

MAGVIT: Masked Generative Video Transformer

Lijun Yu^{††}, Yong Cheng[†], Kihyuk Sohn[†], José Lezama[†], Han Zhang^{†‡}, Huiwen Chang^{†‡}, Alexander G. Hauptmann[‡], Ming-Hsuan Yang[†], Yuan Hao[†], Irfan Essa[†], and Lu Jiang[†]
[†]Carnegie Mellon University, [†]Google Research, [‡]Georgia Institute of Technology

[○]project lead, [▷]technical contribution, [△]advising and method design, [△]advising

lijun@cmu.edu, lujiang@google.com

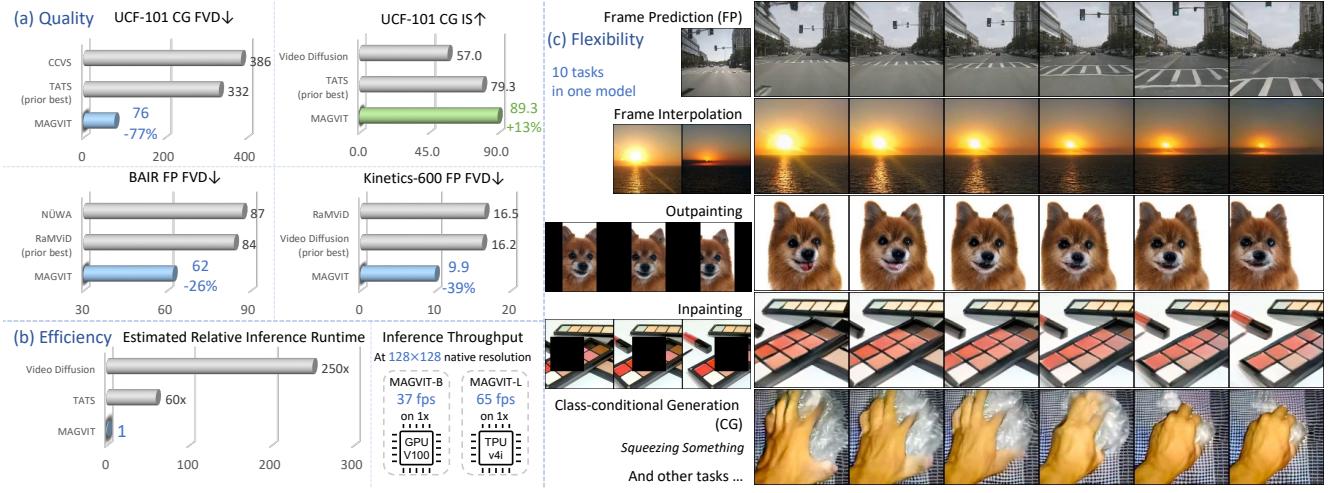


Figure 1. Overview of the video generation **quality**, **efficiency**, and **flexibility** of the proposed MAGVIT model. (a) MAGVIT achieves the state-of-the-art FVD [60] and Inception Score (IS) [49] on two video generation tasks and three benchmarks, in comparison with prior best diffusion models (RaMViD [35], Video Diffusion [33]) and autoregressive models (CCVS [41], TATS [21], NÜWA [69]). (b) It is two orders of magnitude faster than diffusion models and $60\times$ faster than autoregressive models. (c) A single MAGVIT model accommodates different generation tasks, ranging from class-conditional generation to dynamic inpainting of a moving object.

Abstract

We introduce the MAsked Generative VIdeo Transformer, MAGVIT, to tackle various video synthesis tasks with a single model. We introduce a 3D tokenizer to quantize a video into spatial-temporal visual tokens and propose an embedding method for masked video token modeling to facilitate multi-task learning. We conduct extensive experiments to demonstrate the quality, efficiency, and flexibility of MAGVIT. Our experiments show that (i) MAGVIT performs favorably against state-of-the-art approaches and establishes the best-published FVD on three video generation benchmarks, including the challenging Kinetics-600. (ii) MAGVIT outperforms existing methods in inference time by two orders of magnitude against diffusion models and by $60\times$ against autoregressive models. (iii) A single MAGVIT model supports ten diverse generation tasks and generalizes across videos from different visual domains. The source code and trained models will be released to the public at <https://magvit.cs.cmu.edu>.

1. Introduction

Recent years have witnessed significant advances in image and video content creation based on learning frameworks ranging from generative adversarial networks (GANs) [15, 43, 48, 58, 65], diffusion models [25, 33, 35, 47, 64], to vision transformers [44, 45, 68]. Inspired by the recent success of generative image transformers such as DALL-E [46] and other approaches [12, 18, 20, 72], we propose an efficient and effective video generation model by leveraging masked token modeling and multi-task learning.

We introduce the MAsked Generative VIdeo Transformer (*MAGVIT*) for multi-task video generation. Specifically, we build and train a single MAGVIT model to perform a variety of diverse video generation tasks and demonstrate the model’s efficiency, effectiveness, and flexibility against state-of-the-art approaches. Fig. 1(a) shows the quality metrics of MAGVIT on a few benchmarks with efficiency comparisons in (b), and generated examples under different task setups such as frame prediction/interpolation, out/in-painting, and class conditional generation in (c).

MAGVIT models a video as a sequence of visual tokens in the latent space and learns to predict masked tokens with BERT [17]. There are two main modules in the proposed framework. First, we design a 3D quantization model to tokenize a video, with high fidelity, into a low-dimensional spatial-temporal manifold [21, 70]. Second, we propose an effective *masked token modeling* (MTM) scheme for multi-task video generation. Unlike conventional MTM in image understanding [66] or image/video synthesis [12, 26, 28], we present an embedding method to model a video condition using a multivariate mask and show its efficacy in training.

We conduct extensive experiments to demonstrate the quality, efficiency, and flexibility of MAGVIT against state-of-the-art approaches. Specifically, we show that MAGVIT performs favorably on two video generation tasks across three benchmark datasets, including UCF-101 [54], BAIR Robot Pushing [19, 60], and Kinetics-600 [10]. For the class-conditional generation task on UCF-101, MAGVIT reduces state-of-the-art FVD [60] from 332 [21] to 76 ($\downarrow 77\%$). For the frame prediction task, MAGVIT performs best in terms of FVD on BAIR (84 [35] \rightarrow 62, $\downarrow 26\%$) and Kinetics-600 (16 [33] \rightarrow 9.9, $\downarrow 38\%$).

Aside from the visual quality, MAGVIT’s video synthesis is highly efficient. For instance, MAGVIT generates a 16-frame 128×128 video clip in 12 steps, which takes 0.25 seconds on a single TPUv4i [36] device. On a V100 GPU, a base variant of MAGVIT runs at 37 frame-per-second (fps) at 128×128 resolution. When compared at the same resolution, MAGVIT is two orders of magnitude faster than the video diffusion model [33]. In addition, MAGVIT is 60 times faster than the autoregressive video transformer [21] and 4-16 times more efficient than the contemporary non-autoregressive video transformer [26].

We show that MAGVIT is flexible and robust for multiple video generation tasks with a single trained model, including frame interpolation, class-conditional frame prediction, inpainting, and outpainting, etc. In addition, MAGVIT learns to synthesize videos with complex scenes and motion contents from diverse and distinct visual domains, including actions with objects [23], autonomous driving [9], and object-centric videos from multiple views [2].

The main contributions of this work are:

- To the best of our knowledge, we present the first masked multi-task transformer for efficient video generation and manipulation. We show that a trained model can perform ten different tasks at inference time.
- We introduce a spatial-temporal video quantization model design with high reconstruction fidelity.
- We propose an effective embedding method with diverse masks for numerous video generation tasks.
- We show that MAGVIT achieves the best-published fidelity on three widely-used benchmarks, including UCF-101, BAIR Robot Pushing, and Kinetics-600 datasets.

2. Preliminaries: Masked Image Synthesis

The proposed video generation framework is based on a two-stage image synthesis process [20, 46] with non-autoregressive transformers [12, 42]. In the first stage, an image is quantized and flattened into a sequence of discrete tokens by a Vector-Quantized (VQ) auto-encoder [20, 62, 71]. In the second stage, masked token modeling (MTM) is used to train a transformer model [12, 25] on the tokens. Let $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$ be an image and $\mathbf{z} \in \mathbb{Z}^N$ denote the corresponding token sequence of length N .

We take MaskGIT [12] as an example. In the second stage, it applies a binary mask $\mathbf{m}_i \in \{x \rightarrow x, x \rightarrow [\text{MASK}]\}$ to each token to build a corrupted sequence $\bar{\mathbf{z}} = \mathbf{m}(\mathbf{z})$. Condition inputs, such as class labels, are incorporated as the prefix tokens \mathbf{c} . A BERT [17] parameterized by θ is learned to predict the masked tokens in the input sequence $[\mathbf{c}, \bar{\mathbf{z}}]$, where $[\cdot, \cdot]$ concatenates the sequences. The objective is to minimize the cross-entropy between the predicted and the ground-truth token at each masked position:

$$\mathcal{L}_{\text{mask}}(\mathbf{z}; \theta) = \mathbb{E}_{\mathbf{m} \sim p_{\mathcal{U}}} \left[\sum_{\bar{z}_i = [\text{MASK}]} -\log p_{\theta}(\mathbf{z}_i \mid [\mathbf{c}, \bar{\mathbf{z}}]) \right] \quad (1)$$

During training, MaskGIT randomly samples \mathbf{m} from a prior distribution $p_{\mathcal{U}}$ where the mask ratio follows a cosine scheduling function $\gamma(\cdot)$ [12]. Specifically, it first uniformly samples a per-token mask score $s_i \sim \mathcal{U}(0, 1)$ to form a sequence denoted as \mathbf{s} . Then it samples $r \sim \mathcal{U}(0, 1)$ and computes a cut-off threshold s^* as the $\lceil \gamma(r)N \rceil$ th smallest element in \mathbf{s} . Finally, a mask \mathbf{m} is created such that $\mathbf{m}_i(x) = [\text{MASK}]$ if $s_i \leq s^*$ and $\mathbf{m}_i(x) = x$ otherwise.

For inference, the non-autoregressive decoding method [22, 24, 40] is used to synthesize an image [12, 42, 75]. For example, MaskGIT generates an image in $K = 12$ steps [12] from a blank canvas with all visual tokens masked out. At each step, it predicts all tokens in parallel while retaining tokens with the highest prediction scores. The remaining tokens are masked and predicted in the next iteration until all tokens are generated. Similar to the training stage, the mask ratio is computed by the schedule function γ , but with a deterministic input as $\gamma(\frac{t}{K})$, where t is the current step.

3. Masked Generative Video Transformer

Our goal is to design a multi-task video generation model with high quality and inference efficiency. We propose MAsked Generative VIdeo Transformer (MAGVIT), a vision transformer framework that leverages masked token modeling and multi-task learning. MAGVIT generates a video from task-specific condition inputs, such as a frame, a partially-observed video volume, or a class identifier.

The framework consists of two stages. First, we learn a 3D vector-quantized (VQ) autoencoder to quantize a video

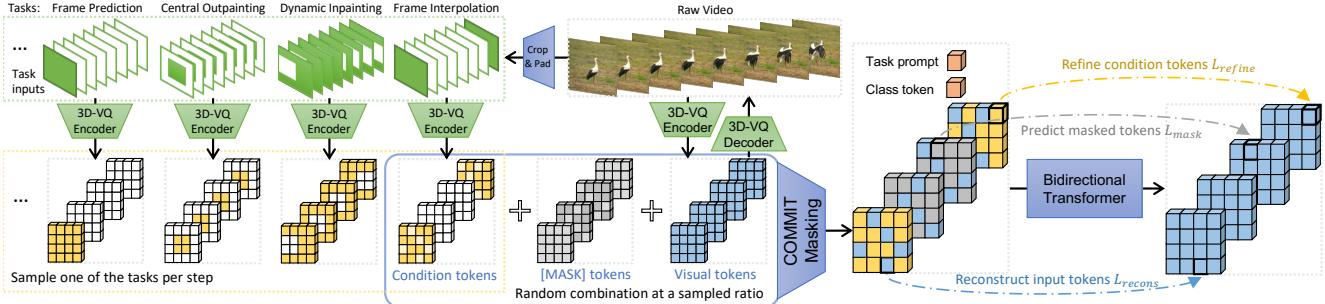


Figure 2. **MAGVIT pipeline overview.** The 3D-VQ encoder quantizes a video into discrete tokens, while the 3D-VQ decoder maps them back to the pixel space. We sample one of the tasks at each training step and build its condition inputs by cropping and padding the raw video, where green denotes valid pixels and white is padding. We quantize the condition inputs with the 3D-VQ encoder and select the non-padding part as condition tokens. The masked token sequence combines condition tokens, [MASK] tokens, and the original visual tokens, with a task prompt and a class token as the prefix. The bidirectional transformer learns to predict the visual tokens through three objectives: refining condition tokens, predicting masked tokens, and reconstructing input tokens.

into discrete tokens. In the second stage, we learn a video transformer by multi-task masked token modeling.

Fig. 2 illustrates the training in the second stage. At each training step, we sample one of the tasks with its prompt token, obtain a task-specific conditional mask, and optimize the transformer to predict all visual tokens given masked inputs. During inference, we adapt the non-autoregressive decoding method to generate tokens conditionally on the task-specific inputs, which will be detailed in Algorithm 1.

3.1. Spatial-Temporal Tokenization

Our video VQ autoencoder is built upon the image VQ-GAN [20]. Let $\mathbf{V} \in \mathbb{R}^{T \times H \times W \times 3}$ be a video clip of T frames. The VQ encoder tokenizes the video as $f_{\mathcal{T}} : \mathbf{V} \rightarrow \mathbf{z} \in \mathbb{Z}^N$, where \mathbb{Z} is the codebook. The decoder $f_{\mathcal{T}}^{-1}$ maps the latent tokens back to video pixels.

