0.982
		✓	43.0(+1.1)	16.6(+0.4)	0.981
MAU [11]	NeurIPS'21		42.2	16.4	0.982
		✓	41.7(-0.5)	16.3(-0.1)	0.982
SimVP [14]	CVPR'22		41.4	16.2	0.982
		✓	43.3(+1.9)	16.8(+0.6)	0.981
PredRNNv2 [30]	TPAMI'23		45.6	16.8	0.980
		✓	45.2(-0.4)	16.8(+0.0)	0.980
Ours			41.2	16.2	0.982
Ours	w/o Inpainter	✓	40.6(-0.6)	16.1(-0.1)	<b>0.983</b>
Ours	w/ Inpainter	✓	<b>40.4</b> (-0.8)	<b>16.0</b> (-0.2)	<b>0.983</b>

TABLE VIII  
COMPARISON RESULTS ON HUMAN3.6M [23] WITH OCCLUSION.

Method	Venue	Occ.	MSE↓	MAE↓	SSIM↑
PhyDNet [31]	CVPR'20		36.9	16.2	0.901
		✓	39.2(+2.3)	17.9(+1.7)	0.870
MAU [11]	NeurIPS'21		33.1	14.9	0.883
		✓	50.7(+17.6)	22.5(+7.6)	0.830
SimVP [14]	CVPR'22		31.6	15.1	0.904
		✓	33.4(+1.8)	16.4(+1.3)	0.897
PredRNNv2 [30]	TPAMI'23		34.8	17.2	0.864
		✓	36.8(+2.0)	18.7(+1.5)	0.842
Ours			23.3	11.9	0.912
Ours	w/o Inpainter	✓	24.7(+1.4)	13.0(+1.1)	0.906
Ours	w/ Inpainter	✓	<b>24.4</b> (+1.1)	<b>12.8</b> (+0.9)	<b>0.907</b>

TABLE IX  
COMPARISON RESULTS ON CALTECH PEDESTRIAN [47] WITH OCCLUSION.

Method	Venue	Occ.	MSE↓	SSIM↑	PSNR↑(db)
CrevNet [16]	ICLR'20		1.55	0.925	29.3
		✓	2.86(+1.31)	0.878	24.7(-4.6)
STMFANet [58]	CVPR'20		1.59	0.927	29.1
		✓	3.18(+1.59)	0.874	25.6(-5.5)
MAU [11]	NeurIPS'21		1.24	0.943	30.1
		✓	2.55(+1.31)	0.898	24.6(-5.5)
SimVP [14]	CVPR'22		1.59	0.927	33.1
		✓	3.23(+1.67)	0.892	30.1(-3.0)
Ours			1.14	0.949	33.1
Ours	w/o Inpainter	✓	2.19(+1.05)	0.917	31.4(-1.7)
Ours	w/ Inpainter	✓	<b>2.08</b> (+0.94)	<b>0.921</b>	<b>32.2</b> (-0.9)

long, so we run 100 epochs in the ablations to save time. The parameters keep the same as in training unless specified.

**FFT Inception block.** We explore the performance of our model using different architectures in the FFT Inception block on Moving MNIST [21]. In particular, we vary the group convolution kernel size from 3 to 11, and insert the Fourier Unit (FU) after the second, the third, and the fourth group convolution branches; the results are shown in Fig. 4(a). From the left figure, we observe that when using more group convolution branches, the performance is unnecessarily getting better. On the contrary, it gets a good trade-off between the performance and the model size, when using two group

TABLE X  
COMPARISON RESULTS ON KTH [25] WITH OCCLUSION.

