

Stochastic Image-to-Video Synthesis using cINNs

Michael Dorckenwald¹ Timo Milbich¹ Andreas Blattmann¹ Robin Rombach¹
Konstantinos G. Derpanis^{2,3,4*} Björn Ommer^{1*}

¹IWR/HCI, Heidelberg University, Germany ²Department of Computer Science, Ryerson University, Canada
³Vector Institute for AI, Canada ⁴Samsung AI Centre Toronto, Canada

Abstract

Video understanding calls for a model to learn the characteristic interplay between static scene content and its dynamics: Given an image, the model must be able to predict a future progression of the portrayed scene and, conversely, a video should be explained in terms of its static image content and all the remaining characteristics not present in the initial frame. This naturally suggests a bijective mapping between the video domain and the static content as well as residual information. In contrast to common stochastic image-to-video synthesis, such a model does not merely generate arbitrary videos progressing the initial image. Given this image, it rather provides a one-to-one mapping between the residual vectors and the video with stochastic outcomes when sampling. The approach is naturally implemented using a conditional invertible neural network (cINN) that can explain videos by independently modelling static and other video characteristics, thus laying the basis for controlled video synthesis. Experiments on four diverse video datasets demonstrate the effectiveness of our approach in terms of both the quality and diversity of the synthesized results. Our project page is available at <https://bit.ly/3dg90fV>.

1. Introduction

Anticipating and predicting what happens next are key features of human intelligence that allow us to understand and deal with the ever-changing environment that governs our everyday life [10]. Consequently, the ability to foresee and hallucinate the future progression of a scene is a cornerstone of artificial visual understanding with applications including autonomous driving [51, 52, 27], medical treatment [7, 18, 8], and robotic planning [20, 24, 12].

Predicting and synthesizing plausible future progressions from a given image requires a deep understanding of how scenes and objects within video are depicted, interplay with each other, and evolve over time. While an

*Indicates equal supervision.

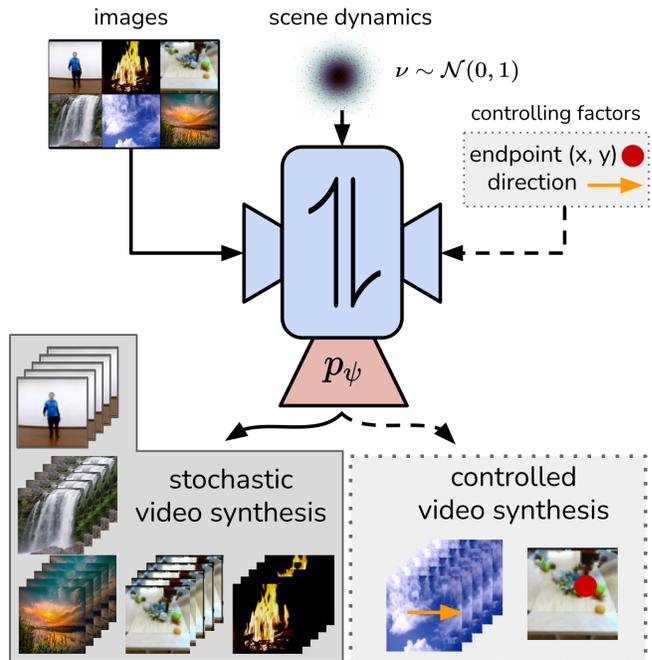


Figure 1. Our approach establishes a bijective mapping between the image and the video domain by introducing a residual representation ν describing the latent scene dynamics. This allows us not only to synthesize diverse videos but also to extend our approach to gain control over the video synthesis task.

image provides information about the observed scene content, such as object appearance and shape, the challenge is to understand the missing information constituting potential futures, such as the scene dynamics setting the scene in motion. Due to the ambiguity and complexity of capturing this information, many works [45, 25, 14, 81] directly focus on predicting likely video continuations, often resorting to simplifying assumptions (e.g., dynamics modelled by optical flow [21, 64]) and side information (e.g., semantic keypoints [61, 78, 49, 22, 5]). However, truly understanding the synthesis problem not only requires to infer such image continuations but, conversely, also demands when observing a video sequence to describe and represent the instantiated scene dynamics animating its initial frame.

Consequently, the image-to-video synthesis task should be modelled as a translation between the image and video domains, ideally by an invertible mapping between them. Since the content information describing an image only accounts for a small fraction of the video information, in particular missing the temporal dimension, learning an invertible mapping requires a dedicated residual representation that captures all missing information. Once learned, given an initial image and an instantiation of the latent residual, we can combine them to synthesize the corresponding future video sequence.

In this paper, we frame image-to-video synthesis as an invertible domain transfer problem and implement it using a conditional invertible neural network (cINN) illustrated in Fig. 1. To account for the domain gap between images and videos, we introduce a dedicated probabilistic residual representation. The bijective nature of our mapping ensures that only information complementary to that in the initial image is captured. Using a probabilistic formulation, the residual representation allows to sample and thus synthesize novel future progressions in video with the same start frame. To reduce the complexity of the learning task, we train a separate conditional variational encoder-decoder architecture to compute a compact, information preserving representation for the video domain. Moreover, our specific framing of learning the residual representation allows to easily incorporating extra conditioning information to exercise control over the image-to-video synthesis process.

Our contributions can be summarized as follows:

- We frame image-to-video synthesis as an invertible domain transfer problem and learn a dedicated residual representation to capture the domain gap.
- Our framework naturally extends to incorporate explicit conditioning factors for exercising control over the synthesis process.
- Extensive evaluations on four diverse video datasets, ranging from structured human motion synthesis to subtle dynamic textures, show strong results demonstrating the effectiveness of our approach.

2. Related Work

Video synthesis. Video synthesis involves a wide range of tasks including video-to-video translation [80], image animation [70, 71], frame interpolation [57, 4], and video prediction. The latter can be divided into unconditional [75, 16] and conditional video generation (the focus of our work). Conditional video generation can be described as finding a future progression given a set of context frames in a deterministic [77, 82, 54, 6] or stochastic manner [45, 25, 14, 3], as pursued here. Several works decrease

the complexity of the synthesis task by using keypoint annotations [55, 61] as conditioning information. A major drawback of this approach is the requirement of semantic keypoint labels which limit consideration to highly structured objects, like humans, and thus exclude the broader range of imagery we consider, e.g., natural scenes. Recent methods aim at improving video prediction quality by use of high capacity architectures with high computational demands, operating in the latent [63] or pixel-space [81], or using attention [16]. In contrast, we propose a model for understanding the image-to-video synthesis process by learning a bijective transformation between the image and video domains using a dedicated residual representation.

Dynamic texture synthesis. Previous work has given special attention to generating dynamic textures. This work can be divided into two groups: (i) methods that exploit the statistics of dynamics textures [74, 84, 86] and (ii) learning-based approaches [87, 90, 48, 21, 85]. To generalize to other video domains, beyond dynamic textures, we introduce a learning-based approach. MDGAN [87] generates landscape videos from a static scene in a deterministic manner. Several methods (e.g., [21, 90]) consider optical flow in their video generation pipeline. The use of optical flow limits application to specific types of imagery, like clouds, at the exclusion of other dynamic textures which grossly violate standard optical flow assumptions [74]. DeepLandscape [48] extends the structure of StyleGAN [37] to animate landscape images. Their model does not attempt to learn full temporal dynamics of videos and works only by a complex optimization scheme for inference, similar to [26] for style transfer. In contrast, our approach allows for efficient feedforward image-to-video synthesis while also maintaining visual quality and temporal coherence.

Invertible Neural Networks. Invertible neural networks (INNs) are bijective functions which makes them attractive for a variety of tasks, such as analyzing inverse problems [1], interpreting neural networks [23], and representation learning [35]. In particular, INNs can be implemented as normalizing flows [65], a special class of likelihood-based generative models which have recently been applied to various tasks, such as image synthesis [40, 2, 62], domain transfer [68, 67, 23, 89], superresolution [50, 89], and video synthesis [43]. In contrast, we use a conditional normalizing flow model to learn a dedicated residual latent, capturing information not contained in the input image. This allows us to both more efficiently learn the bijective mapping and to consider explicit controlling factors.

3. Method

Our goal is to learn the interplay between images and video by explaining video in terms of a single image and the (stochastic) information not captured by the image about the video. Together the deterministic and stochastic content al-

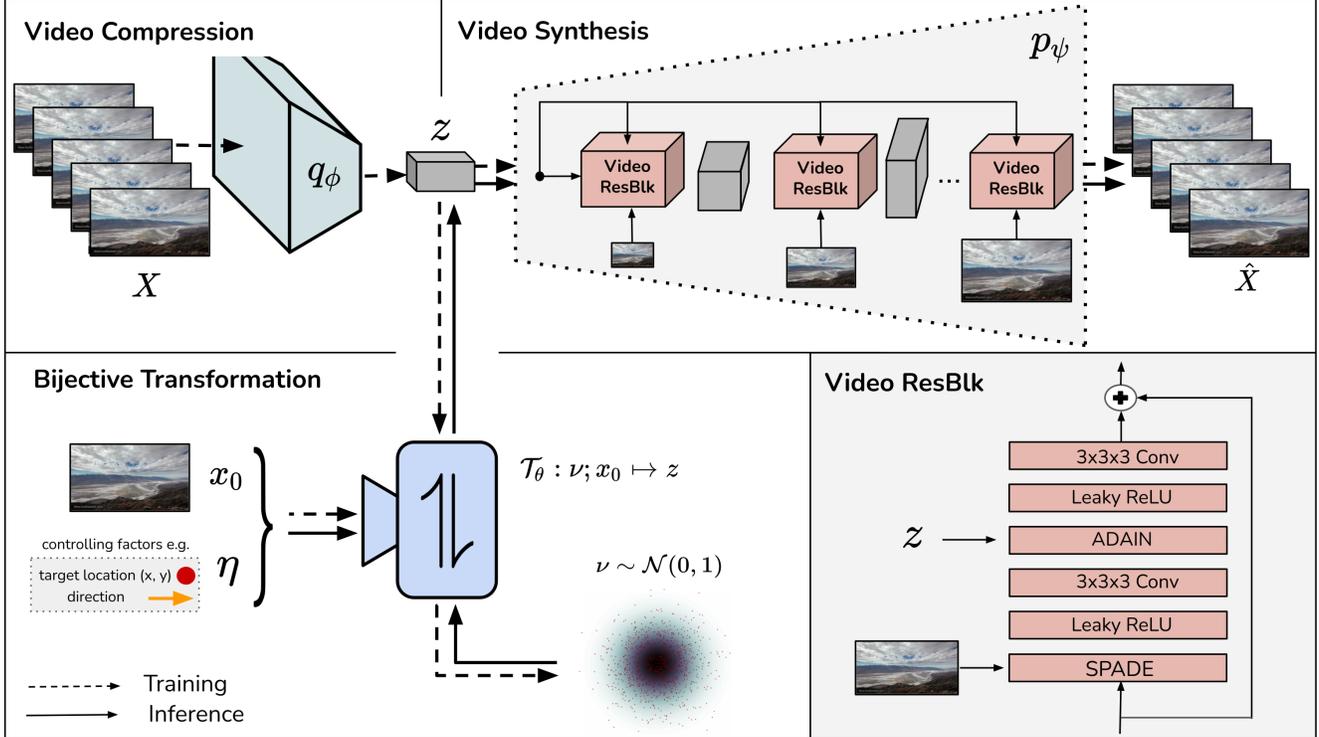


Figure 2. Overview of our proposed framework. We learn an information preserving video representation z using our conditional generative model consisting of an encoder q_ϕ as well as the corresponding decoder p_ψ . The decoder consists of dedicated video residual blocks shown in bottom right. After establishing the video representation, we learn a bijective transformation \mathcal{T} conditioned on the starting frame x_0 and an optionally provided control factor η . During inference, we sample a residual ν , encapsulating the scene dynamics, from the prior distribution and use \mathcal{T}_ϕ to obtain the video representation z . Using our decoder we can then synthesize novel video sequences. Training and inference is indicated by the dotted and solid lines, respectively.