The VQ autoencoder is a crucial module as it not only sets a quality bound for the generation but also determines the token sequence length, hence affecting generation efficiency. Existing methods apply VQ encoders either on each frame independently (2D-VQ) [26, 41] or on a supervoxel (3D-VQ) [21, 70]. We propose different designs that facilitate MAGVIT to perform favorably against other VQ models for video (see Tab. 7).

3D architecture. We design a 3D-VQ network architecture to model the temporal dynamics as follows. The encoder and decoder of VQGAN consist of cascaded residual blocks [29] interleaved by downsampling (average pooling) and upsampling (resizing plus convolution) layers. We expand all 2D convolutions to 3D convolutions with a temporal axis. As the overall downsampling rate is usually different between temporal and spatial dimensions, we use both 3D and 2D downsampling layers, where the 3D ones appear in the shallower layers of the encoder. The decoder mirrors the encoder with 2D upsampling layers in the first few blocks, followed by 3D ones. Appendix A.1 illustrates

the detailed architecture. Note that a token is not only correlated to its corresponding supervoxel but depends on other patches due to the non-local receptive field.

Inflation and padding. We initialize our 3D-VQ with weights from a 2D-VQ in a matching architecture to transfer learned spatial relationships [11], known as 3D inflation. We use inflation on small datasets such as UCF-101 [54]. We use a central inflation method for the convolution layers, where the corresponding 2D kernel fills in the temporally central slice of a zero-filled 3D kernel. The parameters of the other layers are directly copied. To improve token consistency for the same content at different locations [21], we replace the same (zero) padding in the convolution layers with reflect padding, which pads with non-zero values.

Training. We apply the image perceptual loss [20] on each frame. The LeCam regularization [57] is added to the GAN loss to improve the training stability. We adopt the discriminator architecture from StyleGAN [38] and inflate it to 3D. With these components, unlike VQGAN, our model is trained stably with GAN loss from the beginning.

3.2. Multi-Task Masked Token Modeling

In MAGVIT, we adopt various masking schemes to facilitate learning for video generation tasks with different conditions. The conditions can be a spatial region for inpainting/outpainting or a few frames for frame prediction/interpolation. We refer to these partially-observed video conditions as *interior conditions*.

We argue that it is suboptimal to directly unmask the tokens corresponding to the region of the interior condition [12]. As discussed in Section 3.1, the non-local receptive field of the tokenizer can leak the ground-truth information into the unmasked tokens, leading to problematic non-causal masking and poor generalization.

We propose a method, COnditional Masked Modeling by Interior Tokens (or *COMMIT* for short), to embed interior

conditions inside the corrupted visual tokens.

Training. Each training example includes a video \mathbf{V} and the optional class annotation \mathbf{c} . The visual tokens come from the 3D-VQ as $\mathbf{z} = f_{\mathcal{T}}(\mathbf{V})$. At each step, we sample a task prompt ρ , obtain the task-specific interior condition pixels, pad it into $\tilde{\mathbf{V}}$ with the same shape as \mathbf{V} , and get the condition tokens $\tilde{\mathbf{z}} = f_{\mathcal{T}}(\tilde{\mathbf{V}})$. Appendix B.1 lists the padding functions for each task. At a sampled mask ratio, we compute the *multivariate* conditional mask $\mathbf{m}(\cdot | \tilde{\mathbf{z}})$ as

$$m_i(x | \tilde{z}_i) = \begin{cases} \tilde{z}_i & \text{if } s_i \leq s^* \wedge \neg \text{ispad}(\tilde{z}_i) \\ [\text{MASK}] & \text{if } s_i \leq s^* \wedge \text{ispad}(\tilde{z}_i) \\ x & \text{if } s_i > s^* \end{cases} \quad (2)$$

where s_i and s^* are the per-token mask score and the cut-off score introduced in Section 2. $\text{ispad}(\tilde{z}_i)$ returns whether the corresponding supervoxel of \tilde{z}_i in $\tilde{\mathbf{V}}$ only contains padding.

Eq. (2) indicates that COMMIT embeds interior conditions as corrupted visual tokens into the multivariate mask \mathbf{m} , which follows a new distribution $p_{\mathcal{M}}$ instead of the prior $p_{\mathcal{U}}$ for binary masks. With the corrupted token sequence $\bar{\mathbf{z}} = \mathbf{m}(\mathbf{z} | \tilde{\mathbf{z}})$ as input, the *multi-task* training objective is

$$\mathcal{L}(\mathbf{V}; \theta) = \mathbb{E}_{\rho, \tilde{\mathbf{V}}} \mathbb{E}_{\mathbf{m} \sim p_{\mathcal{M}}} \left[\sum_i -\log p_{\theta}(z_i | [\rho, \mathbf{c}, \bar{\mathbf{z}}]) \right] \quad (3)$$

We can decompose the loss in Eq. (3) into three parts according to Eq. (2): $\mathcal{L}_{\text{refine}}$ refines the task-specific condition tokens, $\mathcal{L}_{\text{mask}}$ predicts masked tokens, and $\mathcal{L}_{\text{recons}}$ reconstructs input tokens. Let $\bar{\mathbf{c}} = [\rho, \mathbf{c}, \bar{\mathbf{z}}]$ for simplicity,

$$\begin{aligned} \sum_{i=1}^N -\log p_{\theta}(z_i | [\rho, \mathbf{c}, \bar{\mathbf{z}}]) &= \underbrace{\sum_{\substack{\bar{z}_i = \tilde{z}_i \\ \text{Refine condition tokens } \mathcal{L}_{\text{refine}}}} -\log p_{\theta}(z_i | \bar{\mathbf{c}})}_{\text{Reconstruct input tokens } \mathcal{L}_{\text{recons}}} \\ &+ \underbrace{\sum_{\substack{\bar{z}_i = [\text{MASK}] \\ \text{Predict masked tokens } \mathcal{L}_{\text{mask}}}} -\log p_{\theta}(z_i | \bar{\mathbf{c}})}_{\text{Reconstruct input tokens } \mathcal{L}_{\text{recons}}} + \underbrace{\sum_{\substack{\bar{z}_i = z_i \\ \text{Reconstruct input tokens } \mathcal{L}_{\text{recons}}}} -\log p_{\theta}(z_i | \bar{\mathbf{c}})}_{\text{Reconstruct input tokens } \mathcal{L}_{\text{recons}}} \end{aligned} \quad (4)$$

While $\mathcal{L}_{\text{mask}}$ is the same as the MTM loss in Eq. (1) and $\mathcal{L}_{\text{recons}}$ sometimes is used as a regularizer (e.g., in NLP tasks), $\mathcal{L}_{\text{refine}}$ is a new component introduced by COMMIT.

The COMMIT method facilitates multi-task video generation in three aspects. First, it provides a correct causal masking for all interior conditions. Second, it produces a fixed-length sequence for different conditions of arbitrary regional volume, improving training and memory efficiency since no padding tokens are needed. Third, it achieves state-of-the-art multi-task video generation results (see Tab. 5).

Video generation tasks. We consider *ten* tasks for multi-task video generation where each task has a different interior condition and mask: Frame Prediction (FP), Frame Interpolation (FI), Central Outpainting (OPC), Vertical Outpainting (OPV), Horizontal Outpainting (OPH), Dynamic Outpainting (OPD), Central Inpainting (IPC), and Dynamic Inpainting (IPD), Class-conditional Generation (CG), Class-conditional Frame Prediction (CFP). We provide the detailed definitions in Appendix B.1.

Algorithm 1 Non-autoregressive Decoding by COMMIT

Input: prefix ρ and \mathbf{c} , condition $\tilde{\mathbf{z}}$, number of steps K

Output: predicted visual tokens $\hat{\mathbf{z}}$

```

1:  $s = \mathbf{0}$ ,  $s^* = 1$ ,  $\hat{\mathbf{z}} = \mathbf{0}^N$ 
2: for  $t \leftarrow 0, 1, \dots, K-1$  do
3:    $\bar{\mathbf{z}} \leftarrow \mathbf{m}(\hat{\mathbf{z}} | \tilde{\mathbf{z}}, \mathbf{s}, s^*)$ 
4:    $\hat{\mathbf{z}}_i \sim p_{\theta}(z_i | [\rho, \mathbf{c}, \bar{\mathbf{z}}])$ ,  $\forall i$  where  $s_i \leq s^*$ 
5:    $s_i \leftarrow p_{\theta}(\hat{\mathbf{z}}_i | [\rho, \mathbf{c}, \bar{\mathbf{z}}])$ ,  $\forall i$  where  $s_i \leq s^*$ 
6:    $s^* \leftarrow \text{The } \lceil \gamma(\frac{t+1}{K})N \rceil \text{th smallest value of } \mathbf{s}$ 
7:    $s_i \leftarrow 1$ ,  $\forall i$  where  $s_i > s^*$ 
8: end for
9: return  $\hat{\mathbf{z}} = [\hat{\mathbf{z}}_1, \hat{\mathbf{z}}_2, \dots, \hat{\mathbf{z}}_N]$ 

```

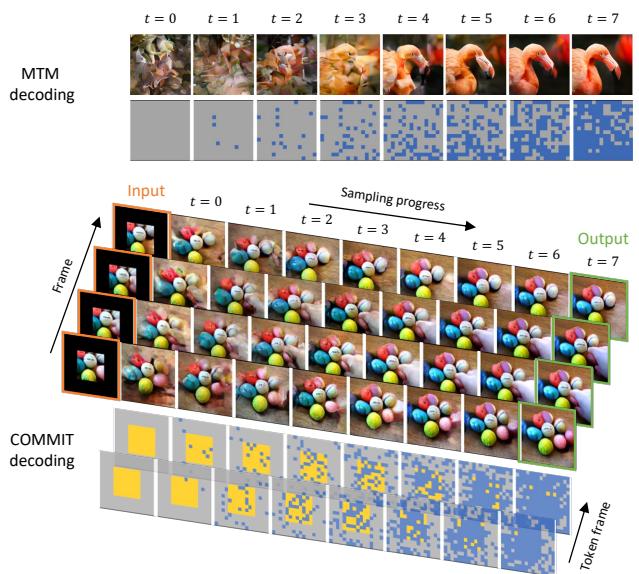


Figure 3. **Comparison between MTM decoding for image [12] and COMMIT decoding for video.** We show the input tokens and the output image/video at each decoding step t , with a central outpainting example for COMMIT. Unlike the MTM denoising decoding from all $[\text{MASK}]$, COMMIT performs a conditional generation process toward the **output tokens** while gradually replacing the **interior condition tokens**. Videos and tokens are temporally down-sampled and stacked for visualization.

Inference. We use a non-autoregressive decoding method to generate video tokens from input conditions in K steps (e.g., 12). Each decoding step follows the COMMIT masking in Eq. (2) with a gradually reduced mask ratio. Algorithm 1 outlines the inference procedure.

Fig. 3 compares the non-autoregressive image decoding [12] and our video decoding procedure. Different from the MTM decoding in [12] which performs denoising from all $[\text{MASK}]$, COMMIT decoding starts from a *multivariate* mask that embeds the **interior conditions**. Guided by this mask, Algorithm 1 performs a conditional transition process toward the **output tokens** by replacing a portion of newly generated tokens at each step. In the end, all tokens are predicted where the interior condition tokens get refined.

Method	Extra Video	Class	FVD \downarrow	IS \uparrow
RaMViD [35]			-	21.71 \pm 0.21
StyleGAN-V* [51]			-	23.94 \pm 0.73
DIGAN [73]			577 \pm 21	32.70 \pm 0.35
DVD-GAN [15]		✓	-	32.97 \pm 1.70
Video Diffusion* [33]			-	57.00 \pm 0.62
TATS [21]			420 \pm 18	57.63 \pm 0.24
CCVS+StyleGAN [41]			386 \pm 15	24.47 \pm 0.13
Make-A-Video* [50]		✓	367	33.00
TATS [21]		✓	332 \pm 18	79.28 \pm 0.38
CogVideo* [34]	✓	✓	626	50.46
Make-A-Video* [50]	✓	✓	81	82.55
MAGVIT-B-CG (ours)		✓	159 \pm 2	83.55 \pm 0.14
MAGVIT-L-CG (ours)		✓	76 \pm 2	89.27 \pm 0.15

Table 1. **Generation performance on the UCF-101 dataset.** Methods in gray are pretrained on additional large video data. Methods with ✓ in the Class column are class-conditional, while the others are unconditional. Methods marked with * use custom resolutions, while the others are at 128×128 . See Appendix C for more comparisons with earlier works.

4. Experimental Results

We conduct extensive experiments to demonstrate the video generation quality (Section 4.2), efficiency (Section 4.3), and flexibility for multi-task generation (Section 4.4). We show a few generation results here, and refer to the web page¹ for more examples.

4.1. Experimental Setups

Datasets. We evaluate the single-task video generation performance of MAGVIT on three standard benchmarks, *i.e.*, class-conditional generation on UCF-101 [54] and frame prediction on BAIR Robot Pushing [19, 60] (1-frame condition) and Kinetics-600 [10] (5-frame condition). For multi-task video generation, we quantitatively evaluate MAGVIT on BAIR and SSv2 [23] on 8-10 tasks. Furthermore, to evaluate model generalizability, we train models with the same learning recipe on three additional video datasets: nuScenes [9], Objectron [2], and 12M Web videos. We show their generated videos in the main paper and quantitative performance in Appendix C.

Evaluation metrics. We use FVD [60] as our primary evaluation metric. Similar to [21, 33], FVD features are extracted with an I3D model trained on Kinetics-400 [11]. We also report the Inception Score (IS) [49] calculated with a C3D [56] model on UCF-101, and PSNR, SSIM [67], and LPIPS [74] on BAIR. We report the mean and standard deviation for each metric calculated over four runs.