Method	Venue	Occ.	SSIM↑	PSNR↑(db)
STMFANet [58]	CVPR'20		0.893	29.85
		✓	0.871(-0.022)	26.27(-3.58)
LMC [35]	CVPR'21		0.894	28.61
		✓	0.879(-0.015)	26.28(-2.33)
SimVP [14]	CVPR'22		0.905	33.72
		✓	0.895(-0.010)	33.00(-0.72)
PredRNNv2 [30]	TPAMI'23		0.858	29.79
		✓	0.827(-0.031)	25.39(-4.40)
Ours			0.912	34.24
Ours	w/o Inpainter	✓	0.904(-0.008)	33.57(-0.67)
Ours	w/ Inpainter	✓	<b>0.905</b> (-0.007)	<b>33.69</b> (-0.55)

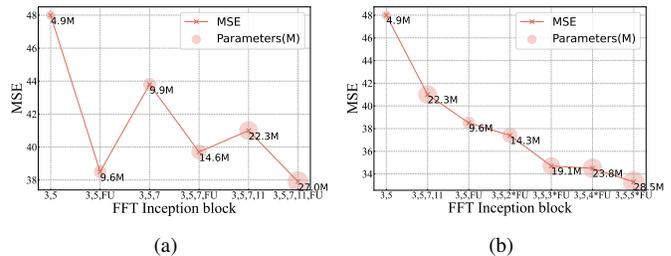


Fig. 4. Ablations of the FFT Inception block on Moving MNIST [21]. (a) Different Fourier Unit (FU) positions; (b) Varying numbers of FUs. We use the comma to separate different branches.

TABLE XI  
MSE OF OUR MODEL WITH DIFFERENT  $M$  FOURIER UNITS.

Dataset	$M = 1$	$M = 2$	$M = 3$
MovingMNIST [21]	38.5	37.4	<b>34.7</b>
TaxiBJ [22]	43.3	<b>41.2</b>	43.4
Human3.6M [23]	23.8	<b>23.3</b>	24.1
CaltechPed. [24]	1.16	<b>1.14</b>	1.21
KTH [25]	<b>25.4</b>	26.3	28.8

convolutions followed by the Fourier Unit. Moreover, we vary the number of FUs from 1 to 5, and depict the results in Fig. 4(b). From the right figure, it can be seen that the prediction performance becomes better when using more FUs, but the model size becomes much larger. Overall, it seems that the performance gets a good balance with two group convolutions and three FUs.

In addition, we vary the number of FUs  $M$  from 1 to 3 on the other datasets and show the results in Table XI. From the table, we see that our FFNet model achieves the best when using two FUs on TaxiBJ [22], Human3.6M [23], and Caltech Pedestrian [24], and one FU on KTH [25]. It suggests a modest number of FUs is enough to achieve promising prediction performance.

**Recovery loss.** The recovery loss estimates the error between the source frame and the recovery frame, and its contribution to the model is governed by the hyper-parameter  $\lambda$ . We vary its value from 0 to 2, and show the results in Table XII. From the table, we observe that our approach improves the video prediction performance with the recovery

TABLE XII  
ABLATIONS ON THE RECOVERY LOSS WITH DIFFERENT  $\lambda$ .

$\lambda$	Moving MNIST [21]			TaxiBJ [22]			Human3.6M [23]			Caltch Pedestrian [47]			KTH [25]		
	MSE↓	MAE↓	SSIM↑	MSE↓	MAE↓	SSIM↑	MSE↓	MAE↓	SSIM↑	MSE↓	MAE↓	SSIM↑	MSE↓	MAE↓	SSIM↑
0.00	42.8	114.6	0.898	40.6	16.1	0.983	25.2	13.5	0.905	2.4	19.1	0.919	42.6	418.5	0.899
0.25	41.4	110.3	0.901	40.6	16.1	0.983	24.9	13.3	0.905	2.3	19.1	0.919	42.4	416.4	0.901
0.50	<b>40.3</b>	<b>105.5</b>	<b>0.908</b>	40.5	16.1	0.983	25.0	13.0	0.905	2.2	18.9	0.920	41.8	415.9	0.903
1.00	42.2	111.0	0.900	<b>40.4</b>	<b>16.0</b>	<b>0.983</b>	<b>24.4</b>	<b>12.8</b>	<b>0.907</b>	<b>2.1</b>	<b>18.8</b>	<b>0.921</b>	<b>40.5</b>	<b>413.7</b>	<b>0.905</b>
2.00	42.6	112.7	0.896	40.6	16.1	0.983	24.6	13.0	0.905	2.4	19.2	0.919	41.9	416.2	<b>0.905</b>