low us to tackle the problem of image-to-video synthesis. In Sec. 3.1, we begin by motivating and introducing our conditional bijective framework for image-to-video mappings and Sec. 3.2 describes the learning process. Sec. 3.3 presents our generative model for video synthesis operating on our learned transformation. Finally, in Sec. 3.4 we extend our model to directly exercise control over factors captured in the residual latent, e.g., direction of motion. Fig. 2 provides an overview of our approach.

3.1. Bijection for Image-to-Video Synthesis

Given an initial image, $x_0 \in \mathbb{R}^{d_x}$, image-to-video synthesis generates a video sequence, $X = [x_1, \dots, x_T]$. This problem is inherently underdetermined with many possible videos conceivable based on x_0 . As a result, we cannot synthesize or explain a video merely with a single frame, but require additional information, ν , such as the scene dynamics. Video synthesis can then be framed as mapping x_0 and a residual ν onto a video X or, equivalently, a representation z thereof,

$$z = \mathcal{T}(\nu; x_0). \quad (1)$$

Commonly, stochastic video prediction methods [45, 25, 55] only focus on synthesizing *arbitrary realistic* videos for

a single initial or a sequence of frames. In contrast, understanding this synthesis process not only demands to explain the missing information ν to be inferred, but also to recover the residual information from video so that it can be modified subsequently. Explaining a video thus requires to estimate this residual information ν , so that x_0 and ν together are isomorphic to the representation z of the video X . Consequently, \mathcal{T} needs to be a conditional bijective mapping between videos and their description in terms of a starting frame x_0 and the remaining residual information ν .

3.2. Inferring an Explicit Residual Representation

Given a single frame x_0 , a multitude of videos are possible with a corresponding z ,

$$z \sim p(z|x_0). \quad (2)$$

Since ν contains all the information of z not captured in x_0 and \mathcal{T} is conditionally bijective, we can invert (1) to obtain the residual

$$\nu = \mathcal{T}^{-1}(z; x_0). \quad (3)$$

Then, by the change-of-variables theorem for probability distributions, \mathcal{T}^{-1} transforms $p(z|x_0)$ as

$$\begin{aligned}
p(z|x_0) &= \frac{p(\nu|x_0)}{|\det J_{\mathcal{T}}(\nu; x_0)|} & (4) \\
&= p(\mathcal{T}^{-1}(z; x_0)|x_0) \cdot |\det J_{\mathcal{T}^{-1}}(z; x_0)|, & (5)
\end{aligned}$$

where $J_{\mathcal{T}}$ denotes the Jacobian of the transformation \mathcal{T} and $|\det[\cdot]|$ the absolute value of the determinant of its input. Using the transformed distribution, $p(z|x_0)$, we can now directly learn our transformation \mathcal{T} and the distribution $p(\nu|x_0)$ by maximum likelihood estimation (MLE). To this end, we need to choose an appropriate prior distribution, which can be analytically evaluated and easily sampled. Since we factorize the residual information ν from the starting frame x_0 , we can assume $p(\nu|x_0) = q(\nu)$ and, thus, resort to the widely used standard normal distribution $q(\nu) = \mathcal{N}(\nu|0, \mathbf{1})$ [41, 88, 28]. Moreover, we parametrize \mathcal{T} as an invertible neural network [58, 17, 42] \mathcal{T}_{θ} with parameters θ which, given the image x_0 , translates between the representations z and ν . Thus, we arrive at the negative log-likelihood minimization problem

$$\min_{\theta \in \Theta} \mathbb{E}_{z, x_0} \left[\log q(\mathcal{T}_{\theta}^{-1}(z; x_0)) - \log |\det J_{\mathcal{T}_{\theta}^{-1}}(z; x_0)| \right]. \quad (6)$$

By simplifying using the standard normal prior and dropping resulting constant terms, we finally arrive at our final objective function

$$\min_{\theta \in \Theta} \mathbb{E}_{z, x_0} \left[\|\mathcal{T}_{\theta}^{-1}(z; x_0)\|_2^2 - \log |\det J_{\mathcal{T}_{\theta}^{-1}}(z; x_0)| \right]. \quad (7)$$

Due to the information-preserving, isomorphic mapping \mathcal{T}_{θ} , ν indeed captures the latent information in X not explained by x_0 .

To generate a video representation z based on an initial frame x_0 , we first sample a residual representation $\nu \sim q(\nu)$ and then apply (1) to obtain $z = \mathcal{T}_{\theta}(\nu; x_0)$.

3.3. Generative Model for Video Synthesis

We now learn a decoding $p(X|z)$ to synthesize video sequences based on z . Since we require z to be a compact, information-preserving video representation, we also need to learn the corresponding encoding $q(z|X)$. Simultaneously learning both is naturally expressed by an autoencoder [41]. Moreover, to optimally enable learning the transformation \mathcal{T}_{θ} , we consider the following modelling constraints: (i) the representation z of the input should be maximally information-preserving to fully capture the residual dynamics information, (ii) we model the residual ν to be a continuous probabilistic model, thus the bijection property of \mathcal{T}_{θ} requires $q(z|X)$ to be a strictly positive density, and (iii) reducing the complexity of the representation z eases the task of learning the bijective mapping \mathcal{T}_{θ} . Thus, while still fully capturing scene dynamics in z , we ideally

exclude all information in the video which is already present in the initial image x_0 .

Learning $p(X|z)$ and $q(z|X)$. Variational latent models [41] are a straightforward choice for stochastic autoencoders. To address (iii) above, we use a conditional variational autoencoder [88] with a parametrized encoder $q_{\phi}(z|X)$ and a parametrized, conditional decoder $p_{\psi}(X|x_0, z)$ with (ϕ, ψ) being their trainable parameters. Such models encourage the distribution of information among latent variables due to the regularization of the capacity of the latent encoding [15, 91, 11]. Thus, using x_0 as a conditioning to represent most of the scene content, the complexity of z can be reduced by forcing the network capacity to focus on capturing the latent information in X . To balance this with maximally preserving the latent residual information in X , we introduce a weighting parameter β to the standard variational lower bound [11],

$$\begin{aligned}
\mathcal{L}_{p_{\psi}, q_{\phi}} &= \mathbb{E}_{z \sim q_{\phi}(z|X)} [\log p_{\psi}(X|x_0, z)] \\
&\quad - \beta D_{\text{KL}}(q_{\phi}(z|X) || q(z)), \quad (8)
\end{aligned}$$

where $q(z)$ denotes a standard normal prior on the encoder q_{ϕ} . The first term optimizes the synthesis quality of the decoding process, thus maximizing information-preservation. While the second term regularizes $q_{\phi}(z|X)$ to match the prior $q(z)$ which constrains its capacity and, thus, encourages the distribution of information among x_0 and z to ease subsequent learning of \mathcal{T}_{θ} . Hence, β allows us to directly balance the informativeness of z and its complexity [15, 91, 11].

Building the video synthesis model. The design of generative architectures significantly influences their synthesis capabilities, especially when dealing with highly complex data. In our conditional model this particularly affects the interplay between information in x_0 and z in $p_{\psi}(X|z, x_0)$. To this end, we construct the conditional decoder p_{ψ} using a sequence of n dedicated video residual blocks operating on increasing spatial and temporal feature resolutions. To optimally facilitate the interplay between z and the content information in x_0 , we combine them both in each block and, thus, at all scales of p_{ψ} . Fig. 2 illustrates the general structure of our video residual blocks used for decoding to a video. The conditioning x_0 is incorporated using a SPADE [59] normalization layer to preserve semantic information throughout the generator. The video representation z is added by means of an ADAIN [37] layer to provide video information at all scales of the decoder. Our encoder q_{ϕ} is implemented as a 3D-ResNet [31] to capture the scene dynamics evolving over time in an input video.

Overall training objective. Following common practice [41], we train our conditional model, (8), using an L_1 reconstruction loss. To emphasize perceptual quality [44] we use a frame-wise perceptual loss ℓ^{ϕ} [19, 36]. Similar to previous work [16, 80], we use a discriminator \mathcal{D}_S applied

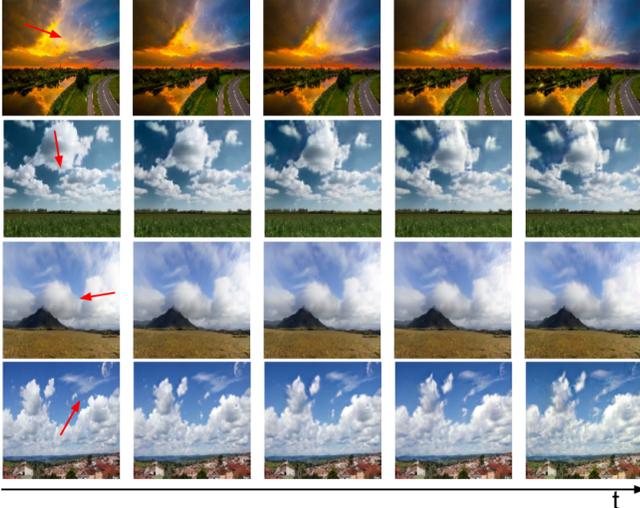


Figure 3. Stochastic video synthesis on Landscape [87] showing subtle motions. Red arrows indicate the direction of motion. Best viewed as video provided in the supplemental.

