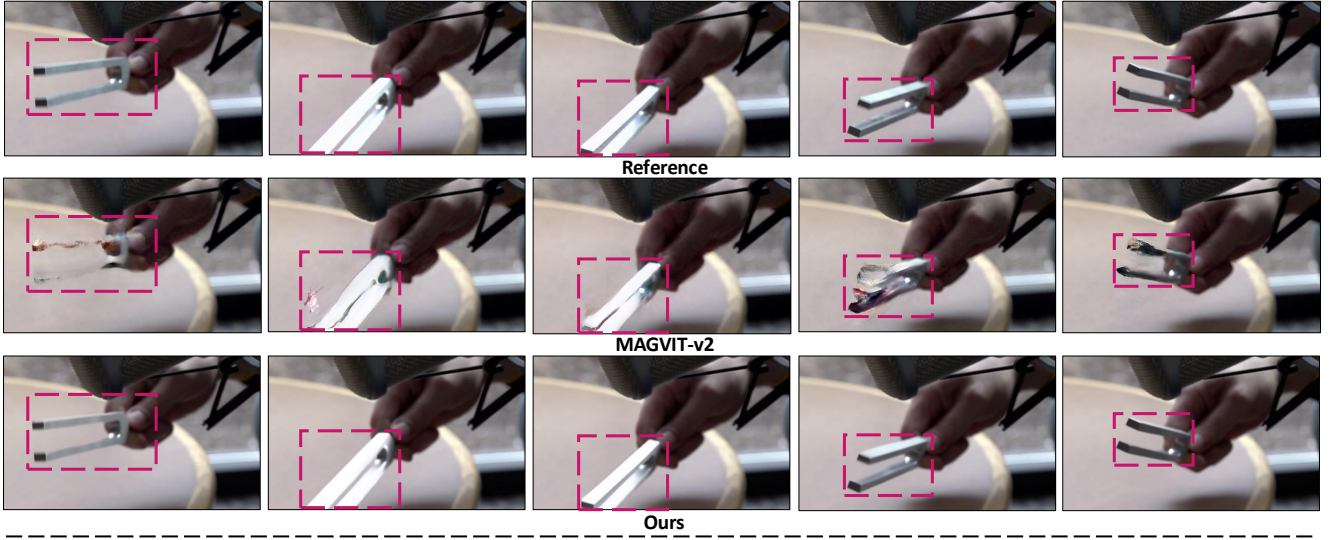


# REGEN: Learning Compact Video Embedding with (Re-)Generative Decoder

Yitian Zhang<sup>1,2†</sup> Long Mai<sup>1</sup> Aniruddha Mahapatra<sup>1</sup> David Bourgin<sup>1</sup>  
Yicong Hong<sup>1</sup> Jonah Casebeer<sup>1</sup> Feng Liu<sup>1</sup> Yun Fu<sup>2</sup>

<sup>1</sup>Adobe Research <sup>2</sup>Northeastern University

<https://bespontaneous.github.io/REGEN/>



Prompt: A sheep behind a fence looking at the camera.

Prompt: Young woman with red hair and leafy crown.

Prompt: Time lapse at the snow land with aurora in the sky.

Figure 1. Reconstruction (top) and text-to-video (T2V) generations (bottom) results at  $32\times$  temporal compression. This figure contains video results (for T2V results), best viewed with Adobe Acrobat Reader. Colored bounding boxes in reconstruction results denote the regions with the most difference where MAGVIT-v2 leads to clear artifacts in the tuning fork.

## Abstract

We present a novel perspective on learning video embedders for generative modeling: rather than requiring an exact reproduction of an input video, an effective embedder should focus on synthesizing visually plausible reconstructions. This relaxed criterion enables substantial improvements in compression ratios without compromising the

quality of downstream generative models. Specifically, we propose replacing the conventional encoder-decoder video embedder with an encoder-generator framework that employs a diffusion transformer (DiT) to synthesize missing details from a compact latent space. Therein, we develop a dedicated latent conditioning module to condition the DiT decoder on the encoded video latent embedding. Our experiments demonstrate that our approach enables superior encoding-decoding performance compared to state-of-

<sup>†</sup>This work was done during an internship at Adobe Research.

the-art methods, particularly as the compression ratio increases. To demonstrate the efficacy of our approach, we report results from our video embedders achieving a temporal compression ratio of up to  $32\times$  ( $8\times$  higher than leading video embedders) and validate the robustness of this ultra-compact latent space for text-to-video generation, providing a significant efficiency boost in latent diffusion model training and inference.

## 1. Introduction

Diffusion models [13, 37] have recently emerged as a major paradigm in generative content creation, achieving breakthroughs across various domains such as images [4, 14, 15, 31, 33], videos [28, 30, 34, 39, 45, 50], and audio [7, 16, 20, 23]. For efficiency, most modern diffusion models compress the input media into a compact latent space before modeling, commonly referred to as latent diffusion models (LDMs) [33]. A core component of the LDM is the embedding model (which we refer to as **embedder** in this paper), which typically follows an encoder-decoder architecture. The encoder transforms the raw data into a compact latent representation, while the diffusion models operate within this latent space and subsequently transform the data back to the raw space using the decoder.

For video generation, state-of-the-art (SOTA) video embedder, exemplified by the continuous-token variants of the MAGVIT-v2 model [46], leverage architectures that enable spatiotemporal processing, offering compression in both spatial and temporal dimensions. Despite these advances, SOTA video embedders typically offer  $8\times$  spatial compression but only  $4\times$  temporal compression.

Increasing compression is inherently challenging due to a fundamental trade-off between representational efficiency and reconstruction fidelity. At high compression ratios, the model struggles to retain sufficient information within the latent space to reproduce high-fidelity details accurately. We verify that adapting existing video embedders to support temporal compression beyond  $4\times$  is highly challenging without significantly degrading reconstruction quality.

This paper presents an alternative perspective on learning video tokenizers tailored for latent diffusion modeling. We propose a generation-oriented view of latent representation learning. *“We argue that, in the context of latent diffusion, the key desired property of the latent space is the ability to generate visually plausible content rather than to faithfully recover an input video.”* We propose transforming the traditional encoder-decoder video embedder into an encoder-generator. This generation-oriented perspective enables a more flexible approach to allowing high-compression in the latent space. By leveraging generative capabilities in the decoder, the encoder can focus only on preserving the essential semantic and structural information. The decoder then synthesizes realistic finer details.

To this end, we introduce **REGEN**, an approach for learning a high-compression video embedder via a diffusion decoder with a resynthesis objective. REGEN utilizes a **diffusion Transformer (DiT)** [31] as the video decoder due to its exceptional video modeling capability compared to diffusion U-Net, as consistently demonstrated by impactful works [21, 28, 30, 50]. Our DiT decoder treats the latent features encoded from an input video sequence as the conditioning signal during the diffusion-based generation process. Both the encoder and the DiT decoder are jointly trained within the diffusion-based training framework to learn to generate the original video content.

To connect the encoded latent features with the DiT decoder, we design a dedicated **latent conditioning module**. This module introduces a novel DiT conditioning mechanism that transforms the latent maps into a content-aware positional embedding, enabling strong conditioning signals for the decoder’s diffusion process. Our proposed latent conditioning mechanism also naturally enables continuous-time decoding that supports not only reconstruction but also interpolation and extrapolation capabilities. More importantly, it also addresses the challenge of encoding and decoding for arbitrary resolutions and aspect ratios, unseen during training, which is difficult to achieve with DiT.

In this work, we aim to enhance the compactness of the latent space beyond the current limits and focus on the axis of increasing the temporal compression of our embedder, keeping the spatial compression fixed at  $8\times$ . Across all the experiments, we follow the design of MAGVIT-v2 [46] to keep the number of latent channels to be 8 for all compression ratios, as we want to isolate the effect of increasing temporal compression on reconstruction quality.

The contributions of the current work are as follows:

1. We present the idea of using a Diffusion Transformer as a decoder in video modeling instead of the conventional VAE [19]-based encoder-decoder. We demonstrate that this modeling strategy allows a re-design of the learning objective in video encoding, allowing us to bypass the fundamental compression-reconstruction trade-off.
2. We introduce a novel latent conditioning module with a content-aware positional encoding formulation. This enables us to effectively turn the encoded latent features into spatiotemporal control signals for the Diffusion Transformer (DiT) decoder. Our decoder not only enables DiT to encode and decode videos of arbitrary aspect ratios and resolutions, but also supports even one-step sampling without utilizing external distillation.
3. We show that REGEN enables superior reconstruction performance compared to extending the VAE-based embedders at high compression ratios. Our method even outperforms even very recent SOTA video embedders at commonly used  $4\times$  temporal compression that have the same number of latent channels as our method. We also

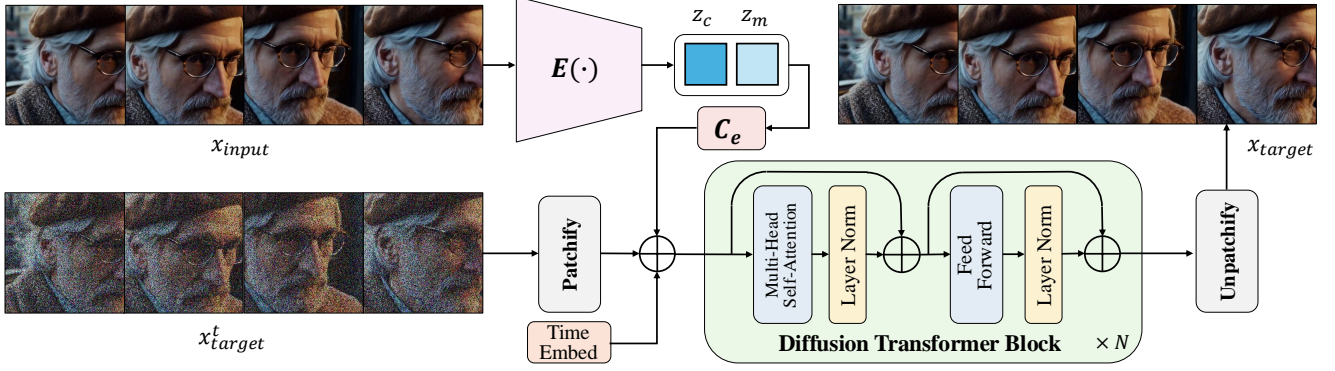


Figure 2. **Overall framework.** Our spatiotemporal video encoder  $E(\cdot)$  encodes the input video  $x_{input}$  into two latent frames, content and motion ( $z_c, z_m$ ). They are processed by the latent expansion module  $C_e$  and serve as conditioning for the generative decoder.

show that our highest-compression ( $32\times$  temporal) latent space is friendly for text-to-video generation.

