
CCVS: Context-aware Controllable Video Synthesis

Guillaume Le Moing^{1,*} Jean Ponce^{1,2} Cordelia Schmid¹

¹Inria and Department of Computer Science, ENS, CNRS, PSL Research University

²Center for Data Science, New York University

Abstract

This presentation introduces a self-supervised learning approach to the synthesis of new video clips from old ones, with several new key elements for improved spatial resolution and realism: It conditions the synthesis process on contextual information for temporal continuity and ancillary information for fine control. The prediction model is doubly autoregressive, in the latent space of an autoencoder for forecasting, and in image space for updating contextual information, which is also used to enforce spatio-temporal consistency through a learnable optical flow module. Adversarial training of the autoencoder in the appearance and temporal domains is used to further improve the realism of its output. A quantizer inserted between the encoder and the transformer in charge of forecasting future frames in latent space (and its inverse inserted between the transformer and the decoder) adds even more flexibility by affording simple mechanisms for handling multimodal ancillary information for controlling the synthesis process (*e.g.*, a few sample frames, an audio track, a trajectory in image space) and taking into account the intrinsically uncertain nature of the future by allowing multiple predictions. Experiments with an implementation of the proposed approach give very good qualitative and quantitative results on multiple tasks and standard benchmarks.

1 Introduction

Feeding machines with extensive video content, and teaching them to create new samples on their own, may deepen their understanding of both the physical and social worlds. Video synthesis has numerous applications from content creation (*e.g.*, deblurring, slow motion) to human-robot interaction (*e.g.*, motion prediction). Despite the photo-realistic results of modern image synthesis models [38], video synthesis is still lagging behind due to the increased complexity of the additional temporal dimension.

An emerging trend is to use autoregressive models, for example transformer architectures [68], for their simplicity, and their ability to model long-range dependencies and learn from large volumes of data [4, 16]. First introduced for natural language processing (NLP) and then successfully applied to visual data [18], the strength of transformers is grounded in a self-attention mechanism which considers all pairwise interactions within the data. The price to pay is a computational complexity which grows quadratically with the data size, which itself depends linearly on the temporal dimension and quadratically on the spatial resolution in the image domain. Although there have been some efforts to reduce the complexity of self-attention [11, 39, 50], using such methods directly on visual data is still limited to low resolutions and impractical without considerable computational power [9, 79].

Some recent works [20, 53] address this problem by using an autoencoder to compress the visual data, and apply the autoregressive model in the latent space. This greatly reduces the memory footprint and computational cost, yet, the greater the compression, the harder it is to faithfully reconstruct frames. The corresponding trade-offs may undermine the practical usability of these approaches. GANs [26] mitigate this issue by “hallucinating” plausible local details in image synthesis [20]. But latent video

*corresponding author: guillaume.le-moing@inria.fr

transformers [53] decode frames independently, which prevents local details from being temporally coherent. Hence, using GANs in this setting may result in flickering effects in the synthesized videos.

We follow the problem decomposition from [20, 53], but introduce a more elaborate compression strategy with CCVS (for *Context-aware Controllable Video Synthesis*), an approach that takes advantage of “context” frames (*i.e.*, both input images and previously synthesized ones) to faithfully reconstruct new ones despite lossy compression. As shown in Figure 1, CCVS relies on optical flow estimation between context and new frames, within temporal skip connections, to let information be shared across timesteps. New content, which cannot be retrieved from context, is synthesized directly from latent compressed features, and adversarial training [26] is used to make up realistic details. Indeed, information propagates in CCVS to new frames as previously synthesized ones become part of the context. Like other video synthesis architectures based on autoregressive models [53, 79], CCVS can be used in many tasks besides future video prediction. Any data which can be expressed in the form of a fixed-size sequence of elements from a finite set (*aka*, tokens) can be processed by a transformer. This applies to video frames (here, via compression and quantization) and to other types of data. Hence, one can easily fuse modalities without having to build complex or task-specific architectures. This idea has been used to control image synthesis [20, 54], and we extend it to a variety of video synthesis tasks by guiding the prediction with different annotations as illustrated in Figure 2. Code, pretrained models, and video samples synthesized by our approach are available at the url <https://16leomoing.github.io/ccvs>. Our main contributions are as follows:

1. an optical flow mechanism within an autoencoder to better reconstruct frames from context,
2. the use of ancillary information to control latent temporal dynamics when synthesizing videos,
3. a performance on par with or better than the state of the art, while being more memory-efficient.

2 Related Work

Video synthesis. In its simplest form, videos are produced without prior information about their content. GAN-based approaches map Gaussian noise into a visually plausible succession of frames. For example, VGAN [71] and ProgressiveVGAN [1], adapt the traditional GAN framework [26] from image to video synthesis by simply using 3D instead of 2D convolutions. These approaches, including recent attempts such as G³AN [76], are computationally expensive, and, by nature, restricted to synthesizing a fixed number of frames due to the constraints of their architecture. Other approaches predict latent motion vectors with a CNN [56], or a recurrent neural network (RNN) [12, 57, 65], and generate frames with individual 2D operations. To avoid the shortcomings of information loss in the sequential processing of RNNs, we use an attention-based autoregressive model instead. Since we forecast temporal dynamics in a compressed space, a large temporal window can be used when predicting new frames, without having to resort to expensive 3D computations. Previous works have attempted to scale video synthesis to higher resolution by using progressive training [1], subsampling [56, 57], reducing the dimension of the discriminator [36], or redefining the task as finding a trajectory in the latent space of a pre-trained image generator [63]. Compression, together with efficient context-aware reconstruction allows us to synthesize videos at high resolution.

Controllable video synthesis. Some of the aforementioned works [12, 56, 57, 76] handle synthesis conditioned on a class label. Another popular control is to use a few priming frames to set off the generation process. This task, known as *future video prediction*, has received a lot of attention recently. Methods based on variational autoencoders (VAE) [2, 15] have been proposed to account for the stochastic nature of future forecasting, *i.e.*, the plurality of possible continuations, disregarded in deterministic predictive models [21, 47, 70]. Yet, their blurry predictions have motivated the incorporation of adversarial training [44], hierarchical architectures [7, 80], fully latent temporal models [22], or normalizing flows [42]. Another line of work infers spatial transformation parameters (*e.g.*, optical flow), and predicts the future by warping past frames (*i.e.*, grid sampling and interpolation as in [33]) in RGB space [21, 24, 28, 72, 81] or in feature space [45], typically using a refinement step to handle occlusions. Lately, autoregressive methods [53, 79] leveraging a self-attention mechanism [68] have also been applied to this task. Our method benefits from the context-efficiency of approaches based on spatial transformation modules and the modeling power of autoregressive networks. In the meantime, other forms of control with interesting applications have emerged. Point-to-point generation [75], a variant of future video prediction, specifies both start and end frames. State or action-conditioned synthesis [28, 48] guides the frame-by-frame evolution with high-level commands.