loss across all datasets. This demonstrates that the inpainting quality of the frames directly influences the future frame prediction. In particular, it achieves the best performance when  $\lambda$  is set to 0.5 for Moving MNIST [21] and 1.0 for the rest.

### F. Qualitative Results

To give an intuitive view on the superiority of our model, we select some challenging cases with occlusion from the five datasets and visualize the predicted frames in Fig. 5 to Fig. 9. Besides, we show the predicted frames without occlusion on Moving MNIST in Fig. 5, where the two digits are seriously overlapped.

In Fig. 5, the two digits can be well generated by our method even when the input digits are heavily occluded. This is because the designed inpainter is able to fill in the missing area of the frame before the translator to capture the temporal dynamics in video. On the contrary, previous methods like PhyDNet [31] and SimVP [14] learns the spatiotemporal features of frames with occlusion, which might bring about some misleading information, leading to inferior predictions.

In Fig. 6, the difference map  $|T - P|$  of the target frame and the predicted frame reflects the quality of the traffic flow prediction. When the dark area becomes larger, it means the prediction is more accurate. From the figure, we see that the dark area produced by our model is much larger that the compared PredRNNv2 [30] that does not consider the occlusion scenario.

Fig. 7 shows a man walking through a room and Fig. 8 shows the street scenario. Here, the masks are randomly placed in different positions. From the figures, we observe that our method generates higher-quality frames compared to others, e.g., the human body and the car are more clearly to see. Fig. 9 show the walking action clip with the occlusion in the middle, and we see that the predicted sequence by ours matches better with the target ones. This is because the developed inpainter is able to recover the occluded area by adopting the fast Fourier convolution, and the designed translator models the temporal dynamics by capturing both the local and the global spatiotemporal features.

## V. CONCLUSION

In this work, we have explored the occluded video prediction problem, which appears frequently in practice but remains untouched yet. To address the occlusion issue, we design the inpainter module which employs the channel-wise

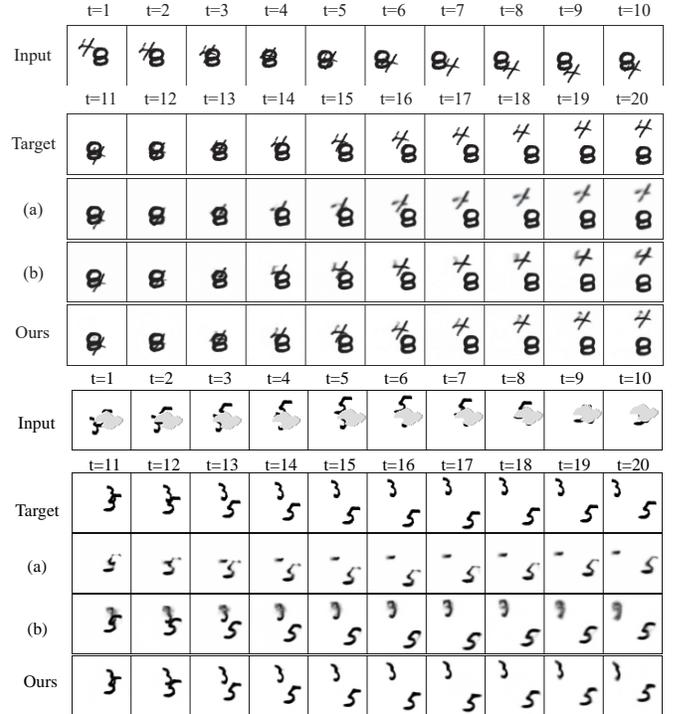


Fig. 5. Predictions on Moving MNIST. (a) PhyDNet [31]; (b) SimVP [14].