## 2. Related work

**Video Embedders.** Early video diffusion models adapted image LDM, using the same image-based embedder for per-frame latent feature extraction [2, 9, 12, 26, 35, 44]. This frame-wise approach neglects temporal relationships, causing inconsistencies and restricting the compactness of the video latent. Recent video diffusion frameworks utilize spatiotemporal embedders for joint processing. The seminal work of MAGVIT-v2 [46] employs causal 3D convolutions for compression across both spatial and temporal dimensions. Originally designed for discrete tokenization, MAGVIT-v2 has been adapted for continuous-token video encoding, facilitating its integration into various prominent video diffusion models [21, 45, 50]. Despite their effectiveness, they typically offer limited temporal compression (e.g.,  $4\times$ ). We aim to enhance the compactness of the latent space beyond the current limits, extending up to  $32\times$  temporal compression. To accomplish it, we offer a novel perspective of using encoder-generator with a (re-)generative diffusion decoding process to escape the reconstruction-compression trade-off.

**Diffusion Autoencoders.** Our work draws inspiration from recent studies that demonstrate the potential of image diffusion models to generate content conditioned on abstract feature vectors extracted from conditioning images [3, 17, 32]. Specifically, [32] shows that an image diffusion model can be semantically guided, akin to style-code manipulation in StyleGAN [18]. Similarly, [3] validates the improvement brought by generative decoders, and [17] highlights the ability to reconstruct images from conditioning feature vectors. We adopt this latent-conditioning diffusion concept to the context of video embedder learning. Complementary to them, to the best of our knowledge, we are the first to explore diffusion transformer autoencoders in the context of

learning highly compact latent spaces for videos.

## 3. Method

Our model consists of two major components. First, the spatiotemporal video encoder that projects the input video sequence into a compact latent space (Section 3.1). The second is the DiT-based generative decoder with that converts the latents back to pixels by taking the latents as a conditioning signal (Section 3.2). The whole model is trained from scratch in an end-to-end fashion with the diffusion objective. Figure 2 illustrates our overall framework.

### 3.1. Spatiotemporal Video Encoder

The goal of our encoder is to encode a video into a compact latent space. Inspired by MAGVIT-v2 [46], we adopt the continuous version of their encoder design due to its ability to encode both images and videos in consistent latent space.

The encoder consists of multiple 3D convolution blocks, which are causal in the temporal dimension. For a video of  $T+1$  frames with dimensions  $H \times W$ , at a spatial downsampling factor of  $m$ , and a temporal downsampling rate of  $k$ , due to the causal nature, the encoder produces a sequence of  $1 + \frac{T}{k}$  latent feature maps each of dimension  $\frac{H}{m} \times \frac{W}{m}$ . To provide flexibility in consistently encoding long videos without going out-of-memory, we follow a chunk-wise encoding scheme that encodes videos of fixed-length frames (a *chunk*). Thus, at a desired temporal compression ratio  $k$ , we encode any input video chunk  $x_{input}$  of length  $k+1$  into two latent frames:

$$z_c, z_m = E(x_{input}), \quad (1)$$

where  $z_c$  and  $z_m$  denote the resulting two latent maps. As the encoder is causal in nature,  $z_c$  only contains information from the first frame. We call this the **content** latent frame.  $z_m$ , contains the compressed motion information of the rest of the frames. We refer to this as the **motion** latent frame. We use 8 latent channels for both  $z_c$  and  $z_m$ .

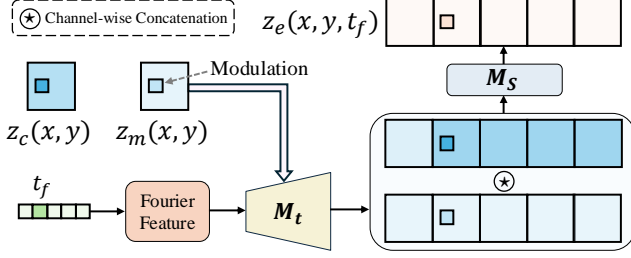


Figure 3. **Latent conditioning module  $C_e$ .** The SIREN network  $M_t$  maps the time coordinate  $t_f$  to a feature vector modulated by the motion latent  $z_m$ . The resulting feature is concatenated with the feature value of  $z_c$  at the corresponding spatial coordinate  $(x, y)$ . The concatenated feature is mapped into the DiT hidden dimension by the projector  $M_s$ . We utilize the first frame prediction from SIREN to replace the first frame of expanded  $z_c$  to ensure consistent representation for both image and video inputs.

### 3.2. Diffusion Transformer Decoder

As opposed to the conventional convolution-based decoder, we use a transformer [41] and model the decoding task as a conditional diffusion process. Given the input video sequence  $x_{input}$  and its corresponding content and motion latents  $[z_c, z_m]$ , the diffusion decoder  $G_d$  is trained to (re-)generate a target video output  $x_{target}$  from a noise map sequence  $N$ , conditioned on the latent maps  $[z_c, z_m]$ . For the reconstruction task,  $x_{target}$  is chosen to be the same as  $x_{input}$ . It’s worth noting that this formulation inherently allows for a more flexible definition of  $x_{target}$  to handle different tasks beyond pure reconstruction. For example, setting  $x_{target}$  to be the temporally upsampled version of  $x_{input}$  corresponds to modeling the temporal interpolation, while setting  $x_{target}$  to be a time-shifted version of  $x_{input}$  leads to the temporal extrapolation.

During training, we corrupt the clean data  $x_{target}$  (denoting as  $x_{target}^0$  for diffusion noise timestep  $t = 0$ ) with the forward diffusion process:

$$q(x_{target}^t | x_{target}^0) = \mathcal{N}(x_{target}^t; \sqrt{\alpha^t} x_{target}^0, (1 - \alpha^t)I), \quad (2)$$

where  $t$  is the diffusion time step, and  $\alpha^t$  is determined via the noise scheduling procedure [24]. The sample at time step  $t$  is obtained with the reparameterization trick:

$$x_{target}^t = \sqrt{\alpha^t} x_{target}^0 + \sqrt{1 - \alpha^t} \epsilon, \quad (3)$$

where  $\epsilon \sim \mathcal{N}(0, I)$ . In the reverse denoising process, the generative decoder learns to invert the forward corruption under the condition  $[z_c, z_m]$ . Specifically, a denoising model  $\epsilon_\theta$  with parameters  $\theta$  is trained to predict the noise  $\epsilon$  at each time step  $t$  given the corresponding noisy version  $x_{target}^t$ , conditioning on  $[z_c, z_m]$  with the simplified learning objective:

$$\mathcal{L}(\theta) = \|\epsilon - \epsilon_\theta(x_{target}^t, [z_c, z_m])\|^2, \quad (4)$$

Under this formulation, we train our video encoder along with the generative decoder in an end-to-end fashion. Our diffusion decoder leverages a DiT architecture for the denoiser backbone due to its superior modeling and scaling capabilities compared to U-Net as pointed to many prior works [21, 28, 30, 50]. Our decoder operates in pixel space, at a patch size of  $p$ . The patch size  $p$  is determined according to the spatial down-sampling ratio in the encoder, which in our experiments is  $8\times$ . Thus, we keep  $p = 8$ .

**Latent Conditioning via Content-Aware Positional Encoding.** The transformer model for video data usually incorporates spatiotemporal positional encoding (PE) of the space-time coordinates as additional inputs. PE is critical not only for addressing the order-invariant nature of transformers but also for capturing the spatiotemporal structure of the generated content. Existing DiT models often adopt a fixed PE scheme where the coordinate-to-embedding mapping. Such a fixed PE scheme struggles with generalizing to input sizes unseen during training [45]. This is particularly problematic in our context, where we leverage DiT as the decoder since the model after training should be applied to decode the latents extracted from inputs of different resolutions and aspect ratios.

We introduce a new conditioning mechanism to address this challenge. The key idea is to generate the positional embedding from the conditioning latent  $[z_c, z_m]$  instead of using a fixed spatial PE. Specifically, we develop a latent expansion module  $C_e$  that takes the condition  $[z_c, z_m]$  as inputs and expand it to the expanded latent  $z_e$  same spatiotemporal dimension as the target output  $x_{target}$ .  $z_e$  is then added to the token embedding and timestep embedding and input into the DiT. In this way,  $C_e$  operates as a mini-decoder that decodes the latent into a full spatiotemporal form, and  $z_e$  operates as a content-aware positional embedding that participates in controlling the spatiotemporal aspects in the synthesized videos.