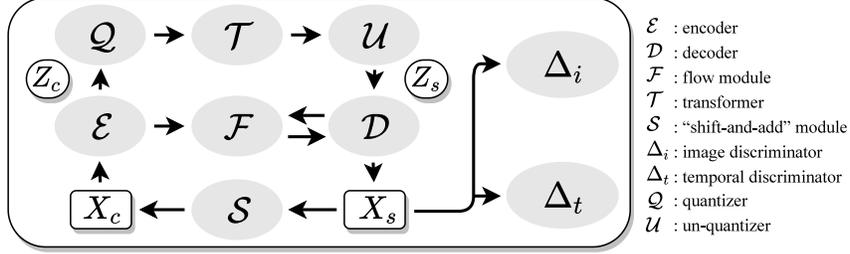


Figure 1: Proposed architecture. Here X_c and X_s respectively stand for the (input) context and (output) synthesized video. Learnable encoder and decoder modules \mathcal{E} and \mathcal{D} are linked by a learnable flow estimation module \mathcal{F} ensuring spatio-temporal consistency between context and synthesized frames. The architecture is doubly autoregressive, with a transformer \mathcal{T} responsible for predicting the features Z_s associated with future frames X_s from the features Z_c associated with context frames X_c , and a simple, parameterless, “shift-and-add” module \mathcal{S} updating X_c as each new frame is generated. The architecture is trained in two steps: The parameters of \mathcal{E} , \mathcal{D} and \mathcal{F} are first estimated (without any future forecasting from the transformer) in an adversarial manner using two discriminators Δ_i and Δ_t to ensure that the frames synthesized are both realistic (Δ_i) and temporally coherent (Δ_t). The two discriminators are then discarded, and the parameters of the transformer \mathcal{T} are estimated with \mathcal{E} , \mathcal{F} and \mathcal{D} frozen. At inference time, the transformer is used only once, latter frames being estimated in an autoregressive manner by the quadruple $\mathcal{D} \mathcal{S} \mathcal{E} \mathcal{F}$. See text for more details.

Additional works consider video synthesis based on one [51] or multiple layouts [46, 74], another video [8], or sound [10, 34, 73]. To account for the variety of potential controls, we leverage the flexibility of transformers, and propose an unifying approach to tackle all of these tasks.

3 Context-aware controllable video synthesis

3.1 Overview of the proposed approach

We consider video data, such that an individual frame x is natively an element of $\mathcal{X} = \mathbb{R}^{H \times W \times 3}$, which can be encoded as some feature z in $\mathcal{Z} = \mathbb{R}^{h \times w \times F}$ with reduced spatial resolution and increased number of channels. We assume that we are given X_c in \mathcal{X}^m corresponding to m successive frames, and our goal is to *synthesize* a plausible representation of the following n frames X_s which lies in \mathcal{X}^n . Given X_c , we can compute the corresponding features Z_c in \mathcal{Z}^m using some encoder $\mathcal{E} : \mathcal{X} \rightarrow \mathcal{Z}$ on each frame individually, then use an autoregressive model (a transformer coupled with a quantization step in our case) to predict features Z_s in \mathcal{Z}^n formed by the n features corresponding to time ticks $m + 1$ to $m + n$. (Note that the values of m and n can be arbitrary, using [by now] traditional masking and sliding window techniques.) These features can finally be converted one by one into the corresponding frames X_s using some decoder $\mathcal{D} : \mathcal{Z} \rightarrow \mathcal{X}$.

The overall approach is illustrated by Figure 1. The use of a quantizer over a learned codebook in our implementation complicates the architecture a bit, but has several advantages, including the reuse of familiar NLP technology [68] and, perhaps more importantly, affording simple mechanisms for handling different types of inputs (from images to sound for example [9, 31]) and the intrinsically uncertain and multimodal nature of future prediction (by sampling different codebook elements according to their likelihood [30]). Concretely, this choice simply amounts to inserting in our architecture, right after the encoder \mathcal{E} , a nearest-neighbor quantizer $\mathcal{Q} : \mathbb{R}^F \rightarrow \llbracket 1, q \rrbracket$ parameterized by an $F \times q$ codebook which, given the embedding z of a frame, returns a $h \times w$ matrix of corresponding *tokens*, that is, the indices of the closest entry in the codebook for feature vector $z_{i,j}$ for all spatial locations (i, j) in $\llbracket 1, h \rrbracket \times \llbracket 1, w \rrbracket$. We abuse notation, and identify \mathcal{Q} with its parameterization by this codebook, so we can optimize over \mathcal{Q} just as we optimize over \mathcal{E} instead of naming explicitly its parameters in the rest of this presentation. An “un-quantizer” $\mathcal{U} : \llbracket 1, q \rrbracket \rightarrow \mathbb{R}^F$, also parameterized (implicitly) by the codebook and associating with each token the corresponding entry of the codebook, is also inserted right before the decoder \mathcal{D} . In this setting, $\mathcal{T} : \llbracket 1, q \rrbracket^{m \times h \times w} \rightarrow \llbracket 1, q \rrbracket^{n \times h \times w}$ takes as input a sequence of tokens, and outputs the tokens for subsequent frames, each one of them chosen among q possibilities as either the one with the highest score, or drawn randomly

from the top- k scores to account for the multimodal nature of future forecasting (here $k \leq q$ is some predefined constant, see [30] for related approaches).

The elements of the architecture described so far are by now (individually) rather classical, with learnable, parametric functions \mathcal{E} , \mathcal{D} , \mathcal{Q} and \mathcal{T} . Besides putting them all together, and as detailed in the rest of this section, we add several original elements: **(a)** The encoding/decoding scheme is improved by the use of two discriminators, Δ_i and Δ_t respectively, trained in an adversarial manner to ensure that the predicted frames are realistic and temporally consistent. **(b)** The context frames X_c are themselves updated each time a new frame is predicted in an autoregressive manner (iteratively fill the sequence up to some predefined capacity, then shift to the left, forgetting the first frame and adding the latest synthetic one on the right). **(c)** The encoder and decoder are linked through a learnable *flow module* \mathcal{F} , allowing the context frames to guide the prediction of the synthetic ones. **(d)** Additional *control variables*, ranging from object trajectories to audio tracks, can be used in the form of sequence- or frame-level annotations to drive the synthesis by adding the corresponding tokens to the ones passed on to the autoregressive model \mathcal{T} .

We detail in this section the concrete components of the approach sketched above, including the autoencoder and quantizer architectures and their training procedure (Section 3.2), and the implementation of the autoregressive model by a transformer [68], illustrated in Figure 2, which we adapt to account for outside control signals (Section 3.3), once again with the corresponding procedure. Further architectural choices are also detailed in Appendix A.

3.2 First stage: training the context-aware autoencoder and the quantizer

Architecture. \mathcal{E} and \mathcal{D} respectively decreases and increases the spatial resolution by using convolutional residual blocks [29], with $(r_k)_{k \in [1:K]}$ the K corresponding resolution levels ($r_k = h_k \times w_k$). It is common practice to augment, as in U-Net [55], the autoencoder with long skip connections between \mathcal{E} and \mathcal{D} to share information across the two models at these intermediate levels and escape the lossy compression bottleneck. Although we cannot apply this directly to our setting since information only flows through \mathcal{D} for predicted timesteps, such skip connections can be established from the encoding stage of a context frame x_c to the decoding stage of a new frame x_s (resulting from features z_s). Similar mechanisms [15, 19] have been proposed in the past for video synthesis but they are only copying static background features from a single context frame. We follow works on semantic segmentation [23] and face frontalization [78] and use a flow module \mathcal{F} to warp features and produce temporally consistent outputs despite motion. We extend this to multi-frame contexts, with significant performance gains and no additional parameter to be learned.

