

Video Prediction by Efficient Transformers

Xi Ye^{a,*}, Guillaume-Alexandre Bilodeau^a

^aLITIV Lab, Polytechnique Montréal, P.O. Box 6079, Station centre-ville, Montreal, H3C3A7, QC, Canada

ARTICLE INFO

Keywords:

Video prediction
Transformers
Video representation learning
Autoregressive generative models
Non-autoregressive generative models

ABSTRACT

Video prediction is a challenging computer vision task that has a wide range of applications. In this work, we present a new family of Transformer-based models for video prediction. Firstly, an efficient local spatial-temporal separation attention mechanism is proposed to reduce the complexity of standard Transformers. Then, a full autoregressive model, a partial autoregressive model and a non-autoregressive model are developed based on the new efficient Transformer. The partial autoregressive model has a similar performance with the full autoregressive model but a faster inference speed. The non-autoregressive model not only achieves a faster inference speed but also mitigates the quality degradation problem of the autoregressive counterparts, but it requires additional parameters and loss function for learning. Given the same attention mechanism, we conducted a comprehensive study to compare the proposed three video prediction variants. Experiments show that the proposed video prediction models are competitive with more complex state-of-the-art convolutional-LSTM based models. The source code is available at <https://github.com/XiYe20/VPTR>.

1. Introduction

Video future frames prediction (VFFP) can be applied to many application areas, for example, autonomous vehicles to predict future traffic outcomes [1], anomaly detection [2] and model-based reinforcement learning [3] to predict users motion. Besides, VFFP is naturally a good self-supervised learning task, and it has drawn much attention in recent years [4, 5]. Indeed, we can train VFFP models to learn good spatio-temporal representations from a large amount of unlabeled videos, and then apply it for many downstream tasks.

In this paper, we focus on the most common video prediction task, i.e. predicting N future frames given L past frames, with L and N being greater than 1. From the perspective of learning a deep neural networks for VFFP, we can formalize the task to be

$$\arg \max_{\theta} p(\hat{x}_{L+N}, \dots, \hat{x}_{L+1} | x_L, \dots, x_1; \theta), \quad (1)$$

where x_t and \hat{x}_t denote the input past frames and predicted future frames respectively, θ denotes the parameters of the neural networks.

Benefiting from the advances in deep learning, the performance of VFFP models is constantly improving. Particularly, the efficient and powerful Convolutional Long Short-Term Memory networks (ConvLSTMs) are the core for almost all the state-of-the-art (SOTA) VFFP models. Nevertheless, ConvLSTMs suffer from some inherent problems typical of recurrent neural networks (RNNs), including slow training and inference speed, error accumulation during inference, vanishing gradient, and predicted frame quality degradation. Researchers keep improving the performance

by developing more and more sophisticated ConvLSTM-based models. For instance, by integrating custom motion-aware units into ConvLSTM [6], or building complex memory modules to store the motion context [7]. Still, VFFP is a challenging task which is far from being solved.

Meanwhile, the Transformer [8] overcomes most of the aforementioned drawbacks of RNNs and made breakthroughs for natural language processing (NLP). Inspired by this, more and more researchers are starting to adapt the Transformer for various computer vision tasks [9, 10, 11], including few recent works for VFFP [12, 13, 14]. However, the Transformer was originally designed for sequence data processing. Given a sequence of elements with length L and feature dimensionality d_{model} , a dot product attention-based Transformer has a complexity of $\mathcal{O}(L^2 d_{model})$. So it is computationally expensive to apply a Transformer to high dimensional visual features. For example, considering a spatiotemporal feature $Z \in \mathbb{R}^{T \times H \times W \times d_{model}}$, where T, H, W denote the time duration, spatial height and spatial width respectively, we need to flatten it to be a sequence $Z \in \mathbb{R}^{(T \cdot H \cdot W) \times d_{model}}$ with length $T \cdot H \cdot W$ along the spatial and temporal dimensions for a standard Transformer. The complexity is therefore $\mathcal{O}((THW)^2 d_{model})$. We still need further research about more efficient and compact visual Transformer, especially for videos. *Therefore, we propose a novel efficient Transformer block with smaller complexity, named VidHRFormer, and we developed a new video prediction Transformer (VPTR) based on it.*

Among the Transformer-based VFFP models [12, 13, 14] that we mentioned earlier, some of them are autoregressive models while some others are non-autoregressive models, and they are based on different attention mechanisms, e.g. a custom convolution multi-head attention (MHA) [12] and standard dot-product MHA [13, 14]. There is no formal comparison of the two typical approaches (autoregressive vs non-autoregressive) in the case of Transformer-based VFFP models so far. *Thus, we developed a non-autoregressive*

*Corresponding author

 xi.ye@polymtl.ca (X. Ye); gabilodeau@polymtl.ca (G. Bilodeau)

 (X. Ye); (G. Bilodeau)

ORCID(s): 0000-0002-1763-6019 (X. Ye); 0000-0003-3227-5060 (G. Bilodeau)

VPTR model (VPTR-NAR), two autoregressive VPTR variants, i.e., the fully autoregressive VPTR model (VPTR-PAR) and the partial autoregressive VPTR model (VPTR-FAR). All VPTR variants share the same attention mechanism and same number of Transformer block layers, which guarantees a fair comparison between the two approaches.

Our main contributions are summarized as:

1) We proposed a new efficient Transformer block, VidHRFormer, for spatio-temporal feature learning by combining spatial local attention and temporal attention in two steps. The new Transformer block successfully reduces the complexity of a standard Transformer block with respect to the same input spatio-temporal feature size, specifically, from $\mathcal{O}((THW)^2 d_{model})$ to $\mathcal{O}(\frac{H^2 W^2}{P^2} + T^2) d_{model}$.

2) Based on this Transformer block, three VPTR models, VPTR-NAR, VPTR-FAR and VPTR-PAR, were developed. We show that the proposed simple attention-based VPTRs are capable of reaching and outperforming more complex SOTA ConvLSTM-based VFFP models.

3) A formal comparison of three VPTR variants was conducted. The results show that VPTR-NAR has a faster inference speed and smaller accumulation of errors during inference, but it is more difficult to train. We solved the training problem of VPTR-NAR by employing a contrastive feature loss that maximizes the mutual information of predicted and ground-truth future frame features.

4) We found that given the same number of Transformer block layers, VPTR-FAR and VPTR-PAR have a worse generalization performance due to the accumulated inference errors, which are introduced by the discrepancy between train and test behaviors. Finally, we compared two different inference methods for VPTR-FAR and VPTR-PAR, the results show that recurrent inference over pixel space introduces less accumulation of errors than recurrent inference over latent space in the case of VPTR-FAR.

This paper is an extension of the work in [15]. We present a new VPTR-PAR variant and we fully analyze auto-regressive vs non-auto-regressive models for our new VidHRFormer.

2. Related work

2.1. Video future frames prediction

Almost all the deep learning-based VFFP models take an encoder to extract the representations of past frames, and then a decoder generates future frame pixels based on those representations. Here, we propose our own taxonomy for VFFP models, which is derived from the different encoding and decoding mechanisms.

Non-decomposing and Decomposing models. Most VFFP models are equipped with an encoder which generates frame-level features. they are classified as non-decomposing models since they assume that the encoders implicitly extract all the necessary information for the prediction of future frames [16, 17, 18, 19, 20]. However, some other VFFP models explicitly decompose the frame visual features, e.g. content and motion, which aims to reduce the difficulty of

prediction by introducing more prior knowledge [21, 22, 23, 24, 25]. There are other decomposing methods, including object keypoint skeleton and appearance decomposition [26, 27], object-centric decomposition [28, 29, 30].

Sequentially-generated and Parallely-generated models. The majority of VFFP models rely on the flexible RNNs for temporal information modeling, which predict future frames recursively. We call these models sequentially-generated models [18, 19, 31, 20, 32]. In contrast, some models take advantage of standard CNNs or 3D-CNNs as the backbones, without the use of RNNs, and generate multiple future frames simultaneously (that is in parallel) [33, 34, 35, 36].

Pixel-direct generation and Transformation-based models. Some VFFP models use decoders to directly synthesize the pixels of future frames based on past frame representations [37, 22, 32, 26, 21], while some other models have no such decoder and generate the future frames by applying spatial transformation operation on the past frames, like warping [30, 38, 21].