As a positioning encoding module,  $C_e$  is designed to map the coordinates  $(x, y, t_f)$  of a token location (after patchification) in DiT input space into an embedding vector where the mapping is conditioned on  $[z_c, z_m]$ . We defined such mapping as:

$$C_e(x, y, t_f | [z_c, z_m]) = M_s(z_c(x, y) \otimes M_t(t_f | z_m(x, y))) \quad (5)$$

where  $\otimes$  denotes the channel-wise concatenation operation. Note that the  $(x, y)$  coordinates in the DiT input space match with those in the feature maps  $[z_c, z_m]$  as we matched the patchification’s patch size with the spatial downsampling factors.  $M_s$  is a projector to align the channel dimension, consisting of a single linear layer and one RMSNorm [47].  $M_t$  is a neural network sub-module that maps the time coordinate  $t_f$  into a feature vector. Inspired by the success of implicit neural representation in modeling video data [6, 27], we implement  $M_t$  with the SIREN [36]. We



Method	Compression	MCL-JCV						DAVIS 2019					
		256×256			512×512			256×256			512×512		
		PSNR	SSIM	rFVD↓	PSNR	SSIM	rFVD↓	PSNR	SSIM	rFVD↓	PSNR	SSIM	rFVD↓
● MAGVIT-v2	8×8×8	26.61	0.718	105.72	29.14	0.771	72.07	22.82	0.602	183.52	24.75	0.660	125.03
● REGEN		<b>28.82</b>	<b>0.785</b>	<b>85.37</b>	<b>32.74</b>	<b>0.846</b>	<b>29.88</b>	<b>26.00</b>	<b>0.711</b>	<b>152.46</b>	<b>29.34</b>	<b>0.778</b>	<b>89.98</b>
● MAGVIT-v2	8×8×16	25.06	0.672	205.75	26.62	0.717	185.69	20.62	0.527	441.24	21.21	0.572	417.43
● REGEN		<b>27.27</b>	<b>0.736</b>	<b>174.29</b>	<b>30.41</b>	<b>0.798</b>	<b>92.48</b>	<b>23.85</b>	<b>0.635</b>	<b>328.83</b>	<b>26.27</b>	<b>0.699</b>	<b>235.13</b>
● MAGVIT-v2	8×8×32	22.97	0.573	536.01	†	†	†	18.23	0.419	1080.15	†	†	†
● REGEN		<b>26.05</b>	<b>0.695</b>	<b>265.96</b>	<b>28.71</b>	<b>0.758</b>	<b>224.56</b>	<b>22.20</b>	<b>0.575</b>	<b>488.89</b>	<b>23.49</b>	<b>0.625</b>	<b>522.20</b>

Table 1. **Reconstruction comparison at high temporal compression.** We compare our method, REGEN, with MAGVIT-v2 at different compression rates on MCL-JCV and DAVIS 2019 datasets. The best results are bold-faced. †MAGVIT-v2 (32×) faces out of memory issue at 512×512, due to the 3D convolution layers in decoder.

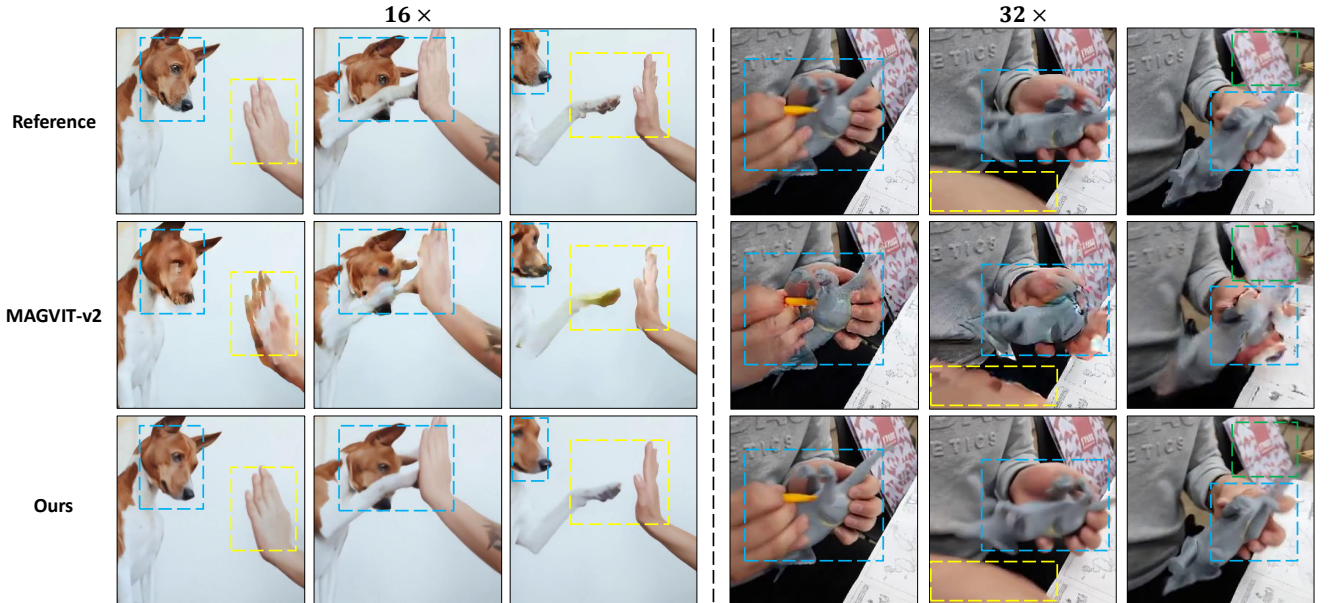


Figure 4. **Effectiveness of REGEN at high temporal compression.** MAGVIT-v2 suffers from strong temporal artifacts in regions of high motion, such as the dog’s face (left) and the toy (right). Areas enclosed in boxes show regions of maximum difference.

condition the mapping in  $M_t$  on the motion feature  $z_m$  by modulating the SIREN network with the feature values of  $z_m$  at the queried  $(x, y)$  location. Thus, the positional information is integrated into the expanded latent so that we can entirely remove the original spatial and temporal positional embedding in DiT, enabling our generative decoder to be generalize across arbitrary resolutions and aspect ratios.

## 4. Experiments

We evaluate our encoder-generative decoder framework in three aspects: (1) effectiveness of generative decoder at various temporal compression ratios, and (2) the compatibility of our compressed latent space for video generation.

### 4.1. Effectiveness of REGEN as Video Embedder

**Implementation Details.** We configure our DiT-based decoder to match the model latency of MAGVIT-v2 [46]: it consists of 24 Transformer blocks, with 16 heads and

a hidden dimension of 2048 using full spatiotemporal self-attention. We set the patch size to be 8 to accommodate the costs. For each compression ratio, we train our model for  $\sim 100K$  iterations under the scenarios of reconstruction, interpolation, and extrapolation objectives. Our encoder shares the same architecture with our implemented MAGVIT-v2. Please refer to the supplementary material for details.

**Evaluation.** Following MAGVIT-v2, we evaluate the reconstruction quality on 2 benchmarks: MCL-JCV [42] and DAVIS 2019 (full resolution) [5]. In both datasets, videos are rescaled and center-cropped to 256x256 and 512x512 resolution for evaluation. We use standard quantitative metrics, PSNR, SSIM, and Fréchet video distance (rFVD) [10, 40], to examine the reconstruction quality. This covers both pixel-based and perceptual quality metrics.

**Baselines.** At 8×8×4 compression, we compare REGEN with SOTA 8-channel video embedders since the number of latent channels greatly affects reconstruction quality. We

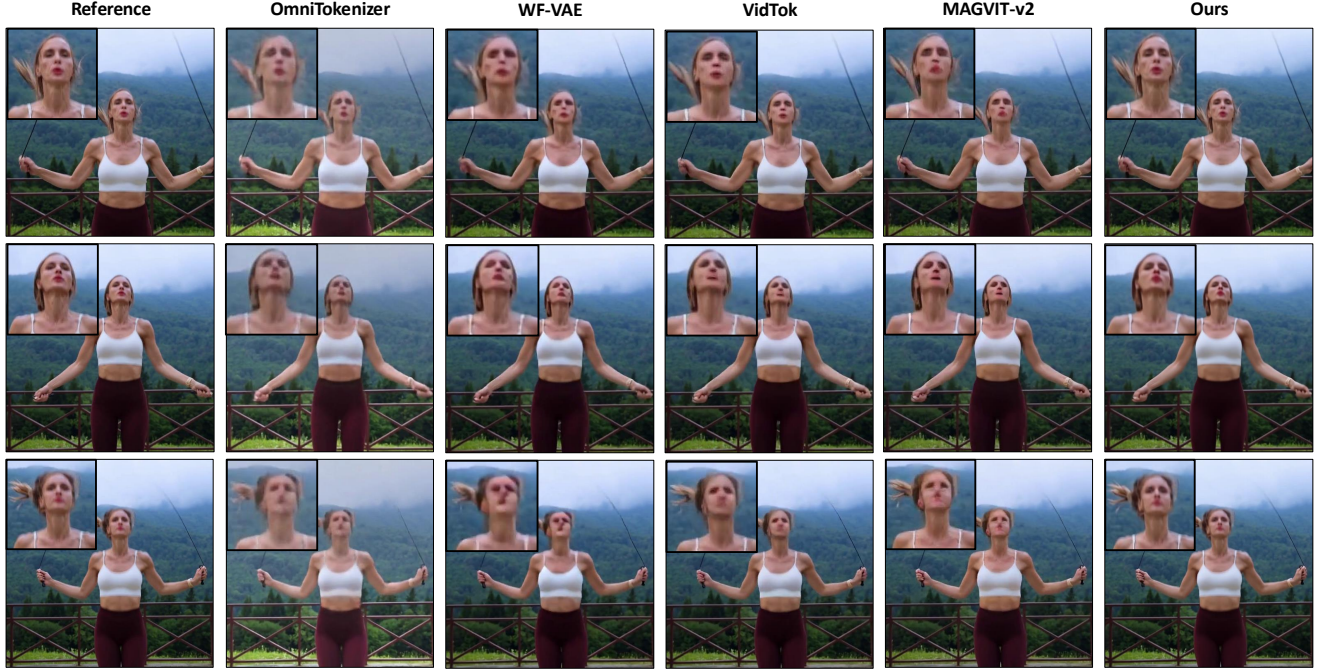


Figure 5. **Effectiveness of REGEN at base  $4\times$  temporal compression.** Current video embedders suffer from ghosting artifacts for videos with large motion, especially in faces (last row). REGEN performs well and retains plausible spatiotemporal structures from the input.