Implementation details. We train MAGVIT to generate 16-frame videos at 128×128 resolution, except for BAIR at 64×64 . The proposed 3D-VQ model quantizes a video into $4 \times 16 \times 16$ visual tokens, where the visual codebook size is 1024. We use the BERT transformer [17] to model the token sequence, which includes 1 task prompt, 1 class token, and

¹<https://magvit.cs.cmu.edu>

Method	K600 FVD \downarrow	BAIR FVD \downarrow
CogVideo [34]	109.2	-
CCVS [41]	55.0 \pm 1.0	99 \pm 2
Phenaki [63]	36.4 \pm 0.2	97
TrIVD-GAN-FP [43]	25.7 \pm 0.7	103
Transframer [44]	25.4	100
MaskViT [26]	-	94
FitVid [4]	-	94
MCVD [64]	-	90
NÜWA [69]	-	87
RaMViD [35]	16.5	84
Video Diffusion [33]	16.2 \pm 0.3	-
MAGVIT-B-FP (ours)	24.5 \pm 0.9	76 \pm 0.1 (48 \pm 0.1)
MAGVIT-L-FP (ours)	9.9 \pm 0.3	62 \pm 0.1 (31 \pm 0.2)

Table 2. **Frame prediction performance on the BAIR and Kinetics-600 datasets.** - marks that the value is unavailable in their paper or incomparable to others. The FVD in parentheses uses a debiased evaluation protocol on BAIR detailed in Appendix B.3. See Appendix C for more comparisons with earlier works.

Method	FVD \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
CCVS [41]	99	-	0.729	-
MCVD [64]	90	16.9	0.780	-
MAGVIT-L-FP (ours)	62	19.3	0.787	0.123

Table 3. **Image quality metrics on BAIR frame prediction.**

1024 visual tokens. Two variants of MAGVIT, *i.e.*, base (B) with 128M parameters and large (L) with 464M, are evaluated. We train both stages with the Adam optimizer [39] in JAX/Flax [5, 30] on TPUs. Appendix B.2 details training configurations.

4.2. Single-Task Video Generation

Class-conditional generation. The model is given a class identifier in this task to generate the full video. Tab. 1 shows that MAGVIT surpasses the previous best published FVD and IS scores. Notably, it outperforms Make-A-Video [50] which is pretrained on additional 10M videos with a text-image prior. In contrast, MAGVIT is just trained on the 9.5K training videos of UCF-101.

Fig. 4 compares the generated videos to baseline models. We can see that CCVS+StyleGAN [41] gets a decent single-frame quality, but yields little or no motion. TATS [21] generates some motion but with artifacts. In contrast, our model produces higher-quality frames with substantial motion.

Frame prediction. The model is given a single or a few frames to generate future frames. In Tab. 2, we compare MAGVIT against highly-competitive baselines. MAGVIT surpasses the previous state-of-the-art FVD on BAIR by a large margin (84 → 62). Inspired by [60], a “debiased” FVD is also reported in the parentheses to overcome the small validation set. See more discussions in Appendix B.3. In Tab. 3, it demonstrates better image quality.

On the large dataset of Kinetics-600, it establishes a new state-of-the-art result, improving the previous best FVD in [33] from 16.2 to 9.9 by a relative 39% improvement.

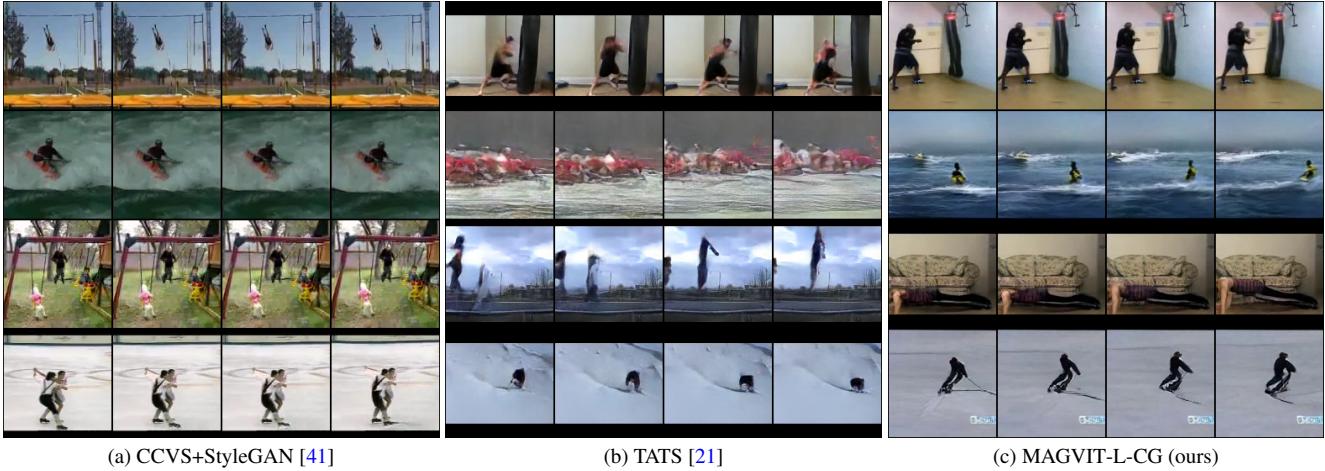


Figure 4. **Comparison of class-conditional generation samples on UCF-101.** 16-frame videos are generated at 128×128 resolution 25 fps and shown at 6.25 fps. Samples for [21, 41] are obtained from their official release². More comparisons are provided in Appendix D.

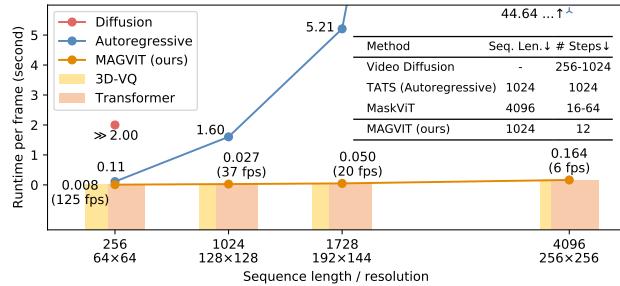


Figure 5. **Inference-time generation efficiency comparison.** The average runtime for generating one frame is measured at different resolutions. The colored bars show the time breakdown between the 3D-VQ and the transformer. The embedded table compares the critical factors of inference efficiency for different methods at 16-frame 128×128 , except for Video Diffusion [33] at 64×64 .

The above results verify MAGVIT’s compelling generation quality, including on the large Kinetics dataset.

4.3. Inference-Time Generation Efficiency

Video generation efficiency is an important metric in many applications. We conduct experiments to validate that MAGVIT offers top speed in video generation. Fig. 5 shows the processing time for each frame on a single V100 GPU at different resolutions. We compare MAGVIT-B with an autoregressive transformer of the same size and a diffusion-based model [33]. At 128×128 resolution, MAGVIT-B runs at 37 frames-per-second (fps). When running on a single TPUs [36], MAGVIT-B runs at 190 fps and MAGVIT-L runs at 65 fps.

Fig. 5 compares the sequence lengths and inference steps of these models. Diffusion models [33] typically require 256-1000 diffusion steps with a 3D U-Net [14]. Autoregressive models, such as TATS [21], decode visual tokens sequentially, which runs 60 times slower than MAGVIT at 128×128 . Compared to the recent non-autoregressive

model MaskViT [26], MAGVIT is 4 to 16 times faster due to more efficient decoding on shorter sequences.

4.4. Multi-task Video Generation

To demonstrate the flexibility in multi-task video synthesis, we train a single MAGVIT model to perform eight tasks on BAIR or ten tasks on SSv2. We do not intend to compare with dedicated models trained on these tasks but to demonstrate a generic model for video synthesis.

Eight tasks on BAIR. We perform a multi-task evaluation on BAIR with eight self-supervised tasks. Tab. 4 lists the “debiased” FVD for each task, where the third column computes the average. We compare the multi-task models (MT) with two single-task baselines trained on unconditional generation (UNC) and frame prediction (FP).

As shown in Tab. 4, the multi-task models achieve better fidelity across all tasks. Single-task models perform considerably worse on the tasks unseen in training (gray values in Tab. 4), especially on the tasks that differ more from the training task. Compared to the single-task models in their training task, MT performs better with a small gain on FP with the same model size.

Ten tasks on SSv2. We evaluate on the large-scale SSv2 dataset, where MAGVIT needs to synthesize 174 basic actions with everyday objects. We evaluate a total of ten tasks, with two of them using class labels (CG and CFP), as shown on the right side of Tab. 4. We observe a pattern consistent with BAIR: multi-task models achieve better average FVD across all tasks. The above results substantiate model generalization trained with the proposed multi-task objective.

4.5. Ablation Study

Conditional MTM. We demonstrate the efficacy of COMMIT by comparing it with conventional MTM meth-

²<https://songweige.github.io/projects/tats/>

Method	Task	BAIR-MT8↓	FP	FI	OPC	OPV	OPH	OPD	IPC	IPD	SSV2-MT10↓	CG	CFP
MAGVIT-B-UNC	Single	150.6	74.0	71.4	119.0	46.7	55.9	389.3	145.0	303.2	258.8	107.7	279.0
MAGVIT-B-FP	Single	201.1	47.7	56.2	247.1	118.5	142.7	366.3	357.3	272.7	402.9	1780.0	59.3
MAGVIT-B-MT	Multi	32.8	47.2	36.0	28.1	29.0	27.8	32.1	31.1	31.0	43.4	94.7	59.3
MAGVIT-L-MT	Multi	22.8	31.4	26.4	21.3	21.2	19.5	20.9	21.3	20.3	27.3	79.1	28.5

Table 4. **Multi-task generation performance on BAIR and SSV2 evaluated by FVD.** Gray values denote unseen tasks during training. We list per-task FVD for all eight tasks on BAIR and the two extra tasks on SSV2 here, and leave the details for SSV2 in Appendix C.

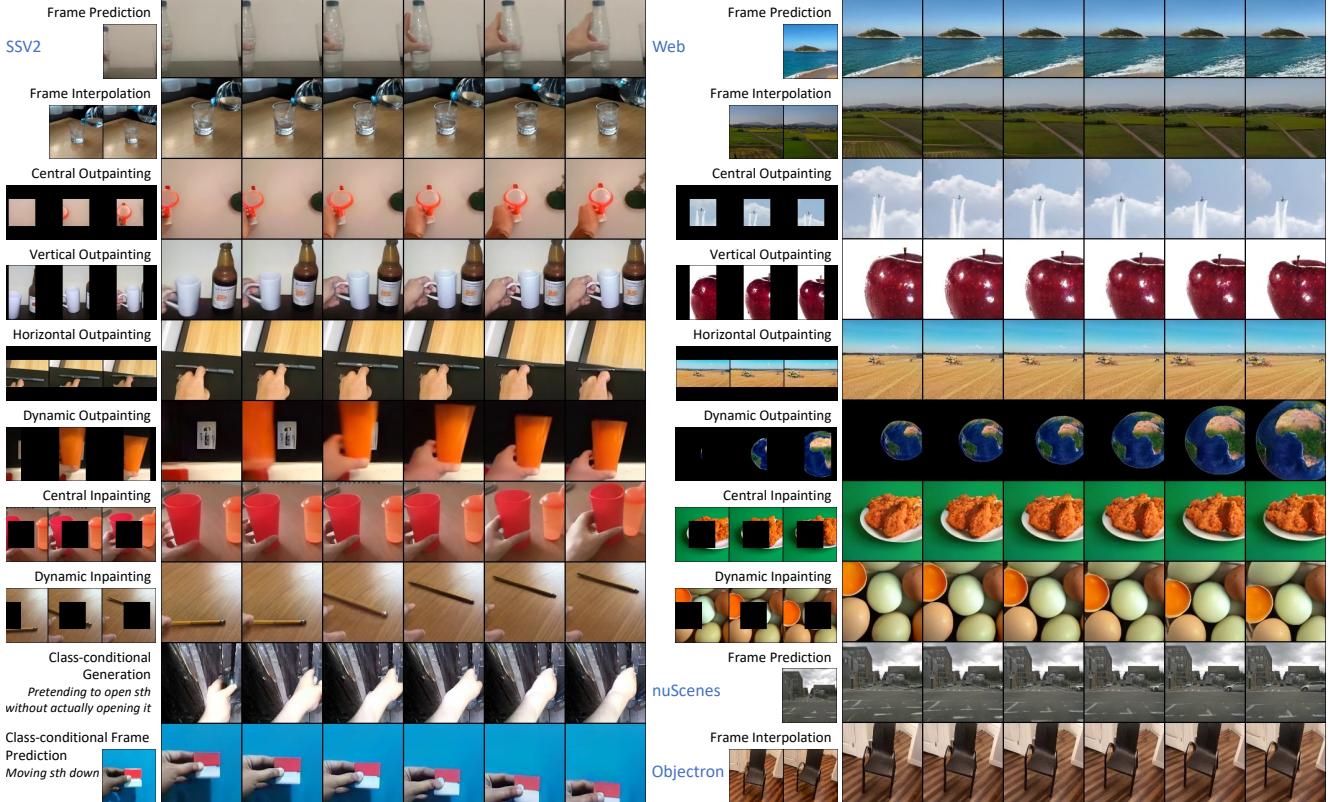


Figure 6. **Multi-task generation samples on four datasets: SSv2 [23], nuScenes [9], Objectron [2], and Web videos.** The left column is from a single ten-task model on SSv2, while the top eight rows on the right are from a single eight-task model on Web data.

Method	Seq. Length	FP FVD↓	MT8 FVD↓
Latent masking in MaskGIT [12]	1024	74	151
Prefix condition	1024-1792	55	-
<i>COMMIT</i> (ours)	$\mathcal{L}_{\text{mask}}$ $\mathcal{L}_{\text{mask}} + \mathcal{L}_{\text{recons}}$ $\mathcal{L}_{\text{mask}} + \mathcal{L}_{\text{recons}} + \mathcal{L}_{\text{refine}}$	388 51 48	143 53 33

Table 5. **Comparison of conditional masked token modeling** on BAIR frame prediction (FP) and eight-task (MT8) benchmarks. - indicates we were not able to train to convergence.

ods, including the latent masking in MaskGIT for image synthesis [12] and the commonly-used prefix condition that prepends cropped condition tokens to the input sequence.