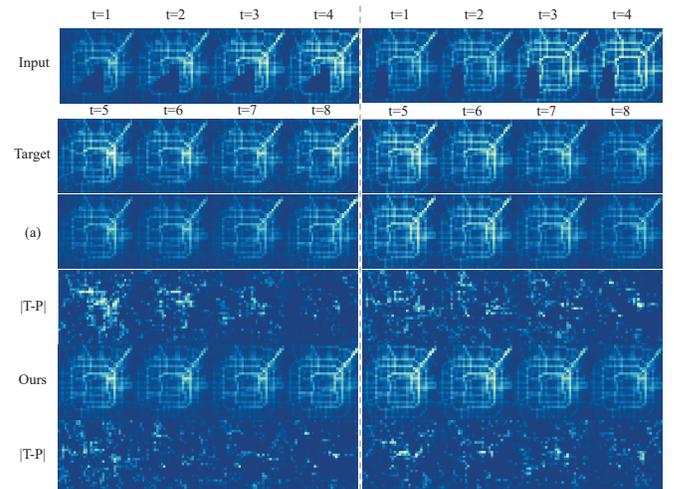


Fig. 6. Predictions on TaxiBJ [22]. (a) PredRNNv2 [30].

fast Fourier convolution to enlarge the receptive field, thus

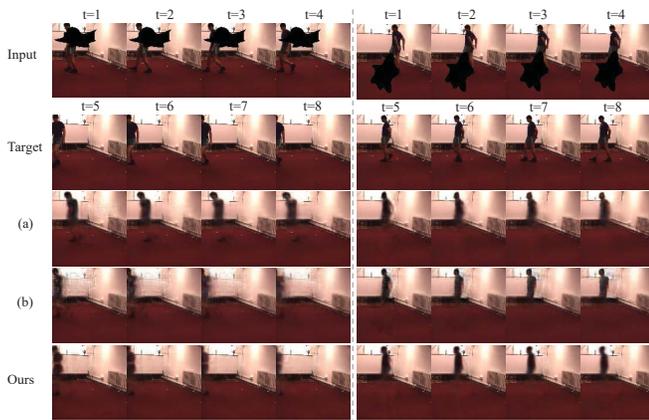


Fig. 7. Predictions on Human3.6M [23]. (a) MAU [11]; (b) SimVP [14].

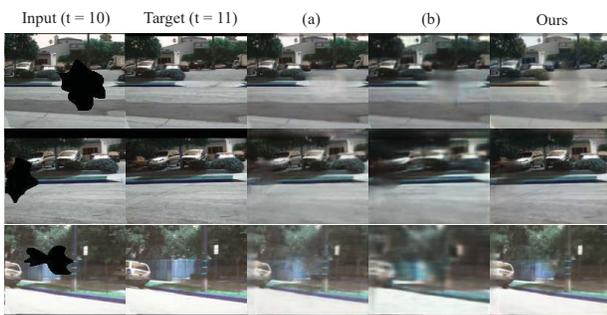


Fig. 8. Predictions on Caltech Ped. [47]. (a) CrevNet [16]; (b) SimVP [14].

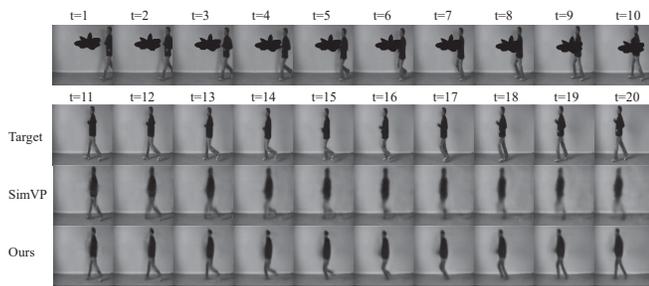


Fig. 9. Predictions on KTH [25].

capturing the global context to recover the missing area in the frame. To model the temporal dynamics, we develop the Fast Fourier Transform Inception block that includes group convolutions and multiple Fourier Units to learn both the local and the global spatiotemporal features, which help to capture the temporal evolution across the video frames. Hence, we proposed the fully-convolutional Fast Fourier Inception Networks, terms FFINet, for occluded video prediction, and conducted comprehensive experiments to verify the effectiveness of the proposed method on several benchmarks.