Method	MCL-JCV			DAVIS 2019		
	PSNR	SSIM	rFVD↓	PSNR	SSIM	rFVD↓
Omni	24.63	0.710	93.35	23.39	0.628	152.01
WF-VAE	31.00	0.804	55.01	27.95	0.737	107.67
VidTok	<b>32.06</b>	<b>0.836</b>	38.85	28.67	<b>0.760</b>	67.24
MAGVIT-v2	31.49	0.829	<u>28.63</u>	28.16	0.758	<u>56.46</u>
REGEN	<b>32.94</b>	<b>0.857</b>	<b>22.40</b>	<b>30.25</b>	<b>0.801</b>	<b>48.38</b>

Table 2. **Reconstruction comparison at base  $4\times$  temporal compression.** We compare REGEN with various SOTA 8-channel video embedders at  $4\times$  temporal compression on MCL-JCV and DAVIS 2019 datasets under  $512\times 512$  inputs. The best results are bold-faced and the second best results are underlined.

compare against OmniTokenizer [43], VidTok [38] and, WF-VAE [22] and MAGVIT-v2 [46]. Since the MAGVIT-v2 weights are not released, we reimplement it and train it on our dataset with over 200K iterations. For higher compression experiments, we adapt MAGVIT-v2’s design to handle different temporal compression ratios. Note that we adjust the location of upsampling layers in MAGVIT-v2 decoder at  $8\times 8\times 32$  to avoid out of GPU memory during training and we follow the original design in other experiments to ensure the optimal reconstruction quality.

**REGEN’s effectiveness at high compression.** We investigate our primary objective, whether a diffusion-based encoder-generator framework can learn a very compact latent space with a high degree of temporal compression that can faithfully reconstruct the input much more effectively than traditional video embedders. As shown in Tab. 1,

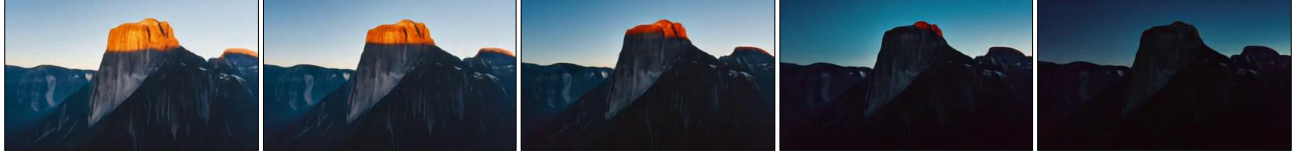
MAGVIT-v2 lags behind our method on all metrics at different temporal compression rates. Moreover, the advantage of REGEN over MAGVIT-v2 widens with the increasing temporal compression ratios, which suggests that naively extending conventional video embedders to high compression rates cannot deliver satisfying results and highlights the generation ability of our decoder. Shown in Fig. 4, MAGVIT-v2 suffers from significant degradation in reconstruction quality and often exhibits visual artifacts at high temporal compressions. In contrast, our method maintains the reconstruction ability to a good extent. Moreover, our decoder is more efficient in terms of GPU memory when decoding long sequences compared to MAGVIT-v2, which is quite important in practical applications.

**Comparison with SOTA video embedders at base  $4\times$  temporal compression.** Although our goal is to strive for higher temporal compression instead of SOTA performance at  $8\times 8\times 4$  compression, we compare REGEN with existing 8-channel video embedders in Tab. 2 to validate the effectiveness of our method. We observe that our implemented MAGVIT-v2 achieves compelling performance compared to recent SOTA methods VidTok [38], WF-VAE [22], and OmniTokenizer [43], which proves the soundness of our experimental setup. Nevertheless, all the approaches lag behind our method in all metrics, and the qualitative comparison. In Fig. 5, we see that REGEN preserves better spatiotemporal structures from the input, especially in the human faces, which other convolution-based embedders, designed explicitly for  $4\times$  compression ratio, fail to recon-

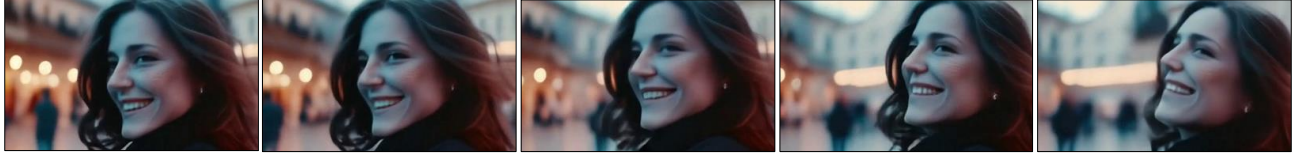




Prompt: A slow cinematic push in on an ostrich standing in a 1980s kitchen.

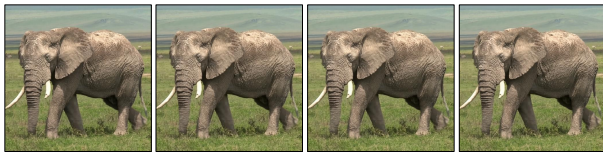


Prompt: A stunning aerial drone footage time lapse of El Capitan in Yosemite National Park at sunset.

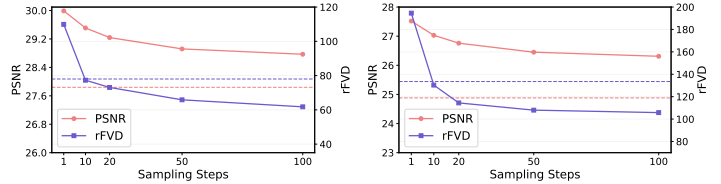


Prompt: A cinematic documentary hand held close up of a woman standing in a busy Italian plaza smirking to herself.

Figure 6. **Text-to-video generation results on our ultra-compact latent space with  $32\times$  temporal compression.** The latent diffusion model generates 132-frame videos with only 8 latent frames on our ultra-compact latent space, offering  $\sim 5\times$  reduction in the number of latent frames compared to current video embedders at  $4\times$  temporal compression.



(a) Example reconstruction results at different sampling steps.



(b) Quantitative results on MCL-JCV.

(c) Quantitative results on DAVIS 2019.

Figure 7. **Reconstruction results at different sampling steps.** Dash lines in (b) and (c) denote the PSNR and rFVD of MAGViT-v2. Our method exhibits strong performance across varying sampling steps and supports few-step and even single-step sampling.

struct accurately. We further provide a holistic comparison with other video embedders that have different latent channels in supplementary material.

## 4.2. Text-to-Video Generation

The previous experiments demonstrate that the compact video embeddings from our methods can achieve promising performance for video reconstruction. We now verify whether our compact latent space is friendly to text-to-video generation, which is our end objective for learning a compact latent space. To validate it, we train a 5B DiT-based latent diffusion model on the ultra-compact latent with  $32\times$  temporal compression for text-to-video generation. The design of the text-to-video model is based on MMDiT [8]. We first train the model with image inputs for 260K iterations, and the model is further trained with mixed image and video inputs for another 80K iterations. We show samples of the generated videos in Fig. 6. Note that due to constraints in computing resources, we trained the model at a relatively small scale with a small number of iterations. Nevertheless, even with small-scale training, we demonstrate that the diffusion model is capable of gen-

erating plausible video content on a much more compact latent space, promising the application of our method in video generation. Notably, the latent diffusion model generates 132-frame videos with only 8 latent frames on our ultra-compact latent space ( $32\times$  temporal compression), offering  $\sim 5\times$  reduction in the number of latent frames compared to current video embedders that operate with  $4\times$  temporal compression, thus reducing the training and inference costs remarkably. Please find more results in the supplementary.

## 5. Discussion

In this section, we address two of the main concerns with the practical usability of our transformer-based diffusion model for embedders: (1) the diffusion model inherently requires many denoising steps at inference, making it very costly, and (2) transformers are difficult to generalize to resolutions and aspect ratios unseen during training. We also provide a method to alleviate the chunking issue by leveraging the extrapolation ability of our decoder.

**Few-step and one-step sampling.** One of the main bottlenecks with the diffusion model is that it requires a num-

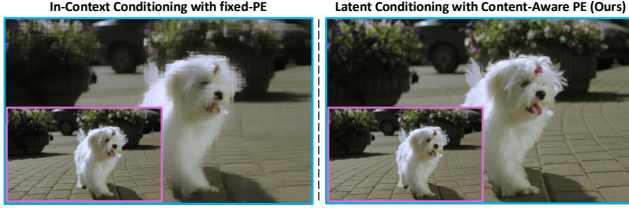


Figure 8. **Video frame reconstruction at different spatial resolutions.** The models are trained under  $192 \times 320$  input and evaluated at the resolution of  $192 \times 320$  (pink bounding boxes) and  $384 \times 640$  (blue bounding boxes). In-context conditioning results in gridding artifacts at a larger resolution, while our method exhibits strong generalization due to the proposed content-aware PE.