Concretely, let e_c^k be features being encoded from x_c , and d_s^k features being decoded from z_s at a given intermediate resolution r_k . We first compute all intermediate context features e_c^k by applying \mathcal{E} to x_c . We then progressively decode features d_s^k for the new frame from low ($k = 1$) to high resolution ($k = K$) by iterating over the following steps: **(a)** apply one decoding sub-module to get d_s^k from d_s^{k-1} , **(b)** use \mathcal{F} to refine the optical flow f_c^k (in $\mathbb{R}^{2 \times r_k}$) which estimates the displacement field from e_c^k to d_s^k (as a proxy to the one from x_c to x_s), and a fusion mask m_c^k (in $\mathbb{R}^{1 \times r_k}$) which indicates the expected similarity between aligned features $e_c^k = W(e_c^k, f_c^k)$ and d_s^k (also as a proxy for the one in image domain) with W corresponding to a standard warping operation, **(c)** use e_c^k and m_c^k to update d_s^k with context information (see update rule (1) below), **(d)** move to resolution level r_{k+1} by going back to **(a)**. We note that \mathcal{F} estimates f_c^k and m_c^k in a coarse-to-fine fashion by refining f_c^{k-1} and m_c^{k-1} (see Appendix A for further details). Temporal skip connections at a given resolution level r_k are defined as the following in-place modification of d_s^k :

$$d_s^k = \sigma(m_c^k) \otimes d_s^k + (\mathbf{1} - \sigma(m_c^k)) \otimes e_c^k, \quad (1)$$

where σ is the Sigmoid function, and \otimes the element-wise product. We note that update rule (1) is quite standard for warping and fusing two streams of spatial information [23, 28, 74]. For concrete implementation of \mathcal{F} , we build upon LiteFlowNet [32], an optical flow estimation model which also combines pyramidal extraction and progressive warping of features. For simple integration into our framework, we use features from \mathcal{E} and \mathcal{D} both in the mask and flow estimation, and in the update (1). This process readily generalizes to multi-frame contextual information (see Appendix D for details). Similar to spatial transformers [33], the warping operation W is differentiable. As a result, gradients from the training losses can backpropagate from \mathcal{D} to \mathcal{E} through \mathcal{F} . This allows end-to-end training of the autoencoder even with information from different timesteps in \mathcal{E} and \mathcal{D} .

Training procedure. The global objective is the linear combination of four auxiliary ones:

$$\mathcal{L} = \lambda_q \mathcal{L}_q + \lambda_r \mathcal{L}_r + \lambda_a \mathcal{L}_a + \lambda_c \mathcal{L}_c, \quad (2)$$

namely a quantization loss (\mathcal{L}_q), a reconstruction loss (\mathcal{L}_r), an adversarial loss (\mathcal{L}_a), and a contextual loss (\mathcal{L}_c), detailed in the next paragraphs.

The codebook is trained by minimizing the reconstruction error between encoded features $z = \mathcal{E}(x)$ and quantized features $z_q = \mathcal{U}(\mathcal{Q}(z))$ (with notations introduced in Section 3.1):

$$\mathcal{L}_q(\mathcal{E}, \mathcal{Q}) = \|\text{sg}(z) - z_q\|_2^2 + \beta \|\text{sg}(z_q) - z\|_2^2, \quad (3)$$

where $\text{sg}(\cdot)$ is the stop gradient operation which constrains its operand to remain constant during backpropagation. The first part moves the codebook entries closer to the encoded features. The second part, known as the *commitment loss* [67], reverses the roles played by the two variables.

In regions with complex textures and high frequency details, local patterns shifted by a few pixels in x and its reconstruction $\hat{x} = \mathcal{D}(z_q)$ may result in large pixel-to-pixel errors, while being visually satisfactory. We thus define the recovery objective as the L_1 loss in both RGB space and between features from a VGG network [59] pretrained on ImageNet [14]:

$$\mathcal{L}_r(\mathcal{E}, \mathcal{Q}, \mathcal{F}, \mathcal{D}) = \|x - \hat{x}\|_1 + \|\text{VGG}(x) - \text{VGG}(\hat{x})\|_1. \quad (4)$$

To tackle cases where information cannot be recovered from context due to occlusion, and the compressed features are insufficient to create plausible reconstructions due to lossy compression, we supplement our architecture with an image discriminator Δ_i , made of downsampling residual blocks [29], to encourage realistic outputs. Δ_i tries to distinguish real images from reconstructed ones (\mathcal{L}_d), while \mathcal{E} and \mathcal{D} fools Δ_i into assuming reconstructed images are as good as real ones (\mathcal{L}_a):

$$\mathcal{L}_d(\Delta_i) = \ln(1 + e^{\Delta_i(x)}) + \ln(1 + e^{-\Delta_i(\hat{x})}), \quad (5) \quad \mathcal{L}_a(\mathcal{E}, \mathcal{Q}, \mathcal{F}, \mathcal{D}) = \ln(1 + e^{\Delta_i(\hat{x})}). \quad (6)$$

We employ a similar strategy on sequences of consecutive frames to improve the temporal consistency using a 3D temporal discriminator Δ_t , a direct extension of 2D image discriminator Δ_i .

The success of our method relies on accurate motion estimation in \mathcal{F} , a difficult task which benefits from self-supervision [35]. Therefore, we train the autoencoder with augmented views of the input frames as context. Custom augmentations functions $A : \mathcal{X} \rightarrow \mathcal{X}$ include: rotation, scaling, translation, elastic deformation, and combinations of these. Augmented views are obtained by warping x by the suitable flow a_c : $A(x) = W(x, a_c)$, and the inverted flow f_c from $A(x)$ to x can be approximated.¹ We resort to flow inversion because directly reconstructing distorted views may encourage similar defects during inference. Moreover, we balance between self-recovery and context-recovery objectives by additionally applying a blurring function $B : \mathcal{X} \rightarrow \mathcal{X}$ and an occlusion mask o_c to the augmented frames $x_c = o_c \otimes B(A(x))$. This augmentation strategy is illustrated in Appendix F. We define the contextual loss as:

$$\mathcal{L}_c(\mathcal{E}, \mathcal{Q}, \mathcal{F}, \mathcal{D}) = \|f_c - \hat{f}_c\|_2^2 + \|o'_c - \sigma(\hat{m}_c)\|_2^2, \quad (7)$$

where $o'_c = W(o_c, f_c)$, and \hat{f}_c and \hat{m}_c are the flow and mask estimated by \mathcal{F} . In practice, this loss is applied at intermediate resolution levels r_k for improved training.