Method	LPIPS ↓	FID ↓	DTFVD ↓	FVD ↓	DIV VGG ↑	DIV I3D ↑
MDGAN ² [87]	0.49	68.9	2.35	385.1	–	–
DTVNet ² [90]	0.35	74.5	2.78	693.4	0.00	0.00
DL ^{2,†} [48]	0.41	41.1	1.73	351.5	–	–
AL ² [21]	0.26	16.4	1.24	307.0	0.97	0.71
Ours	0.23	10.5	0.59	134.4	0.71	1.22

Table 1. Quantitative evaluation of video synthesis quality and diversity on Landscape [87]. Numeric superscripts indicate the source of the results, cf. Sec. 4.3. The diversity score based on the I3D [73] trained on DTDB [30] can be found in the supplemental. [†] provided pretrained model from DL [48] was trained on their unreleased dataset.

to each frame and \mathcal{D}_T on the temporal level. Both discriminators are optimized using the hinge formulation [46, 9]. Thus, the overall training objective can be summarized as

$$\mathcal{L} = \mathcal{L}_{p_\psi, q_\phi} + \mathcal{L}_{\mathcal{D}_T} + \mathcal{L}_{\mathcal{D}_S}. \quad (9)$$

Please see the supplemental for further details of our loss.

3.4. Controllable Video Synthesis

There are many factors comprising the latent residual ν . Understanding the image-to-video process allows us to directly exercise control over such factors and thus over the progression of the depicted scene in the input image x_0 . Assuming $\eta \in \mathbb{R}^{d_\eta}$ represents such a factor, e.g., the target location of a moving object, we can explicitly model it while learning our bijective mapping \mathcal{T}_θ as $\mathcal{T}_\theta(\nu; x_0, \eta)$. Note, now ν constitutes the residual latent information to *both* x_0 and η . Since such individual factors are typically low in information themselves, in general there is no benefit in considering them when learning the conditional decoder p_ψ in contrast to the richer information in x_0 . Image-to-video synthesis now extends to additionally manually adjusting η to a fixed value η^* to infer a video representation

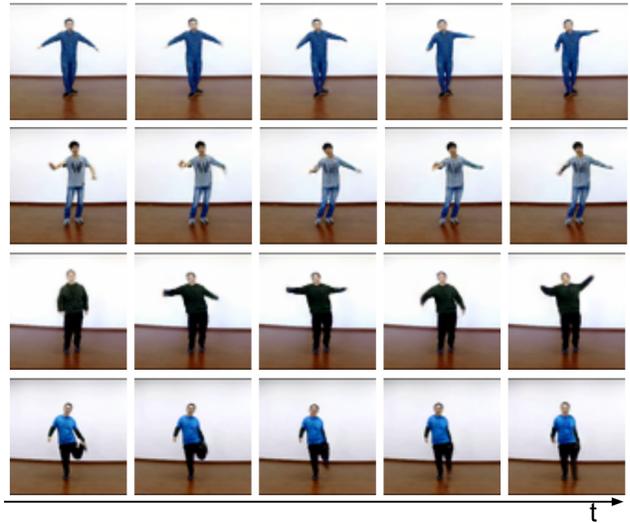


Figure 4. Stochastic video synthesis on iPER [47] showing structured, diverse human motion. Best viewed as video provided in the supplemental.

Method	FVD ↓	DIV VGG ↑	DIV I3D ↑
SAVP ³ [45]	368.6	0.00*	0.01*
SRVP ³ [25]	336.3	0.34	1.01
IVRNN ³ [14]	191.4	0.23	0.57
Ours	132.9	0.50	1.63
Ours w/o cINN	180.6	0.32	1.21
Ours w/o x_0	381.5	0.73	2.15
Ours w/o ADAIN	156.7	0.48	1.60

Table 2. Quantitative evaluation of video synthesis quality and diversity on iPER [47]. Numeric superscripts indicate the source of the results, cf. Sec. 4.3. * SAVP experienced mode collapse due to training instabilities originating from the two involved discriminators.

$z = \mathcal{T}_\theta(\nu; x_0, \eta^*)$ which is then used to synthesize a video sequence using p_ψ .

4. Experiments

We evaluate the efficacy of our video synthesis method on a diverse set of four video datasets which range from human motion to stochastic dynamics as encompassed by natural landscape scenery. Video prediction results and comparisons are best viewed as videos which are available in the supplemental and on our project page¹. Implementation details can be found in the supplemental. Our PyTorch [60] implementation can be found on our GitHub page². Unless otherwise stated, we generate 16 frame predictions.

4.1. Datasets

Here, we summarize the four diverse datasets used in our evaluation. We train all models on a sequence length of 16.

¹<https://bit.ly/3dg90fV>

²<https://bit.ly/3t66bnU>

Method	Fire					Vegetation					Waterfall					Clouds				
	LPIPS ↓	FID ↓	FVD ↓	DTFVD ↓	DIV ↑	LPIPS ↓	FID ↓	FVD ↓	DTFVD ↓	DIV ↑	LPIPS ↓	FID ↓	FVD ↓	DTFVD ↓	DIV ↑	LPIPS ↓	FID ↓	FVD ↓	DTFVD ↓	DIV ↑
DG ³ [84]	0.18	29.4	361.3	0.40	–	0.22	71.6	290.3	0.86	–	0.25	143.4	1680.6	2.41	–	0.17	73.5	217.5	0.40	–
AL ³ [21]	0.28	48.4	1475.9	11.42	0.74	0.28	48.9	271.0	1.48	0.93	0.32	124.3	1847.8	5.94	0.98	0.27	38.7	142.1	0.76	1.52
Ours	0.23	24.2	376.8	0.79	1.10	0.21	18.2	123.8	0.52	0.86	0.25	66.8	1126.5	2.52	0.61	0.25	18.3	179.3	0.73	0.98

Table 3. Quantitative evaluation of video synthesis quality and diversity (based on VGG [72]) on DTDB [30]. The diversity score based on the I3D [73] trained on DTDB [30] can be found in the supplemental. Note, DG [84] directly optimizes on test samples. Numeric superscripts indicate the source of the results, cf. Sec. 4.3.

A detailed description of the evaluation protocol for each dataset can be found in the supplemental.

Landscape [87] consists of ~ 3000 time-lapse videos of dynamic sky scenes, e.g., cloudy skies and night scenes with moving stars. This dataset contains a wide range of sky appearances and motion speeds. Following previous work [87, 90], we evaluate on a sequence length of 32 frames. We compare with recent work on landscape synthesis [87, 21, 90, 48]. To generate sequences of length 32 we apply our model sequentially, meaning we use the last predicted frame from the last generated 16 frame block as input for the next set of 16 frames.

Dynamic Texture DataBase (DTDB) [30] contains more than 10,000 dynamic texture videos. For evaluation, we focus on the following classes: fire, clouds, vegetation, and waterfall. Each texture class consists of 150 to 300 videos. We train one model for each texture (same as for [21, 84]) on a sequence length of 16 on a resolution of 128×128 .

BAIR Robot Pushing [20] consists of a randomly moving robotic arm that pushes and grasps objects in a box. It contains around 40k training and 256 test videos. This dataset is used by prior work as a benchmark due its stochastic nature and the real-world application. We follow the standard protocol [81, 76, 16, 63] and evaluate on a sequence length of 16 frames on a resolution of 64×64 .

Impersonator (iPER) [47] is a recent dataset that contains humans with diverse styles of clothing executing various random actions. The entire dataset contains 206 videos with a total of 241,564 frames. We follow the train/test split defined in [47] which leads to training set and test sets containing 180k and 49k frames, respectively. We evaluate our model on a sequence length of 16 on a 64×64 resolution.

4.2. Evaluation Metrics

Synthesis quality. We evaluate the video synthesis quality using the Fréchet Video Distance (FVD) [76] which is sensitive to both perceptual quality and temporal coherence. This metric represents the spatiotemporal counterpart to the Fréchet Inception Distance (FID) [32] which is based on an I3D network [73] trained on Kinetics [38], a large-scale human action dataset. To evaluate dynamic textures, we introduce the Dynamic Texture Fréchet Video Distance (DTFVD) by replacing the pre-trained network with one we trained on DTDB for classification [30]. The motivation behind introducing DTFVD is that we seek a metric that is sensitive to the types of dynamics encapsulated by dynamic textures, rather than human action-related motions

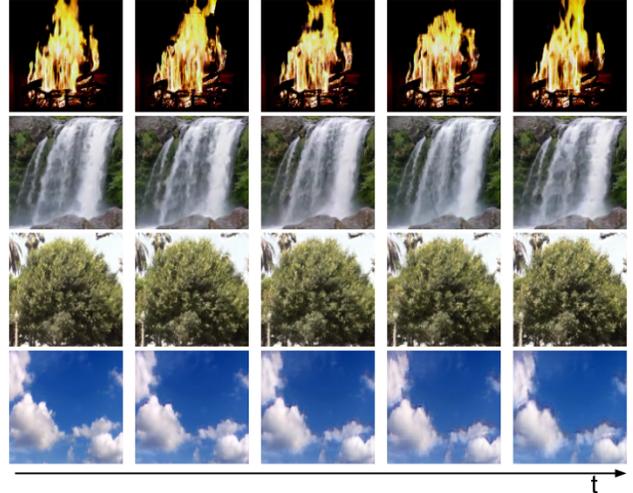


Figure 5. Stochastic video synthesis on DTDB [30] for diverse texture categories. Best viewed as video provided in the supplemental.

as captured by FVD. To further evaluate dynamic textures, we also evaluate perceptual quality in terms of the FID [32] and the Learned Perceptual Image Patch Similarity (LPIPS) [19, 36] metrics.

Diversity. Photorealism and plausible dynamics are not the only factors we are interested in. In addition, our model is capable of stochastically generating plausible videos from a single image. Following previous work [45] on video synthesis, we measure the diversity between video sequence predictions given an initial frame x_0 as their average mutual distance in the feature space of a VGG-16 network [72] pre-trained on ImageNet [69]. In contrast to [45], we use the Euclidean distance instead of the Cosine distance. Moreover, we also report diversity on pre-trained I3D [73] models (similar to above) which is sensitive to both appearance and motion instead of comparing samples frame-wise. We discuss and compare our chosen diversity measures in the supplemental.

4.3. Quantitative Evaluation

For comparison, we use reported performance from the corresponding paper (marked by ¹), where possible, otherwise we report numbers based on pretrained models (marked by ²) or retrained models using the official code (marked by ³) provided by the author.