Method	192×320			384×640		
	PSNR	SSIM	rFVD↓	PSNR	SSIM	rFVD↓
In-context	25.71	0.709	135.89	23.39	0.587	441.98
Ours	<b>26.04</b>	<b>0.720</b>	<b>128.80</b>	<b>29.41</b>	<b>0.785</b>	<b>57.01</b>

Table 3. **Ablation of conditioning mechanism on DAVIS 2019 dataset.** In-context conditioning exhibits inferior performance compared to ours at the **training resolution** and cannot generalize well to **unseen resolutions**. The best results are bold-faced.

ber of denoising steps, even at inference, to achieve good-quality results, which causes a lot of overhead for practical use. However, our decoder is in fact tasked with an easier generation problem compared to text-to-video generation as it is equipped with a very strong conditioning signal provided by our latent conditioning module. Thus, we perform inference on the same model with different numbers of steps to inspect how significantly the reconstruction quality degrades as the number of sampling steps reduces. Fig. 7a shows reconstruction examples obtained at different sampling steps using our  $4 \times$  compression model. Interestingly, we observe that it is possible to reduce the number of sampling steps significantly (even down to 1-step sampling) without noticeable visual quality degradation. We further evaluate our model at different sampling steps on all test videos from DAVIS 2019 and MCL-JCV datasets. Fig. 7b and Fig. 7c suggest that more sampling steps lead to slightly worse PSNR and better rFVD scores, which can be attributed to the increase in sharpness in the reconstructed videos. To conclude, our diffusion decoder offers a flexible decoding scheme and can even act like a feed-forward model without utilizing external distillation, promising its application for practical purposes.

**Ablation of conditioning mechanism.** A significant challenge with using a transformer as a decoder is that conventional transformers with fixed positional encoding (PE) can not generalize to unseen resolutions at inference. In contrast, our content-aware positional encoding makes REGEN inherently flexible to varying aspect ratios and resolutions at inference, unseen during training. We verify that with



Figure 9. **Alleviating the chunking issue with latent extension.** The x-t slice is obtained by extracting a short segment (shown as the red line in the video frame) from 2 chunk frames and latent extension offers a smoother transition across the chunk boundary.

an experiment where we replace our condition module with the in-context conditioning design in DiT [31]. This conditioning scheme retains the fixed-PE scheme like other conventional conditioning designs and performs conditioning via concatenating the conditioning tokens (*i.e.* the encoded latents  $[z_c, z_m]$  in our context) with the input tokens along the sequence dimension. Fig. 8 shows that our method generalizes well to different resolutions. In contrast, in-context conditioning leads to strong gridding artifacts when evaluating at higher resolutions not seen during training, severely impairing its functionality as a video decoder. The quantitative results in Tab. 3 further validate this.

**Alleviating the chunking issue.** Like other conventional video embedders, REGEN sometimes exhibits slight jumps across the intersection of two chunks due to the chunk-wise encoding scheme, which is widely adapted to process arbitrary video lengths. To mitigate this issue, we leverage the extrapolation capability of the decoder to generate the last chunk frames based on the latent of the current chunk. Following SDEdit [29], we utilize the predicted last frame from the previous chunk to guide the generation of the next chunk so that we can better align consecutive chunks and reduce the jumping in an autoregressive manner. Figure 9 demonstrates that our latent extension strategy mitigates the jumping and offers a smoother transition across the chunk boundary. Note the visible horizontal line in the middle time-slice visualization in the Vanilla decoding, which becomes much smoother when the latent extension is applied.

## 6. Conclusion and Limitation

In this work, we provide a new perspective on constructing the latent space for video generation. We propose that a good latent space needs to be able to generate visually plausible content that faithfully recovers an input video without needing to match the input to the pixel-perfect level of detail. To this end, we present an encoder-generator framework, instead of the commonly-used encoder-decoder framework trained with VAE objective, where we leverage a diffusion transformer as the decoder to ensure the latent space only captures essential semantic and structural features of the input video and the details are synthesized during decoding by the generative decoder. We show that this



approach allows us to achieve a temporal compression ratio as high as  $32\times$  with reconstruction quality far surpassing existing state-of-the-art conventional video embedders trained for such an aggressive level of compactness.

**Limitation.** While our method enables highly compact video embedding, there still remain some challenges and future directions: (1) While one-step (or few-step) sampling is efficient, the training of diffusion-based decoders remains computationally expensive. Due to limited resources, we could not conduct extensive ablation on architecture configuration with the number of transformer blocks and other architectural details. Investigating this is a future direction. (2) Although our method of latent extension could mitigate the transition challenge between chunks to great extent, it cannot fully remove the jumping issue in such a training-free manner. (3) Our decoder operates in the pixel space, requiring a large patch size for efficiency (8 in our case), which can degrade generation quality [31] compared to the commonly used patch size of 1 or 2.

## References

- [1] Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chattopadhyay, Yongxin Chen, Yin Cui, Yifan Ding, et al. Cosmos world foundation model platform for physical ai. *arXiv preprint arXiv:2501.03575*, 2025. 11
- [2] Jie An, Songyang Zhang, Harry Yang, Sonal Gupta, Jia-Bin Huang, Jiebo Luo, and Xi Yin. Latent-shift: Latent diffusion with temporal shift for efficient text-to-video generation. *arXiv preprint arXiv:2304.08477*, 2023. 3
- [3] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023. 3
- [4] blackforestlabs. Flux. <https://blackforestlabs.ai/>, 2024. Accessed: 2024-XX-XX. 2
- [5] Sergi Caelles, Jordi Pont-Tuset, Federico Perazzi, Alberto Montes, Kevis-Kokitsi Maninis, and Luc Van Gool. The 2019 davis challenge on vos: Unsupervised multi-object segmentation. *arXiv preprint arXiv:1905.00737*, 2019. 5
- [6] Zeyuan Chen, Yinbo Chen, Jingwen Liu, Xingqian Xu, Vidit Goel, Zhangyang Wang, Humphrey Shi, and Xiaolong Wang. Videoir: Learning video implicit neural representation for continuous space-time super-resolution. In *CVPR*, 2022. 4
- [7] Jade Copet, Felix Kreuk, Itai Gat, Tal Remez, David Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre Défossez. Simple and controllable music generation. *NeurIPS*, 2024. 2
- [8] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024. 7
- [9] Songwei Ge, Seungjun Nah, Guilin Liu, Tyler Poon, Andrew Tao, Bryan Catanzaro, David Jacobs, Jia-Bin Huang, Ming-Yu Liu, and Yogesh Balaji. Preserve your own correlation: A noise prior for video diffusion models. In *ICCV*, 2023. 3
- [10] Songwei Ge, Aniruddha Mahapatra, Gaurav Parmar, Jun-Yan Zhu, and Jia-Bin Huang. On the content bias in fréchet video distance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7277–7288, 2024. 5
- [11] Yoav HaCohen, Nisan Chiprut, Benny Brazowski, Daniel Shalem, Dudu Moshe, Eitan Richardson, Eran Levin, Guy Shiran, Nir Zabari, Ori Gordon, et al. Ltx-video: Realtime video latent diffusion. *arXiv preprint arXiv:2501.00103*, 2024. 11
- [12] Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffusion models for high-fidelity long video generation. *arXiv preprint arXiv:2211.13221*, 2022. 3
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 2020. 2
- [14] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 2
- [15] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *NeurIPS*, 2022. 2
- [16] Qingqing Huang, Daniel S Park, Tao Wang, Timo I Denk, Andy Ly, Nanxin Chen, Zhengdong Zhang, Zhishuai Zhang, Jiahui Yu, Christian Frank, et al. Noise2music: Text-conditioned music generation with diffusion models, march 2023a. *arXiv preprint arXiv:2302.03917*. 2
- [17] Drew A Hudson, Daniel Zoran, Mateusz Malinowski, Andrew K Lampinen, Andrew Jaegle, James L McClelland, Loic Matthey, Felix Hill, and Alexander Lerchner. Soda: Bottleneck diffusion models for representation learning. In *CVPR*, 2024. 3
- [18] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019. 3
- [19] Diederik P Kingma, Max Welling, et al. Auto-encoding variational bayes, 2013. 2
- [20] Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi Parikh, Yaniv Taigman, and Yossi Adi. Audiogen: Textually guided audio generation. *arXiv preprint arXiv:2209.15352*, 2022. 2
- [21] PKU-Yuan Lab and Tuzhan AI etc. Open-sora-plan, 2024. 2, 3, 4, 12
- [22] Zongjian Li, Bin Lin, Yang Ye, Liuhan Chen, Xinhua Cheng, Shenghai Yuan, and Li Yuan. Wf-vae: Enhancing video vae by wavelet-driven energy flow for latent video diffusion model. *arXiv preprint arXiv:2411.17459*, 2024. 6, 12
- [23] Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. Audioldm: Text-to-audio generation with latent diffusion models. *arXiv preprint arXiv:2301.12503*, 2023. 2