## REFERENCES

- [1] X. Shi, Z. Chen, H. Wang, D. Yeung, W. Wong, and W. Woo, “Convolutional LSTM network: A machine learning approach for precipitation nowcasting,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2015, pp. 802–810.
- [2] H. Wu, Z. Yao, J. Wang, and M. Long, “Motionrnn: A flexible model for video prediction with spacetime-varying motions,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 15 435–15 444.
- [3] X. Bei, Y. Yang, and S. Soatto, “Learning semantic-aware dynamics for video prediction,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 902–912.
- [4] L. Castrejón, N. Ballas, and A. C. Courville, “Improved conditional vrnnns for video prediction,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 7607–7616.
- [5] M. Babaeizadeh, C. Finn, D. Erhan, R. H. Campbell, and S. Levine, “Stochastic variational video prediction,” in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- [6] E. Denton and R. Fergus, “Stochastic video generation with a learned prior,” in *Proceedings of the International Conference on Machine Learning (ICML)*, 2018, pp. 1182–1191.
- [7] J. L. Elman, “Finding structure in time,” *Cognitive Science*, vol. 14, no. 2, pp. 179–211, 1990.
- [8] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proc. IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [9] Y. Wang, M. Long, J. Wang, Z. Gao, and P. S. Yu, “Predrnn: Recurrent neural networks for predictive learning using spatiotemporal lstms,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 879–888.
- [10] Y. Wang, L. Jiang, M. Yang, L. Li, M. Long, and L. Fei-Fei, “Eidetic 3d LSTM: A model for video prediction and beyond,” in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- [11] Z. Chang, X. Zhang, S. Wang, S. Ma, Y. Ye, X. Xiang, and W. Gao, “Mau: A motion-aware unit for video prediction and beyond,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, 2021.
- [12] Z. Liu, R. A. Yeh, X. Tang, Y. Liu, and A. Agarwala, “Video frame synthesis using deep voxel flow,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 4473–4481.
- [13] H. Gao, H. Xu, Q. Cai, R. Wang, F. Yu, and T. Darrell, “Disentangling propagation and generation for video prediction,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 9005–9014.
- [14] Z. Gao, C. Tan, L. Wu, and S. Z. Li, “Simvp: Simpler yet better video prediction,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 3160–3170.
- [15] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, J. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, T. Uszkoreit, and N. Houlsby, “An image is worth 16x16 words: transformers for image recognition at scale,” in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- [16] W. Yu, Y. Lu, S. Easterbrook, and S. Fidler, “Efficient and information-preserving future frame prediction and beyond,” in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- [17] Z. Wang, Z. Liu, G. Li, Y. Wang, T. Zhang, L. Xu, and J. Wang, “Spatio-temporal self-attention network for video saliency prediction,” *IEEE Transactions on Multimedia (TMM)*, vol. 25, pp. 1161–1174, 2023.
- [18] L. Chi, B. Jiang, and Y. Mu, “Fast fourier convolution,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [19] Y. Katznelson, “An introduction to harmonic analysis,” *The American Mathematical Monthly*, vol. 77, no. 4, 2005.
- [20] R. Suvorov, E. Logacheva, A. Mashikhin, A. Remizova, A. Ashukha, A. Silvestrov, N. Kong, H. Goka, K. Park, and V. Lempitsky, “Resolution-robust large mask inpainting with fourier convolutions,” in *Proceedings of the IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2022, pp. 3172–3182.
- [21] N. Srivastava, E. Mansimov, and R. Salakhutdinov, “Unsupervised learning of video representations using lstms,” in *Proceedings of the International Conference on Machine Learning (ICML)*, 2015, pp. 843–852.
- [22] J. Zhang, Y. Zheng, and D. Qi, “Deep spatio-temporal residual networks for citywide crowd flows prediction,” in *Proceedings of the AAAI conference on artificial intelligence(AAAI)*, 2017, pp. 1655–1661.
- [23] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu, “Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments,” *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 36, no. 7, pp. 1325–1339, 2014.
- [24] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, “Vision meets robotics: The KITTI dataset,” *International Journal of Robotics Research (IJRR)*, vol. 32, no. 11, pp. 1231–1237, 2013.