3.3 Second stage: predicting temporal dynamics with transformers

Architecture. We follow [20] and adopt an architecture similar to Image-GPT [9] for \mathcal{T} . Instead of modeling a single annotated frame as in [54], we design our model to account for sequences of N such frames to allow prediction of temporal dynamics controlled by ancillary information in the form of video- and frame-level annotations. We have shown, in Section 3.2, how to represent a frame as a sequence of $h \times w$ tokens (indices in $\llbracket 1, q \rrbracket$) through encoding and quantization. Similar strategies can be applied to cater to other types of data, with or without compression depending on their complexity, and, thereby, turn ancillary information into tokens as well. The capacity of the transformer (the maximum sequence length it can process) is $L = l_v + N * (l_f + h * w)$, where l_v and l_f are the size of video- and frame-level annotations respectively. The final layer of the model predicts, for every i in $\llbracket 1, L - 1 \rrbracket$, a vector \hat{o}_i (of size q , the number of possible tokens) which scores the likelihood of the $i + 1^{\text{th}}$ token given all preceding ones.

¹Our implementation of flow inversion approximation is detailed in Appendix C.

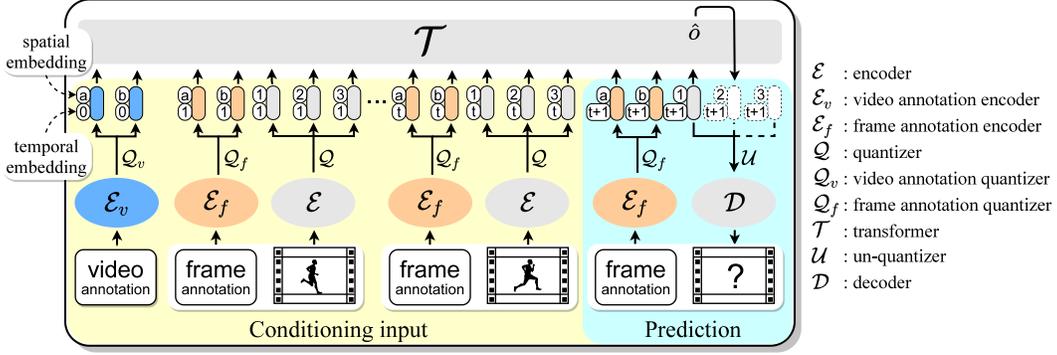


Figure 2: Illustration of the transformer to predict latent temporal dynamics with control-related data. For clarity, we omit the flow module \mathcal{F} between \mathcal{E} and \mathcal{D} . Input video-level annotation, frame-level annotations, and conditioning frames are encoded and quantized to form the initial token sequence, and mapped to corresponding embeddings. Subsequent frame tokens are obtained autoregressively with \mathcal{T} , while the ones corresponding to video- and frame-level annotations guide the prediction.

Training procedure. To learn the parameters of \mathcal{T} , we load a complete sequence (N frames corresponding to L tokens), and try to predict the $L - 1$ last tokens based on the $L - 1$ first ones. This is done by maximizing the log-likelihood of the data using the cross-entropy loss:

$$\mathcal{L}(\mathcal{T}) = - \sum_{i=2}^L \log \left(\frac{\exp(\langle \hat{\delta}_{i-1}, e_{\tau(i)} \rangle)}{\sum_j \exp(\langle \hat{\delta}_{i-1}, e_j \rangle)} \right), \quad (8)$$

where $\tau(i)$ is the i^{th} ground-truth token, and e_j a vector of 0's with a 1 in its j^{th} coordinate.

Inference. During inference we use \mathcal{T} to complete autoregressively an input sequence of tokens. To avoid known pitfalls (*e.g.*, repetitive synthesis) and allow diverse outcomes, we use top- k sampling whereby the next token is randomly chosen among the k most likely ones, weighted by their scores $\hat{\delta}$. Although \mathcal{T} processes sequences of N annotated frames, we predict ones of arbitrary length by using a temporal sliding window. Tokens corresponding to video- and frame-level annotations are given as input to \mathcal{T} (not predicted, even for future frames) to guide the synthesis. We show in the experiments how this control can translate into a variety of interesting tasks.

4 Experiments

We assess the merits of C CVS in the light of extensive experiments on various video synthesis tasks.

Datasets. BAIR Robot Pushing [19] consists of 43k training and 256 test videos of a robotic arm interacting with objects from a fixed viewpoint. A high-resolution version has recently been released. We manually annotate, in 500 frames, the (x, y) location of the arm in image space to train a position estimator, which we use for state-conditioned synthesis. To account for real world scenarios, we evaluate on Kinetics-600 [5], a large and diverse action-recognition dataset with approximately 500K videos. We also test our method on AudioSet-Drums [25] for sound-conditioned synthesis on music performance, containing 6k and 1k video clips in train and test splits respectively. Other datasets and tasks are covered in Appendix E.

Metrics. We use the Fréchet video distance (FVD) [66] which measures the distribution gap between real and synthetic videos in the feature space of an Inception3D network [6] pretrained on Kinetics-400 [40]. It estimates the visual quality and temporal consistency of samples as well as the diversity in unconditioned scenarios. In conditioned ones, we use another metric for diversity (DIV) which is the mean pixel-wise distance among synthetic trajectories conditioned on the same input. For near-deterministic motions (*e.g.*, in reconstructions, or constrained tasks), there is a one-to-one mapping between real video frames and synthetic ones, and we include pairwise image quality assessments: the structural similarity index measure (SSIM) [77] which evaluates a per-frame

Table 1: Ablation study of the autoencoder on BAIR (256×256). We evaluate self- and context-recovery modules in different scenarios: synthesizing 16-frame videos from known compressed features (“Reconstruction”), by inferring compressed features with \mathcal{T} given the real trajectory of the robotic arm (“State-conditioned”), or without the trajectory (“Pred.” and “Unc.”). The first real frame is used as initial context in all cases, except for “Unc.” where it is synthesized by StyleGAN2 [38].

Self-recovery				Ctxt.-recovery			Reconstruction		State-conditioned		Pred.	Unc.	
RGB	VGG	Δ_i	Δ_t	\mathcal{F}	Sup.	Ctxt.	FVD↓	PSNR↑	FVD↓	PSNR↑	FVD↓	FVD↓	
✓							0	1200 \pm 6	20.0	1238 \pm 24	17.8	1265 \pm 22	1321 \pm 10
✓	✓						0	700 \pm 11	17.9	714 \pm 5	17.4	704 \pm 7	765 \pm 12
✓	✓	✓					0	323 \pm 3	17.8	355 \pm 5	16.6	377 \pm 11	566 \pm 22
✓	✓	✓	✓				0	389 \pm 12	18.2	401 \pm 4	16.9	407 \pm 10	627 \pm 9
✓	✓	✓	✓	✓			1	98 \pm 2	22.7	106 \pm 2	22.2	142 \pm 6	350 \pm 11
✓	✓	✓	✓	✓	✓		1	87 \pm 3	24.4	97 \pm 4	22.1	128 \pm 4	301 \pm 10
✓	✓	✓	✓	✓	✓		8	62 \pm 1	25.4	76 \pm 3	22.7	109 \pm 6	299 \pm 4
✓	✓	✓	✓	✓	✓		15	60 \pm 1	25.6	75 \pm 2	22.8	110 \pm 3	297 \pm 7
<i>Training longer</i>				<i>(num. epochs \times 3)</i>			45\pm1	26.8	67\pm1	22.3	100\pm2	293\pm7	

“Sup.”: self-supervision of \mathcal{F} ; “Ctxt.”: number of context frames taken into account (in \mathcal{F}) when decoding current frame (in \mathcal{D}).