- [24] Xingchao Liu, Xiwen Zhang, Jianzhu Ma, Jian Peng, et al. InstafLOW: One step is enough for high-quality diffusion-based text-to-image generation. In *ICLR*, 2023. 4
- [25] I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 11
- [26] Zhengxiong Luo, Dayou Chen, Yingya Zhang, Yan Huang, Liang Wang, Yujun Shen, Deli Zhao, Jingren Zhou, and Tieniu Tan. Videofusion: Decomposed diffusion models for high-quality video generation. *arXiv preprint arXiv:2303.08320*, 2023. 3
- [27] Long Mai and Feng Liu. Motion-adjustable neural implicit video representation. In *CVPR*, 2022. 4
- [28] Willi Menapace, Aliaksandr Siarohin, Ivan Skorokhodov, Ekaterina Deyneka, Tsai-Shien Chen, Anil Kag, Yuwei Fang, Aleksei Stoliar, Elisa Ricci, Jian Ren, et al. Snap video: Scaled spatiotemporal transformers for text-to-video synthesis. In *CVPR*, 2024. 2, 4
- [29] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021. 8, 11
- [30] OpenAI. Sora. <https://openai.com/index/video-generation-models-as-world-simulators/>, 2024. Accessed: 2024-XX-XX. 2, 4
- [31] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *ICCV*, 2023. 2, 8, 9
- [32] Konpat Preechakul, Nattanat Chatthee, Suttisak Wizatwongsa, and Supasorn Suwajanakorn. Diffusion autoencoders: Toward a meaningful and decodable representation. In *CVPR*, 2022. 3
- [33] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 2
- [34] Runway. Runway gen-3: Next-generation video generation model. <https://runwayml.com>, 2023. Accessed: 2023-XX-XX. 2
- [35] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiuyan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022. 3
- [36] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *NeurIPS*, 2020. 4
- [37] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, 2015. 2
- [38] Anni Tang, Tianyu He, Junliang Guo, Xinle Cheng, Li Song, and Jiang Bian. Vidtok: A versatile and open-source video tokenizer. *arXiv preprint arXiv:2412.13061*, 2024. 6, 12
- [39] Kuaishou AI Team. Kling. <https://kling.kuaishou.com/en>, 2024. Accessed: 2024-XX-XX. 2
- [40] 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 arXiv:1812.01717*, 2018. 5
- [41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 4
- [42] Haiqiang Wang, Weihao Gan, Sudeng Hu, Joe Yuchieh Lin, Lina Jin, Longguang Song, Ping Wang, Ioannis Katsavounidis, Anne Aaron, and C-C Jay Kuo. Mcl-jcv: a jnd-based h. 264/avc video quality assessment dataset. In *ICIP*, 2016. 5
- [43] Junke Wang, Yi Jiang, Zehuan Yuan, Bingyue Peng, Zuxuan Wu, and Yu-Gang Jiang. Omnitokenizer: A joint image-video tokenizer for visual generation. *NeurIPS*, 2025. 6, 12
- [44] Wenjing Wang, Huan Yang, Zixi Tuo, Huiguo He, Junchen Zhu, Jianlong Fu, and Jiaying Liu. Videofactory: Swap attention in spatiotemporal diffusions for text-to-video generation. 2023. 3
- [45] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024. 2, 3, 4
- [46] Lijun Yu, José Lezama, Nitesh B Gundavarapu, Luca Versari, Kihyuk Sohn, David Minnen, Yong Cheng, Vighnesh Birodkar, Agrim Gupta, Xiuye Gu, et al. Language model beats diffusion—tokenizer is key to visual generation. *arXiv preprint arXiv:2310.05737*, 2023. 2, 3, 5, 6, 11, 12, 13
- [47] Biao Zhang and Rico Sennrich. Root mean square layer normalization. *NeurIPS*, 2019. 4
- [48] Guozhen Zhang, Yuhan Zhu, Haonan Wang, Youxin Chen, Gangshan Wu, and Limin Wang. Extracting motion and appearance via inter-frame attention for efficient video frame interpolation. In *CVPR*, 2023. 14
- [49] Sijie Zhao, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Muyao Niu, Xiaoyu Li, Wenbo Hu, and Ying Shan. Cv-vae: A compatible video vae for latent generative video models. *NeurIPS*, 2025. 12
- [50] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all, 2024. 2, 3, 4, 11, 12

## A. Additional Video Examples

Due to the space limit in the main paper, we provide additional video examples to validate our method and please check our [project website](#) for all the results displayed in a web UI.

## B. Implementation Detail

### B.1. Dataset

For video embedder training, the dataset is composed of 15 million videos and 300 million images. For text-to-video generation, we keep the same image resources and leverage 1 million videos to construct the dataset. The images are at the resolution of  $256\times$  with various aspect ratios, while all videos are at the resolution of  $192\times 320$ .

### B.2. Training

#### B.2.1. MAGVIT-v2

We reimplement MAGVIT-v2 [46] at the temporal compression ratio of  $4\times$ ,  $8\times$ ,  $16\times$ ,  $32\times$ , and train the models at the resolution of  $128\times 128$  on our dataset with over 200K iterations. For  $4\times$  temporal compression rate, we adopt the original 4-stage design and encode 17 frames into 5 latent frames at each chunk. For higher compression rate experiments, we add an additional stage to both encoder and decoder, and encode 17 frames into 3, 2 latent frames for temporal compression rate  $8\times$ ,  $16\times$ , respectively. As for the extreme  $32\times$  temporal compression rate, every 33 frames will be encoded into 2 latent frames. We train both MAGVIT-v2 and its discriminator with the AdamW [25] optimizer and we feed mixed images and videos to the network based on the pre-defined ratios (Image: 20%, Video: 80%). All models are trained with 32 A100 GPUs and each GPU holds either 6 image inputs or 2 video sequences.

#### B.2.2. REGEN

Our method is trained at the temporal compression ratio of  $4\times$ ,  $8\times$ ,  $16\times$ ,  $32\times$  on our dataset with over 100K iterations. We train our method under the setting of reconstruction, interpolation and extrapolation, and each scenario will be randomly sampled at every iteration based on pre-defined probabilities. Specifically, we tend more towards reconstruction in the early training stage, and progressively increase the probabilities of interpolation and extrapolation with more training iterations. Our spatiotemporal video encoder adopts the same architectural design as our implemented MAGVIT-v2 but encodes each chunk into 2 latent frames with varying compression rates. We train the spatiotemporal video encoder and the generative decoder in an end-to-end fashion with the AdamW optimizer. Similar with MAGVIT-v2 training, the model will be trained with mixed images and videos based on the pre-defined ratios (Image: 20%, Video: 80%). All models are trained using 32

A100 GPUs, but we leverage larger batch sizes compared to MAGVIT-v2 as we did not utilize 3D convolution layers in the decoder and the GPU can process more samples. The image batch size is set to be 28 for all experiments and the video batch size is 6, 4, 3, 1 for temporal compression ratio of  $4\times$ ,  $8\times$ ,  $16\times$ ,  $32\times$ , respectively.

### B.3. Architecture Details

We illustrate the architecture detail of our method at  $4\times$  temporal compression rate in Fig.10. The spatiotemporal video encoder consists of four stages and it will encode every 5 raw frames into 2 latent frames with the spatial compression rate of  $8\times$ . The base channels are 128 and the channel multiplier for different stages is 1, 2, 4, 6, respectively. The generative decoder comprises 24 diffusion transformer blocks and takes the output of encoder as the conditioning signal. For higher compression rate experiments, we add an additional stage to the encoder while keeping the same number of diffusion transformer blocks in the decoder. The channel multiplier is adjusted to 1, 2, 4, 6, 6, accordingly.

### B.4. Alleviating Chunking Issue with Latent Extension

Following the idea of SDEdit [29], we utilize the last frame prediction from the previous chunk to guide the generation of the next chunk in an auto-regressive manner. Specifically, given a set of latent frames, we decode the first chunk in the reconstruction setting and preserve the intermediate results of the last frame at each sampling step. When decoding later chunks, we extend the latent by extrapolation and control the decoder to generate  $T_c + 1$  frames with the time shift of -1, where  $T_c$  represents the length of each chunk. In this manner, we can ensure there is one overlapped frame between consecutive chunks. During the generation of later chunks, we leverage the intermediate results of the last frame from the previous chunk, which is also the first frame of the current chunk, to guide the generation of the current chunk so that we could mitigate the jumping issue in an auto-regressive way. We provide an additional example in Fig. 11 to demonstrate the effect.

## C. Additional Comparisons with SOTA Video Embedders

As many existing video embedders have different latent channel dimensions. Some video embedders, like Open-SORA (OS) [50] and Open-SORA-Plan (OSP) only have 4 latent channels. For fairness of comparison, we compare against these SOTA embedders by calculating the compression factor following LTX-Video [11]:  $r = \frac{C \times H \times W \times (T-1)}{c \times h \times w \times t}$ . Shown in Fig. 12, we further compare REGEN with Cosmos-Tokenizer [1], LTX-Video [11], Om-