Landscape. Tab. 1 provides a summary of our evaluation on Landscape in terms of perceptual quality and temporal coherence. As can be seen, we generally outper-

Method	FVD ↓	DIV VGG ↑	DIV I3D ↑
Video Flow ¹ [43]	131.0	–	–
SRVP ² [25]	141.7	0.93	1.65
IVRNN ³ [14]	121.3	0.69	1.13
SAVP ^{1,2} [45]	116.4	0.98	1.70
LVT ¹ [63]	125.8	–	–
DVD-GAN ¹ [16]	109.8	–	–
Video Transformer ¹ [81]	94.0	–	–
Ours	99.3	0.98	1.93
Ours w/o cINN	134.5	0.59	0.94
Ours w/o x_0	272.6	2.40 [†]	2.48 [†]
Ours w/o ADAIN	131.2	0.78	1.73

Table 4. Quantitative evaluation and ablation study of generation quality and diversity on BAIR [20]. Numeric superscripts indicate the source of the results, cf. Sec. 4.3. [†] denotes high diversity due to artifacts.

form all methods across all metrics. Animating Landscape (AL) [21] stores the motion embeddings of all training instances in their codebook and uses them to generate videos during inference. In this way, AL is able to reproduce the diversity of the training videos. DTVNet [90] does not enforce a distribution over their representation and consequently is limited to deterministic video generations. DeepLandscape [48] (DL) does not learn dynamics from videos, but rather uses a manually constructed set of homographies. The pretrained model provided by DL was trained on their unreleased dataset. In contrast, we explicitly model and learn the dynamics distribution and by that, are able to synthesize *novel* dynamics to set scenes in motion.

DTDB. We observe similar results on DTDB (Tab. 3) on nearly all dynamic textures (fire, waterfall, and vegetation) across all perceptual quality and coherence metrics. For the clouds, AL achieves better results due to the fact that this motion can be faithfully described by optical flow. Here, we also consider results from methods dedicated to dynamic texture synthesis [84, 86] as strong baselines. These methods are not exactly comparable as they directly optimize on *test samples*. We only present results for DG [84], as Xie et al. [86] did not converge when trained on all test samples.

BAIR. We achieve strong results in terms of video quality (Tab. 4), even when compared with the computationally expensive transformer based approach [81]. In terms of diversity, we are on par with the state-of-the-art stochastic video prediction approaches.

iPER. The evaluation of articulated human motion is presented in Tab. 2. We achieve superior results to recent approaches for video prediction [25, 14, 45] in terms of FVD and diversity. Note, that we only condition on one frame in comparison to the baselines which use two [45, 14] and eight context frames [25].

4.4. Qualitative Evaluation

Image-to-video synthesis. We provide samples for all datasets. On Landscape [87] we see that our model is able to

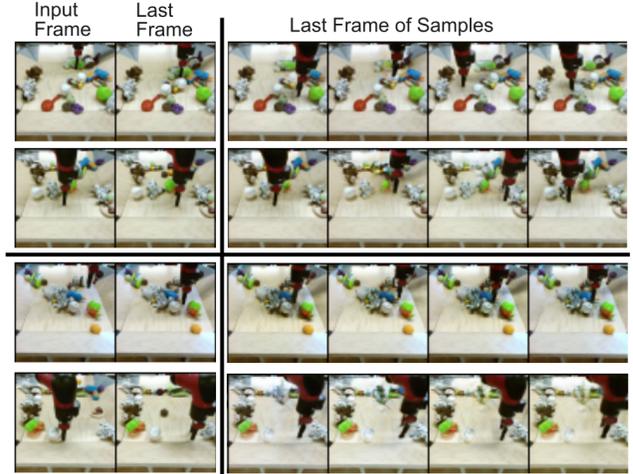


Figure 6. Qualitative evaluation of diversity on BAIR [20]. **Stochastic video synthesis:** (top two rows, left-to-right) input frame and last frame from a BAIR sequence and four frames representing the last frames from sampled videos generated using the input frame alone. The generated frames show a high degree of stochasticity in terms of the end effector position, as desired. **Controlled video synthesis:** (bottom two rows, left-to-right) input frame and last frame from a BAIR sequence, and four frames representing the last frames from sampled videos generated using both the input frame and the 3D end effector position in the last frame. The end effector position in the last frame is in close agreement with the position control input, as desired.

synthesize realistic samples (see Fig. 3) from diverse, complex scenes captured in the input image. In Fig. 4, we show samples on iPER [47] which illustrates the complexity of motion in the dataset. In Fig. 5, we visualize one sample per DTDB class which shows the variety of dynamic textures used for evaluation. Lastly, Fig. 6 (top two rows) show the diversity in our video samples by way of the differences across the last generated frame per sample on BAIR [20].

Controllable video synthesis. A strength of our model is the ability to exert explicit control over the synthesis process. As described in Sec. 3.4, we control this process by introducing a factor η . Here, we consider two different factors for controllable video synthesis on BAIR [20] and DTDB [30]. On BAIR we condition the synthesis process on the 3D location of the robot arm’s end effector in the last frame; we use the location provided in the groundtruth. Fig. 6 (bottom two rows) shows several samples of the last frame of each sequence of our controllable synthesis. It can clearly be seen that the last frames of our samples match closely to the groundtruth end frame. As a second example, this time on DTDB [30], we condition the video synthesis of clouds based on the 2D direction of motion, again through manipulating η . This is visualized in Fig. 7 where four different directions are considered. To aid in the visualization, we also include the optical flow fields, estimated with [33], to show the consistency between the motion direction used

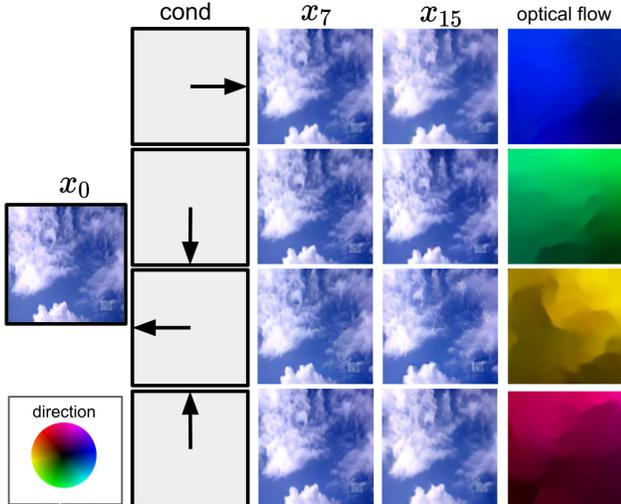


Figure 7. Examples of controlling video synthesis of clouds in DTDB [30] starting at frame x_0 using motion direction inputs (indicated by arrows). We show intermediate frames x_7 and x_{15} . The color wheel indicates flow direction.

for conditioning and the direction realized in the generated videos. As can be seen, the conditioning and generated motion directions are in close agreement. For results on controlled video-to-video synthesis (cf. Sec. 3.4), please refer to the supplemental.

Motion transfer. Finally, we illustrate the capability of our model to transfer a motion contained in one sequence to a set of initial frames for video synthesis. Fig. 8 illustrates this process using Landscape [87], where the top row contains the motion to be transferred and the bottom three rows show the generated video sequences realized by combining the transferred motion and the initial frames. As can be clearly seen, the original motion is successfully transferred to each of the scenes.

4.5. Ablation study

To evaluate the design choices of our approach, we now perform ablation studies on BAIR [20] and iPER [47]: (Ours w/o x_0) represents implementing our video generator, p_ψ , without conditioning on the input image, x_0 , thus z also captures the full scene content information, (Ours w/o ADAIN) similarly denotes removing the ADAIN input of z in our proposed Video ResBlk, i.e., p_ψ only has access to z via the bottleneck and (Ours w/o cINN) stands for removing the cINN resulting in a cVAE framework.

In Tab. 2 and Tab. 4, we observe significant performance drops for all ablations compared to our full model (Ours). In particular removing the conditioning image, x_0 , from the generator, p_ψ , greatly affects the synthesis quality. This is due to the generator not having direct access to the static information depicted in the initial frame x_0 . When removing the ADAIN input of z from our Video ResBlk, the infor-

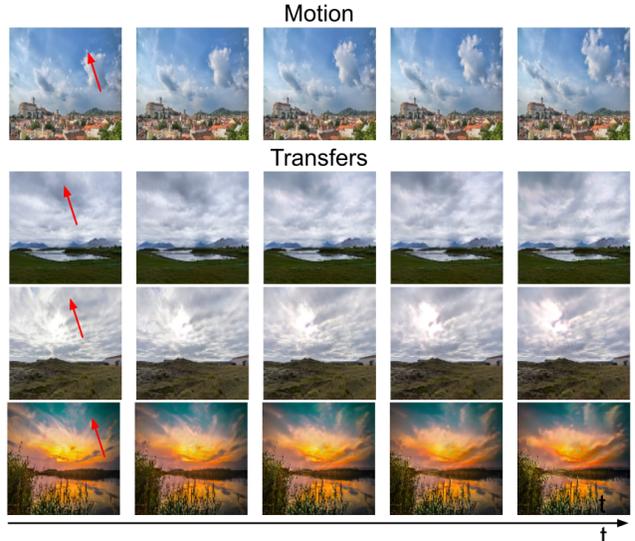


Figure 8. Transferring motion across videos on Landscape [87]. (top row, left-to-right) source video for target motion. (bottom three rows, left-to-right) animating different starting frames by transferring motion from source video. Red arrows indicate the 2D direction of motion. Best viewed as video provided in the supplemental.

mation of z is now only available at the lowest scale of p_ψ , in contrast to the multi-scale information flow in our full model. Moreover, the cVAE-only model (w/o cINN) results in worse performance both in quality and diversity, which can be explained by the trade-off between synthesis quality and capacity regularization, as discussed in Sec. 3.3.

5. Conclusion

In summary, we introduced a novel model for understanding image-to-video synthesis based on a bijective transformation, instantiated as a cINN, between the video and image domains plus residual information. The probabilistic residual representation allows to sample and synthesize novel, plausible progressions in video with the same initial frame. Moreover, our framework allows for incorporating additional controlling factors to guide the image-to-video synthesis process. Our empirical evaluation and comparison to strong baselines on four diverse video datasets demonstrated the efficacy of our stochastic image-to-video synthesis approach.

Acknowledgement

This work was started as part of M.D.’s internship at Ryerson University and was supported in part by the DAAD scholarship, the NSERC Discovery Grant program (K.G.D.), the German Research Foundation (DFG) within project 421703927 (B.O.) and the BW Stiftung (B.O.). K.G.D. contributed to this work in his capacity as an Associate Professor at Ryerson University.