- [25] C. Schüldt, I. Laptev, and B. Caputo, "Recognizing human actions: A local SVM approach," in *Proceedings of the Computational Vision and Active Perception Laboratory (CVAP)*, 2004, pp. 32–36.
- [26] Y. Wang, Z. Gao, M. Long, J. Wang, and P. S. Yu, "Predrnn++: Towards a resolution of the deep-in-time dilemma in spatiotemporal predictive learning," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2018, pp. 5110–5119.
- [27] Y. Wang, J. Zhang, H. Zhu, M. Long, J. Wang, and P. S. Yu, "Memory in memory: A predictive neural network for learning higher-order non-stationarity from spatiotemporal dynamics," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 9154–9162.
- [28] M. Oliu, J. Selva, and S. Escalera, "Folded recurrent neural networks for future video prediction," in *Proceedings of the European Conference on Computer Vision (ECCV)*, vol. 11218, 2018, pp. 745–761.
- [29] J. Su, W. Byeon, J. Kossaiji, F. Huang, J. Kautz, and A. Anandkumar, "Convolutional tensor-train LSTM for spatio-temporal learning," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [30] Y. Wang, H. Wu, J. Zhang, Z. Gao, J. Wang, P. S. Yu, and M. Long, "Predrnn: A recurrent neural network for spatiotemporal predictive learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 45, no. 2, pp. 2208–2225, 2023.
- [31] V. L. Guen and N. Thome, "Disentangling physical dynamics from unknown factors for unsupervised video prediction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 11 474–11 484.
- [32] S. Park, K. Kim, J. Lee, J. Choo, J. Lee, S. Kim, and E. Choi, "Vid-ode: Continuous-time video generation with neural ordinary differential equation," in *Proceedings of the AAAI conference on artificial intelligence (AAAI)*, 2021, pp. 2412–2422.
- [33] N. Kim and J. Kang, "Dynamic motion estimation and evolution video prediction network," *IEEE Transactions on Multimedia (TMM)*, vol. 23, pp. 3986–3998, 2021.
- [34] Z. Ye, M. Xia, R. Yi, J. Zhang, Y.-K. Lai, X. Huang, G. Zhang, and Y.-J. Liu, "Audio-driven talking face video generation with dynamic convolution kernels," *IEEE Transactions on Multimedia (TMM)*, vol. 25, pp. 2033–2046, 2023.
- [35] S. Lee, H. G. Kim, D. H. Choi, H. Kim, and Y. M. Ro, "Video prediction recalling long-term motion context via memory alignment learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 3054–3063.
- [36] Z. Chang, X. Zhang, S. Wang, S. Ma, and W. Gao, "STRPM: A spatiotemporal residual predictive model for high-resolution video prediction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 13 926–13 935.
- [37] G. Chen, W. Zhang, H. Lu, S. Gao, Y. Wang, M. Long, and X. Yang, "Continual predictive learning from videos," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 10 718–10 727.
- [38] W. Yu, W. Chen, S. Yin, S. Easterbrook, and A. Garg, "Modular action concept grounding in semantic video prediction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 3595–3604.
- [39] D. Weissenborn, O. Täckström, and J. Uszkoreit, "Scaling autoregressive video models," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- [40] R. Rakhimov, D. Volkhonskiy, A. Artemov, D. Zorin, and E. Burnaev, "Latent video transformer," in *Proceedings of the 16th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications*, 2021, pp. 101–112.
- [41] Z. Yang, X. Yang, and Q. Lin, "Tctn: A 3d-temporal convolutional transformer network for spatiotemporal predictive learning," *arXiv preprint arXiv:2112.01085*, 2021.