Deterministic and Stochastic models. The future prediction is by nature multimodal [39], i.e. stochastic, even though we are given multiple past frames as context. However, it is difficult to model the stochasticity of the future and thus most VFFP models ignore it, even though it affects the predicted image quality [16, 36, 37]. Stochastic models normally make uncertainty estimation based on VAE or GAN, such as stochastic variational video prediction (SV2P) model [19], stochastic video generation with a learned prior (SVG-LP) [20], and improved conditional variational RNNs (VRNNs) [31].

Note that the categories described above are not mutually exclusive as a model can be a combination of arbitrary different encoders and decoders. For example, our proposed VPTR belongs to the deterministic, non-decomposing, pixel-direct generated model categories. The VPTR-NAR variant is considered to be a parallely-generated model, while VPTR-PAR and VPTR-FAR variants are sequentially-generated models.

Recently, many work achieve SOTA performance by integrating attention mechanism or memory augmented modules in the ConvLSTM models for a better long-term motion information learning [40, 7, 6]. For example, LMC-Memory model [7] stores the long-term motion context by a novel memory alignment learning, and the motion information is recalled during test to facilitate the long-term prediction. Specifically, LMC-Memory takes the difference image between adjacent frames as motion information, same as MCnet [21], to learn the long-term motion context embeddings as an external memory during training. During test, the recalled motion context memory feature is concatenated with spatial feature of current time step to predict the next frame. VPEG [40] shares a similar idea with LMC-Memory. VPEG also learns a pool of motion features, which are called examples, of the whole dataset, but VPEG aims to approximate the stochastic predictions by those examples and thus

bypassing the variational inference of other aforementioned stochastic video prediction models.

Zhang et al. [6] proposed an attention-based motion-aware unit (MAU) to increase the temporal receptive field of RNNs. The proposed MAU is composed of an attention module and a fusion module, where the attention module considers the correlation between the current step features and different history states within a temporal receptive field as an attention score, and the fusion module is responsible to aggregate those features. In this way, the MAU is able to capture better motion information and receive a broader temporal receptive field than the vanilla ConvLSTMs. However, those attention or memory augmented ConvLSTM architectures tend to be complex and it is hard to understand and follow the spatio-temporal information flow in the network. In this paper, we demonstrate that a highly modulated custom Transformer block is capable of capturing good spatio-temporal information for video prediction with a simpler network architecture.

2.2. Transformers in computer vision

The architecture of the proposed VPTR-NAR model follows the architecture of Detection Transformer (DETR) [41], which relaxes the dependence on complex region proposal and non-maximal suppression in object detection by using Transformers. DETR follows the classical neural machine translation (NMT) Transformer architecture. The 2D image features extracted by a CNN are flattened into a sequence and fed into a standard Transformer encoder. The output features of the Transformer encoder, normally referred as memories, together with some learned object queries are fed into a Transformer decoder. The decoder output features, corresponding to each object query, are considered as the representation of each potential object, and finally a small feed forward neural network is used to do the classification or bounding box detection based on the extracted potential object features.

Efficient visual Transformers. The bottleneck of a basic Transformer comes from the quadratic complexity of the attention score calculation. There are mainly two types of models to reduce the computation cost for visual Transformers. The first class of models reduces the flattened sequence length by different methods. ViT and many successive works [9, 42, 11] divide high resolution input into local patches, 2D or 3D, and then concatenate along the channel dimension to tokenize the patch. Pooling can be an alternative to sequence length reduction [43]. The second class of models introduces sparse attention to reduce the complexity, e.g. restricting the attention over a local region [44, 45, 46], or decomposing the global attention into a series of axial-attention [47, 48, 11]. HRFormer [46] is an example of local region attention-based Transformer, which is designed for image classification and dense prediction. Essentially, HRFormer replaces the convolution layers of HRNet [49] by self-attention. The multi-resolution parallel architecture is the same as HRNet.

Our VidHRFormer block is inspired by the HRFormer block. Specifically, a HRFormer block is composed by a

local-window multi-head self attention layer and a depth-wise convolution feed-forward network. The input feature maps $Z \in \mathbb{R}^{H \times W \times C}$ is firstly evenly divided into P non-overlapping local patches, each patch is $Z_p \in \mathbb{R}^{\frac{H}{P} \times \frac{W}{P} \times C}$. Then a multi-head self-attention is performed for each patch. Finally, the depth-wise convolution is used to exchange information of different local patches. HRFormer is similar to the Swin Transformer [45]. The Swin Transformer takes a cyclic shift window partitioning procedure instead of a depth-wise convolution layer to exchange the information of local patches.

2.3. Transformers for VFFP

A recent ConvTransformer [12] model follows the architecture of DETR [50], which also inspired the development of our VPTR-NAR. Despite the similarities, our VPTR-NAR is different from ConvTransformer with respect to the fundamental attention mechanism. Specifically, ConvTransformer proposed a custom hybrid multi-head attention module that replaces both the linear projection and dot-product attention operation by a convolution, but our VPTR-NAR uses the standard multi-head attention module. The ConvTransformer attention map calculation has a complexity of $\mathcal{O}(HWTk_c^2d_{model}^2)$, where k_c is the size of convolutional kernels. Given our experimental configurations, i.e., $H = 8$, $W = 8$, $P = 4$, $T = 20$, and $d_{model} = 512$, our efficient Transformer block (with a complexity of $\mathcal{O}(\frac{H^2W^2}{P^2} + T^2)d_{model}$) is more efficient. Another more recent VideoGPT [13] model is capable of both single image animation and VFFP. VideoGPT takes a 3D CNN as backbone to encode video clips into spatial-temporal features, which are then flattened to be a sequence to train a standard Transformer with the GPT manner. VideoGPT shares a similar architecture and train/test behaviours as our VPTR-FAR, as both of them fully decompose a joint distribution into a product of conditional distribution, like GPT. But VideoGPT performs the attention along spatial and temporal jointly while our VPTR-FAR performs the attention along spatial and temporal separately with a smaller computational complexity. More importantly, VideoGPT takes a 3D CNN to downsample the time dimension of input videos to facilitate the temporal information modeling. In our case, we do not downsample. Our VPTR models solely depend on attention for a full temporal information modeling, without downsampling and loss of information. Besides, VideoGPT is a stochastic model based on VQ-VAE [51], instead of a deterministic model as our VPTR. Another recent work NÜWA [14] shares a similar idea as VideoGPT.

3. The proposed VPTR models

3.1. Overall framework of VPTR

Our overall video prediction framework is illustrated in Fig. 1. It consists of an autoencoder and the VPTR module itself. Specifically, a CNN encoder shared by all the past frames extracts the visual features of each frame, then a VPTR, based on VidHRFormer blocks, is applied to predict

the visual features of each future frame based on the past frame features. Finally, we reconstruct the pixels of each future frame with the CNN decoder. The detail architectures of the autoencoder and three different VPTR variants are described in the following subsections.

In contrast to most ConvLSTM-based models, here we disentangle the visual feature extraction and the prediction process, similar to VideoGPT [13] and VPEG [40]. The benefit of the disentanglement is that we are able to derive the explicit representations of the input past video clips, and use them for downstream tasks, like action recognition, early action prediction etc. In other words, because of the disentanglement, it is easier to extend VPTR to solve self-supervised learning tasks, in contrast to most ConvLSTM-based models, especially in the case of the Transformation-based models, due to the lack of explicit representations. We naturally introduce a two-stages training strategy because of the feature extraction and prediction disentanglement, which increases the flexibility and reduces the learning burden of VPTR. For simplicity, Fig. 1 only shows the inference behavior, the two-stages training strategy is described in detail at the end of this section.

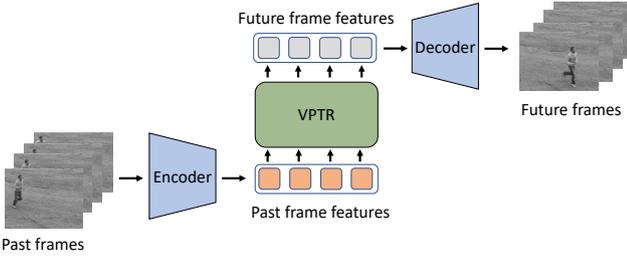


Figure 1: Overall framework of the proposed VFFP model. VPTR predicts over latent feature space.