References

- [1] Lynton Ardizzone, Jakob Kruse, Carsten Rother, and Ullrich Köthe. Analyzing inverse problems with invertible neural networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019. [2](#)
- [2] Lynton Ardizzone, Carsten Lüth, Jakob Kruse, Carsten Rother, and Ullrich Köthe. Guided image generation with conditional invertible neural networks. *CoRR*, 2019. [2](#), [15](#)
- [3] Mohammad Babaeizadeh, Chelsea Finn, Dumitru Erhan, Roy H. Campbell, and Sergey Levine. Stochastic variational video prediction. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018. [2](#)
- [4] Wenbo Bao, Wei-Sheng Lai, Chao Ma, Xiaoyun Zhang, Zhiyong Gao, and Ming-Hsuan Yang. Depth-aware video frame interpolation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3703–3712, 2019. [2](#)
- [5] Andreas Blattmann, Timo Milbich, Michael Dorcenwald, and Bjorn Ommer. Behavior-driven synthesis of human dynamics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12236–12246, June 2021. [1](#)
- [6] Andreas Blattmann, Timo Milbich, Michael Dorcenwald, and Bjorn Ommer. Understanding object dynamics for interactive image-to-video synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5171–5181, June 2021. [2](#)
- [7] Biagio Brattoli, Uta Büchler, Michael Dorcenwald, Philipp Reiser, Linard Filli, Fritjof Helmchen, Anna-Sophia Wahl, and Björn Ommer. Unsupervised behaviour analysis and magnification (ubam) using deep learning. *Nature Machine Intelligence*, 2021. [1](#)
- [8] Biagio Brattoli, Uta Büchler, Anna-Sophia Wahl, Martin E. Schwab, and Björn Ommer. LSTM self-supervision for detailed behavior analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3747–3756, 2017. [1](#)
- [9] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019. [5](#)
- [10] Andreja Bubic, D Yves von Cramon, and Ricarda I Schubotz. Prediction, cognition and the brain. *Frontiers in Human Neuroscience*, 4:25, 2010. [1](#)
- [11] Christopher P. Burgess, Irina Higgins, Arka Pal, Loïc Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. Understanding disentangling in β -VAE. *CoRR*, 2018. [4](#)
- [12] Arunkumar Byravan, Felix Leeb, Franziska Meier, and Dieter Fox. SE3-Pose-Nets: structured deep dynamics models for visuomotor planning and control. *CoRR*, 2017. [1](#)
- [13] João Carreira and Andrew Zisserman. Quo vadis, action recognition? A new model and the kinetics dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4724–4733, 2017. [16](#)
- [14] Lluís Castrejón, Nicolas Ballas, and Aaron C. Courville. Improved conditional vrns for video prediction. In *Proceedings of the International Conference on Computer Vision (ICCV)*, pages 7607–7616, 2019. [1](#), [2](#), [5](#), [7](#), [13](#), [16](#)
- [15] Xi Chen, Diederik P. Kingma, Tim Salimans, Yan Duan, Prfulla Dhariwal, John Schulman, Ilya Sutskever, and Pieter Abbeel. Variational lossy autoencoder. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017. [4](#)
- [16] Aidan Clark, Jeff Donahue, and Karen Simonyan. Adversarial video generation on complex datasets. *CoRR*, 2019. [2](#), [4](#), [6](#), [7](#)
- [17] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017. [4](#), [15](#)
- [18] Michael Dorcenwald, Uta Büchler, and Björn Ommer. Unsupervised magnification of posture deviations across subjects. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8253–8263, 2020. [1](#)
- [19] Alexey Dosovitskiy and Thomas Brox. Generating images with perceptual similarity metrics based on deep networks. In *Neural Information Processing Systems (NeurIPS)*, pages 658–666, 2016. [4](#), [6](#)
- [20] Frederik Ebert, Chelsea Finn, Alex X. Lee, and Sergey Levine. Self-supervised visual planning with temporal skip connections. In *Conference on Robot Learning (CoRL)*, pages 344–356, 2017. [1](#), [6](#), [7](#), [8](#), [13](#), [16](#)
- [21] Yuki Endo, Yoshihiro Kanamori, and Shigeru Kuriyama. Animating landscape: Self-supervised learning of decoupled motion and appearance for single-image video synthesis. *ACM Transactions on Graphics*, pages 175:1–175:19, 2019. [1](#), [2](#), [5](#), [6](#), [7](#), [13](#), [16](#), [17](#)
- [22] Patrick Esser, Johannes Haux, Timo Milbich, and Björn Ommer. Towards learning a realistic rendering of human behavior. In *ECCV Workshops*, pages 409–425, 2018. [1](#)
- [23] Patrick Esser, Robin Rombach, and Björn Ommer. A disentangling invertible interpretation network for explaining latent representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9220–9229, 2020. [2](#)
- [24] Chelsea Finn and Sergey Levine. Deep visual foresight for planning robot motion. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pages 2786–2793, 2017. [1](#)
- [25] Jean-Yves Franceschi, Edouard Delasalles, Mickaël Chen, Sylvain Lamprier, and Patrick Gallinari. Stochastic latent residual video prediction. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 3233–3246, 2020. [1](#), [2](#), [3](#), [5](#), [7](#), [16](#)
- [26] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. A neural algorithm of artistic style. *CoRR*, 2015. [2](#)
- [27] Patrick Gebert, Alina Roitberg, Monica Haurilet, and Rainer Stiefelhagen. End-to-end prediction of driver intention using 3D convolutional neural networks. In *Proceedings of the IEEE Intelligent Vehicles Symposium (IV)*, pages 969–974, 2019. [1](#)
- [28] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville,

- and Yoshua Bengio. Generative adversarial networks. *CoRR*, 2014. 4
- [29] Ishaan Gulrajani, Faruk Ahmed, Martín Arjovsky, Vincent Dumoulin, and Aaron C. Courville. Improved training of wasserstein gans. In *Neural Information Processing Systems (NeurIPS)*, pages 5767–5777, 2017. 15
- [30] Isma Hadji and Richard P. Wildes. A new large scale dynamic texture dataset with application to convnet understanding. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 334–351, 2018. 5, 6, 7, 8, 13, 14, 16, 17, 18
- [31] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 4, 14, 15
- [32] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. In *Neural Information Processing Systems (NeurIPS)*, pages 6626–6637, 2017. 6
- [33] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. FlowNet 2.0: Evolution of optical flow estimation with deep networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1647–1655, 2017. 7
- [34] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5967–5976, 2017. 15
- [35] Jörn-Henrik Jacobsen, Arnold W. M. Smeulders, and Edouard Oyallon. i-revnet: Deep invertible networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018. 2
- [36] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 694–711, 2016. 4, 6
- [37] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4401–4410, 2019. 2, 4, 15
- [38] Will Kay, João 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. *CoRR*, 2017. 6, 16
- [39] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015. 15, 16
- [40] Diederik P. Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In *Neural Information Processing Systems (NeurIPS)*, pages 10236–10245, 2018. 2, 15
- [41] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2014. 4, 15
- [42] Ivan Kobyzev, Simon Prince, and Marcus Brubaker. Normalizing flows: An introduction and review of current methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 4
- [43] 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 *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020. 2, 7
- [44] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 1558–1566, 2016. 4
- [45] Alex X. Lee, Richard Zhang, Frederik Ebert, Pieter Abbeel, Chelsea Finn, and Sergey Levine. Stochastic adversarial video prediction. *CoRR*, 2018. 1, 2, 3, 5, 6, 7, 16
- [46] Jae Hyun Lim and Jong Chul Ye. Geometric GAN. *CoRR*, 2017. 5
- [47] Wen Liu, Zhixin Piao, Jie Min, Wenhan Luo, Lin Ma, and Shenghua Gao. Liquid warping GAN: A unified framework for human motion imitation, appearance transfer and novel view synthesis. In *Proceedings of the International Conference on Computer Vision (ICCV)*, pages 5903–5912, 2019. 5, 6, 7, 8, 13, 14, 16
- [48] Elizaveta Logacheva, Roman Suvorov, Oleg Khomenko, Anton Mashikhin, and Victor Lempitsky. DeepLandscape: Adversarial Modeling of Landscape Videos. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 256–272, 2020. 2, 5, 6, 7
- [49] Dominik Lorenz, Leonard Bereska, Timo Milbich, and Björn Ommer. Unsupervised part-based disentangling of object shape and appearance. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10955–10964, 2019. 1
- [50] Andreas Lugmayr, Martin Danelljan, Luc Van Gool, and Radu Timofte. SRFlow: Learning the super-resolution space with normalizing flow. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 715–732, 2020. 2
- [51] Manuel Martin, Alina Roitberg, Monica Haurilet, Matthias Horne, Simon Reiß, Michael Voit, and Rainer Stiefelhagen. Drive&act: A multi-modal dataset for fine-grained driver behavior recognition in autonomous vehicles. In *Proceedings of the International Conference on Computer Vision (ICCV)*, pages 2801–2810, 2019. 1
- [52] Sarfaraz Masood, Abhinav Rai, Aakash Aggarwal, Mohammad Najam Doja, and Musheer Ahmad. Detecting distraction of drivers using convolutional neural network. *Pattern Recognition Letters*, pages 79–85, 2020. 1
- [53] Lars M. Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for gans do actually converge? In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 3478–3487, 2018. 15