Tab. 5 compares these methods on the BAIR dataset where the same 3D-VQ tokenizer is used in all approaches. As discussed in Section 3.2, latent masking in [12], which

directly unmasks tokens of the condition region at inference time, leads to poor generalization, especially for the multi-task setup. Prefix condition produces a long sequence of variable length, making it less tractable for multi-task learning. In contrast, COMMIT yields a fixed-length sequence and better generalizability for both single- and multi-task setups.

Training losses. The bottom section of Tab. 5 shows the contribution of the training loss components in Eq. (4).

Decoding methods. Tab. 6 compares Algorithm 1 with existing autoregressive (AR) and non-autoregressive (NAR) decoding methods. We consider two NAR baselines, *i.e.*, MaskGIT [12] for image and MaskViT [26] for video synthesis. We use the same 3D-VQ tokenizer for MaskGIT, AR, and MAGVIT. As shown, the proposed decoding algo-

Decoding Method	Tokenizer	Type	Param.	Seq. Len. \downarrow	# Steps \downarrow	FVD \downarrow
MaskGIT [12]	2D-VQ	NAR	53M+87M	4096	12	222 (177)
	3D-VQ	NAR	41M+87M	1024	12	122 (74)
MaskViT [26]	2D-VQ	NAR	53M+189M	4096	18	94*
AR	3D-VQ	AR	41M+87M	1024	1024	91 (56)
MAGVIT (ours)	3D-VQ	NAR	41M+87M	1024	12	76 (48)

Table 6. **Comparison of decoding methods** on BAIR frame prediction benchmark. The number of parameters is broken down as VQ + Transformer. NAR is non-autoregressive and AR is autoregressive. FVD and debiased FVD (in parentheses) are reported. * marks the quoted number from their paper.

Tokenizer	From Scratch		ImageNet [16]		Initialization		
	FVD \downarrow	IS \uparrow	FVD \downarrow	IS \uparrow	FVD \downarrow	IS \uparrow	
MaskGIT [12] 2D-VQ	240	80.9	216	82.6	-	-	
TATS [21] 3D-VQ	162	80.6	-	-	-	-	
		Average		Central			
MAGVIT 3D-VQ-B (ours)	127	82.1	103	84.8	58	87.0	
MAGVIT 3D-VQ-L (ours)	45	87.1	35	88.3	25	88.9	

Table 7. **Comparison of tokenizer architectures and initialization methods** on UCF-101 training set reconstruction results. The 2D-VQ compresses by 8×8 spatially and the 3D-VQ compresses by $4 \times 8 \times 8$ spatial-temporally.

rithm produces the best quality with the 3D-VQ and has a $4 \times$ shorter sequence than the 2D-VQ. While the AR transformer obtains a reasonable FVD, it takes over $85 \times$ more steps at inference time.

VQ architecture and training techniques. We evaluate the design options of our 3D-VQ model in MAGVIT. Tab. 7 lists the reconstruction FVD and IS metrics on the UCF-101 training set, which are different from the generation metrics as they measure the intermediate quantization. Nevertheless, reconstruction quality bounds the generation quality.

Tab. 7 compares the proposed 3D architecture with existing 2D [12] and 3D [21] VQ architectures. We train the MaskGIT [12] 2D-VQ and our 3D-VQ with the same protocol and evaluate the official TATS [21] 3D-VQ model. We compare two inflation methods for our 3D-VQ model, *i.e.*, average [11] and central inflation.

The results show the following. First, 3D-VQ models, despite producing a higher compression rate, show better video reconstruction quality than 2D-VQ, even with fewer parameters. Second, the proposed VQ performs favorably against baseline architectures with a similar size and gets much better with a larger model. Third, ImageNet [16] initialization boosts the performance for 2D and 3D models, where the central inflation outperforms the average inflation. The results demonstrate the excellent reconstruction fidelity of our tokenizer design.

5. Related Work

GAN-based approaches. Early success in video synthesis has been made by GAN models [1, 6, 7, 15, 27, 37, 48, 51,

55, 58, 65, 73]. Training instability and lack of generation diversity [12] are known issues of GAN models.

Autoregressive transformers. Inspired by the success of GPT [8], autoregressive transformers have been adapted for image [13, 18, 20, 46, 72] and video generation [4, 34, 68, 69]. A focus for video is autoregressive modeling of visual dynamics. Studies have switched from modeling the raw pixels [13, 61] to the discrete codes in a latent space [45, 70]. The state-of-the-art model TATS [21] uses two hierarchical transformers to reduce the computation for long video generation, with tokens learned by a 3D-VQGAN [20]. Unlike prior works, we introduce a non-autoregressive transformer with higher efficiency and flexibility.

Non-autoregressive transformers. Concurrently, a few methods use non-autoregressive transformers for image synthesis [12, 42, 75]. Section 2 reviews a state-of-the-art model called MaskGIT [12]. Compared with these approaches [26, 28], we present an embedding mask to model multi-task video conditions with better quality.

Diffusion models. Diffusion models have recently received much attention for image synthesis. For example, the state-of-the-art video diffusion model [33] extends the image denoising diffusion model [3, 32, 52, 53, 59] by incorporating 3D U-Net [14] architectures and joint training on both images and videos. Despite its high-quality, sampling speed is a bottleneck hindering the application of diffusion models in video synthesis. We show a different solution to train a highly-efficient model that offers compelling quality.

Multi-task video synthesis. Multi-task video synthesis [28, 44, 69] is yet to be well-studied. Transframer [44] is closest to our work, which adopts an image-level representation for autoregressive modeling of tasks based on frame prediction. We present an efficient non-autoregressive transformer for multi-task generation, and verify the quality and efficiency on ten video generation tasks.

Text-to-video. All of our models are trained only on public benchmarks, except the Web video model. We leave the text-to-video task as our future work. As shown in [31, 50, 63], training such models requires large, and sometimes non-public, datasets of paired texts and images.

6. Conclusion

In this paper, we propose MAGVIT, a generic and efficient mask-based video generation model. We introduce a high-quality 3D-VQ tokenizer to quantize a video and design COMMIT for multi-task conditional masked token modeling. We conduct extensive experiments to demonstrate the video generation quality, efficiency, and flexibility for multi-task generation. Notably, MAGVIT establishes a new state-of-the-art quality for class conditional generation on UCF-101 and frame prediction on BAIR Robot Pushing and Kinetics-600 datasets.

Acknowledgements

The authors would like to thank Tom Duerig, Victor Gomes, Paul Natsev along with the Multipod committee for sponsoring the computing resources. We appreciate valuable feedback and leadership support from David Salesin, Jay Yagnikm, Tomas Izo, and Rahul Sukthankar throughout the project. Special thanks to Wolfgang Macherey for supporting the project. We thank David Alexander Ross and Yu-Chuan Su for many helpful comments for improving the paper. We also give thanks to Sarah Laszlo and Hugh Williams for creating the MAGVIT model card, Bryan Seybold and Albert Shaw for extending the features, Jonathan Ho and Tim Salimans for providing the JAX code pointer for FVD computation, and the Scenic team for the infrastructure support. We are thankful to Wenhe Liu, Xinyu Yao, Mingzhi Cai, Yizhi Zhang, and Zhao Jin for proof reading the paper.

References

- [1] Dinesh Acharya, Zhiwu Huang, Danda Pani Paudel, and Luc Van Gool. Towards high resolution video generation with progressive growing of sliced wasserstein gans. *arXiv:1810.02419*, 2018. 8
- [2] Adel Ahmadyan, Liangkai Zhang, Artsiom Ablavatski, Jianing Wei, and Matthias Grundmann. Objectron: A large scale dataset of object-centric videos in the wild with pose annotations. In *CVPR*, 2021. 2, 5, 7
- [3] Jacob Austin, Daniel Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. Structured denoising diffusion models in discrete state-spaces. In *NeurIPS*, 2021. 8
- [4] Mohammad Babaeizadeh, Mohammad Taghi Saffar, Suraj Nair, Sergey Levine, Chelsea Finn, and Dumitru Erhan. Fitvid: Overfitting in pixel-level video prediction. *arXiv:2106.13195*, 2021. 5, 8
- [5] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. 5
- [6] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. In *ICLR*, 2018. 8
- [7] Tim Brooks, Janne Hellsten, Miika Aittala, Ting-Chun Wang, Timo Aila, Jaakko Lehtinen, Ming-Yu Liu, Alexei A Efros, and Tero Karras. Generating long videos of dynamic scenes. *arXiv:2206.03429*, 2022. 8
- [8] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *NeurIPS*, 2020. 8
- [9] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liang, Qiang Xu, Anush Krishnan, Yu Pan, Giacarlo Baldan, and Oscar Beijbom. nuScenes: A multimodal dataset for autonomous driving. In *CVPR*, 2020. 2, 5, 7
- [10] Joao Carreira, Eric Noland, Andras Banki-Horvath, Chloe Hillier, and Andrew Zisserman. A short note about Kinetics-600. *arXiv:1808.01340*, 2018. 2, 5
- [11] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the Kinetics dataset. In *CVPR*, 2017. 3, 5, 8
- [12] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. MaskGIT: Masked generative image transformer. In *CVPR*, 2022. 1, 2, 3, 4, 7, 8
- [13] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pre-training from pixels. In *ICML*, 2020. 8
- [14] Özgün Çiçek, Ahmed Abdulkadir, Soeren S Lienkamp, Thomas Brox, and Olaf Ronneberger. 3D U-Net: learning dense volumetric segmentation from sparse annotation. In *MICCAI*, 2016. 6, 8
- [15] Aidan Clark, Jeff Donahue, and Karen Simonyan. Adversarial video generation on complex datasets. *arXiv:1907.06571*, 2019. 1, 5, 8
- [16] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009. 8
- [17] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2019. 2, 5
- [18] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. CogView: Mastering text-to-image generation via transformers. In *NeurIPS*, 2021. 1, 8
- [19] Frederik Ebert, Chelsea Finn, Alex X Lee, and Sergey Levine. Self-supervised visual planning with temporal skip connections. In *CoRL*, 2017. 2, 5
- [20] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *CVPR*, 2021. 1, 2, 3, 8
- [21] Songwei Ge, Thomas Hayes, Harry Yang, Xi Yin, Guan Pang, David Jacobs, Jia-Bin Huang, and Devi Parikh. Long video generation with time-agnostic VQGAN and time-sensitive transformer. In *ECCV*, 2022. 1, 2, 3, 5, 6, 8
- [22] Marjan Ghazvininejad, Omer Levy, Yinhan Liu, and Luke Zettlemoyer. Mask-Predict: Parallel decoding of conditional masked language models. In *EMNLP-IJCNLP*, 2019. 2
- [23] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The “something something” video database for learning and evaluating visual common sense. In *ICCV*, 2017. 2, 5, 7
- [24] Jiatao Gu and Xiang Kong. Fully non-autoregressive neural machine translation: Tricks of the trade. In *ACL-IJCNLP Findings*, 2021. 2
- [25] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *CVPR*, 2022. 1, 2

- [26] Agrim Gupta, Stephen Tian, Yunzhi Zhang, Jiajun Wu, Roberto Martín-Martín, and Li Fei-Fei. MaskViT: Masked visual pre-training for video prediction. *arXiv:2206.11894*, 2022. 2, 3, 5, 6, 7, 8
- [27] Sonam Gupta, Arti Keshari, and Sukhendu Das. RV-GAN: Recurrent GAN for unconditional video generation. In *CVPRW*, 2022. 8
- [28] Ligong Han, Jian Ren, Hsin-Ying Lee, Francesco Barbieri, Kyle Olszewski, Shervin Minaee, Dimitris Metaxas, and Sergey Tulyakov. Show me what and tell me how: Video synthesis via multimodal conditioning. In *CVPR*, 2022. 2, 8
- [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3
- [30] Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2020. 5
- [31] 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:2210.02303*, 2022. 8
- [32] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 8
- [33] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. In *ICLR Workshops*, 2022. 1, 2, 5, 6, 8
- [34] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. CogVideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv:2205.15868*, 2022. 5, 8
- [35] Tobias Höppe, Arash Mehrjou, Stefan Bauer, Didrik Nielsen, and Andrea Dittadi. Diffusion models for video prediction and infilling. *arXiv:2206.07696*, 2022. 1, 2, 5
- [36] Norman P Jouppi, Doe Hyun Yoon, Matthew Ashcraft, Mark Gottscho, Thomas B Jablin, George Kurian, James Laudon, Sheng Li, Peter Ma, Xiaoyu Ma, et al. Ten lessons from three generations shaped google’s TPUs. In *ISCA*, 2021. 2, 6
- [37] Emmanuel Kahembwe and Subramanian Ramamoorthy. Lower dimensional kernels for video discriminators. *Neural Networks*, 132:506–520, 2020. 8
- [38] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019. 3
- [39] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2014. 5
- [40] Xiang Kong, Lu Jiang, Huiwen Chang, Han Zhang, Yuan Hao, Haifeng Gong, and Irfan Essa. BLT: Bidirectional layout transformer for controllable layout generation. In *ECCV*, 2022. 2
- [41] Guillaume Le Moing, Jean Ponce, and Cordelia Schmid. CCVS: Context-aware controllable video synthesis. In *NeurIPS*, 2021. 1, 3, 5, 6
- [42] José Lezama, Huiwen Chang, Lu Jiang, and Irfan Essa. Improved masked image generation with Token-Critic. In *ECCV*, 2022. 2, 8
- [43] Pauline Luc, Aidan Clark, Sander Dieleman, Diego de Las Casas, Yotam Doron, Albin Cassirer, and Karen Simonyan. Transformation-based adversarial video prediction on large-scale data. *arXiv:2003.04035*, 2020. 1, 5
- [44] Charlie Nash, João Carreira, Jacob Walker, Iain Barr, Andrew Jaegle, Mateusz Malinowski, and Peter Battaglia. Transframer: Arbitrary frame prediction with generative models. *arXiv:2203.09494*, 2022. 1, 5, 8
- [45] Ruslan Rakhimov, Denis Volkonskiy, Alexey Artemov, Denis Zorin, and Evgeny Burnaev. Latent video transformer. In *VISIGRAPP (5: VISAPP)*, 2021. 1, 8
- [46] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML*, 2021. 1, 2, 8
- [47] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 1
- [48] Masaki Saito, Eiichi Matsumoto, and Shunta Saito. Temporal generative adversarial nets with singular value clipping. In *ICCV*, 2017. 1, 8
- [49] Masaki Saito, Shunta Saito, Masanori Koyama, and So-suke Kobayashi. Train sparsely, generate densely: Memory-efficient unsupervised training of high-resolution temporal gan. *IJCV*, 128(10):2586–2606, 2020. 1, 5
- [50] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv:2209.14792*, 2022. 5, 8
- [51] Ivan Skorokhodov, Sergey Tulyakov, and Mohamed Elhosiny. StyleGAN-V: A continuous video generator with the price, image quality and perks of StyleGAN2. In *CVPR*, 2022. 5, 8
- [52] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, 2015. 8
- [53] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *NeurIPS*, 2019. 8
- [54] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. *arXiv:1212.0402*, 2012. 2, 3, 5
- [55] Yu Tian, Jian Ren, Menglei Chai, Kyle Olszewski, Xi Peng, Dimitris N Metaxas, and Sergey Tulyakov. A good image generator is what you need for high-resolution video synthesis. In *ICLR*, 2021. 8
- [56] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3D convolutional networks. In *ICCV*, 2015. 5
- [57] Hung-Yu Tseng, Lu Jiang, Ce Liu, Ming-Hsuan Yang, and Weilong Yang. Regularizing generative adversarial networks under limited data. In *CVPR*, 2021. 3
- [58] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. MoCoGAN: Decomposing motion and content for video generation. In *CVPR*, 2018. 1, 8
- [59] Belinda Tzen and Maxim Raginsky. Neural stochastic differential equations: Deep latent gaussian models in the diffusion limit. *arXiv:1905.09883*, 2019. 8