3.2. Encoder and decoder

We adapted the ResNet-based autoencoder from the Pix2Pix model [52]. In order to match with the VPTR model feature dimensionality d_{model} , the only modification is that we change the encoder output feature channels and decoder input feature channels to be of size d_{model} . The reason to use Pix2Pix autoencoder is that there are no skip connections between the encoder and decoder. Normally, the encoder-decoder skip connections are good for a higher reconstruction image quality, like the famous U-Net architecture [53]. However, it is incompatible with our strategy of disentangling feature extraction and prediction. We chose this strategy because we hypothesize that the learned visual features by the encoder will include all the information of the input and thus would be better for future self-supervised learning tasks, while skip connections would enable some information to bypass the bottleneck encodings. More importantly, we aim to predict future frames instead of simply reconstructing past frames. The encoder-decoder skip connections from past frame features to future frame features make the two-stages training impractical because future image transformations tend to be ignored and predictions

are not learned. Training jointly does not solve this problem. This was confirmed by preliminary experiments. We failed to conduct a successful joint training of the Transformer and the autoencoder.

The loss function to train the encoder and decoder includes three terms and is given by,

$$\mathcal{L}_{rec} = \mathcal{L}_2(X, \hat{X}) + \mathcal{L}_{gdl}(X, \hat{X}) + \lambda_1 \arg \min_G \max_D \mathcal{L}_{GAN}(G, D), \quad (2)$$

where \mathcal{L}_2 denotes the MSE loss (Eq. 3), \mathcal{L}_{gdl} denotes image gradient difference loss [33] (Eq. 4) and \mathcal{L}_{GAN} denotes the GAN loss (Eq. 5). Both \mathcal{L}_{gdl} and \mathcal{L}_{GAN} are utilized to learn the details in the images and thus provide a higher visual quality. X and \hat{X} denote the original frames and reconstructed frames respectively, x_i denotes a single frame, λ_1 and α are hyperparameters. In Eq. 5, D denotes a discriminator, which is not shown in Fig. 1, and the combination of the encoder and decoder is considered to be a generator G . We train \mathcal{L}_{GAN} with the PatchGAN [52] manner.

$$\mathcal{L}_2(X, \hat{X}) = \sum_{i=1}^n \|x_i - \hat{x}_i\|_2^2 \quad (3)$$

$$\mathcal{L}_{gdl}(X, \hat{X}) = \sum_{i=1}^n \sum_{i,j} \left| |x_{i,j} - x_{i-1,j}| - |\hat{x}_{i,j} - \hat{x}_{i-1,j}| \right|^\alpha + \left| |x_{i,j-1} - x_{i,j}| - |\hat{x}_{i,j-1} - \hat{x}_{i,j}| \right|^\alpha \quad (4)$$

$$\mathcal{L}_{GAN}(G, D) = \mathbb{E}_X[\log D(X)] + \mathbb{E}_{\hat{X}}[\log(1 - D(G(X)))] \quad (5)$$

3.3. VidHRFormer Block

In order to reduce the complexity of a standard Transformer and to make it practical for high-dimensional video representation learning, our solution is to apply attention only over local spatial patch and separate the spatial and temporal attention. The proposed new video representation learning Transformer block is named as VidHRFormer block, see the blue area of Fig. 3a for the detail architecture. Essentially, we extend the HRFormer block [46] by integrating a temporal multi-head attention layer, together with some other necessary feed-forward and normalization layers.

Local spatial multi-head self-attention (MHSA). Considering a batch of video feature maps $Z \in \mathbb{R}^{N \times T \times H \times W \times d_{model}}$, the local spatial MHSA is shared by frames at different time steps. So Z is firstly reshaped and evenly divided into $P = \frac{HW}{K^2}$ local patches $\{Z_1, Z_2, \dots, Z_p\}$ along the H and W dimensions, where each local patch is of size $K \times K$, therefore $Z_p \in \mathbb{R}^{(NT) \times K^2 \times d_{model}}$. The multi-head self-attention is formulated as $MHSA(Z_p) = Concat[head(Z_p)_1, \dots, head(Z_p)_h]$, where $Concat$ denotes the concatenation operation and $head(Z_p)_i \in \mathbb{R}^{K^2 \times \frac{d_{model}}{h}}$ is calculated by

$$head(Z_p)_i = softmax\left[\frac{((Z_p^O W_i^O)(Z_p^K W_i^K))}{\sqrt{d_{model}/h}}\right] Z_p W_i^V, \quad (6)$$

where W_i^Q, W_i^K, W_i^V are linear projection matrices for the query, key and value of each head i respectively, Z_p^Q and Z_p^K denote the key and query for attention. The complexity of local spatial MHSA is $\mathcal{O}(\frac{H^2W^2}{p^2}d_{model})$.

Another critical component for a Transformer is the positional encoding, which enables the Transformer to output different representation of the same input element given different context. Both fixed absolute 2D positional encoding [41] and relative positional encoding (RPE) [54] are good candidates for the local patch to get Z_p^Q and Z_p^K . These two different positional encodings are compared in the experiments.

Convolutional feed-forward neural network (Conv FFN). Conv FFN is also shared by frames at different time step of the video. The outputs of the local spatial MHSA, $\{Z_1, Z_2, \dots, Z_p\}$ are assembled back to be $Z \in \mathbb{R}^{(NT) \times H \times W \times d_{model}}$. The Conv FFN layer includes a 3×3 depth-wise convolution and two point-wise MLPs. Different from the original HRFormer block, all the normalization layers in Conv FFN are layer normalization, instead of batch normalization.

Temporal MHSA. A temporal MHSA is placed on top of the previously described local spatial MHSA and Conv FFN to model the temporal dependency between frames. It is shared by features across different spatial locations. In other words, the input feature map $Z \in \mathbb{R}^{(NT) \times H \times W \times d_{model}}$ is reshaped to be $Z \in \mathbb{R}^{(NHW) \times T \times d_{model}}$. The temporal MHSA uses the same standard multi-head self-attention as the local spatial MHSA and thus the complexity of temporal MHSA is $\mathcal{O}(T^2d_{model})$. However, here, for simplicity, we selected a fixed absolute 1D positional encoding for the time steps. Same as the standard Transformer, we place a MLP feed-forward neural network after the temporal MHSA. Finally, the output feature map is reshaped back to be $Z \in \mathbb{R}^{N \times T \times H \times W \times d_{model}}$ to be used by the next VidHRFormer layer.

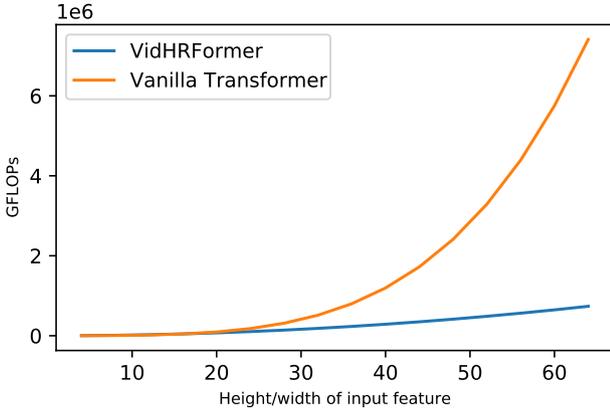


Figure 2: GFLOPs comparison of a VidHRFormer block and a vanilla Transformer block.

To sum up, the total complexity of VidHRFormer block is the combination of spatial local window MHSA complexity and temporal MHSA complexity, i.e., $\mathcal{O}(\frac{H^2W^2}{p^2} + T^2)d_{model}$. This results in a better efficiency as compared to a standard Transformer that has a complexity of $\mathcal{O}((THW)^2d_{model})$. For a better demonstration that the proposed VidHRFormer block is much more efficient than a vanilla Transformer block, we plotted the GFLOPs curves w.r.t an increasing H or W (T is fixed) in Figure 2. As the height or width of the input feature increases, the GFLOPs of a vanilla Transformer quickly increase and makes it infeasible for application. We would observe a similar phenomenon if we fix the H/W and vary the video feature length T .

In the following sections, we describe three different VPTR models that are based on the VidHRFormer block.