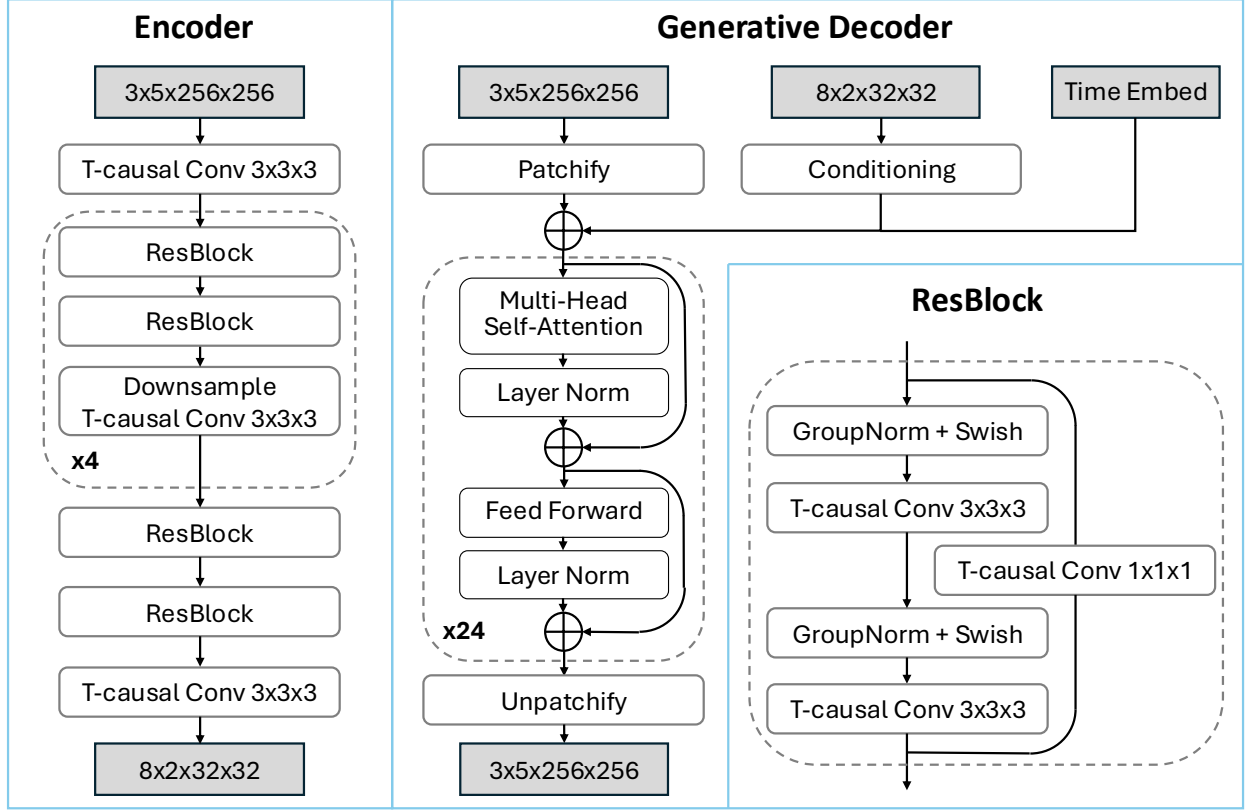


Figure 10. Architecture details of our method at  $4\times$  temporal compression rate. T-causal Conv stands for temporal causal convolution where  $\text{stride} = 1 \times 1 \times 1$ . Downsample T-causal Conv stands for temporal causal convolution where stride will be selected from  $\{1 \times 1 \times 1, 1 \times 2 \times 2, 2 \times 1 \times 1, 2 \times 2 \times 2\}$ , depending on the target compression ratios. The spatiotemporal video encoder is composed of 4 stages and the generative decoder is made up of 24 diffusion transformer blocks.

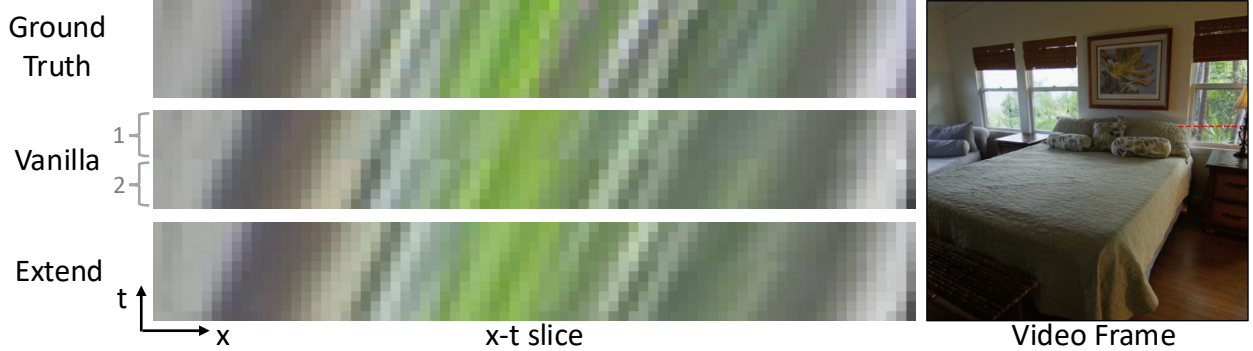


Figure 11. Alleviating the chunking issue with latent extension. The x-t slice is obtained by extracting a short segment (shown as the red line in the video frame) from 2 chunk frames.

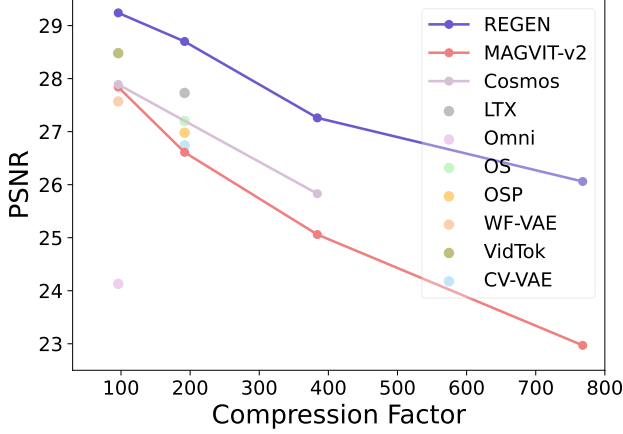
niTokenizer [43], Open-SORA (OS) [50], Open-SORA-Plan (OSP) [21], WF-VAE [22], VidTok [38] and CV-VAE [49]. One can observe that our method exhibits better performance at various compression factors on both datasets, demonstrating the effectiveness of our generative decoder and the soundness of our experimental setup.

#### D. Efficiency Analysis of REGEN

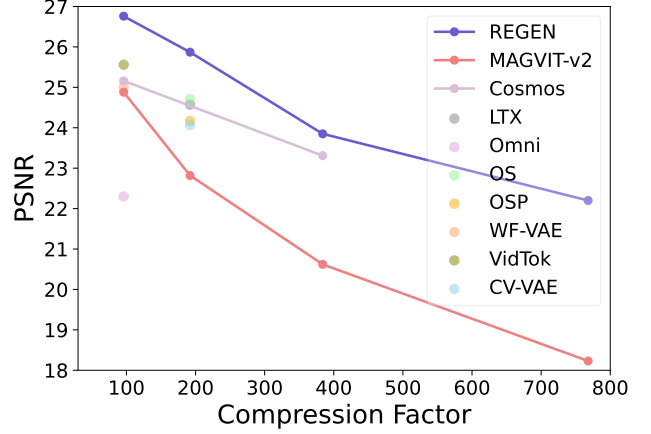
For the training efficiency of latent generative models, we have the same efficiency with MAGVIT-v2 [46] as we leverage the same encoder with it and decoders will not affect the training of latent generative models.

Although operating on the noise input, REGEN is equipped with a large patch size to accommodate the costs and ensure efficient inference. Since REGEN supports one-





(a) Comparisons with state-of-the-art video embedders on MCL-JCV dataset under  $256 \times 256$  inputs.



(b) Comparisons with state-of-the-art video embedders on DAVIS 2019 dataset under  $256 \times 256$  inputs.

Figure 12. Holistic comparison with state-of-the-art video embedders with different latent channel dimensions. The compression factors is calculated by  $r = \frac{C \times H \times W \times (T-1)}{c \times h \times w \times t}$ . REGEN obtains better PSNR compared to existing video embedders at various compression factors on different datasets.

step sampling, we measure the running time (milliseconds) of REGEN on one A100 GPU with one-step decoding and compare with MAGVIT-v2 at various compression rates under  $256 \times 256$  inputs. One can observe from Tab. 4 that REGEN costs similar or fewer running times except at the compression rate of  $8 \times 8 \times 32$ . The reason is that we have to adjust the location of upsampling layers in MAGVIT-v2 decoder at  $8 \times 8 \times 32$  and shift all upsampling layers towards the end blocks to avoid out of GPU memory during training due to the causal 3D convolutional layers, even on 80GB A100s. Note that we avoid doing this for other compression ratios of MAGVIT-v2 as it is not an optimal design and lowers reconstruction quality, as also mentioned in MAGVIT-v2 work [46]. For  $32 \times$  temporal compression variant, this design choice is unavoidable to ensure the decoding process fits within the available GPU memory. Consequently, the computational complexity of this model is lower than that of the ideal design. Even so, this version of  $32 \times$  MAGVIT-v2 cannot decode videos of even  $512 \times 512$  resolution which highlights the scalability of our transformer-based decoder.

Method	Latency (ms)			
	$8 \times 8 \times 4$	$8 \times 8 \times 8$	$8 \times 8 \times 16$	$8 \times 8 \times 32$
MAGVIT-v2	88	317	343	153*
REGEN	89	159	295	548

Table 4. Comparison of decoder latency at various compression rates on one A100 GPU under  $256 \times 256$  inputs. \*MAGVIT-v2  $32 \times$  has small latency because we have a different decoder design where we move all upsampling blocks toward the end layers. This is not optimal from the perspective of reconstruction quality for the decoder, but this design is unavoidable otherwise the model gives Out Of GPU Memory on 80GB A100 GPUs.

## E. Effectiveness of Scaling Convolution-Based Video Embedders

While we have shown that REGEN exhibits significant advantage over MAGVIT-v2 at high temporal compression rates, e.g.,  $16 \times$  and  $32 \times$ , we further scale up the MAGVIT-v2 decoder to explore whether larger model size could help to mitigate the performance gap. Specifically, we scale up the MAGVIT-v2 decoder through the width dimension until it matches the parameter of REGEN or reaches the GPU memory. Tab. 5 shows that simply scaling up the model size do not always translate to better results at high compression rates, e.g., the expanded MAGVIT-v2 has worse SSIM and rFVD compared to the original version at  $16 \times$  compression. Although scaling up model size has shown some improvement at  $32 \times$  compression, it still lags behind REGEN obviously, suggesting that the idea of using generative decoders to escape the compression-reconstruction trade-off is effective.