Table 2: Ablation study of the transformer on BAIR. We adopt notations and evaluation setups from Table 1.

Architecture		Top-k	State-conditioned		Pred.	Unc.		
Layer	Head	Dec.	Frame	State	FVD↓	PSNR↑	FVD↓	FVD↓
6	4		1	1	73 \pm 2	23.0	281 \pm 7	474 \pm 17
12	8		1	1	73 \pm 3	23.1	262 \pm 7	435 \pm 8
24	16		1	1	70 \pm 3	23.3	331 \pm 9	521 \pm 19
24	16	✓	1	1	69 \pm 2	23.2	321 \pm 9	479 \pm 34
24	16	✓	10	1	65\pm2	22.4	127 \pm 7	308 \pm 19
24	16	✓	100	1	67 \pm 1	22.3	121 \pm 2	314 \pm 12
24	16	✓	100	10	67 \pm 1	22.3	100\pm2	293\pm7

“Dec.”: spatio-temporal decomposition of positional embeddings.

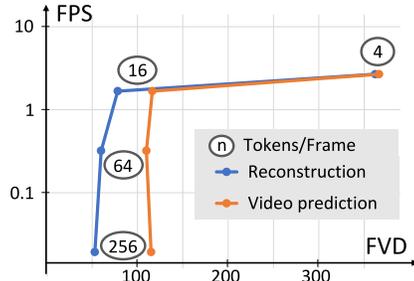


Figure 3: Quality and speed of synthesis vs. compression on a Nvidia V100 GPU.

conformity (combination of luminance, contrast and structure), and the peak signal-to-noise ratio (PSNR) which is directly related to the root mean squared error. For each metric, we compute the mean and standard deviation (std) over 5 evaluation runs (2 for Kinetics due to its voluminous 50k video test set). For clarity, the std value is shown only when it is greater than the reported precision.

Training details. All our models are trained on 4 Nvidia V100 GPUs (32GB VRAM each), using ADAM [41] optimizer, for multiple 20 hour runs. We adapt the batch size to fill the memory available. We use a learning rate of 0.02 to train the autoencoder, and exponential moving average [83] to obtain its final parameters. We use weighting factors (1, 10, 1, 1) and 0.25 for $(\lambda_q, \lambda_r, \lambda_a, \lambda_c)$ and β in Equations (2) and (3) respectively. We use a learning rate of 10^{-5} to train the transformer.

4.1 Ablation study

We conduct an ablation study of CCVS to show the individual contribution of the key components of the proposed autoencoder (Table 1), transformer (Table 2), and the effect of compression (Figure 3).

First, we fix the transformer and observe the incremental improvements in synthesis quality when adding self- and context-recovery modules to the autoencoder (Table 1). In particular, L_1 loss in the feature space of a VGG net [59] produces sharper videos than using the same loss in the RGB space alone. Predicted frames using image and temporal discriminators (Δ_i and Δ_t) display greater realism and temporal consistency. The flow module \mathcal{F} significantly improves the performance on all metrics by allowing context frames to guide the reconstruction of synthetic ones. The self-supervision of \mathcal{F} , the use of larger context windows, and longer training times, further improve the quality of the synthesis. Δ_t seems to deteriorate the FVD at first, but when all modules are combined it improves FVD by almost a factor 2 as it encourages better temporal consistency in the presence of context.

Table 3: Future video prediction on BAIR (64×64), synthesizing 16-frame videos given a few conditioning frames (“Cond.”). We include some extensions of our method at higher resolution.

Method	Cond.	FVD ↓	Code Avail.	Memory, compute	
MoCoGAN [65]	4	503	✓	16GB, 23h*	
SVG-FP [15]	2	315	✓	12GB, 6 to 24h*	
CDNA [21]	2	297	✓	10GB, 20h*	
SV2P [2]	2	263	✓	16GB, 24 to 48h*	
SRVP [22]	2	181	✓	36GB, 168h*	
VideoFlow [42]	3	131		128GB, 336h*	
LVT [53]	1	126±3	✓	128GB, 48h	
SAVP [44]	2	116	✓	32GB, 144h	
DVD-GAN-FP [12]	1	110		2TB, 24h*	
Video Transformer (S) [79]	1	106±3		256GB, 33h*	
TriVD-GAN-FP [45]	1	103		1TB, 280h*	
Low res. CCVS (ours)	1	99±2	✓	128GB, 40h	
Video Transformer (L) [79]	1	94±2		512GB, 336h*	
			SSIM ↑ (t = 8)	SSIM ↑ (t = 15)	
High res. ** CCVS (ours)	1	80±3	0.729	0.683	
+ end frame	2	81±2	0.766	0.839	
+ state	1	50±1	0.885	0.863	

*: value confirmed by authors;

** : training / inference / SSIM at 256×256 , FVD at 64×64 .

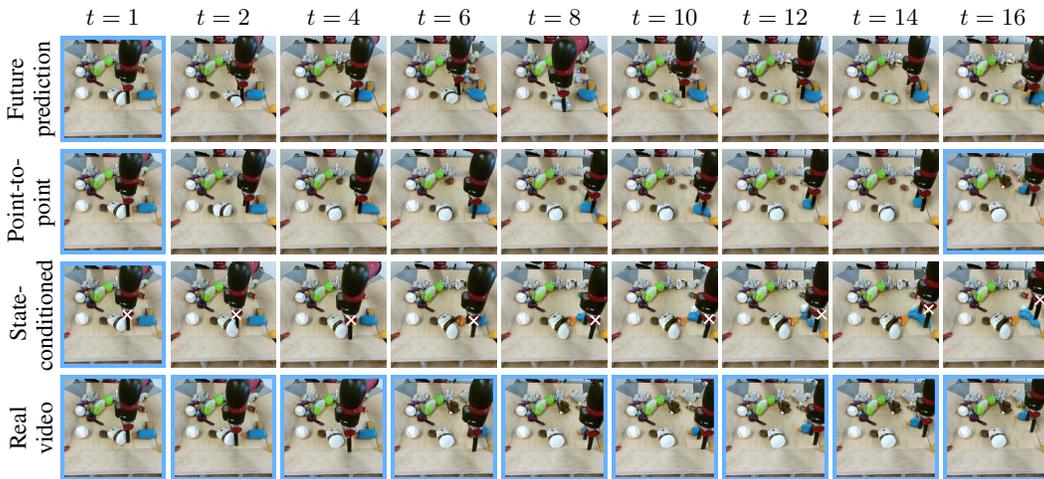


Figure 4: Qualitative samples for different types of control on BAIR (256×256). Zoom in for details.