3.4. VPTR-FAR (fully autoregressive)

Together with a well-trained CNN autoencoder, the VPTR-FAR parameterizes the following distribution:

$$p(x_1, \dots, x_L, \dots, x_{L+N}) = \prod_{t=1}^{L+N} p(x_t | x_{t-1}, \dots, x_1) \quad (7)$$

In other words, VPTR-FAR predicts the next frame conditioned on all previous frames. This is the most common paradigm for most SOTA VFFP models. To enforce the causal relationship between the next frame and previous frames, an attention mask is applied to the temporal MHSA module. The fully autoregressive VPTR model is simply composed of a stack of multiple VidHRFormer blocks, see Fig. 3a.

For training, we feed the ground-truth frame feature sequence $\{z_1, \dots, z_{L+N-1}\}$ generated by the encoder into VPTR-FAR, which then predicts the future feature sequence $\{\hat{z}_2, \dots, \hat{z}_{L+N}\}$ for the decoder to generate frames $\{\hat{x}_2, \dots, \hat{x}_{L+N}\}$. We define the training loss of VPTR-FAR as follows (see Eq. 3 and 4):

$$\mathcal{L}_{FAR} = \sum_{t=2}^{L+N} \mathcal{L}_2(x_t, \hat{x}_t) + \sum_{t=2}^{L+N} \mathcal{L}_{gdl}(x_t, \hat{x}_t) \quad (8)$$

For test, given the ground-truth feature sequence $\{z_1, \dots, z_L\}$ of all past frames, there are two different options for the recurrent prediction of the future frames. The first one is only taking the VPTR module to recurrently predicting all the future frame features, i.e. $\hat{z}_t = \mathcal{T}(z_1, \dots, z_{t-1}), t \in [L+1, \dots, L+N]$, where \mathcal{T} denotes the VPTR-FAR module. Then we get $\hat{x}_t = Dec(\hat{z}_t), t \in [L+1, \dots, L+N]$, where Dec denotes the CNN frame decoder. The second option includes two additional steps. We firstly decode one future feature to be frame \hat{x}_t , and then encode the frame back into a latent feature before the prediction of next future frame feature, i.e., $\hat{z}_t = Enc(Dec(\mathcal{T}(z_1, \dots, z_{t-1}))), t \in [L+1, \dots, L+N]$, where Enc denotes the CNN frame encoder. The second method significantly reduces the accumulated test error during inference, and we analyze the reasons in the experiments section.

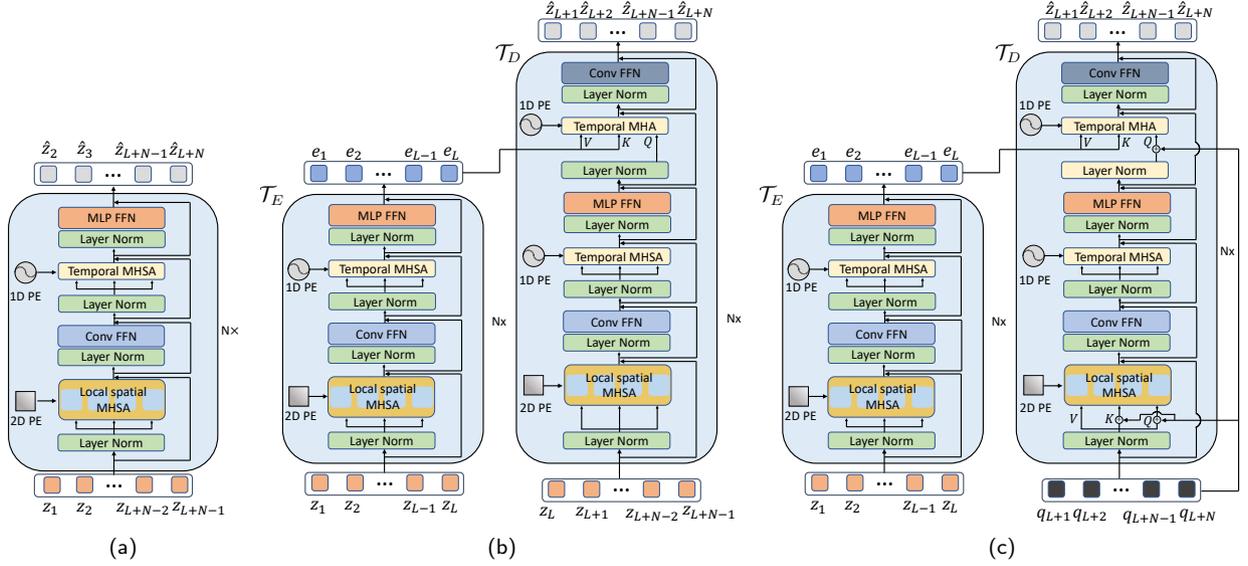


Figure 3: (a) VPTR-FAR. The blue area indicates the proposed basic VidHRFormer block. A temporal attention mask is applied to the Temporal MHSA module for VPTR-FAR. (b) VPTR-PAR. There is an attention mask for the Temporal MHA layer of the autoregressive Transformer decoder \mathcal{T}_D . (c) VPTR-NAR. The Transformer decoder \mathcal{T}_D is non-autoregressive.

3.5. VPTR-PAR (partially autoregressive)

The partially autoregressive variant is illustrated in Fig. 3b. It consists a Transformer encoder and decoder, where the encoder \mathcal{T}_E encodes all past frame features $z_t, t \in [1, L]$ to be memories $e_t, t \in [1, L]$. The architecture of \mathcal{T}_E , left part of Fig. 3b, is the same as the VPTR-FAR, except that there is no temporal attention mask for the temporal MHSA module. The autoregressive Transformer decoder \mathcal{T}_D of VPTR-PAR, right part of Fig. 3b, includes two more layers compared with \mathcal{T}_E , a temporal multi-head attention (MHA) layer and another output Conv FFN layer. The Temporal MHA layer is also called the encoder-decoder attention layer, which takes the memories as value and key, while the query is derived from the $\{z_L, \dots, z_{L+N-1}\}$. A temporal attention mask is applied to Temporal MHSA layer of \mathcal{T}_D to achieve the autoregressive modeling. Theoretically, VPTR-PAR models the following distribution:

$$p(x_{L+N}, \dots, x_{L+1} | x_L, \dots, x_1) = \prod_{t=L+1}^{L+N} p(x_t | x_{t-1}, \dots, x_{L+1}, x_L, \dots, x_1) \quad (9)$$

Compared with VPTR-FAR, we only decompose the probability distribution of future frames to be the product of a series of conditional distributions, so this model is named as partially autoregressive VPTR. During training, the frame features $\{z_1, \dots, z_{L+N-1}\}$ are divided into two parts, as shown in Fig. 3b, and fed into \mathcal{T}_E and \mathcal{T}_D respectively. Following the convention of NMT research, the input of \mathcal{T}_E is called

source sequence, and input of \mathcal{T}_D is called target sequence. The training loss of VPTR-PAR is formulated as

$$\mathcal{L}_{PAR} = \sum_{t=L+1}^{L+N} \mathcal{L}_2(x_t, \hat{x}_t) + \sum_{t=L+1}^{L+N} \mathcal{L}_{gdl}(x_t, \hat{x}_t), \quad (10)$$

where \mathcal{L}_{gdl} and \mathcal{L}_2 are defined in Eq. 4 and Eq. 3 respectively.

Same as VPTR-FAR, we can use the same two different recurrent inference methods for VPTR-PAR.