- [60] 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:1812.01717*, 2018. 1, 2, 5
- [61] Aäron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *ICML*, 2016. 8
- [62] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. In *NeurIPS*, 2017. 2
- [63] Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang, Mohammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. Phenaki: Variable length video generation from open domain textual description. *arXiv:2210.02399*, 2022. 5, 8
- [64] Vikram Voleti, Alexia Jolicoeur-Martineau, and Christopher Pal. Masked conditional video diffusion for prediction, generation, and interpolation. In *NeurIPS*, 2022. 1, 5
- [65] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. In *NeurIPS*, 2016. 1, 8
- [66] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhajit Som, et al. Image as a foreign language: BEiT pretraining for all vision and vision-language tasks. *arXiv:2208.10442*, 2022. 2
- [67] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE TIP*, 13(4):600–612, 2004. 5
- [68] Dirk Weissenborn, Oscar Täckström, and Jakob Uszkoreit. Scaling autoregressive video models. In *ICLR*, 2019. 1, 8
- [69] Chenfei Wu, Jian Liang, Lei Ji, Fan Yang, Yuejian Fang, Daxin Jiang, and Nan Duan. NÜWA: Visual synthesis pre-training for neural visual world creation. In *ECCV*, 2022. 1, 5, 8
- [70] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. VideoGPT: Video generation using vq-vae and transformers. *arXiv:2104.10157*, 2021. 2, 3, 8
- [71] Jiahui Yu, Xin Li, Jing Yu Koh, Han Zhang, Ruoming Pang, James Qin, Alexander Ku, Yuanzhong Xu, Jason Baldridge, and Yonghui Wu. Vector-quantized image modeling with improved VQGAN. In *ICLR*, 2022. 2
- [72] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv:2206.10789*, 2022. 1, 8
- [73] Sihyun Yu, Jihoon Tack, Sangwoo Mo, Hyunsu Kim, Junho Kim, Jung-Woo Ha, and Jinwoo Shin. Generating videos with dynamics-aware implicit generative adversarial networks. In *ICLR*, 2022. 5, 8
- [74] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 5
- [75] Zhu Zhang, Jianxin Ma, Chang Zhou, Rui Men, Zhikang Li, Ming Ding, Jie Tang, Jingren Zhou, and Hongxia Yang. M6-UFC: Unifying multi-modal controls for conditional image synthesis. *arXiv:2105.14211*, 2021. 2, 8

MAGVIT: Masked Generative Video Transformer

Supplementary Materials

Overview

This supplementary document provides additional details to support our main manuscript, organized as follows:

- Appendix A presents the 3D-VQ architectures and the transformer models in MAGVIT.
- Appendix B includes additional implementation details in training and evaluation.
- Appendix C provides more quantitative evaluation results, which include:
 - Comparisons to more published results on the three benchmarks in the paper: UCF-101 [31], BAIR [12, 35], and Kinetics-600 [6].
 - Multi-task results on Something-Something-v2 (SSv2) [14].
 - Results on three additional datasets: NuScenes [5], Objectron [3] and Web video datasets.
- Appendix D shows more qualitative examples of the generated videos.

We attach a demo video for MAGVIT and show more generated examples on this web page¹.

A. MAGVIT Model Architecture

A.1. 3D-VQ Tokenizer

Fig. 1 shows the architectures of the MAGVIT 3D-VQ module and compares it with the 3D-VQ module in TATS [13] which held the previous state-of-the-art for video generation. Compared with TATS, the major design choices in MAGVIT 3D-VQ are listed below.

- Average pooling, instead of strided convolution, is used for down-sampling.
- Nearest resizing and convolution are used for up-sampling.
- We use spatial down- and up-sampling layers near the latent space and spatial-temporal down- and up-sampling layers near the pixel space, resulting in mirrored encoder-decoder architecture.
- A single deeper 3D discriminator is designed rather than two shallow discriminators for 2D and 3D separately.
- We quantize into a much smaller vocabulary of 1,024 as compared to 16,384.

¹<https://magvit.cs.cmu.edu>

- We use group normalization [42] instead of batch normalization [20] and Swish [26] activation function instead of SiLU [16].

The quantitative comparison of the 3D-VQ from TATS and MAGVIT were presented in Table 6 of the main paper. In addition, Fig. 3 below qualitatively compares their reconstruction quality on UCF-101. Figs. 4 and 5 show MAGVIT’s high-quality reconstruction on example YouTube videos.

We design two variants of the MAGVIT 3D-VQ module, *i.e.*, the base (B) with 41M parameters and the large (L) with 158M parameters, excluding the discriminators.

A.2. Transformer

MAGVIT uses the BERT transformer architecture [10] adapted from the Flaxformer implementation². Following the transformer configurations in ViT [11], we use two variants of transformers, *i.e.*, base (B) with 87M parameters and large (L) with 306M in all our experiments. Tab. 1 lists the detailed configurations for each variant. A huge (H) transformer is only used to train on the large Web video dataset and generate demo videos.

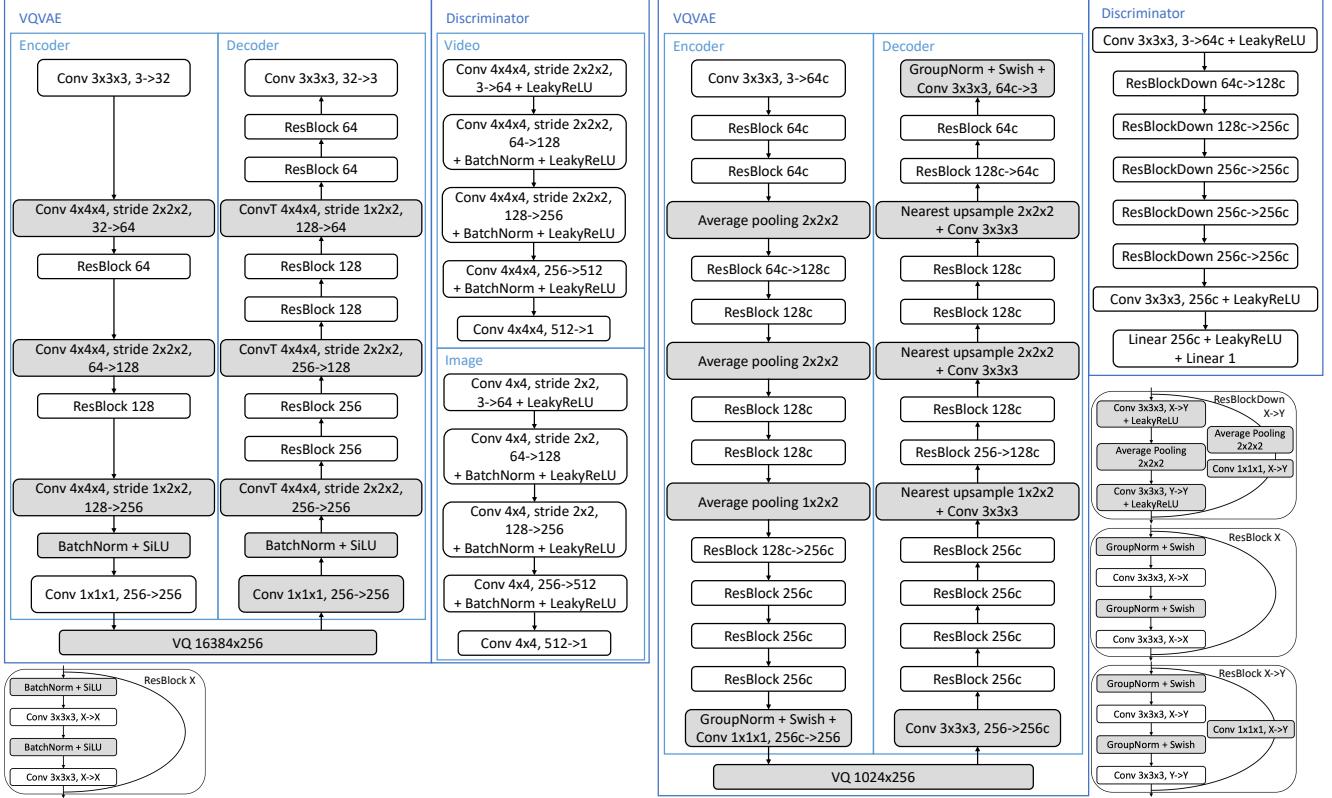
B. Implementation Details

B.1. Task Definitions

We employ a total of ten tasks for multi-task video generation. Each task is characterized by a few adjustable settings such as interior condition shape, padding function, and optionally prefix condition. Fig. 2 illustrates the interior condition regions for each task under the above setup. Given a video of shape $T \times H \times W$, we define the tasks as following:

- Frame Prediction (FP)
 - Interior condition: t frames at the beginning; $t = 1$.
 - Padding: replicate the last given frame.
- Frame Interpolation (FI)
 - Interior condition: t_1 frames at the beginning and t_2 frames at the end; $t_1 = 1, t_2 = 1$.
 - Padding: linear interpolate between the last given frame at the beginning and the first given frame at the end.
- Central Outpainting (OPC)

²<https://github.com/google/flaxformer>



(a) TATS [13] 3D-VQ. (32M+14M parameters)

(b) MAGVIT (ours) 3D-VQ.
(41M+15M parameters at $c = 1$ (B), 158M+61M at $c = 2$ (L))

Figure 1. **Comparison of 3D-VQ model architectures between MAGVIT and the TATS [13].** We highlight the blocks with major differences in gray background and detail their design differences in Appendix A.1. We train the models to quantize 16-frame clips of 128×128 resolution into $4 \times 16 \times 16$ tokens. The number of parameters in parentheses are broken down between VQVAE and discriminators.

Model	Param.	# heads	# layers	Hidden size	MLP dim
MAGVIT-B	87 M	12	12	768	3072
MAGVIT-L	305 M	16	24	1024	4096
MAGVIT-H	634 M	16	32	1280	5120

Table 1. Transformer architecture configurations used in MAGVIT.

- Interior condition: a rectangle at the center with height h and width w ; $h = 0.5H$, $w = 0.5W$.
- Padding: pad the nearest pixel for each location (edge padding).
- Vertical Outpainting (OPV)
 - Interior condition: a centered vertical strip with width w ; $w = 0.5W$.
 - Padding: edge padding.
- Horizontal Outpainting (OPH)
 - Interior condition: a centered horizontal strip with height h ; $h = 0.5H$.
 - Padding: edge padding.

- Dynamic Outpainting (OPD)
 - Interior condition: a moving vertical strip with width w ; $w = 0.5W$.
 - Direction of movement: left to right.
 - Padding: zero padding.
- Central Inpainting (IPC)
 - Interior condition: everything but a rectangle at the center with height h and width w ; $h = 0.5H$, $w = 0.5W$.
 - Padding: zero padding.
- Dynamic Inpainting (IPD)
 - Interior condition: everything but a vertically centered

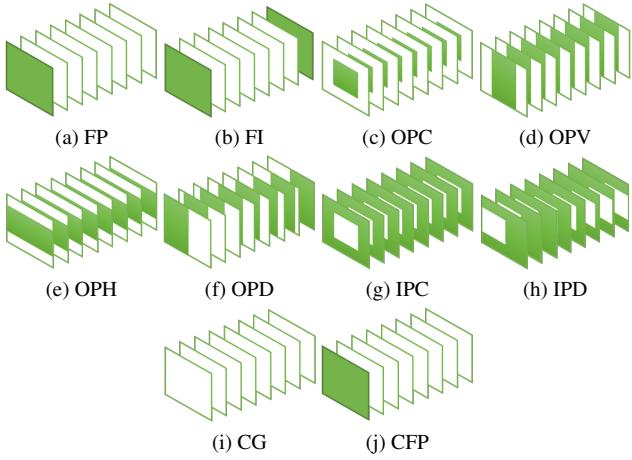


Figure 2. **Interior condition regions for each task**, where green denotes valid pixels and white pixels denote the task-specific paddings discussed in Appendix B.1. The tasks are Frame Prediction (FP), Frame Interpolation (FI), Central Outpainting (OPC), Vertical Outpainting (OPV), Horizontal Outpainting (OPH), Dynamic Outpainting (OPD), Central Inpainting (IPC), Dynamic Inpainting (IPD), Class-conditional Generation (CG), and Class-conditional Frame Prediction (CFP).

moving rectangle with height h and width w ; $h = 0.5H, w = 0.5W$.