Method	$8 \times 8 \times 16$			$8 \times 8 \times 32$		
	PSNR	SSIM	rFVD↓	PSNR	SSIM	rFVD↓
MAGVIT-v2	20.62	0.527	441.24	18.23	0.419	1080.15
MAGVIT-v2◇	20.83	0.508	486.48	18.98	0.437	1020.80
REGEN	<b>23.85</b>	<b>0.635</b>	<b>328.83</b>	<b>22.20</b>	<b>0.575</b>	<b>488.89</b>

Table 5. Comparisons of expanded MAGVIT-v2 with REGEN on DAVIS 2019 dataset under  $256 \times 256$  inputs. We expand the MAGVIT-v2 decoder by scaling up the width dimension and ◇ denotes the expanded version. The best results are bold-faced.

## F. Latent Interpolation and Extrapolation

As described in the Method, our INR-based latent conditioning module not only supports reconstruction, but can

generalize to interpolation and extrapolation as well with a uniform design. To examine the interpolation ability of our method, we compare it with two baselines: (1) Frame Averaging: it averages the ground truth to get the interpolated frames; (2) Ours + External Interpolation: it applies off-the-shelf interpolation model [48] on our reconstructed frames. Fig. 13 shows that simply averaging the frames results in clear artifacts on the interpolated frames, while the results of our model and EMA display a smoother transition and align well with the ground truth. Apart from interpolation, our design supports extrapolation as well, where the model is going to predict the previous or future frames based on the given input. Shown in Fig. 14, our approach forecasts future motion based on prior frame sequences and the results demonstrate strong alignment with the ground truth, which highlights the generation ability of our decoder. Note that our method is able to predict the past frames as well which is shown in Fig. 15, promising the application of our method for chunk-free generation to mitigate the jumping issue.



Figure 13.  $2\times$  interpolation results. Given input frames with purple bounding boxes, the model is asked to conduct interpolation to predict the frame with blue bounding box and we compare our method with frame averaging and external interpolation model.

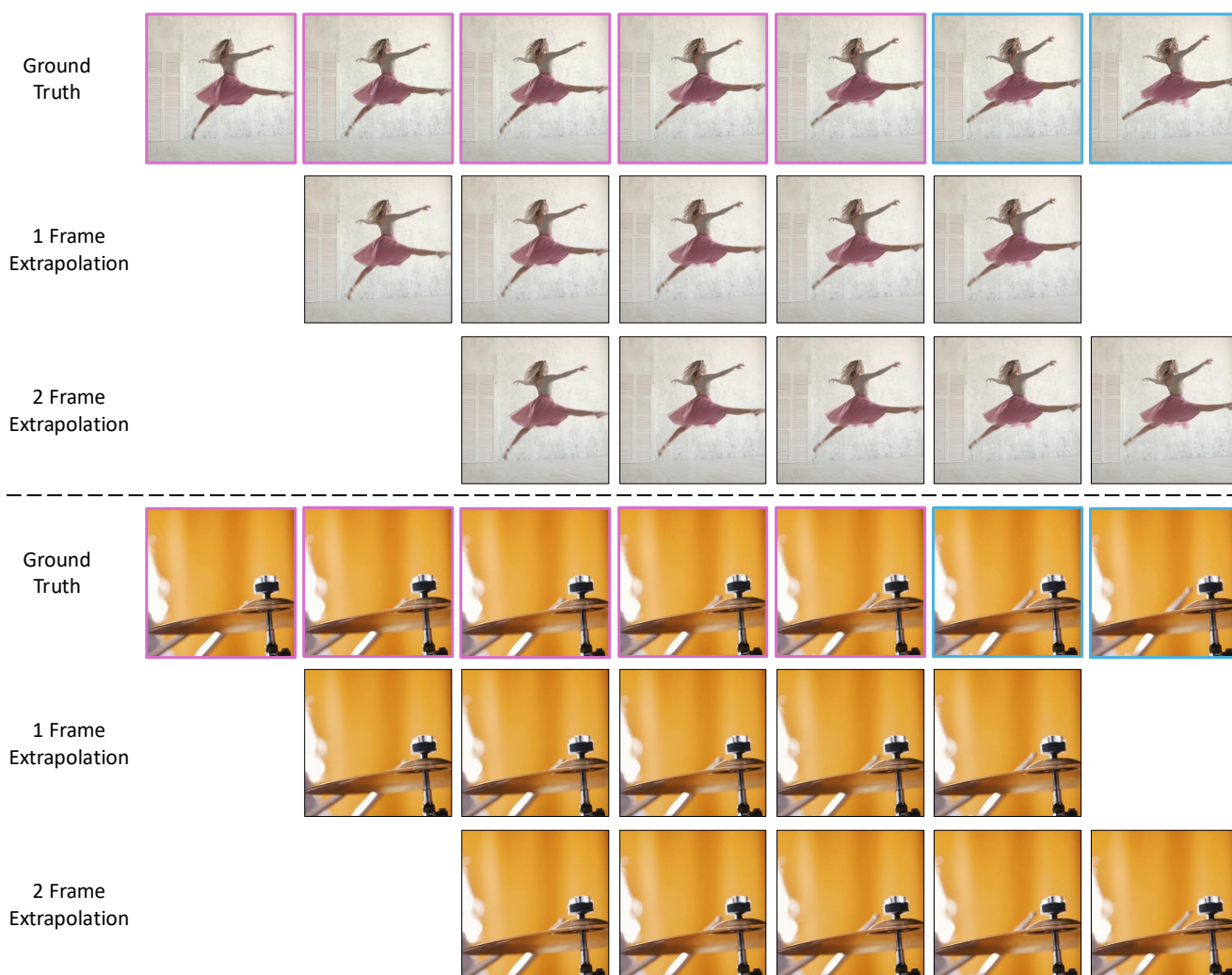


Figure 14. Forward latent extrapolation results. Given input frames with purple bounding boxes, the model is asked to conduct extrapolation to predict the future frame with blue bounding box.



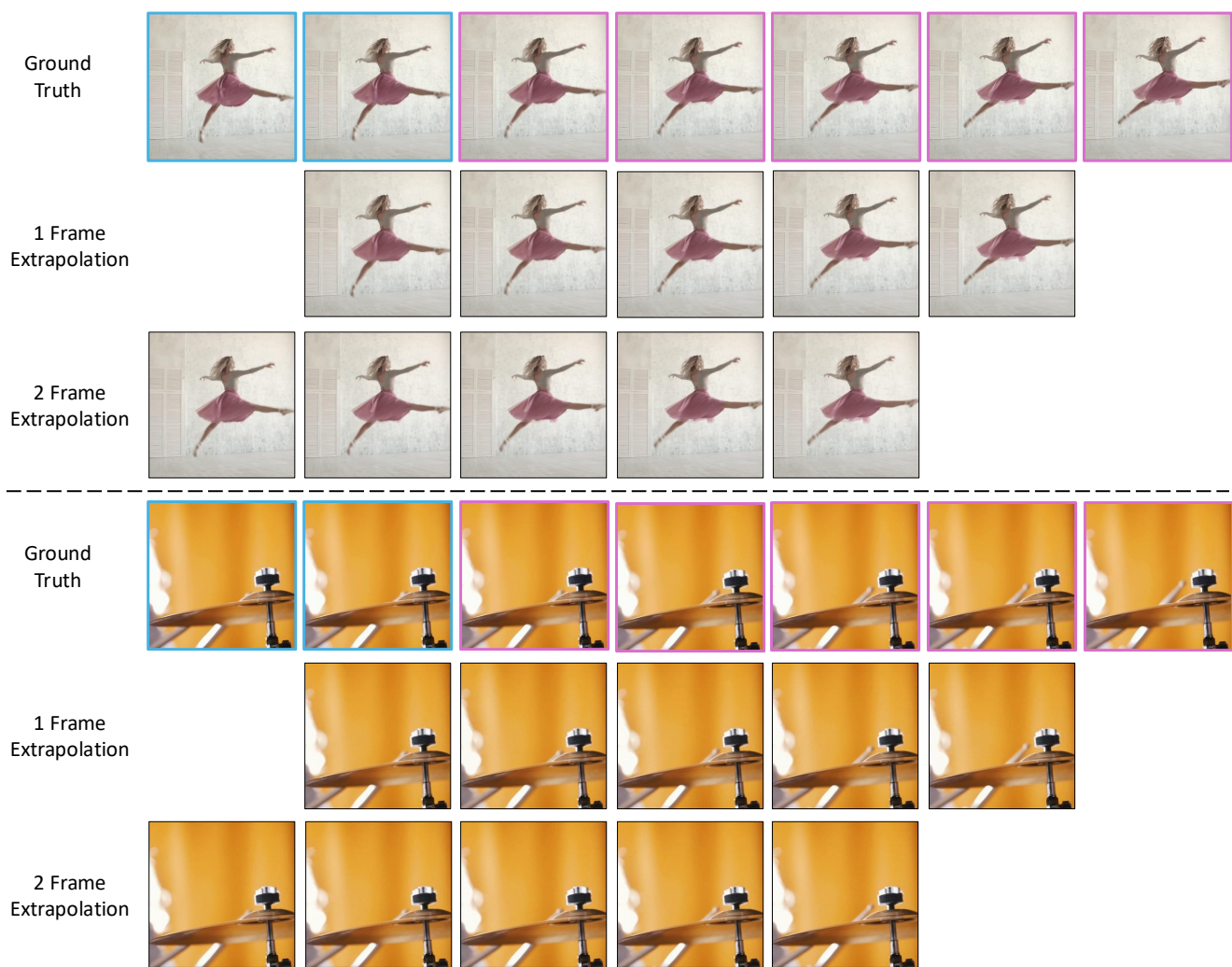


Figure 15. Backward latent extrapolation results. Given input frames with purple bounding boxes, the model is asked to conduct extrapolation to predict the past frame with blue bounding box.