3.6. VPTR-NAR (non-autoregressive)

In order to reduce the predicted frames accumulated error and increase the inference speed of the two autoregressive counterparts, we propose a non-autoregressive variant (VPTR-NAR), which is inspired by the architecture of DETR [41]. VPTR-NAR is illustrated in Fig. 3c. VPTR-NAR shares the same \mathcal{T}_E as VPTR-PAR, while the \mathcal{T}_D of them are slightly different. Firstly, target sequence for \mathcal{T}_D is substituted with zero [41], instead of $\{z_1, \dots, z_{L+N-1}\}$ generated by the CNN encoder, and a future frames query sequence $\{q_{L+1}, \dots, q_{L+N}\}$ is fed into two different sublayers of \mathcal{T}_D , where $q_t \in \mathbb{R}^{H \times W \times C}, t \in [L+1, L+N]$. The future frame query sequence is randomly initialized and updated during training. Secondly, there is no temporal attention mask for any temporal attention layer of VPTR-NAR, because we do not need to impose conditional dependency among each future frame query. Theoretically, VPTR-NAR directly models the following conditional distribution:

$$p(x_{L+N}, \dots, x_{L+1} | x_L, \dots, x_1) \quad (11)$$

Contrastive feature loss for VPTR-NAR. Unfortunately, a combination loss of MSE and GDL, i.e., $\mathcal{L} = \sum_{t=L+1}^{L+N} \mathcal{L}_2(x_t, \hat{x}_t) + \mathcal{L}_{gdl}(x_t, \hat{x}_t)$, is not enough to train VPTR-NAR. Specifically, we observe that all the predicted future frames of one video clip are somewhat similar to each other when VPR-NAR is trained with the same loss function as VPTR-FAR and VPTR-NAR, which indicates that VPTR-NAR cannot learn good motion information in this case. In NLP, a similar phenomenon is also observed for some non-autoregressive NMT models, where the Transformer decoder frequently generates repeated tokens [55]. This is because autoregressive models make the estimation of joint distribution to be tractable and thus they are easier to train, even though slower. To deal with this problem, we maximize the mutual information between predicted future frame feature \hat{z}_t and the future frame feature z_t (ground-truth) generated by the CNN encoder by adapting another contrastive feature loss \mathcal{L}_c [56], where $t \in [L+1, L+N]$. \mathcal{L}_c is defined by

$$\mathcal{L}_c(z_t, \hat{z}_t) = \frac{1}{2} \sum_{s=1}^{S_t} l_c(\hat{v}_s, v_s, sg(\bar{v}_s)) + l_c(v_s, \hat{v}_s, sg(\hat{\bar{v}}_s)), \quad (12)$$

where $v_s \in \mathbb{R}^{d_{model}}$ denotes a feature vector at spatial location s of z_t , $\bar{v}_s \in \mathbb{R}^{(S_t-1) \times d_{model}}$ denotes the collection of feature vectors at all other spatial locations of z_t . \hat{v}_s and $\hat{\bar{v}}_s$ of \hat{z}_t are defined in the same way. The total number of spatial locations in a feature map is $S_t = H \times W$, and sg is the stop gradient operation. l_c is the infoNCE-based contrastive loss formulated as follows,

$$l_c(v, v^+, v^-) = -\log \frac{\exp(s(v, v^+))}{\exp(s(v, v^+)) + \sum_{m=1}^M \exp(s(v, v^-))}. \quad (13)$$

Similar to other contrastive learning objectives, (v, v^+) denotes a positive pair of examples and (v, v^-) denotes a negative pair of examples. In detail, $v^+ \in \mathbb{R}^{d_{model}}$ is the spatially-corresponding ground-truth feature vector of a feature vector $v \in \mathbb{R}^{d_{model}}$, while $v^- \in \mathbb{R}^{M \times d_{model}}$ denotes the M other ground-truth feature vectors at different spatial location. $s(v1, v2)$ is the feature dot-product similarity operation between feature vectors $v1$ and $v2$. We can finally compose the loss function of VPTR-NAR as

$$\mathcal{L}_{NAR} = \sum_{t=L+1}^{L+N} \mathcal{L}_2(x_t, \hat{x}_t) + \mathcal{L}_{gdl}(x_t, \hat{x}_t) + \lambda_2 \mathcal{L}_c(z_t, \hat{z}_t). \quad (14)$$

Different from VPTR-FAR and VPTR-PAR, VPTR-NAR predicts N future frames simultaneously instead of recurrently. Besides, the testing behavior and training behavior of VPTR-NAR are the same. For the case that we need to predict more than N future frames during test, we can use VPTR-NAR as a block-wise autoregressive model,

i.e., taking the N predicted future frames as past frames and feeding them back into the encoder for predicting the next N future frames.

3.7. Training strategy

We divide the VFFP model training process into two stages. In stage one, the VPTR module is ignored and we only train the encoder and decoder as a normal autoencoder with the loss function of Eq. 2, which aims to reconstruct all the frames of the whole training set perfectly. In stage two, we freeze the well-trained encoder and decoder and only optimize the parameters of the VPTR module. VPTR-FAR, VPTR-PAR and VPTR-NAR are trained with the loss functions defined in Eq. 8, Eq. 10 and Eq. 14 respectively. As we mentioned before, a two-stage training strategy is naturally compatible with the disentanglement of feature extraction and prediction process. Moreover, an end-to-end training of Transformers and the CNN autoencoder is much more difficult and it requires a much bigger memory for computations. Besides, we observe that a final joint fine-tuning of autoencoder and VPTR after two-stage training is not beneficial. Furthermore, the two-stage training strategy gives us a more flexible framework since we are allowed to test different VPTR variants without repetitive training of the encoder and decoder.

4. Experiments

4.1. Datasets and Metrics

The proposed VPTR models are evaluated over three datasets: BAIR [57], KTH [58], and MovingMNIST [59].

BAIR showcases video clips of a randomly moving robot arm that pushes different objects on a table. The training and testing sets of BAIR are predefined by the dataset authors. The frame size of BAIR is 64×64 . We normalize it for training, but without any data augmentation.

KTH includes grayscale video clips of different human actions. Following the experiments setup of previous work [21, 16], the frames are center cropped to be square and then resize to 64×64 . Besides, frames without human are removed based on the results of a person detector. Video clips of persons 1-16 are used for training, and persons 17-25 are used for testing. The dataset pixels are normalized before training. Random horizontal and vertical flips of each video clip are utilized as data augmentation.

MovingMNIST is a synthetic dataset where two MNIST characters randomly move in a square. We utilize the same MovingMNIST dataset as E3D-LSTM [60] instead of randomly generating frames on the fly. The frame size of MovingMNIST is 64×64 . The training set contains 10000 clips, test set contains 3000 clips and valid set contains 2000 clips. The same data augmentation method as KTH is applied to MovingMNIST dataset.

Following the experimental protocols of previous work, our VPTR models are trained to predict 10 future frames given 2 past frames for the BAIR dataset, and predict 10 future frames given 10 past frames for KTH and MovingMNIST.

Metrics. Learned Perceptual Image Patch Similarity (LPIPS) and Structural Similarity Index Measure (SSIM) are used to evaluate all the three datasets. Peak Signal-to-Noise Ratio (PSNR) is used to evaluate the KTH and BAIR dataset, and Mean Square Error (MSE) is used to evaluate the MovingMNIST dataset. All the LPIPS values are presented in 10^{-3} scale. Smaller values are better for LPIPS and MSE, while larger values are better for PSNR and SSIM.

4.2. Implementation

Training stage one. For KTH and BAIR, the output layer of the decoder is Tanh. For MovingMNIST, the output layer of decoder is Sigmoid. We set $\lambda_1 = 0.01$ for KTH and MovingMNIST, $\lambda_1 = 0$ for BAIR dataset. The optimizer is Adam, with betas of (0.5, 0.999) and a learning rate of $2e^{-4}$.

Training stage two. We define the visual features dimension of each frame as $H = 8, W = 8, d_{model} = 528$. For local spatial MHSA, the local patch size is $K = 4$. VPTR-FAR includes 12 layers of VidHRFormer blocks. For VPTR-NAR and VPTR-PAR, the number of layers of \mathcal{T}_E is 4, and the number of layers of \mathcal{T}_D is 8. All Transformers are optimized by AdamW with a learning rate of $1e^{-4}$. Gradient clipping is employed to stabilize the training. $\lambda_2 = 0.1$ for the loss function of VPTR-NAR (Eq. 14).

4.3. Results and discussion

Results on KTH. We present the best results of all three VPTR variants in Table 1. The reported average metrics are calculated over 20 predicted frames, which follows the same evaluation protocol of previous works. Compared with some SOTA models, all three VPTR variants outperform them in terms of LPIPS by a large margin. For PSNR and SSIM, our proposed VPTR models reach competitive performances.

Fig. 4 shows some prediction examples for a qualitative comparison. Comparing LMC-Memory with VPTR-NAR and VPTR-FAR, we observe that VPTR-NAR is better than VPTR-FAR and that both our VPTR models generate predictions that are more aligned with the ground-truth. It indicates that the VPTR-NAR and VPTR-FAR learns a better cyclic hand waving movements that is only condition on the past frames. We suspect that LMC-Memory retrieves some plausible but inaccurate motion from the learned memory bank and thus the prediction deviates from the ground-truth. Besides, due to a bigger accumulation of errors during inference, the last few predicted frames of both VPTR-PAR and VPTR-FAR are worse than VPTR-NAR. We analyzed the reasons in detail at section 4.5 and 4.6.