- Direction of movement: left to right.
- Padding: zero padding.
- Class-conditional Generation (CG)
 - Prefix condition: class label.
- Class-conditional Frame Prediction (CFP)
 - Prefix condition: class label.
 - Interior condition: t frames at the beginning; $t = 1$.
 - Padding: replicate the last given frame.

B.2. Training

MAGVIT is trained in two stages where we first train the 3D-VQ tokenizer and then train the transformer with a frozen tokenizer. We follow the same learning recipe across all datasets, with the only variation in the number of training epochs. Here are the training details for both stages:

- 3D-VQ:
 - Video: 16 frames, frame stride 1, 128×128 resolution. (64×64 resolution for BAIR)
 - Base channels: 64 for B, 128 for L.
 - VQVAE channel multipliers: 1, 2, 2, 4. (1, 2, 4 for 64×64 resolution).
 - Discriminator channel multipliers: 2, 4, 4, 4, 4. (2, 4, 4, 4 for 64×64 resolution)
 - Latent shape: $4 \times 16 \times 16$.
 - Vocabulary size: 1,024.
 - Embedding dimension: 256.
 - Initialization: central inflation from a 2D-VQ trained

Dataset	3D-VQ		Transformer	
	B	L	B	L
UCF-101	500	2000	2000	2000
BAIR	400	800	400	800
BAIR-MT	400	800	1200	1600
Kinetics-600	45	180	180	360
SSv2	135	400	720	1440
nuScenes	1280	5120	2560	10240
Objectron	1000	2000	1000	2000
Web	5	20	10	20

Table 2. Training epochs for each dataset.

on ImageNet with this setup.

- Peak learning rate: 10^{-4} .
- Learning rate schedule: linear warm up and cosine decay.
- Optimizer: Adam with $\beta_1 = 0$ and $\beta_2 = 0.99$.
- Generator loss type: Non-saturating.
- Generator adversarial loss weight: 0.1.
- Perceptual loss weight: 0.1.
- Discriminator gradient penalty: r1 with cost 10.
- EMA model decay rate: 0.999.
- Batch size: 128 for B, 256 for L.
- Speed: 0.41 steps/sec on 16 TPU-v2 chips for B, 0.56 steps/sec on 32 TPU-v4 chips for L.
- Transformer:
 - Sequence length: 1026.
 - Hidden dropout rate: 0.1.
 - Attention dropout rate: 0.1.
 - Mask rate schedule: cosine.
 - Peak learning rate: 10^{-4} .
 - Learning rate schedule: linear warm up and cosine decay.
 - Optimizer: Adam with $\beta_1 = 0.9$ and $\beta_2 = 0.96$.
 - Weight decay 0.045.
 - Label smoothing: 10^{-4} .
 - Max gradient norm: 1.
 - Batch size: 256.
 - Speed: 1.24 steps/sec on 16 TPU-v2 chips for B, 2.70 steps/sec on 32 TPU-v4 chips for L.

Using more hardware resources can speed up the training. We train MAGVIT models for each dataset separately. The training epochs for each dataset are listed in Tab. 2.

B.3. Evaluation

Evaluation metrics. The FVD [35] is used as the primary evaluation metric. We follow the official implementation³ in extracting video features with an I3D model trained on

³https://github.com/google-research/google-research/tree/master/frechet_video_distance

Kinetics-400 [7]. We report Inception Score (IS) [28]⁴ on the UCF-101 dataset which is calculated with a C3D [33] model trained on UCF-101. We further include image quality metrics: PSNR, SSIM [39] and LPIPS [45] (computed by the VGG features) on the BAIR dataset.

Sampling protocols. We follow the sampling protocols from previous works [9, 13] when evaluating on the standard benchmarks, *i.e.* UCF-101, BAIR, and Kinetics-600. We sample 16-frame clips from each dataset without replacement to form the real distribution in FVD and extract condition inputs from them to feed to the model. We continuously run through all the samples required (*e.g.*, 40,000 for UCF-101) with a single data loader and compute the mean and standard deviation for 4 folds. When evaluating on other datasets, due to the lack of prior works, we adapt the above protocol based on the dataset size to ensure sample diversity.

For our MAGVIT model, we use the following COMMIT decoding hyperparameters by default: cosine schedule, 12 steps, temperature 4.5. Below are detailed setups for each dataset:

- UCF-101:
 - Dataset: 9.5K videos for training, 101 classes.
 - Number of samples: $10,000 \times 4$.
 - Resolution: 128×128 .
 - Real distribution: random clips from the training videos.
- BAIR:
 - Dataset: 43K videos for training and 256 videos for evaluation.
 - Number of samples: $25,600 \times 4$.
 - Resolution: 64×64 .
 - Real distribution: the first 16-frame clip from each evaluation video.
 - COMMIT decoding: exponential schedule, temperature 400.
- Kinetics-600:
 - Dataset: 384K videos for training and 29K videos for evaluation.
 - Number of samples: $50,000 \times 4$.
 - Generation resolution: 128×128 .
 - Evaluation resolution: 64×64 , via central crop and bilinear resize.
 - Real distribution: 6 sampled clips (2 temporal windows and 3 spatial crops) from each evaluation video.
 - COMMIT decoding: uniform schedule, temperature 7.5.
- SSv2:
 - Dataset: 169K videos for training and 24K videos for evaluation, 174 classes.
 - Number of samples: $50,000 \times 4$.

⁴<https://github.com/pfnet-research/tgan2>

- Resolution: 128×128 .
- Real distribution for the CG task: random clips from the training videos.
- Real distribution for the other tasks: 2 sampled clips (2 temporal windows and central crop) from each evaluation video.
- nuScenes:
 - Dataset: 5.4K videos for training and 0.6K videos for evaluation, front camera only, 32 frames per video.
 - Number of samples: $50,000 \times 4$.
 - Resolution: 128×128 .
 - Real distribution: 48 sampled clips (16 temporal windows and 3 spatial crops) from each evaluation video.
- Objectron:
 - Dataset: 14.4K videos for training and 3.6K videos for evaluation.
 - Number of samples: $50,000 \times 4$.
 - Resolution: 128×128 .
 - Real distribution: 5 sampled clips (5 temporal windows and central crop) from each evaluation video.
- Web videos:
 - Dataset: $\sim 12M$ videos for training and 26K videos for evaluation.
 - Number of samples: $50,000 \times 4$.
 - Resolution: 128×128 .
 - Real distribution: randomly sampled clips from evaluation videos.

For the “random clips” above, we refer to the combination of a random temporal window and a random spatial crop on a random video. For the fixed number of “temporal windows” or “spatial crops”, deterministic uniform sampling is used.

For the image quality metrics on BAIR in Table 3 of the main paper, CCVS [22] generates at 256×256 while the others are at 64×64 . When calculating PSNR and SSIM, we follow [37] in using the best value from 100 trials for each evaluation video.

Debiased FVD on BAIR Computing FVD is difficult on the BAIR dataset due to its small evaluation target of only 256 16-frame clips. Following the standard evaluation protocol, we generate 100 predictions for each clip to create 256,00 samples [4].

The real distribution to compute FVD in this way is highly biased with the insufficient evaluation videos [35]. We can see this by a simple experiment where we compute the training FVD with only 256 training videos. We observe that this 256-sample training FVD (64) is far worse than the regular training FVD with all 43K videos (13), showing the biased FVD computation.

To bridge the gap, we use uniformly sampled 16-frame clips from the 256 30-frame evaluation videos, which results in $256 \times 15 = 3840$ clips. The uniform sampling

Method	Extra Video	Class	FVD↓	IS↑
VGAN [38]		✓	-	8.31 ± 0.09
TGAN [27]			-	11.85 ± 0.07
MoCoGAN* [34]		✓	-	12.42 ± 0.07
ProgressiveVGAN [2]		✓	-	14.56 ± 0.05
TGAN [27]		✓	-	15.83 ± 0.18
RaMViD [19]			-	21.71 ± 0.21
LDVD-GAN [21]			-	22.91 ± 0.19
StyleGAN-V*# [30]			-	23.94 ± 0.73
VideoGPT [43]			-	24.69 ± 0.30
TGANv2 [28]		✓	1209 ± 28	28.87 ± 0.67
MoCoGAN-HD# [32]			838	32.36
DIGAN [44]			655 ± 22	29.71 ± 0.53
DIGAN# [44]			577 ± 21	32.70 ± 0.35
DVD-GAN# [9]		✓	-	32.97 ± 1.70
Video Diffusion*# [17]			-	57.00 ± 0.62
TATS [13]			420 ± 18	57.63 ± 0.24
CCVS+StyleGAN# [22]			386 ± 15	24.47 ± 0.13
Make-A-Video* [29]		✓	367	33.00
TATS [13]		✓	332 ± 18	79.28 ± 0.38
CogVideo* [18]		✓	✓	626
Make-A-Video* [29]		✓	✓	81
MAGVIT-B-CG (ours)		✓	159 ± 2	83.55 ± 0.14
MAGVIT-L-CG (ours)		✓	76 ± 2	89.27 ± 0.15

Table 3. **Generation performance on the UCF-101 dataset.** Methods in gray are pretrained on additional large video data. Methods with ✓ in the Class column are class-conditional, while the others are unconditional. Methods marked with * use custom resolutions, while the others are at 128×128 . Methods marked with # additionally used the test set in training.

yields a better representation of the evaluation set. Under this new protocol, MAGVIT-L-FP achieves FVD 31 instead of 62, which is more aligned with its training set performance (FVD=8).

We report this “debiased FVD” in addition to the standard FVD computation on the BAIR dataset, with the default COMMIT decoding hyperparameters. We also use it for BAIR multi-task evaluation and ablation studies on BAIR .

C. Additional Quantitative Evaluation

Class-conditional generation. Tab. 3 shows a detailed comparison with the previously published results on the UCF-101 [31] class-conditional video generation benchmark, where the numbers are quoted from the cited papers. Note that CogVideo [18] and Make-A-Video [29] are pretrained on additional 5-10M videos before finetuning on UCF-101, where Make-A-Video further uses a text-image

Method	K600 FVD↓	BAIR FVD↓
LVT [25]	224.7	126 ± 3
Video Transformer [40]	170.0 ± 5.0	94 ± 2
CogVideo* [18]	109.2	-
DVD-GAN-FP [9]	69.1 ± 1.2	110
CCVS [22]	55.0 ± 1.0	99 ± 2
Phenaki [36]	36.4 ± 0.2	97
VideoGPT [43]	-	103
TrIVD-GAN-FP [23]	25.7 ± 0.7	103
Transframer [24]	25.4	100
MaskViT [15]	-	94
FitVid [4]	-	94
MCVD [37]	-	90
NÜWA [41]	-	87
RaMViD [19]	16.5	84
Video Diffusion [17]	16.2 ± 0.3	-
MAGVIT-B-FP (ours)	24.5 ± 0.9	$76 \pm 0.1 (47 \pm 0.1)$
MAGVIT-L-FP (ours)	9.9 ± 0.3	$62 \pm 0.1 (31 \pm 0.2)$

Table 4. **Frame prediction performance on the BAIR and Kinetics-600 datasets.** - marks that the value is unavailable in their paper or incomparable to others. The FVD in parentheses uses a debiased evaluation protocol on BAIR detailed in Appendix B.3. Methods marked with * is pretrained on additional large video data.

prior trained on a billion text-image pairs. The remaining models, including MAGVIT, are only trained on 9.5K training videos of UCF-101, or 13.3K training and testing videos of UCF-101 for those marked with #. Fig. 6 provides a visual comparison to the baseline methods.

As shown, even the smaller MAGVIT-B performs favorably against previous state-of-the-art model TATS [13] by a large margin. MAGVIT-L pushes both the FVD ($332 \rightarrow 76$, $\downarrow 77\%$) and IS ($79.28 \rightarrow 89.27$, $\uparrow 13\%$) to a new level, while outperforming the contemporary work Make-A-Video [29] which is pretrained on significantly large extra training data.

Frame prediction. For the frame prediction task on BAIR Robot Pushing [12, 35] (1-frame condition) and Kinetics-600 [6] (5-frame condition), Tab. 4 provides a detailed comparison with previously published results. We use “-” to mark the FVDs that either is unavailable in their paper or incomparable to others. For example, Video Diffusion [17]’s FVD reported in their paper was on a different camera angle (top-down view `image_main`⁵) and is hence incomparable to others.

MAGVIT achieves state-of-the-art quality in terms of

⁵https://www.tensorflow.org/datasets/catalog/bair_robot_pushing_small

Method	Task	Avg↓	FP	FI	OPC	OPV	OPH	OPD	IPC	IPD	CG	CFP
MAGVIT-B-UNC	Single	258.8	278.8	91.0	67.5	27.3	36.2	711.5	319.3	669.8	107.7	279.0
MAGVIT-B-FP	Single	402.9	59.3	76.2	213.2	81.2	86.3	632.7	343.1	697.9	1780.0	59.3
MAGVIT-B-MT	Multi	43.4	71.5	38.0	38.8	23.3	26.1	33.4	23.3	25.3	94.7	59.3
MAGVIT-L-MT	Multi	27.3	33.8	25.0	21.1	16.8	17.0	23.5	13.5	15.0	79.1	28.5

Table 5. Multi-task generation performance on Something-Something-V2 evaluated by FVD. Gray values denote unseen tasks during training.

Method	nuScenes-FP	Objectron-FI	Web-MT8	FP	FI	OPC	OPV	OPH	OPD	IPC	IPD
MAGVIT-B	29.3	-	33.0	84.9	33.9	34.4	21.5	22.1	26.0	20.7	20.4
MAGVIT-L	20.6	26.7	21.6	45.5	30.9	19.9	15.3	14.5	20.2	12.0	14.7

Table 6. Generation performance on NuScenes, Objectron, and Web videos evaluated by FVD.