- [54] Timo Milbich, Miguel Ángel Bautista, Ekaterina Sutter, and Björn Ommer. Unsupervised video understanding by reconciliation of posture similarities. In *Proceedings of the International Conference on Computer Vision (ICCV)*, pages 4404–4414, 2017. [2](#)
- [55] Matthias Minderer, Chen Sun, Ruben Villegas, Forrester Cole, Kevin P. Murphy, and Honglak Lee. Unsupervised learning of object structure and dynamics from videos. In *Neural Information Processing Systems (NeurIPS)*, pages 92–102, 2019. [2](#), [3](#)
- [56] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018. [15](#)
- [57] Simon Niklaus and Feng Liu. Context-aware synthesis for video frame interpolation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1701–1710, 2018. [2](#)
- [58] George Papamakarios, Eric T. Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *CoRR*, 2019. [4](#)
- [59] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2337–2346, 2019. [4](#), [15](#)
- [60] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Neural Information Processing Systems (NeurIPS)*, pages 8024–8035, 2019. [5](#), [14](#)
- [61] Karl Pertsch, Oleh Rybkin, Jingyun Yang, Shenghao Zhou, Konstantinos G. Derpanis, Kostas Daniilidis, Joseph J. Lim, and Andrew Jaegle. Keyframing the future: Keyframe discovery for visual prediction and planning. In *LADC, Proceedings of Machine Learning Research*, pages 969–979, 2020. [1](#), [2](#)
- [62] Albert Pumarola, Stefan Popov, Francesc Moreno-Noguer, and Vittorio Ferrari. C-Flow: Conditional Generative Flow Models for Images and 3D Point Clouds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7946–7955, 2020. [2](#)
- [63] Ruslan Rakhimov, Denis Volkhonskiy, Alexey Artemov, Denis Zorin, and Evgeny Burnaev. Latent video transformer. In *International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications, VISIGRAPP*, pages 101–112, 2021. [2](#), [6](#), [7](#)
- [64] Fitsum A. Reda, Guilin Liu, Kevin J. Shih, Robert Kirby, Jon Barker, David Tarjan, Andrew Tao, and Bryan Catanzaro. Sdcnet: Video prediction using spatially-displaced convolution. *CoRR*, 2018. [1](#)
- [65] Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 1530–1538, 2015. [2](#)
- [66] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 1278–1286, 2014. [15](#)
- [67] Robin Rombach, Patrick Esser, and Björn Ommer. Making Sense of CNNs: Interpreting Deep Representations and Their Invariances with INNs. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 647–664, 2020. [2](#), [15](#)
- [68] Robin Rombach, Patrick Esser, and Björn Ommer. Network-to-network translation with conditional invertible neural networks. In *Neural Information Processing Systems (NeurIPS)*, 2020. [2](#)
- [69] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Fei-Fei Li. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision (IJCV)*, 115(3), 2015. [6](#), [16](#)
- [70] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. Animating arbitrary objects via deep motion transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2377–2386, 2019. [2](#)
- [71] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. In *Neural Information Processing Systems (NeurIPS)*, pages 7135–7145, 2019. [2](#)
- [72] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015. [6](#), [16](#)
- [73] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2826, 2016. [5](#), [6](#), [16](#), [17](#)
- [74] Matthew Tesfaldet, Marcus A. Brubaker, and Konstantinos G. Derpanis. Two-stream convolutional networks for dynamic texture synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6703–6712, 2018. [2](#)
- [75] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. MoCoGAN: Decomposing Motion and Content for Video Generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1526–1535, 2018. [2](#)
- [76] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphaël Marinier, Marcin Michalski, and Sylvain Gelly. Towards Accurate Generative Models of Video: A New Metric & Challenges. *CoRR*, 2018. [6](#), [16](#)
- [77] Ruben Villegas, Jimei Yang, Seunghoon Hong, Xunyu Lin, and Honglak Lee. Decomposing motion and content for natural video sequence prediction. In *Proceedings of the In-*

- ternational Conference on Learning Representations (ICLR), 2017. 2
- [78] Jacob Walker, Kenneth Marino, Abhinav Gupta, and Martial Hebert. The pose knows: Video forecasting by generating pose futures. In *Proceedings of the International Conference on Computer Vision (ICCV)*, pages 3352–3361, 2017. 1
- [79] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional GANs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8798–8807, 2018. 15
- [80] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Nikolai Yakovenko, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. In *Neural Information Processing Systems (NeurIPS)*, pages 1152–1164, 2018. 2, 4
- [81] Dirk Weissenborn, Oscar Täckström, and Jakob Uszkoreit. Scaling autoregressive video models. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020. 1, 2, 6, 7
- [82] Nevan Wichers, Ruben Villegas, Dumitru Erhan, and Honglak Lee. Hierarchical long-term video prediction without supervision. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 6033–6041, 2018. 2
- [83] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–19, 2018. 14, 15
- [84] Jianwen Xie, Ruiqi Gao, Zilong Zheng, Song-Chun Zhu, and Ying Nian Wu. Learning Dynamic Generator Model by Alternating Back-Propagation through Time. In *Proceedings of the National Conference on Artificial Intelligence (AAAI)*, pages 5498–5507, 2019. 2, 6, 7, 13
- [85] Jianwen Xie, Yang Lu, Ruiqi Gao, Song-Chun Zhu, and Ying Nian Wu. Cooperative training of descriptor and generator networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, pages 27–45, 2020. 2
- [86] Jianwen Xie, Song-Chun Zhu, and Ying Nian Wu. Synthesizing dynamic patterns by spatial-temporal generative convnet. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1061–1069, 2017. 2, 7
- [87] Wei Xiong, Wenhan Luo, Lin Ma, Wei Liu, and Jiebo Luo. Learning to generate time-lapse videos using multi-stage dynamic generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2364–2373, 2018. 2, 5, 6, 7, 8, 13, 14, 16
- [88] Xinchen Yan, Jimei Yang, Kihyuk Sohn, and Honglak Lee. Attribute2image: Conditional image generation from visual attributes. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 776–791, 2016. 4
- [89] Jason J. Yu, Konstantinos G. Derpanis, and Marcus A. Brubaker. Wavelet flow: Fast training of high resolution normalizing flows. In *Neural Information Processing Systems (NeurIPS)*, 2020. 2
- [90] Jiangning Zhang, Chao Xu, Liang Liu, Mengmeng Wang, Xia Wu, Yong Liu, and Yunliang Jiang. DTVNet: Dynamic Time-Lapse Video Generation via Single Still Image. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 300–315, 2020. 2, 5, 6, 7, 13
- [91] Shengjia Zhao, Jiaming Song, and Stefano Ermon. InfoVAE: Information maximizing variational autoencoders. *CoRR*, 2017. 4
- [92] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A. Efros, Oliver Wang, and Eli Shechtman. Toward multimodal image-to-image translation. In *Neural Information Processing Systems (NeurIPS)*, pages 465–476, 2017. 16

Supplemental

A. Additional Visualizations

For each of our experiments conducted in the main paper, we provide additional video material, consisting of 17 videos in total. To further highlight the benefits of our proposed framework, in the course of our supplemental video material, we compare to *five* approaches. **Due to the collective large size of the videos, the supplemental with the corresponding videos is provided on our [project page](#).** For each video, multiple cycles are shown (indicated left-bottom) as well as the corresponding video playback rate in frames-per-second (FPS) (right-bottom). The file structure of our provided video material is as follows:

```
supplemental_material_222
|
+---A1-Landscape
|
+---A2-iPER
|
+---A3-DTDB
|
+---A4-BAIR
|
+---A5-Controllable_Video_Synthesis
|
+---A6-Failure_Cases
```

We next discuss the video material for each experiment individually. Each subsection matches its corresponding file (e.g., ‘A.1.Landscape’ corresponds to ‘...--A1-Landscape’) which contains the discussed video sequences.

A.1. Landscape

For the Landscape dataset [87], we provide the corresponding video (`Landscape_samples.mp4`) to the samples depicted in Fig. 3 in the main paper. Additionally, we show a qualitative comparison to previous work, i.e., AL [21], DTVNet [90], and MDGAN [87] in `Landscape_comparison.mp4`, with ‘GT’ denoting the ground-truth. We clearly observe that our model synthesizes more appealing and realistic video sequences compared to the competing methods. Both MDGAN [87] and DTVNet [90] produce blurry videos when using the officially provided pretrained weights and code from the respective webpages. While AL produces decent animations in the presence of small motion, when animating fast motions, however, warping artifacts are present, cf. e.g., row 3. These artifacts become even more evident when AL is applied to DTDB (Sec. A.3). In contrast, our method produces realistic looking results in the case of both small and

large motions. Next, we evaluate the diversity of the generated samples in `Landscape_diversity.mp4`. The video contains multiple future progressions for a given starting frame, x_0 . It can be seen that our approach produces diverse samples capturing a broad range of motion directions, as well as speeds. Moreover, we demonstrate in `Landscape_longer_duration.mp4` the capability of our model to synthesize longer sequences (48 frames) by sequentially applying our model on the last frame of the previously predicted video sequence.

A.2. iPER

For the iPER dataset [47], we provide the corresponding video (`iPER_samples.mp4`) to the samples depicted in Fig. 4 in the main paper. We further provide a qualitative comparison to the best performing method IVRNN [14] on iPER in `iPER_comparison.mp4` with ‘GT’ denoting the ground-truth. Our method produces more natural motions, e.g., row 3, compared to [14]. Note, that both methods suffer from artifacts due to the low image resolution of 64×64 , such as vanishing hands in motion.

A.3. DTDB

For each dynamic texture from DTDB [30] used in our main paper, we provide examples (`Clouds.mp4`, `Fire.mp4`, `vegetation.mp4`, `Waterfall.mp4`) for stochastic image-to-video synthesis for random starting frames, x_0 , comparing our proposed approach to AL [21] and DG [84]. As described in the main paper, DG [84] is directly optimized on test samples, thus overfitting directly to the test distribution. Consequently, we observe that their generations almost perfectly reproduce the ground-truth motion which is most evident for the clouds texture. However, their method suffers from blurring due to optimization using an L2 pixel loss. Similar to the comparisons on the Landscape dataset (Sec. A.1), AL [21] has problems with learning and reproducing the motion of dynamic textures exhibiting rapid motion changes, such as fire. This is explained by the susceptibility of optical flow to inaccuracies when capturing very fast motion, as well as dynamic patterns outside the scope of optical flow, e.g., flicker. Moreover, in the clouds examples (last row) AL wrongly sets the landscape into motion. Our model, on the other hand, produces sharp video sequences with realistic looking motions for *all* textures.

A.4. BAIR

In `BAIR_comparison.mp4`, we provide a qualitative comparison to a strong baseline, IVRNN [14], on the BAIR dataset [20]. While both approaches are able to render the robot’s end effector and the visible environment well, we observe significant differences when it comes to the effector interacting with or occluding background objects. An

example of this difficulty can be seen when interacting with the object in the middle of the scene in row 2. IVRNN is unable to depict the object structure and texture during the interaction which results in heavy blur due to averaging over all possible future states. In contrast, this interaction looks much more natural in the video sequence predicted by our model (also row 2). Moreover, the last row (back of the scene, right) illustrates a problem of IVRNN which sometimes occurs in the presence of object occlusions. Specifically, the object which is occluded at the beginning is eventually revealed and is synthesized as a blurry texture, by that, averaging over all possible realizations. Again, our model does not suffer from this problem and correctly handles object occlusions. Additionally, `BAIR_diversity.mp4` qualitatively illustrates the prediction diversity of our model by animating a fixed starting frame x_0 multiple times. Again, ‘GT’ denotes ground-truth. Our model synthesizes diverse samples by broadly covering motions in the x , y , and z directions.