Results on MovingMNIST. The results on the MovingMNIST dataset are shown in the right part of Table 1. In terms of SSIM, the performance of our VPTR models is close to the SOTA, but there are large gaps in terms of MSE and LPIPS, especially for VPTR-FAR and VPTR-PAR. After inspection of some prediction examples, we find that our VPTRs make poor predictions for the overlapping characters.

Results on BAIR. Because the robot arm motion in BAIR dataset is random and we only condition on two past

Table 1

Results on KTH and MovingMNIST. \uparrow : higher is better, \downarrow : lower is better. **Boldface**: best results.

Methods	KTH				Moving MNIST		
	10 \rightarrow 20						
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	MSE \downarrow	SSIM \uparrow	LPIPS \downarrow	
MCNET [21]	25.95	0.804	-	-	-	-	
PredRNN++ [61]	28.47	0.865	228.9	46.5	0.898	59.5	
E3D-LSTM [60]	29.31	0.879	-	41.3	0.910	-	
STMFA net [16]	29.85	0.893	118.1	-	-	-	
Conv-TT-LSTM [62]	28.36	0.907	133.4	53.0	0.915	40.5	
LMC-Memory [7]	28.61	0.894	133.3	41.5	0.924	46.9	
VPTR-NAR	26.96	0.879	86.1	63.6	0.882	107.5	
VPTR-PAR	25.40	0.836	84.8	93.2	0.859	138.4	
VPTR-FAR	26.13	0.859	79.6	107.2	0.844	157.8	

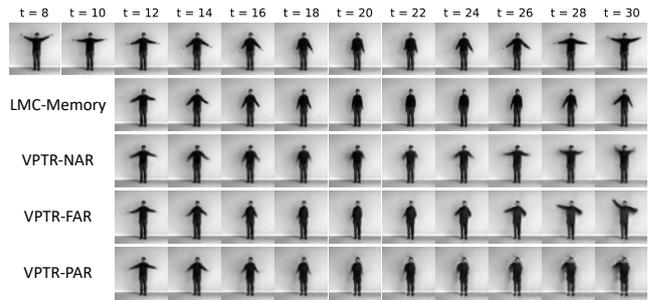


Figure 4: Qualitative results on KTH dataset. The first row is ground-truth. For the past frames, $t \in [1, 10]$. For the future frames, $t \in [11, 30]$.

frames to prediction ten future frames, BAIR is a more challenging dataset than KTH and MovingMNIST. The test results on BAIR dataset are listed in Table 2. We find that the performance of VPTR-NAR is close to the reference models in terms of all three metrics. Particularly, VPTR-NAR outperforms two MCVD [63] models in terms of both SSIM and PSNR, it also outperforms SVG-LP with regard to PSNR. We note that the predicted robot arm becomes blurry after the first few frames due to the deterministic nature of VPTRs, see Fig. 5. Our VPTRs could be extended to be stochastic models easily, and we expect that the stochastic version of VPTRs would achieve better performance on the BAIR dataset. The autoregressive variants have worse performances compared to the non-autoregressive variant because of the accumulation of errors, see the ablation study for more detailed analysis.

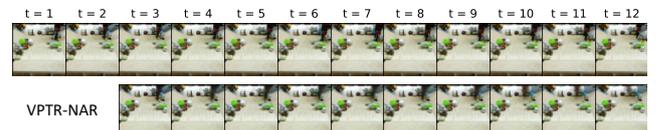


Figure 5: Qualitative results on BAIR dataset. The first row is ground-truth. For the past frames, $t \in [1, 2]$. For the future frames, $t \in [3, 12]$.

Table 2

Results on BAIR. \uparrow : higher is better, \downarrow : lower is better. **Boldface**: best results.

Methods	2 \rightarrow 28		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
SV2P [19]	20.36	0.817	91.4
MCVD-concat [63]	17.70	0.797	-
MCVD-spatin [63]	17.70	0.789	-
SVG-LP [20]	17.72	0.815	60.3
Improved VRNN [31]	-	0.822	55.0
STMFANet [16]	21.02	0.844	93.6
VPTR-NAR	17.77	0.813	70.0
VPTR-PAR	15.94	0.745	104.8
VPTR-FAR	15.76	0.724	110.7

4.4. Running time comparison

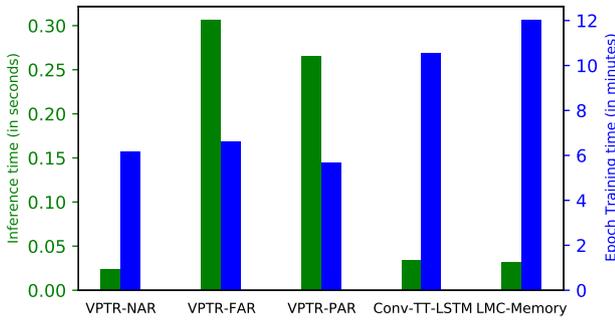


Figure 6: Comparison of inference time and epoch training time.

In order to show that VPTR is less complex and more efficient compared with ConvLSTM-based methods, we evaluated the inference time and epoch training time on the KTH dataset of all our VPTR variants and two SOTA ConvLSTM-based methods, Conv-TT-LSTM and LMC-Memory. Specifically, we measured the average inference time for predicting 10 future frames given 10 past frames. The results show that VPTR-NAR is the fastest one because of the non-autoregressive prediction. VPTR-FAR and VPTR-PAR are significantly slower than all other models because autoregressive prediction by Transformer is expensive. The results of average epoch training time indicate that all three VPTR variants have a similar training speed and are almost two times faster than the ConvLSTM-based models, thanks to the benefit from the parallel computation of Transformers. On the contrary, ConvLSTM-based models need to slowly backpropagate through each time step during training. All models are tested on the same NVidia RTX3090 GPU with the same batch size.

4.5. Ablation Study

We conducted a thorough ablation study to investigate the influence of positional encodings, separated spatial and temporal attention and autoregressive inference methods. The ablation study results are summarized in Table 3. Besides, we report the number of parameters and inference FLOPs of each different model in 4.

Table 3

Ablation study on KTH and MovingMNIST. \uparrow : higher is better, \downarrow : lower is better. **Boldface**: best results.

Methods	KTH			Moving MNIST		
	10 \rightarrow 20					
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	MSE \downarrow	SSIM \uparrow	LPIPS \downarrow
VPTR-NAR-BASE	26.92	0.881	94.6	64.2	0.880	114.2
VPTR-NAR-RPE	26.96	0.879	86.1	63.6	0.882	107.5
VPTR-NAR-FSTA	26.25	0.872	101.1	68.0	0.872	128.7
VPTR-PAR-BASE	25.04	0.832	86.0	93.2	0.859	138.4
VPTR-PAR-RPE	25.40	0.836	84.8	104.2	0.848	139.4
VPTR-PAR-FSTA	24.33	0.814	89.9	81.6	0.873	105.9
VPTR-PAR-RIL	23.14	0.694	193.6	194.1	0.362	516.1
VPTR-FAR-BASE	25.71	0.816	79.5	108.3	0.843	157.3
VPTR-FAR-RPE	26.13	0.859	79.6	107.2	0.844	157.8
VPTR-FAR-RIL	21.61	0.678	192.7	138.2	0.821	445.7

RPE. We take the VPTRs with fixed absolute positional encodings as the base models, i.e. VPTR-NAR-BASE, VPTR-PAR-BASE and VPTR-FAR-BASE. To investigate the influence of relative positional encodings for all VPTR variants, we substituted the 2D absolute positional encoding of all local spatial MHSA module with a learned 2D RPE, which give us VPTR-NAR-RPE, VPTR-PAR-RPE and VPTR-FAR-RPE. Our experiments show that in general, RPE is beneficial for the performance on both KTH and MovingMNIST datasets because the RPE models outperform the base models with regard to most metrics. Therefore, we can conclude that RPE is better than absolute positional encodings for VPTR. However, RPE would also introduce additional computational cost, see Table 4, which shows that FLOPs of RPE variants are slightly larger than the base models.