FVD on both datasets, with a 39% relative improvement on the large-scale Kinetics benchmark than the highly-competitive Video Diffusion baseline [17]. Fig. 7 and Fig. 8 below provide visual comparisons to the baseline methods on BAIR and Kinetics-600, respectively.

Multi-task video generation. Having verified single-task video generation, Tab. 5 shows per-task performance of the ten tasks on the large-scale Something-Something-v2 (SSv2) [14] dataset. SSv2 is a challenging dataset commonly used for action recognition, whereas this work benchmarks video generation on it for the first time. On this dataset, a model needs to synthesize 174 basic actions with everyday objects. Fig. 9 shows examples of generated videos for each task on this dataset.

We compare the multi-task models (MT) with two single-task baselines trained on unconditional generation (UNC) and frame prediction (FP). The multi-task models show consistently better average FVD across all tasks compared with the single-task baselines.

Results on nuScenes, Objectron, and 12M Web Videos. Tab. 6 shows the generation performance on three additional datasets, *i.e.*, nuScenes [5], Objectron [3], and 12M Web videos which contains 12 million videos we collected from the web. We evaluate our model on the frame prediction task on nuScenes, the frame interpolation task on Objectron, and the 8-task suite on the Web videos. Fig. 10 shows examples of generated videos for each task. The results substantiate the generalization performance of MAGVIT on videos from distinct visual domains and the multi-task learning recipe on large-scale data.

D. Qualitative Examples

D.1. High-Fidelity Tokenization

Comparison of tokenizers. Fig. 3 compares the reconstruction quality of the three VQ tokenizers on the UCF-101, including the 2D-VQ from MaskGIT [8], the 3D-VQ from TATS [13], and MAGVIT 3D-VQ, where the videos are taken from the UCF-101 training set. We obtain the TATS model from their official release ⁶. We train the MaskGIT 2D-VQ and MAGVIT 3D-VQ using the same protocol on the UCF-101 dataset.

We can see that the MaskGIT 2D-VQ produces a reasonable image quality, but falls short of frame consistency which causes significant flickering when played as a video (*e.g.*, the curtain color in the first row and the wall color in the third row). TATS 3D-VQ has a better temporal consistency but loses details for moving objects (*e.g.*, the woman’s belly in the second row). In contrast, our 3D VQ produces consistent frames with greater details reconstructed for both static and moving pixels.

Scalable tokenization. Since the tokenizers are trained in an unsupervised manner, they exhibit remarkable generalization performances and can be scaled to big data as no labels are required. To demonstrate this, we train a large MAGVIT 3D-VQ on the large YouTube-8M [1] dataset while ignoring the labels, and use the model to quantize randomly sampled videos on YouTube.

Figs. 4 and 5 show the original and reconstructed videos from YouTube at 240p (240 × 432) resolution with arbitrary lengths (*e.g.* 4,096 frames). Although the tokenizer is only trained with 16-frame 128×128 videos, it produces high reconstruction fidelity for high spatial-temporal resolutions that are unseen in training. Our 3D-VQ model com-

⁶<https://songweige.github.io/projects/tats/>

presses the video by a factor of 4 temporally, by 8×8 spatially, and by 2.4 (24 bits \rightarrow 10 bits) per element, yielding a $614.4 \times$ compression rate. Despite such high compression, the reconstructed results show stunning details and are almost indistinguishable from the real videos.

D.2. Single-Task Generation Examples

Fig. 6 compares the generated samples from CCVS+StyleGAN [22], the prior state-of-the-art TATS [13], and MAGVIT on the UCF-101 class-conditional generation benchmark. As shown in Fig. 6, CCVS+StyleGAN [22] gets a decent single-frame quality attributing to the pretrained StyleGAN, but yields little or no motion. TATS [13] generates some motion but with clear artifacts. In contrast, our model produces higher-quality frames with substantial motion.

Fig. 7 compares the generated samples between the state-of-the-art RaMViD [19] and MAGVIT on the BAIR frame prediction benchmark given 1-frame condition. As shown, the clips produced by MAGVIT maintaining a better visual consistency and spatial-temporal dynamics.

Fig. 8 compares the generated samples from RaMViD [19] and MAGVIT on the Kinetics-600 frame prediction benchmark given 5-frame condition. Note that RaMViD generates video in 64×64 and MAGVIT in 128×128 where the standard evaluation is carried out on 64×64 . As shown, given the conditioned frames, MAGVIT generates plausible actions with greater details.

D.3. Multi-Task Generation Examples

Fig. 9 shows multi-task generation results on 10 different tasks from a single model trained on SSv2. Fig. 10 shows multi-task samples for three other models trained on nuScenes, Objectron, and Web videos. These results substantiate the multi-task flexibility of MAGVIT.

The diverse video generation tasks that MAGVIT is capable of can enable many useful applications. For example, Figs. 11 and 12 show a few unrawide inpainting samples by repeatedly performing the vertical inpainting task. MAGVIT can easily generate nice large panorama videos given a small condition.



(a) MaskGIT [8] 2D-VQ



(b) TATS [13] 3D-VQ



(c) MAGVIT 3D-VQ-L (ours)



(d) Real

Figure 3. **Comparison of tokenizers on UCF-101 training set reconstruction.** Videos are reconstructed at 16 frames 64×64 resolution 25 fps and shown at 12.5 fps, with the ground truth in (d). MaskGIT 2D-VQ produces a reasonable image quality, but falls short of frame consistency which causes significant flickering when played as a video (e.g., the curtain color in the first row and the wall color in the third row). TATS 3D-VQ has a better temporal consistency but loses details for moving objects (e.g., the woman’s belly in the second row). In contrast, our 3D VQ produces consistent frames with greater details reconstructed for both static and moving pixels.

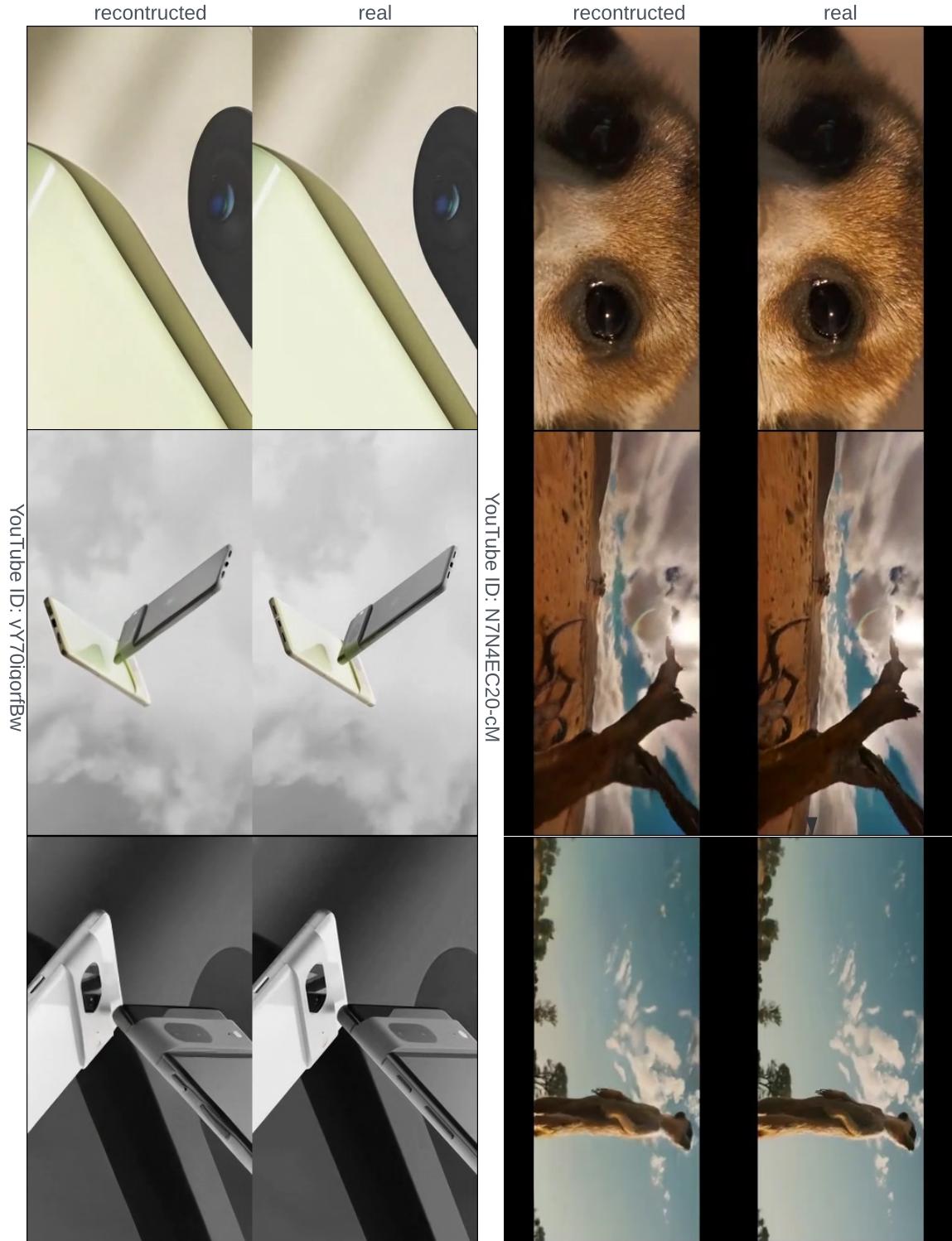


Figure 4. Our 3D-VQ model produces high reconstruction fidelity with scalable spatial-temporal resolution. For each group, the top row contains real YouTube videos and the bottom row shows the reconstructed videos from the discrete tokens. The original videos are in 240p (240×432) resolution with N frames. Our 3D-VQ model represents the video as $\frac{N}{4} \times 30 \times 54$ discrete tokens with a codebook of size 1024, representing a total compression rate of 614.4. Despite such high compression, the reconstructed results show stunning details and are almost indistinguishable from the real videos.

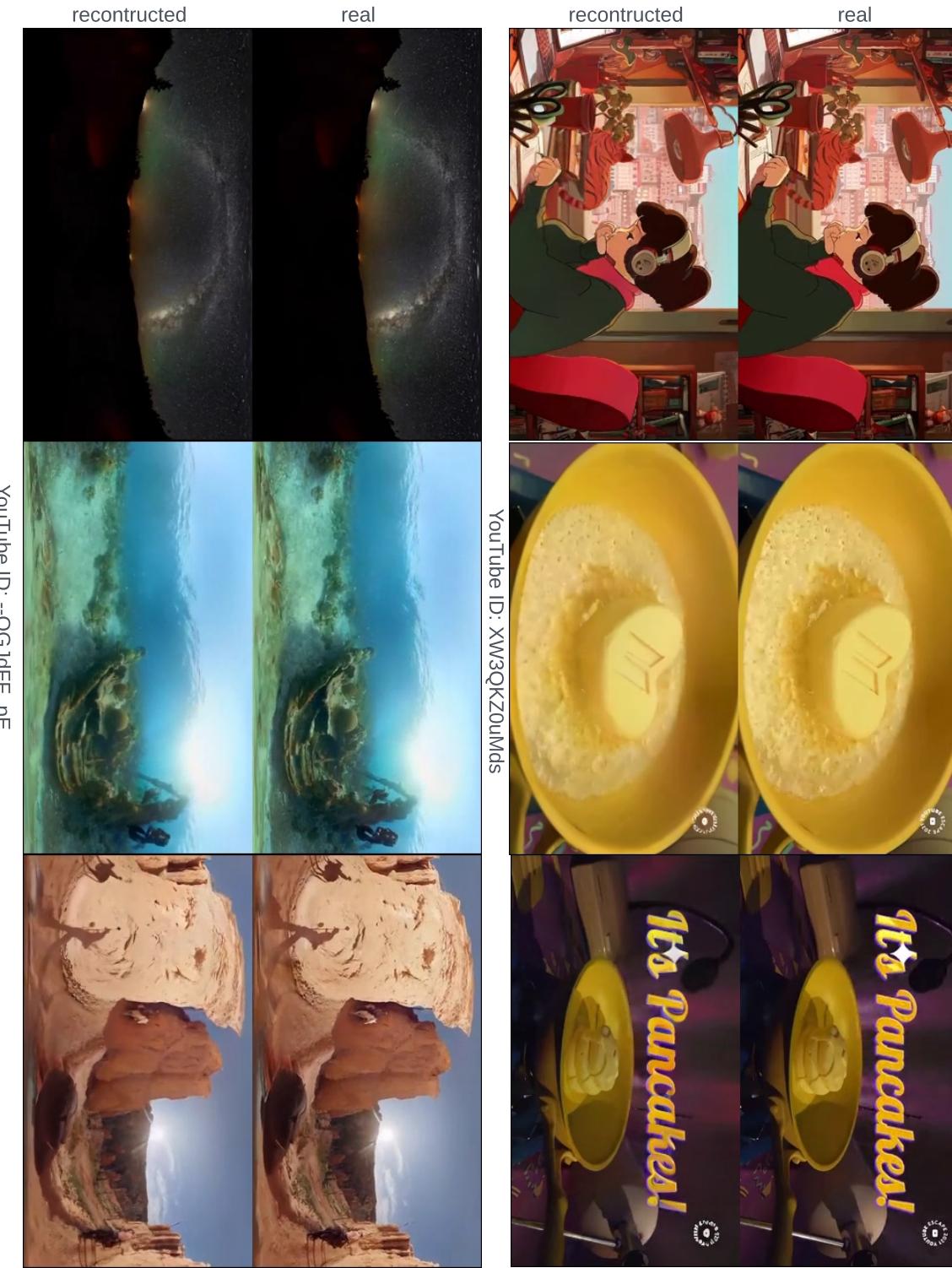


Figure 5. Our 3D-VQ model produces high reconstruction fidelity with scalable spatial-temporal resolution. For each group, the top row contains real YouTube videos and the bottom row shows the reconstructed videos from the discrete tokens. The original videos are in 240p (240×432) resolution with N frames. Our 3D-VQ model represents the video as $\frac{N}{4} \times 30 \times 54$ discrete tokens with a codebook of size 1024, representing a total compression rate of 614.4. Despite such high compression, the reconstructed results show stunning details and are almost indistinguishable from the real videos.