We fix the autoencoder, and compare different architectures and sampling strategies for the transformer (Table 2). Increasing the model capacity by adding layers and increasing expressivity (number of attention heads) along with a spatio-temporal decomposition of positional embeddings (shown in Figure 2, detailed in supplementary material) yields small improvements. Top- k sampling is beneficial for stochastic tasks (video prediction and unconditional synthesis), whereas always selecting the most probable tokens results in little to no motion. Guiding temporal dynamics with state-conditioned synthesis reduces the advantage of sampling by narrowing down possible outcomes.

Finally, we explore the effect of compression (in terms of the number of tokens for each frame) on synthesis speed (FPS) and quality (FVD) of reconstructed and predicted videos (Figure 3). We use a compression of 64 tokens in our default setup since it gives the best FVD while retaining a reasonable speed. The critical drop of FPS for low compression ratios is due to the pairwise consideration of all input tokens in \mathcal{T} . Note that using a smaller temporal window in \mathcal{T} may allow additional speed-ups.

Table 4: Future video prediction on Kinetics (64×64) of 16-frame videos from 5 consecutive input frames.

Method	FVD ↓	GPU/TPU Mem.
LVT [53]	225	128GB, 48h
Video Transformer [79]	170 \pm 5	2TB, 336h*
DVD-GAN-FP [12]	69 \pm 1	2TB, 144h*
CCVS (<i>ours</i>)	55 \pm 1	128GB, 300h
TriVD-GAN-FP [45]	26 \pm 1	16TB, 160h*

*: value confirmed by authors.

Table 5: Codebook size on Kinetics.

Size	Reconstruction			Pred.
	FVD ↓	SSIM ↑	PSNR ↑	FVD ↓
1024	58 \pm 1	0.907	31.2	66 \pm 2
4096	54 \pm 1	0.917	31.6	64 \pm 1
16384	49 \pm 1	0.923	32.0	55 \pm 1
65536	45 \pm 1	0.928	32.2	61 \pm 1
∞	42 \pm 1	0.963	34.5	229 \pm 1

Table 6: Sound-conditional video synthesis on AudioSet-Drums (64×64).

Method	Cond. Audio	SSIM ↑			PSNR ↑			
		$t = 16$	$t = 30$	$t = 45$	$t = 16$	$t = 30$	$t = 45$	
SVG-LP [15]	15	0.971 \pm 0.017	0.661 \pm 0.010	0.510 \pm 0.008	30.0 \pm 1.1	16.6 \pm 0.3	13.5 \pm 0.1	
Vougioukas <i>et al.</i> [73]	15	✓	0.940 \pm 0.017	0.904 \pm 0.007	0.896 \pm 0.015	26.2 \pm 1.0	23.8 \pm 0.2	23.3 \pm 0.3
Sound2Sight [10]	15	✓	0.984 \pm 0.009	0.954 \pm 0.007	0.947 \pm 0.007	33.2 \pm 0.1	27.9 \pm 0.5	27.0 \pm 0.3
CCVS (<i>ours</i>)	15	✓	0.987 \pm 0.001	0.956 \pm 0.006	0.945 \pm 0.008	33.7 \pm 0.4	28.4 \pm 0.6	27.3 \pm 0.5

4.2 Quantitative and qualitative studies

BAIR. For future video prediction on BAIR (Table 3), CCVS trained at 64×64 resolution (*low res.*) is on par with the best method (L-size version of [79]), but requires much less computing resources, and outperforms [79] under similar resources. We also propose *high res.* CCVS which is not strictly comparable to the prior arts as we use 256×256 image resolution for training and test, and resize the synthesized frames to 64×64 for computing FVD. However, using this variant demonstrates the performance gains that can arise by scaling CCVS. We additionally address point-to-point synthesis (with the end frame as a video-level annotation) and state-conditioned synthesis (with the estimated 2D position of the arm as a frame-level annotation). Point-to-point synthesis is more difficult than video prediction: Not only does it require realistic video continuations, but also ones which explain the end position of all visible objects. Hence, FVD score is constant despite the additional input. Still, this yields better SSIM for mid-point and one-before-last frames. State-conditioned synthesis improves on FVD and mid-synthesis SSIM as motion becomes near-deterministic. Some synthetic frames for these tasks are shown in Figure 4. CCVS creates plausible high-quality videos in various settings, and true interactions with objects compared to previous attempts [48] at the same resolution.

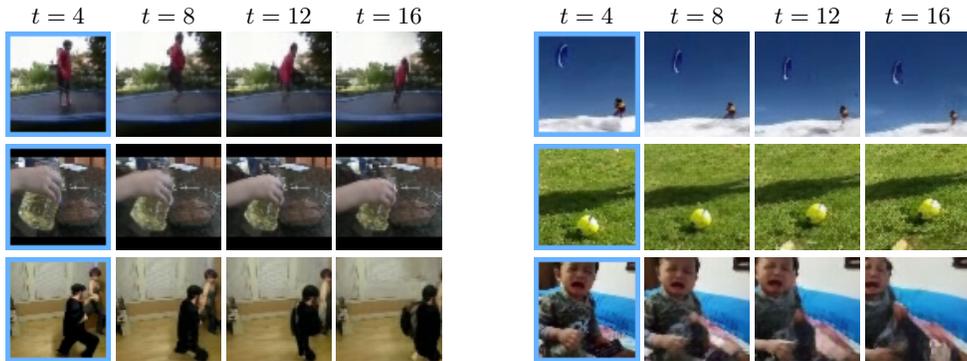
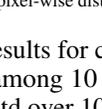
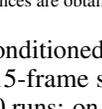
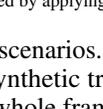
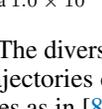
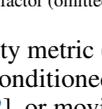
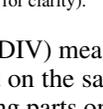
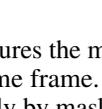
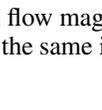
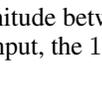
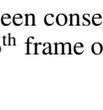
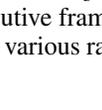
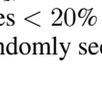
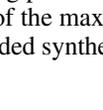
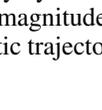


Figure 5: Qualitative samples on future video prediction on Kinetics (64×64). Zoom in for details.

Kinetics. CCVS ranks second on Kinetics video prediction benchmark. Kinetics contains more diversity than BAIR, and the reconstruction is thus more difficult. A solution [54] is to increase the codebook size (Table 5) but it stops translating into better prediction FVD at some point. We also try removing the quantization step (equivalent to an infinite codebook), and directly regressing latent features with \mathcal{T} (instead of ranking the likelihood of possible tokens). It allows better reconstructions, yet the prediction FVD is high. Figure 5 shows examples of synthetic continuations on

	DIV* \uparrow	Input $t = 1$	Seed 1 $t = 15$	Seed 2 $t = 15$	Seed 3 $t = 15$	Seed 4 $t = 15$	Seed 5 $t = 15$	Seed 6 $t = 15$
BAIR	Whole frames							
	19.53 \pm 2.60							
Kinetics	Moving parts							
	20.85 \pm 2.90							
AudioSet-Drums	Moving parts							
	48.16 \pm 13.83							
AudioSet-Drums	Whole frames							
	2.25 \pm 0.26							
	Moving parts							
	65.45 \pm 6.00							

*: the pixel-wise distances are obtained by applying a 1.0×10^{-3} factor (omitted for clarity).