Spatial-temporal separated attention. We propose to separate the spatial and temporal attention since we aim to reduce the complexity of the standard Transformer for video feature learning. However, this separated attention mechanism means that a feature at one location only attends to partial locations of the whole spatio-temporal space. In order to analyze the impact of the separated attention, the encoder-decoder attention layers of VPTR-NAR and VPTR-FAR are replaced with a full spatio-temporal attention (FSTA), which has a complexity of $\mathcal{O}(\frac{H^2W^2T^2}{p^2}d_{model})$. As we only replace the encoder-decoder attention layers, the increased computation cost is affordable. The FLOPs of FSTA variants shown in Table 4 validate our arguments. We find that FSTA is not beneficial by comparing the VPTR-NAR-FSTA and VPTR-PAR-FSTA with their base counterparts. Consequently, it is safe to conclude that the alternate stacking of multiple spatial and temporal attention layers, as proposed with our VidHRFormer block, is capable of propagating global information from past frames to future frames.

Autoregressive inference methods. For two autoregressive variants, VPTR-FAR and VPTR-PAR, we can perform recurrently inference over latent space (RIL) or recurrently inference over pixel space (RIP) as we have mentioned in

Table 4

Comparison of FLOPs and number of parameters for different models in the ablation study.

Methods	Moving MNIST 10 → 10	
	FLOPs	#Params
VPTR-NAR-BASE	197.00G	165.83M
VPTR-NAR-RPE	201.13G	162.48M
VPTR-NAR-FSTA	197.81G	165.83M
VPTR-PAR-BASE	664.44G	164.71M
VPTR-PAR-RPE	679.59G	164.87M
VPTR-PAR-FSTA	668.93G	164.71M
VPTR-PAR-RIL	601.71G	164.71M
VPTR-FAR-BASE	1.40T	136.37M
VPTR-FAR-RPE	1.45T	133.02M
VPTR-FAR-RIL	1.34T	136.37M

Section 3.4. VPTR-FAR-BASE and VPTR-PAR-BASE are evaluated by RIP. The FLOPs presented in Table 4 indicate that RIL variants require fewer FLOPs than the RIP models, i.e., it is a little faster than RIP. However, we observe that VPTR-FAR-RIL and VPTR-PAR-RIL are outperformed by a large margin due to the severe accumulation of errors of RIL.

Accumulation of errors exists for any inference process with autoregressive models that are trained with a teacher-force manner, and we analyze the detailed reasons in section 4.6. RIL gets a larger accumulated error than RIP because the VPTR-FAR and VPTR-PAR receive only supervision from the pixel space during training. According to the loss functions in Eq. 8 and Eq. 10, it is clear that there is no direct penalty on the distance between the feature space generated by the CNN encoder and the feature space predicted by the Transformer. Furthermore, the pixel space dimensionality is smaller than the latent space dimensionality of the autoencoder, which is a common case for VFFP, as there are no skip connections from encoder to decoder and good reconstruction visual quality is expected. Therefore, we can hypothesize that recurrent inference solely by the Transformer predictor would make the predicted features quickly deviate from the ground-truth features (learned by the autoencoder during stage one). However, decoding the feature firstly and then encoding it back into latent space by the CNN encoder restricts the deviation to some degree.

4.6. Comparison of VPTR variants

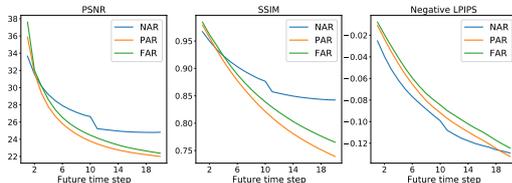


Figure 7: Results of VPTR variants on KTH for increasing prediction steps.

For better visualization of the difference between the three VPTR variants, i.e., VPTR-NAR-BASE, VPTR-FAR-BASE and VPTR-PAR-BASE, the curves of the metrics with respect to the predicted future time steps are plotted in Fig. 7. Recall that for LPIPS, a smaller value is better. The performance of VPTR-FAR and VPTR-PAR are close to each other, but VPTR-FAR is slightly better. Comparing the loss function of VPTR-FAR (Eq. 8) and VPTR-PAR (Eq. 10), the loss function of VPTR-FAR is calculated from $t = 2$ to $t = L + N$, while VPTR-PAR is only trained by the loss that is calculated from $t = L + 1$ to $t = L + N$. We believe the better performance of VPTR-FAR can be attributed to the larger supervision it receives during training.

Comparing the autoregressive variants with VPTR-NAR, VPTR-FAR and VPTR-PAR achieve a better PSNR and SSIM than VPTR-NAR during the first few prediction steps, but their performance drops quickly due to the accumulation of errors introduced by the recurrent inference. The PSNR and SSIM curves demonstrate that VPTR-NAR has a smaller quality degradation. For the last 10 steps of the LPIPS curve, VPTR-NAR also has a smaller slope than VPTR-FAR.

The accumulated inference errors of autoregressive VPTR variants are mainly due to the discrepancy between training and testing behaviors. Specifically, the previously predicted frames are used during inference instead of the ground-truth as during training, which leads to a worse generalization ability for VPTR-FAR or VPTR-PAR given the same number of VidHRFormer layers as VPTR-NAR. On the contrary, the training and testing behaviors of VPTR-NAR are the same. However, it is more difficult for the VPTR-NAR to directly estimate the joint distribution, so an additional contrastive feature loss and more parameters (the learnable future frame queries) are required.

VPTR-NAR has another advantage, i.e., a faster inference speed. Consider predicting N future frames given L past frames by the three VPTR variants. Then VPTR-NAR has a complexity of $\mathcal{O}(N^2)$, but the complexity for VPTR-FAR and VPTR-PAR is $\mathcal{O}(\sum_{n=1}^N n^2)$. Even though VPTR-PAR has the same complexity as VPTR-FAR in terms of the predicted future frames length, VPTR-PAR has a faster inference speed. Indeed, in VPTR-PAR, the past frames are processed by 4 layers of VidHRFormer block (\mathcal{T}_E), and each future frame are generated by passing through another 8 layers of VidHRFormer block (\mathcal{T}_D). However, in VPTR-FAR, each future frame is generated by passing all previous frames through 12 layers of VidHRFormer block, which costs more time. For simplicity, in this assessment, we ignored the H , W and d_{model} of video features, the computation cost of processing past frames, and we supposed that the predicted future frames length during test is the same as of the training. We tested the inference speed of three models on an NVidia RTX3090, the results in Figure 6 show that VPTR-NAR is 12.75 times faster than VPTR-FAR, and VPTR-PAR is 1.2 times faster than VPTR-FAR. The FLOPs results summarized in Table 4 also validate our analysis.

5. Conclusion

In this paper, an efficient VidHRFormer block is proposed for spatio-temporal representation learning, and three different VFFP models are developed based on it. Our proposed VidHRFormer block could be applied to many other video processing tasks as a backbone. We compared the performance of the proposed VPTRs with SOTA models on various datasets, and our proposed methods are competitive with more complex ConvLSTM-based models. Finally, we conducted a through ablation study to analyze the influence of different modules for three VPTR variants, and we observed that VPTR-NAR achieves a better performance than VPTR-FAR and VPTR-PAR.