A.5. Controllable Video Synthesis

In this section, we present qualitative experiments for the following controlled video prediction task: *controlled image-to-video synthesis*, *motion transfer*, and *controlled video-to-video synthesis*.

Controlled image-to-video synthesis. The video `Endpoint_BAIR.mp4` illustrates several image-to-video generations while controlling $\eta = (x, y, z)$, the 3D end effector position, similar to Fig. 6 in our main paper. It shows that, while in each example the effector approximately stops at the provided end position (end frame of GT), its movements between the starting and end frame, which are inferred by the sampled residual representations $\nu \sim q(\nu)$, exhibit significantly varying and natural progressions. Moreover, in `Direction_Clouds1.mp4` we provide additional video examples for controlling the direction of cloud movements with η , similar to Fig. 7 in our main paper. We observe that our model renders crisp future progressions (row 2-5) of a given starting frame x_0 , while following our provided movement control (top row).

Motion transfer. Next, we analyze the application of our model for the task of directly transferring a query motion extracted from a given landscape video \tilde{X} to a random starting frame x_0 . To this end, we extract the residual representation $\tilde{\nu}$ of \tilde{X}_0 by first obtaining its video representation $\tilde{z} = q_\phi(z|\tilde{X})$ and corresponding residual $\tilde{\nu} = \mathcal{T}_\theta^{-1}(\tilde{z}; \tilde{x}_0)$ with \tilde{x}_0 being the starting frame of \tilde{X} . We use $\tilde{\nu}$ to animate the starting frame x_0 . `Transfer_Landscape.mp4` shows that our model accurately transfers the query motion, e.g., as the corresponding direction and speed of the clouds, to the target landscape images (rows 1-3, left-to-right).

Controlled video-to-video synthesis. In controlled video-to-video synthesis, we explicitly adjust the initial factor

$\tilde{\eta}$ of an observed video sequence \tilde{X} . To this end, we first obtain its video representation $\tilde{z} = q_\phi(z|\tilde{X})$ followed by extracting the corresponding residual information $\tilde{\nu} = \mathcal{T}_\theta^{-1}(\tilde{z}; \tilde{x}_0, \tilde{\eta})$. Subsequently, to generate the video sequence depicting our controlled adjustment of \tilde{X} , we simply choose a new value $\tilde{\eta} = \tilde{\eta}^*$ and perform the image-to-sequence inference process. This can be seen in the video `Direction_Clouds2.mp4` using cloud video sequences from DTDB [30]. In each example (second row), the motion direction of the query video (leftmost) is adjusted by the provided control (top row). To highlight that the residual representations ν in these cases actually correspond to the query video, we additionally animate the initial image of the query videos by sampling a new residual representation $\nu \sim q(\nu)$ and apply the same controls (bottom rows). We observe that, while the directions of the synthesized videos are identical, their speeds are significantly different, as desired. In the case of video-to-video synthesis, the movement speed remains the same, in contrast to the image-to-video case, where the movement speed has changed due to the changed residual representation.

A.6. Failure Cases

We highlight two types of failure cases we observed which are visualized in the video `Failure_cases.mp4`:

- When the starting frame depicts a complex posture (e.g., folded arms or a leg in the air) on iPER [47] the model has difficulty synthesizing realistic continuations.
- While the Landscape dataset [87] mainly covers naturally progressing cloud motions, there is also a small subset of fast timelapse videos. Due to the underrepresentation of such examples in the dataset, our model struggles to correctly capture fast paced timelapse data without explicitly resorting to data-balancing techniques during training.

B. Implementation Details

Here, we provide a detailed overview of our network architecture as well as the training procedure. The PyTorch [60] implementation can be found on our GitHub page³.

B.1. Network Details

Encoder. The encoder $q_\phi(z|X)$ follows the structure of a 3D ResNet-18 [31] using GroupNorm [83] as a normalization layer. Two convolutions with a kernel size of 4×4 are used to obtain an one-dimensional latent representation for representing the mean μ and log variance $\log \sigma^2$.

³<https://bit.ly/3t66bnU>

During training, we sample from $q_\phi(z|X)$ using the reparametrization trick [41, 66].

Decoder. The decoder $p_\psi(X|x_0, z)$ consists of $n = 6$ video residual blocks, with each block followed by nearest-neighbor upsampling to upscale the feature map in space and time (except the last one). This structure is illustrated in Fig. 9. The video representation, z , is inserted into the generator using a fully connected layer matching the initial feature map. The channel factor, ch_f , defines the number of channels and by that, the depth of the model. For BAIR and iPER, we set ch_f to 64, otherwise we set it to 32. Depending on the dataset, time length, and resolution, the last two up-scaling layers needs to be adjusted. The video representation z is inserted to the decoder using a fully connected layer matching the initial feature map. We use GroupNorm [83] in SPADE [59] and instance normalization in the ADAIN [37] layer. If the input and the output channels do not match, a 1×1 convolution is used to adjust the channel dimensions. For matching the output channels, we use a 3D convolution followed by a Tanh activation function. Moreover, spectral norm [56] is used in the decoder.

Bijection Transformation. The bijective transformation, \mathcal{T}_θ , is realized as a normalizing flow consisting of a stacked sequence of n_f invertible neural networks (INNs) operating on the video representation, z . We use $n_f = 20$ invertible blocks for all datasets. Each block consists 8 normalizing flows consisting of actnorm [40], affine coupling layers [17], and fixed shuffling layers, following previous work [67]. Each affine coupling layer is parameterized by two fully connected layers. In every affine coupling layer, we additionally insert the conditioning information following previous work [2, 67]. The feature representation for the starting frame x_0 is obtained by a pretrained Autoencoder optimized for reconstructing images.

Discriminators. For the static discriminator, a patch discriminator [34] is used and for the temporal discriminator a 3D ResNet [31].

B.2. Training Details

The loss objective for the generative model of a video sequence $X = [x_1, \dots, x_T] \sim p_X(X) \in \mathbb{R}^{d_x}$ with the corresponding starting frame $x_0 \in \mathbb{R}^{d_x}$ and a video representation $z \sim q_\phi(z|X) \in \mathbb{R}^{d_z}$ can be written as

$$\begin{aligned} \mathcal{L}_{p_\psi, q_\phi} = & \mathbb{E}_{\substack{X \sim p_X(X) \\ z \sim q_\phi(z|X)}} \left[\lambda \| \| X - p_\psi(X|x_0, z) \|_1 \right. \\ & + \ell^\phi(X, p_\psi(X|x_0, z)) - \mathcal{D}_T(p_\psi(X|x_0, z)) \\ & \left. - \mathcal{D}_S(p_\psi(X|x_0, z)) + \lambda_F \ell_F(X, p_\psi(X|x_0, z)) \right] \\ & + \beta D_{\text{KL}}(q_\phi(z|X) \| q(z)), \end{aligned} \quad (10)$$

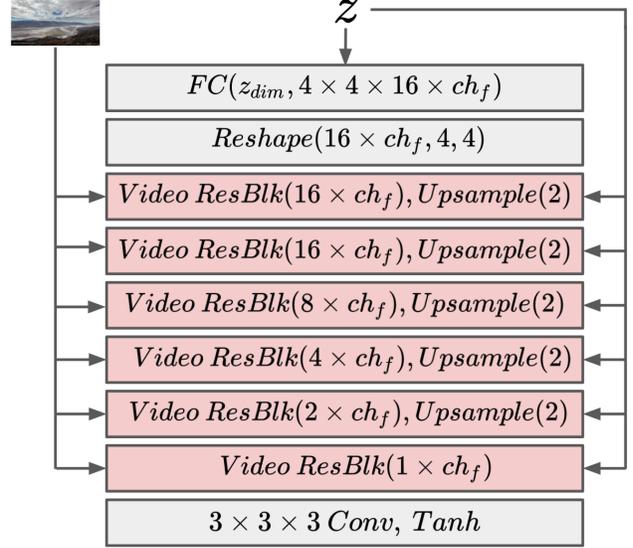


Figure 9. Overview of the decoder structure.

where ℓ_F denotes the feature matching loss [79] to stabilize the training. The hyperparameters λ and λ_F are both set to 10.

The loss objective for the temporal discriminator can be written as

$$\begin{aligned} \mathcal{L}_{\mathcal{D}_T} = & \mathbb{E}_{X \sim p_X(X)} [\rho(1 - \mathcal{D}_T(X)) + \lambda_{GP} \| \nabla \mathcal{D}_T(X) \|_2^2] \\ & + \mathbb{E}_{\substack{X \sim p_X(X) \\ z \sim q_\phi(z|X)}} [\rho(1 + \mathcal{D}_T(p_\psi(X|x_0, z))], \end{aligned} \quad (11)$$

where $\| \nabla \mathcal{D}_T(X) \|_2^2$ denotes the gradient penalty [53, 29] to stabilize the discriminator training and ρ the ReLU activation function. The weighting factor λ_{GP} was set to 10.

For the spatial discriminator, the objective can be formulated as

$$\begin{aligned} \mathcal{L}_{\mathcal{D}_S} = & \mathbb{E}_{X \sim p_X(X)} [\rho(1 - \mathcal{D}_S(X)) \\ & + \mathbb{E}_{\substack{X \sim p_X(X) \\ z \sim q_\phi(z|X)}} [\rho(1 + \mathcal{D}_S(p_\psi(X|x_0, z)))]]. \end{aligned} \quad (12)$$

The overall loss objective can be summarized as

$$\mathcal{L} = \mathcal{L}_{p_\psi, q_\phi} + \mathcal{L}_{\mathcal{D}_T} + \mathcal{L}_{\mathcal{D}_S}. \quad (13)$$

Our video synthesis model is trained using Adam [39] with a learning rate of $2 \cdot 10^{-4}$, $\beta_1 = 0.5$, $\beta_2 = 0.9$, weight decay of 10^{-5} , and exponential learning rate decay. The dimension of z is set to $d_z = 64$ for all datasets. The weighting term β of the Kullback-Leibler divergence loss D_{KL} is set to $\beta = 1 \cdot 10^{-5}$. For the controllable video synthesis task, we discretize the conditioning ν_1 to one-hot vectors. For the 3D end effector position, the x, y and z axis is discretized into 10 bins. For the clouds, the motion direction is discretized into 36 bins. The 3D end effector position was

Method	Landscape	Fire	Vegetation	Waterfall	Clouds
AL[21]	1.49	0.36	0.30	0.80	1.22
Ours	3.41	1.42	0.98	1.11	1.51

Table 5. Diversity scores based on the I3D [73] trained on DTDB [30]. The average difference between ground-truth samples are a factor of ~ 1000 smaller for the I3D [73] network trained on DTDB [30] as the one trained Kinetics [38]. For presentation purposes, the numbers in the table have been multiplied by a factor of 1000.

provided by [20] and for the clouds [30] we manually labelled the direction. The normalizing flow, \mathcal{T}_θ , was trained using Adam [39] with a learning rate of $1 \cdot 10^{-5}$ and linear learning rate decay.