Figure 6: Diversity results for conditioned scenarios. The diversity metric (DIV) measures the mean pixel-wise distance among 10 15-frame synthetic trajectories conditioned on the same frame. We report the mean and std over 100 runs: on whole frames as in [82], or moving parts only by masking static regions (optical flow magnitude between consecutive frames $< 20\%$ of the max magnitude) as in [69]. We show, for the same input, the 15th frame of various randomly seeded synthetic trajectories.

5-frame unseen test sequences. CCVS produces realistic and temporally coherent outputs which display various types of motion (*e.g.*, body, hand, camera).

AudioSet-Drums. CCVS achieves top performance on sound-conditioned video synthesis on AudioSet Drums (Table 6). Figure 6 shows quantitative and qualitative insights on the diversity of the synthetic trajectories conditioned on the same input for the three datasets. The diversity metrics (DIV) computed on whole frames is lower on AudioSet Drums than on the other two datasets. This is explained by the fact that motion is quite repetitive and involves a limited portion of the frame. The same metric on moving parts only, and the end position of the drummer’s hand and upper left cymbal in qualitative samples, show the diversity of synthetic trajectories. An ablation of CCVS with/without audio guidance as well as more qualitative results on diversity can be found in Appendix E.

5 Discussion

CCVS is on par or better than the state-of-the-art on standard benchmarks, uses less computational resources, and scales to high resolution. Training neural networks is environmentally costly, due to the carbon footprint to power processing hardware [17, 61]. Methods sparing GPU-hours like ours are crucial to make AI less polluting [43, 61], and move from a “Red” to a “Green” AI [58]. Future work will include exploring new codebook strategies and synthesis guided by textual information.

Limitations. CCVS uses a complex architecture and a two-stage training strategy. Simplification of both is an interesting direction for improving the method. Moreover, CCVS lacks global regularization of motion (flow computed on pairs of timesteps), and its efficiency relies on recycling context information such that synthesizing content from scratch (*i.e.*, no input frame given) remains difficult.

Broader impact. The increased accessibility and the many controls CCVS offers could accelerate the emergence of questionable applications, such as “deepfakes” (*e.g.*, a video created from someone’s picture and an arbitrary audio) which could lead to harassment, defamation, or dissemination of fake news. On top of current efforts to automate their detection [49], it remains our responsibility to grow awareness of these possible misuses. Despite these worrying aspects, our contribution has plenty of positive applications which outweigh the potential ethical harms. Our efficient compression scheme is a step in the direction of real-time solutions: *e.g.*, enhancing human-robot interactions, or improving the safety of self-driving cars by predicting the trajectories of people and vehicles nearby.

Acknowledgements

This work was granted access to the HPC resources of IDRIS under the allocation 2020-AD011012227 made by GENCI. It was funded in part by the French government under management of Agence Nationale de la Recherche as part of the “Investissements d’avenir” program, reference ANR-19-P3IA-0001 (PRAIRIE 3IA Institute). JP was supported in part by the Louis Vuitton/ENS chair in artificial intelligence and the Inria/NYU collaboration. We thank the reviewers for useful comments.

References

- [1] Dinesh Acharya, Zhiwu Huang, Danda Pani Paudel, and Luc Van Gool. Towards high resolution video generation with progressive growing of sliced wasserstein GANs. *arXiv preprint*, 2018.
- [2] Mohammad Babaeizadeh, Chelsea Finn, Dumitru Erhan, Roy H. Campbell, and Sergey Levine. Stochastic variational video prediction. In *ICLR*, 2018.
- [3] Yoshua Bengio, Nicholas Léonard, and Aaron C. Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint*, 2013.
- [4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *NeurIPS*, 2020.
- [5] Joao Carreira, Eric Noland, Andras Banki-Horvath, Chloe Hillier, and Andrew Zisserman. A short note about Kinetics-600. *arXiv preprint*, 2018.
- [6] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*, 2017.
- [7] Lluís Castrejon, Nicolas Ballas, and Aaron Courville. Improved conditional VRNNs for video prediction. In *ICCV*, 2019.
- [8] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A. Efros. Everybody dance now. In *ICCV*, 2019.
- [9] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pretraining from pixels. In *ICML*, 2020.
- [10] Anoop Cherian, Moitrey Chatterjee, and Narendra Ahuja. Sound2Sight: Generating visual dynamics from sound and context. In *ECCV*, 2020.
- [11] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, David Belanger, Lucy Colwell, and Adrian Weller. Rethinking attention with performers. In *ICLR*, 2021.
- [12] Aidan Clark, Jeff Donahue, and Karen Simonyan. Adversarial video generation on complex datasets. *arXiv preprint*, 2019.
- [13] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016.
- [14] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [15] Emily Denton and Rob Fergus. Stochastic video generation with a learned prior. In *ICML*, 2018.
- [16] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *ACL*, 2019.
- [17] Payal Dhar. The carbon impact of artificial intelligence. *Nature Machine Intelligence*, 2020.
- [18] 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.

- [19] Frederik Ebert, Chelsea Finn, Alex X. Lee, and Sergey Levine. Self-supervised visual planning with temporal skip connections. In *CoRL*, 2017.
- [20] Patrick Esser, Robin Rombach, and Björn Ommer. Taming transformers for high-resolution image synthesis. In *CVPR*, 2021.
- [21] Chelsea Finn, Ian Goodfellow, and Sergey Levine. Unsupervised learning for physical interaction through video prediction. In *NeurIPS*, 2016.
- [22] Jean-Yves Franceschi, Edouard Delasalles, Mickaël Chen, Sylvain Lamprier, and Patrick Gallinari. Stochastic latent residual video prediction. In *ICML*, 2020.
- [23] Raghudeep Gadde, Varun Jampani, and Peter V. Gehler. Semantic video cnns through representation warping. In *ICCV*, 2017.
- [24] Hang Gao, Huazhe Xu, Qi-Zhi Cai, Ruth Wang, Fisher Yu, and Trevor Darrell. Disentangling propagation and generation for video prediction. In *ICCV*, 2019.
- [25] Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. In *ICASSP*, 2017.
- [26] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, 2014.
- [27] Lena Gorelick, Moshe Blank, Eli Shechtman, Michal Irani, and Ronen Basri. Actions as space-time shapes. *TPAMI*, 2007.
- [28] Zekun Hao, Xun Huang, and Serge Belongie. Controllable video generation with sparse trajectories. In *CVPR*, 2018.
- [29] Kaifeng He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [30] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *ICLR*, 2020.
- [31] Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Ian Simon, Curtis Hawthorne, Noam Shazeer, Andrew M. Dai, Matthew D. Hoffman, Monica Dinulescu, and Douglas Eck. Music transformer. In *ICLR*, 2019.
- [32] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. LiteFlowNet: A lightweight convolutional neural network for optical flow estimation. In *CVPR*, 2018.
- [33] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and Koray Kavukcuoglu. Spatial transformer networks. In *NeurIPS*, 2015.
- [34] Amir Jamaludin, Joon Son Chung, and Andrew Zisserman. You said that?: Synthesising talking faces from audio. *IJCV*, 2019.
- [35] Rico Jonschkowski, Austin Stone, Jonathan T. Barron, Ariel Gordon, Kurt Konolige, and Anelia Angelova. What matters in unsupervised optical flow. In *ECCV*, 2020.
- [36] Emmanuel Kahembwe and Subramanian Ramamoorthy. Lower dimensional kernels for video discriminators. *Neural Networks*, 2020.
- [37] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *CVPR*, 2014.
- [38] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *CVPR*, 2020.
- [39] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are RNNs: Fast autoregressive transformers with linear attention. In *ICML*, 2020.
- [40] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The kinetics human action video dataset. *arXiv preprint*, 2017.
- [41] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.