- [42] Y. Wu, Q. Wen, and Q. Chen, "Optimizing video prediction via video frame interpolation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 17 793–17 802.
- [43] J. Huang, Y. Jin, K. M. Yi, and L. Sigal, "Layered controllable video generation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, vol. 13676, 2022, pp. 546–564.
- [44] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2014, pp. 1724–1734.
- [45] Y. Kwon and M. Park, "Predicting future frames using retrospective cycle GAN," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 1811–1820.
- [46] J. Xu, B. Ni, and X. Yang, "Progressive multi-granularity analysis for video prediction," *International Journal of Computer Vision (IJCV)*, vol. 129, no. 3, pp. 601–618, 2021.
- [47] P. Dollár, C. Wojek, B. Schiele, and P. Perona, "Pedestrian detection: A benchmark," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009, pp. 304–311.
- [48] R. Villegas, J. Yang, S. Hong, X. Lin, and H. Lee, "Decomposing motion and content for natural video sequence prediction," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- [49] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Transactions on Image Processing (TIP)*, vol. 13, no. 4, pp. 600–612, 2004.
- [50] Z. Chang, X. Zhang, S. Wang, S. Ma, and W. Gao, "Stam: A spatiotemporal attention based memory for video prediction," *IEEE Transactions on Multimedia (TMM)*, vol. 25, pp. 2354–2367, 2023.
- [51] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015.
- [52] L. N. Smith and N. Topin, "Super-convergence: Very fast training of neural networks using large learning rates," in *Proceedings of the International Society for Optics and Photonics*, 2019.
- [53] Z. Li, C. Lu, J. Qin, C. Guo, and M. Cheng, "Towards an end-to-end framework for flow-guided video inpainting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 17 541–17 550.
- [54] X. Liang, L. Lee, W. Dai, and E. P. Xing, "Dual motion GAN for future-flow embedded video prediction," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 1762–1770.
- [55] W. Lotter, G. Kreiman, and D. D. Cox, "Deep predictive coding networks for video prediction and unsupervised learning," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- [56] Z. Hao, X. Huang, and S. J. Belongie, "Controllable video generation with sparse trajectories," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 7854–7863.
- [57] W. Byeon, Q. Wang, R. K. Srivastava, and P. Koumoutsakos, "Contextvp: Fully context-aware video prediction," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 781–797.
- [58] B. Jin, Y. Hu, Q. Tang, J. Niu, Z. Shi, Y. Han, and X. Li, "Exploring spatial-temporal multi-frequency analysis for high-fidelity and temporal-consistency video prediction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 4553–4562.
- [59] D. Geng, M. Hamilton, and A. Owens, "Comparing correspondences: Video prediction with correspondence-wise losses," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 3355–3366.
- [60] X. Jia, B. D. Brabandere, T. Tuytelaars, and L. V. Gool, "Dynamic filter networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2016, pp. 667–675.
- [61] A. X. Lee, R. Zhang, F. Ebert, P. Abbeel, C. Finn, and S. Levine, "Stochastic adversarial video prediction," in *Proceedings of the 7th International Conference on Learning Representations (ICLR)*, 2019.
- [62] X. Gao, Y. Jin, Q. Dou, C. Fu, and P. Heng, "Accurate grid keypoint learning for efficient video prediction," in *Proceedings of the International Conference on Intelligent Robots and Systems (IROS)*, 2021, pp. 5908–5915.