C. Evaluation Details

C.1. Diversity Metric

Besides synthesis quality, diversity is the main criteria we use to evaluate and compare stochastic video synthesis approaches. The assessment of diversity is typically based on measures utilizing feature representations of pretrained models [45, 92]. For instance, SAVP [45] uses a VGG network [72] trained for classification on ImageNet [69] to yield frame-wise representations of video sequences. Based on these representations, videos are compared based on their frame-wise differences measured using a given distance metric. The guiding intuition is that more diverse sample sets should exhibit larger feature differences on average. To this end, SAVP [45] uses the Cosine distance. We argue that this evaluation distance has a major drawback: the Cosine distance only measures the angle between feature vectors, thus discarding crucial information represented by the vector norms. For instance, two data points may lie approximately on a line (i.e., a Cosine distance of 0) but still are located far from each other. Hence, diversity is measured based on incomplete information.

To circumvent this issue, we replace the Cosine distance with the Euclidean distance which also takes the magnitude of a vector into account. Moreover, to explicitly capture temporal information, we also investigate replacing the frame-based VGG feature extractor with an I3D model [73] which directly yields representations that capture the appearance and dynamics of the entire video sequence. Tab. 6 compares the discussed diversity measures. It can be seen that independent of the diversity measure, the order of the approaches is the same. We employ both VGG MSE and I3D MSE measures in our experiments. Note that the I3D feature extractors have been trained on similar datasets as the videos to be evaluated, i.e., Kinetics [38] for human motion [47] and DTDB [30] for Landscape [87]. Moreover, we report the missing diversity scores based on the I3D [73] from the main paper on Landscape [21] and DTDB [30] in Tab. 5.

Method	VGG Cosine	VGG MSE	I3D MSE
SAVP ^{†,3} [45]	0.000	0.00	0.01
SRVP ³ [25]	0.040	0.34	1.01
IVRNN ³ [14]	0.023	0.23	0.57
Ours	0.076	0.50	1.63

Table 6. Comparison of different diversity metrics on iPER [47]. [†] SAVP experienced mode collapse due to training instabilities originating from the two involved discriminators. The VGG based feature extractors have been pretrained on ImageNet [69]. The I3D feature extractor has been pretrained on Kinetics [13]. ³ denotes models trained using the official code from their corresponding webpages.

C.2. Evaluation Protocol

For comparisons on each dataset, we use the reported numbers from the corresponding paper, where possible, otherwise we use pretrained models or train models from scratch using the code from the official webpage⁴. Here, we list the evaluation protocol for each dataset.

BAIR [20]. We follow the standard protocol [76] for computing the FVD score by evaluating videos on a sequence length of 16 on a resolution of 64×64 using all 256 test videos. Diversity is measured by predicting five future progression given the starting frames from all 256 test sequences and computing the Euclidean distance in the VGG-16 [72] as well as in the I3D [73] feature space between the corresponding generated videos.

iPER [47]. For evaluating the FVD score, we use 1000 randomly sampled sequences from the test set as well as the corresponding generations. Note, for a fair comparison, we concatenate the *last* conditioning frame to the generated rather than all conditioning frames since previous work condition on up to eight frames. This results in a sequence of length 17 for computing the FVD score. For computing the diversity, we predict five future progression for each of the 1000 test sequences and measure the diversity based on that.

Landscape [87]. We create an evaluation set by randomly sampling six times sequences of length 32 from each test video with length over 32 resulting in 918 videos. Based on these sequences, FVD, DTFVD, LPIPS, and FID are computed. As explained in the main paper, our model is trained

⁴

<https://github.com/edouardelasalles/srvp>
https://github.com/facebookresearch/improved_vrnn
https://github.com/alexlee-gk/video_prediction
https://github.com/jianwen-xie/Dynamic_generator
<https://github.com/zilongzheng/STGConvNet>
<https://github.com/endo-yuki-U/Animating-Landscape>
<https://github.com/zhangzjn/DTVNet>
<https://github.com/weixiong-ur/mdgan>

on a sequence length of 16 but applied two times by using the last predicted frame as input for the next prediction. For diversity, we again generate five future progressions for each sequence of the 918 evaluation sequences and use the same procedure described for BAIR.

DTDB [30]. We create an evaluation set by using five sequences of length 16 from each test video resulting in between 90 and 385 test sequences depending on the texture. Based on these sequences, the FVD, DTFVD, LPIPS, and FID are computed. This evaluation procedure is the same for each texture. We train one model for AL [21] as well as for our approach on each texture. For diversity, we again generate five future progressions for each sequence of the evaluation set and use the same procedure described for BAIR.

C.3. Dynamic Texture FVD (DTFVD)

In Sec. 4.3 of our main paper, we introduced a dedicated FVD metric for the domain of dynamics textures, the Dynamic Texture Fréchet Video Distance (DTFVD). To this end, we trained a network on DTDB [30] for the task of dynamic texture classification. The motivation behind introducing DTFVD is to provide an additional metric which is sensitive to the types of appearances and dynamics encapsulated by dynamic textures, rather than human action-related motions, as captured by FVD. For the DTFVD network, we use the same architecture as used for the FVD model, i.e., an I3D network [73]. At convergence (cf. Fig. 11), the DTFVD model achieved 81.7% training accuracy, while achieving 84.0% test accuracy, thus indicating that the model yields well generalizing features capturing the appearance and dynamics in DTDB. A similar conclusion can be drawn by looking at the confusion matrix in Fig. 10 computed for the test set of DTDB, which shows a dominant diagonal structure. Note, we used dropout with a probability of $p = 0.5$ to avoid overfitting, which explains why the classification performance is higher on the test set than on the training set. To evaluate sequences with lengths of 16 as well as 32 we train two separate networks.

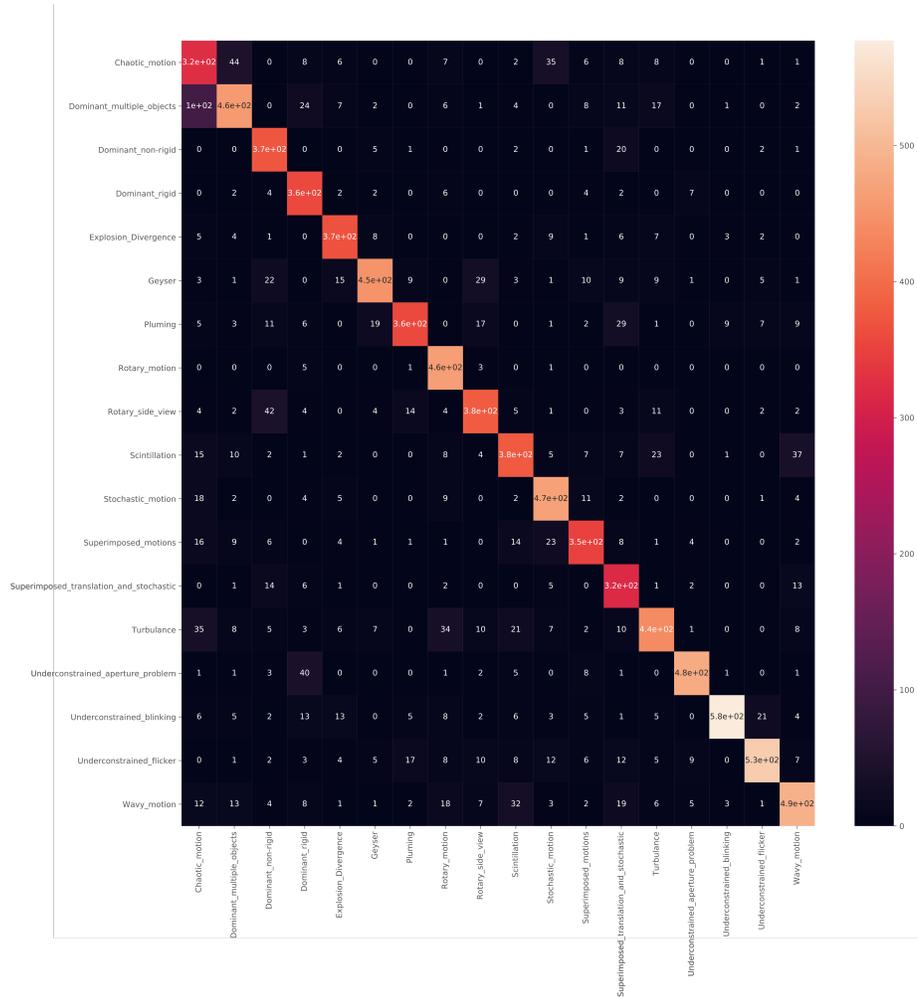


Figure 10. Confusion matrix on the test set of DTDB [30] computed from our DTFVD backbone model.

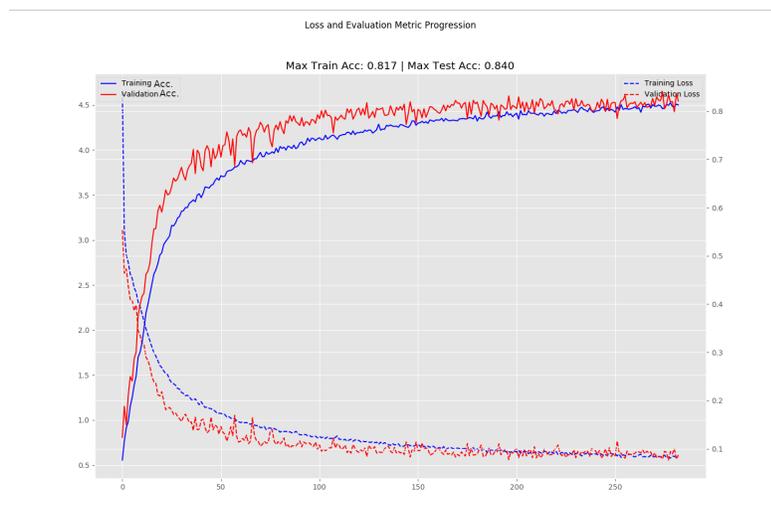


Figure 11. Training and validation loss while optimizing our DTFVD backbone network on a sequence length of 32. Similar accuracy on both dataset splits indicate a well-generalizing model.