- [42] Manoj Kumar, Mohammad Babaeizadeh, Dumitru Erhan, Chelsea Finn, Sergey Levine, Laurent Dinh, and Durk Kingma. VideoFlow: A conditional flow-based model for stochastic video generation. In *ICLR*, 2020.
- [43] Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. Quantifying the carbon emissions of machine learning. In *NeurIPS workshop*, 2019.
- [44] Alex X. Lee, Richard Zhang, Frederik Ebert, Pieter Abbeel, Chelsea Finn, and Sergey Levine. Stochastic adversarial video prediction. *arXiv preprint*, 2018.
- [45] Pauline Luc, Aidan Clark, Sander Dieleman, Diego de Las Casas, Yotam Doron, Albin Cassirer, and Karen Simonyan. Transformation-based adversarial video prediction on large-scale data. *arXiv preprint*, 2020.
- [46] Arun Mallya, Ting-Chun Wang, Karan Sapra, and Ming-Yu Liu. World-consistent video-to-video synthesis. In *ECCV*, 2020.
- [47] Michael Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond mean square error. In *ICLR*, 2016.
- [48] Willi Menapace, Stéphane Lathuilière, Sergey Tulyakov, Aliaksandr Siarohin, and Elisa Ricci. Playable video generation. In *CVPR*, 2021.
- [49] Thanh Thi Nguyen, Quoc Viet Hung Nguyen, Cuong M. Nguyen, Dung Nguyen, Duc Thanh Nguyen, and Saeid Nahavandi. Deep learning for deepfakes creation and detection: A survey. *arXiv preprint*, 2019.
- [50] Anselm Levskaya Nikita Kitaev, Lukasz Kaiser. Reformer: The efficient transformer. In *ICLR*, 2020.
- [51] Junting Pan, Chengyu Wang, Xu Jia, Jing Shao, Lu Sheng, Junjie Yan, and Xiaogang Wang. Video generation from single semantic label map. In *CVPR*, 2019.
- [52] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *CVPR*, 2019.
- [53] Ruslan Rakhimov, Denis Volkhonskiy, Alexey Artemov, Denis Zorin, and Evgeny Burnaev. Latent video transformer. *arXiv preprint*, 2020.
- [54] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML*, 2021.
- [55] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015.
- [56] Masaki Saito, Eiichi Matsumoto, and Shunta Saito. Temporal generative adversarial nets with singular value clipping. In *ICCV*, 2017.
- [57] Masaki Saito, Shunta Saito, Masanori Koyama, and Sosuke Kobayashi. Train sparsely, generate densely: Memory-efficient unsupervised training of high-resolution temporal GAN. *IJCV*, 2020.
- [58] Roy Schwartz, Jesse Dodge, Noah A. Smith, and Oren Etzioni. Green AI. *ACM*, 2020.
- [59] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [60] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint*, 2012.
- [61] Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and policy considerations for deep learning in NLP. In *ACL*, 2019.
- [62] Andrew Tao, Karan Sapra, and Bryan Catanzaro. Hierarchical multi-scale attention for semantic segmentation. *arXiv preprint*, 2020.
- [63] Yu Tian, Jian Ren, Menglei Chai, Kyle Olszewski, Xi Peng, Dimitris N. Metaxas, and Sergey Tulyakov. A good image generator is what you need for high-resolution video synthesis. In *ICLR*, 2021.
- [64] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, , and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *ICCV*, 2015.
- [65] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. MoCoGAN: Decomposing motion and content for video generation. In *CVPR*, 2018.

- [66] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges. *arXiv preprint*, 2018.
- [67] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. In *NeurIPS*, 2017.
- [68] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- [69] Ruben Villegas, Jimei Yang, Seunghoon Hong, Xunyu Lin, and Honglak Lee. Decomposing motion and content for natural video sequence prediction. In *ICLR*, 2017.
- [70] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Anticipating visual representations from unlabeled video. In *CVPR*, 2016.
- [71] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. In *NeurIPS*, 2016.
- [72] Carl Vondrick and Antonio Torralba. Generating the future with adversarial transformers. In *CVPR*, 2017.
- [73] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. End-to-end speech-driven facial animation with temporal GANs. In *BMVC*, 2018.
- [74] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. In *NeurIPS*, 2018.
- [75] Tsun-Hsuan Wang, Yen-Chi Cheng, Chieh Hubert Lin, Hwann-Tzong Chen, and Min Sun. Point-to-point video generation. In *ICCV*, 2019.
- [76] Yaohui Wang, Piotr Bilinski, Francois Bremond, and Antitza Dantcheva. G3AN: Disentangling appearance and motion for video generation. In *CVPR*, 2020.
- [77] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *TIP*, 2004.
- [78] Yuxiang Wei, Ming Liu, Haolin Wang, Ruifeng Zhu, Guosheng Hu, and Wangmeng Zuo. Learning flow-based feature warping for face frontalization with illumination inconsistent supervision. In *ECCV*, 2020.
- [79] Dirk Weissenborn, Oscar Täckström, and Jakob Uszkoreit. Scaling autoregressive video models. In *ICLR*, 2020.
- [80] Bohan Wu, Suraj Nair, Roberto Martin-Martin, Li Fei-Fei, and Chelsea Finn. Greedy hierarchical variational autoencoders for large-scale video prediction. In *CVPR*, 2021.
- [81] Yue Wu, Rongrong Gao, Jaesik Park, and Qifeng Chen. Future video synthesis with object motion prediction. In *CVPR*, 2020.
- [82] Dingdong Yang, Seunghoon Hong, Yunseok Jang, Tianchen Zhao, and Honglak Lee. Diversity-sensitive conditional generative adversarial networks. In *ICLR*, 2019.
- [83] Yasin Yaz, Chuan-Sheng Foo, Stefan Winkler, Kim-Hui Yap, Georgios Piliouras, and Vijay Chandrasekhar. The unusual effectiveness of averaging in GAN training. In *ICLR*, 2019.