References

- [1] Jan-Aike Bolte, Andreas Bar, Daniel Lipinski, and Tim Fingscheidt. Towards Corner Case Detection for Autonomous Driving. In *2019 IEEE Intelligent Vehicles Symposium (IV)*, pages 438–445, June 2019. ISSN: 2642-7214.
- [2] Wen Liu, Weixin Luo, Dongze Lian, and Shenghua Gao. Future Frame Prediction for Anomaly Detection – A New Baseline. In *CVPR*, 2018.
- [3] Felix Leibfried, Nate Kushman, and Katja Hofmann. A Deep Learning Approach for Joint Video Frame and Reward Prediction in Atari Games. In *ICML Workshop on Principled Approaches to Deep Learning*, 2016.
- [4] Y. Bengio, A. Courville, and P. Vincent. Representation Learning: A Review and New Perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, August 2013.
- [5] Xiaolong Wang and Abhinav Gupta. Unsupervised Learning of Visual Representations Using Videos. In *ICCV*, 2015.
- [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.
- [7] Sangmin Lee, Hak Gu Kim, Dae Hwi Choi, Hyung-II Kim, and Yong Man Ro. Video Prediction Recalling Long-Term Motion Context via Memory Alignment Learning. In *CVPR*, 2021.
- [8] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017.
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- [10] Patrick Esser, Robin Rombach, and Björn Ommer. Taming Transformers for High-Resolution Image Synthesis. In *CVPR*, 2021.
- [11] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. ViViT: A video vision transformer. In *ICCV*, 2021.
- [12] 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]*, November 2020.
- [13] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. VideoGPT: Video Generation using VQ-VAE and Transformers. In *arXiv:2104.10157 [cs]*, September 2021. arXiv: 2104.10157.
- [14] Chenfei Wu, Jian Liang, Lei Ji, Fan Yang, Yuejian Fang, Daxin Jiang, and Nan Duan. NVUWA: Visual Synthesis Pre-training for Neural visUal World creation. *arXiv:2111.12417 [cs]*, November 2021.
- [15] Xi Ye and Guillaume-Alexandre Bilodeau. VPTR: Efficient Transformers for Video Prediction. In *26th International Conference on Pattern Recognition (ICPR)*, 2022.
- [16] 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.
- [17] Y. Wang, J. Wu, M. Long, and J. B. Tenenbaum. Probabilistic Video Prediction From Noisy Data With a Posterior Confidence. In *CVPR*, 2020.
- [18] Dinesh Jayaraman, Frederik Ebert, Alyosha Efros, and Sergey Levine. Time-agnostic prediction: Predicting predictable video frames. In *ICLR*, 2019.
- [19] Mohammad Babaeizadeh, Chelsea Finn, Dumitru Erhan, Roy Campbell, and Sergey Levine. Stochastic variational video prediction. In *ICLR*, 2018.
- [20] Emily Denton and Rob Fergus. Stochastic Video Generation with a Learned Prior. In *ICML*, 2018.
- [21] Ruben Villegas, Jimei Yang, Seunghoon Hong, Xunyu Lin, and Honglak Lee. Decomposing motion and content for natural video sequence prediction. In *ICLR*, 2017.
- [22] Jean-Yves Franceschi, Edouard Delasalles, Mickael Chen, Sylvain Lamprier, and Patrick Gallinari. Stochastic Latent Residual Video Prediction. In *ICLR*, 2020.
- [23] Yunseok Jang, Gunhee Kim, and Yale Song. Video Prediction with Appearance and Motion Conditions. In *ICML*. PMLR, 2018.
- [24] S. Tulyakov, M. Liu, X. Yang, and J. Kautz. MoCoGAN: Decomposing Motion and Content for Video Generation. In *CVPR*, 2018.
- [25] Emily L. Denton and Vighnesh Birodkar. Unsupervised Learning of Disentangled Representations from Video. In *NIPS*, volume 30, 2017.
- [26] Haoye Cai, Chunyan Bai, Yu Wing Tai, and Chi Keung Tang. Deep video generation, prediction and completion of human action sequences. In *ECCV*, 2018.
- [27] Naoya Fushishita, Antonio Tejero-de Pablos, Yusuke Mukuta, and Tatsuya Harada. Long-Term Human Video Generation of Multiple Futures Using Poses. In *Computer Vision – ECCV 2020 Workshops*, pages 596–612, 2020.
- [28] Jun Ting Hsieh, Bingbin Liu, De An Huang, Li Fei-Fei, and Juan Carlos Niebles. Learning to decompose and disentangle representations for video prediction. In *NIPS*, 2018.
- [29] Adam Kosior, Hyunjik Kim, Yee Whye Teh, and Ingmar Posner. Sequential attend, infer, repeat: Generative modelling of moving objects. In *NIPS*, 2018.
- [30] Xionghao 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.
- [31] L. Castrejon, N. Ballas, and A. Courville. Improved Conditional VRNNs for Video Prediction. In *ICCV*, October 2019.
- [32] Y. Kwon and M. Park. Predicting Future Frames Using Retrospective Cycle GAN. In *CVPR*, 2019.
- [33] Michael Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond mean square error. In *ICLR*, 2016.
- [34] Baoyang Chen, Wenmin Wang, and Jinzhuo Wang. Video Imagination from a Single Image with Transformation Generation. In *Proceedings of the on Thematic Workshops of ACM Multimedia*, 2017.
- [35] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. In *NIPS*, 2016.
- [36] Y. Wu, R. Gao, J. Park, and Q. Chen. Future Video Synthesis With Object Motion Prediction. In *CVPR*, 2020.
- [37] 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.
- [38] B. Jin, Y. Hu, Y. Zeng, Q. Tang, S. Liu, and J. Ye. VarNet: Exploring Variations for Unsupervised Video Prediction. In *IROS*, 2018.
- [39] Katerina Fragkiadaki, Jonathan Huang, Alex Alemi, Sudheendra Vijayanarasimhan, Susanna Ricco, and Rahul Sukthankar. Motion Prediction Under Multimodality with Conditional Stochastic Networks. In *arXiv:1705.02082 [cs]*, 2017.

- [40] Jingwei Xu, Huazhe Xu, Bingbing Ni, Xiaokang Yang, and Trevor Darrell. Video Prediction via Example Guidance. In *ICML*, 2020.
- [41] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-End Object Detection with Transformers. In *ECCV*, 2020.
- [42] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid Vision Transformer: A Versatile Backbone for Dense Prediction Without Convolutions. In *ICCV*, 2021.
- [43] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *ICCV*, 2021.
- [44] Pengchuan Zhang, Xiyang Dai, Jianwei Yang, Bin Xiao, Lu Yuan, Lei Zhang, and Jianfeng Gao. Multi-Scale Vision Longformer: A New Vision Transformer for High-Resolution Image Encoding. In *ICCV*, 2021.
- [45] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021.
- [46] Yuhui Yuan, Rao Fu, Lang Huang, Weihong Lin, Chao Zhang, Xilin Chen, and Jingdong Wang. HRFormer: High-resolution transformer for dense prediction. In *NeurIPS*, 2021.
- [47] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. CCNet: Criss-Cross Attention for Semantic Segmentation. In *ICCV*, 2019.
- [48] Huiyu Wang, Yukun Zhu, Bradley Green, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Axial-DeepLab: Stand-Alone Axial-Attention for Panoptic Segmentation. In *ECCV*, 2020.
- [49] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep High-Resolution Representation Learning for Human Pose Estimation. In *CVPR*, 2019.
- [50] Tim Meinhardt, Alexander Kirillov, Laura Leal-Taixe, and Christoph Feichtenhofer. TrackFormer: Multi-Object Tracking with Transformers. In *arXiv:2101.02702 [cs]*, January 2021. arXiv: 2101.02702.
- [51] Aaron van den Oord, Oriol Vinyals, and koray kavukcuoglu. Neural discrete representation learning. In *NIPS*, 2017.
- [52] Phillip Isola, Jun Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017.
- [53] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation. In *MICCAI*, 2015.
- [54] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. In *NAACL*, 2018.
- [55] Yiren Wang, Fei Tian, Di He, Tao Qin, ChengXiang Zhai, and Tie-Yan Liu. Non-autoregressive machine translation with auxiliary regularization. In *AAAI*, volume 33, pages 5377–5384, 2019. Number: 01.
- [56] Alex Andonian, Taesung Park, Bryan Russell, Phillip Isola, Jun-Yan Zhu, and Richard Zhang. Contrastive feature loss for image prediction. In *ICCVW*, 2021.
- [57] Frederik Ebert, Chelsea Finn, Alex X. Lee, and Sergey Levine. Self-supervised visual planning with temporal skip connections. In *CoRL*, 2017.
- [58] C. Schuldt, I. Laptev, and B. Caputo. Recognizing human actions: a local SVM approach. In *ICPR*, 2004.
- [59] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhudinov. Unsupervised Learning of Video Representations using LSTMs. In *ICML*, 2015.
- [60] 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.
- [61] Yunbo Wang, Zhifeng Gao, Mingsheng Long, Jianmin Wang, and Phillip S. Yu. PredRNN++: Towards A Resolution of the Deep-in-Time Dilemma in Spatiotemporal Predictive Learning. In *ICML*, 2018.
- [62] Jiahao Su, Wonmin Byeon, Jean Kossaifi, Furong Huang, Jan Kautz, and Animashree Anandkumar. Convolutional Tensor-Train LSTM for Spatio-temporal Learning. In *NeurIPS*, 2020.
- [63] 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.