(a) CCVS+StyleGAN [22]



(b) TATS [13]

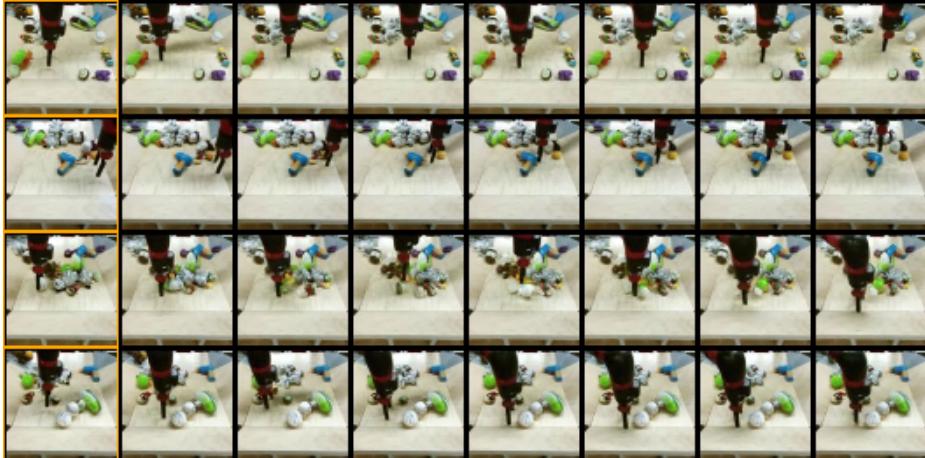


(c) MAGVIT-L-CG (ours)

Figure 6. **Comparison of class-conditional generation samples on UCF-101.** 16-frame videos are generated at 128×128 resolution 25 fps and shown at 12.5 fps. Samples for [13, 22] are obtained from their official release (<https://songweige.github.io/projects/tats/>). CCVS+StyleGAN gets a decent single-frame quality attributing to the pretrained StyleGAN, but yields little or no motion. TATS generates some motion but with clear artifacts. In contrast, our model produces higher-quality frames with substantial motion.



(a) RaMViD [19]

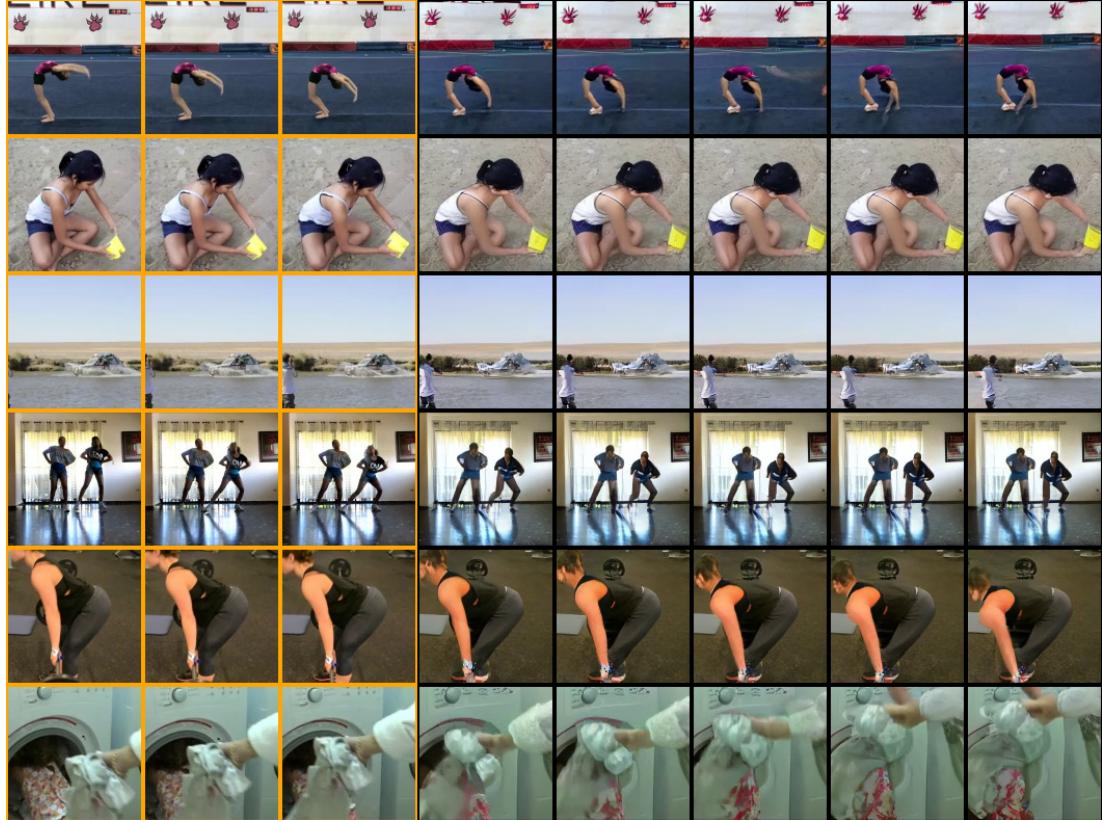


(b) MAGViT-L-FP (ours)

Figure 7. Comparison of frame prediction samples on BAIR unseen evaluation set. 16-frame videos are generated at 64×64 resolution 10 fps given the first frame as condition and shown at 5 fps where condition frames are marked in orange. Samples for [19] are obtained from their official release (<https://sites.google.com/view/video-diffusion-prediction>). As shown, the clips produced by MAGViT maintaining a better visual consistency and spatial-temporal dynamics.



(a) RaMViD [19] at 64×64 resolution, condition information is unavailable.



(b) MAGViT-L-FP (ours) at 128×128 resolution, condition frames are marked in orange.

Figure 8. Comparison of frame prediction samples on Kinetics-600 unseen evaluation set. 16-frame videos are generated at 25 fps given 5-frame condition. Samples for [19] are obtained from their official release (<https://sites.google.com/view/video-diffusion-prediction>). As shown, given the conditioned frames, MAGViT generates plausible actions with greater details.

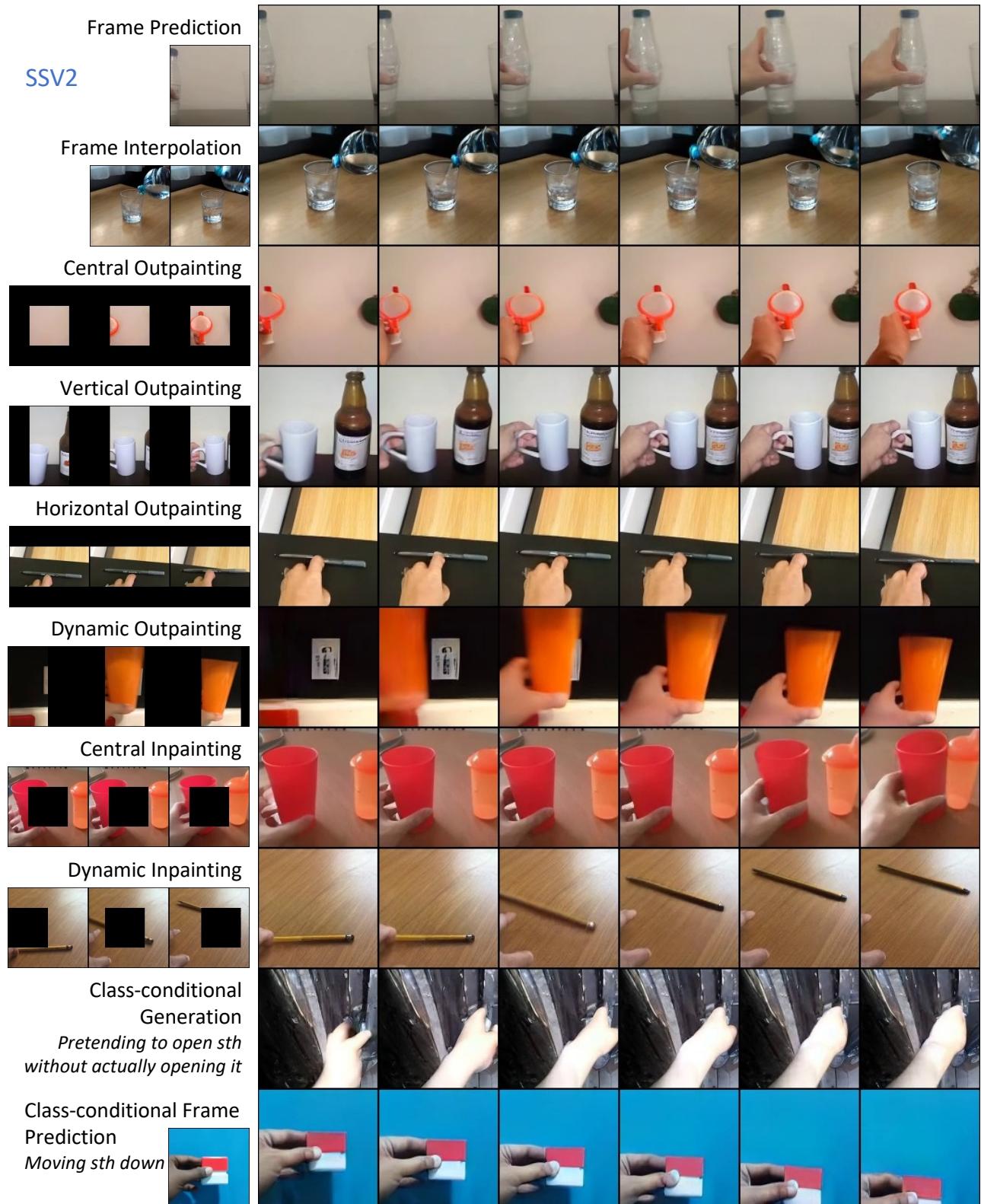


Figure 9. **Multi-task generation results** for the model only trained on the Something-Something-V2 dataset [14]. The condition used to generate the shown videos are taken from the Something-Something-V2 evaluation videos.

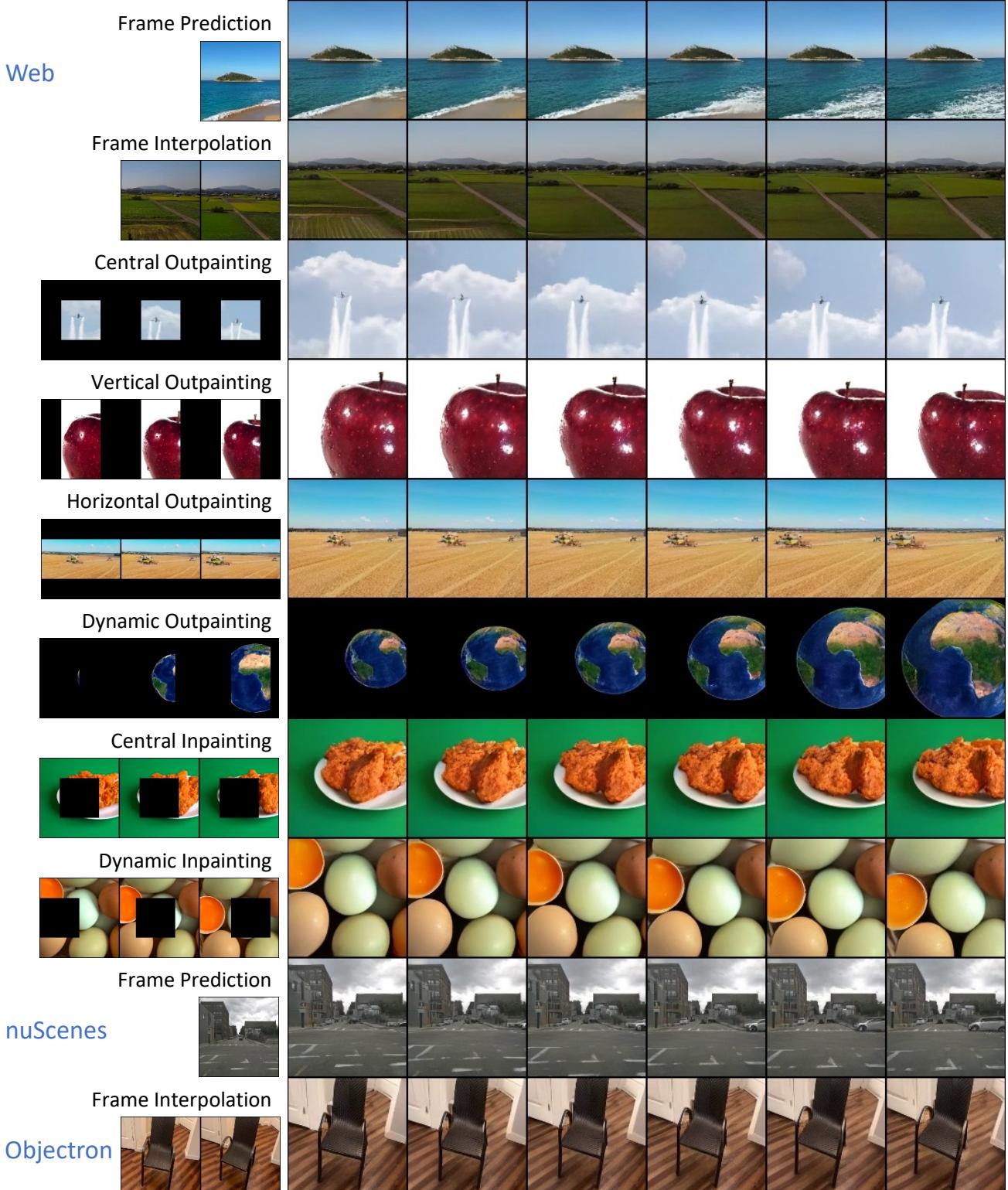


Figure 10. **Multi-task generation results** for three models trained on nuScenes [5], Objectron [3], and 12M Web videos, respectively. The condition used to generate the shown videos are taken from the evaluation set.

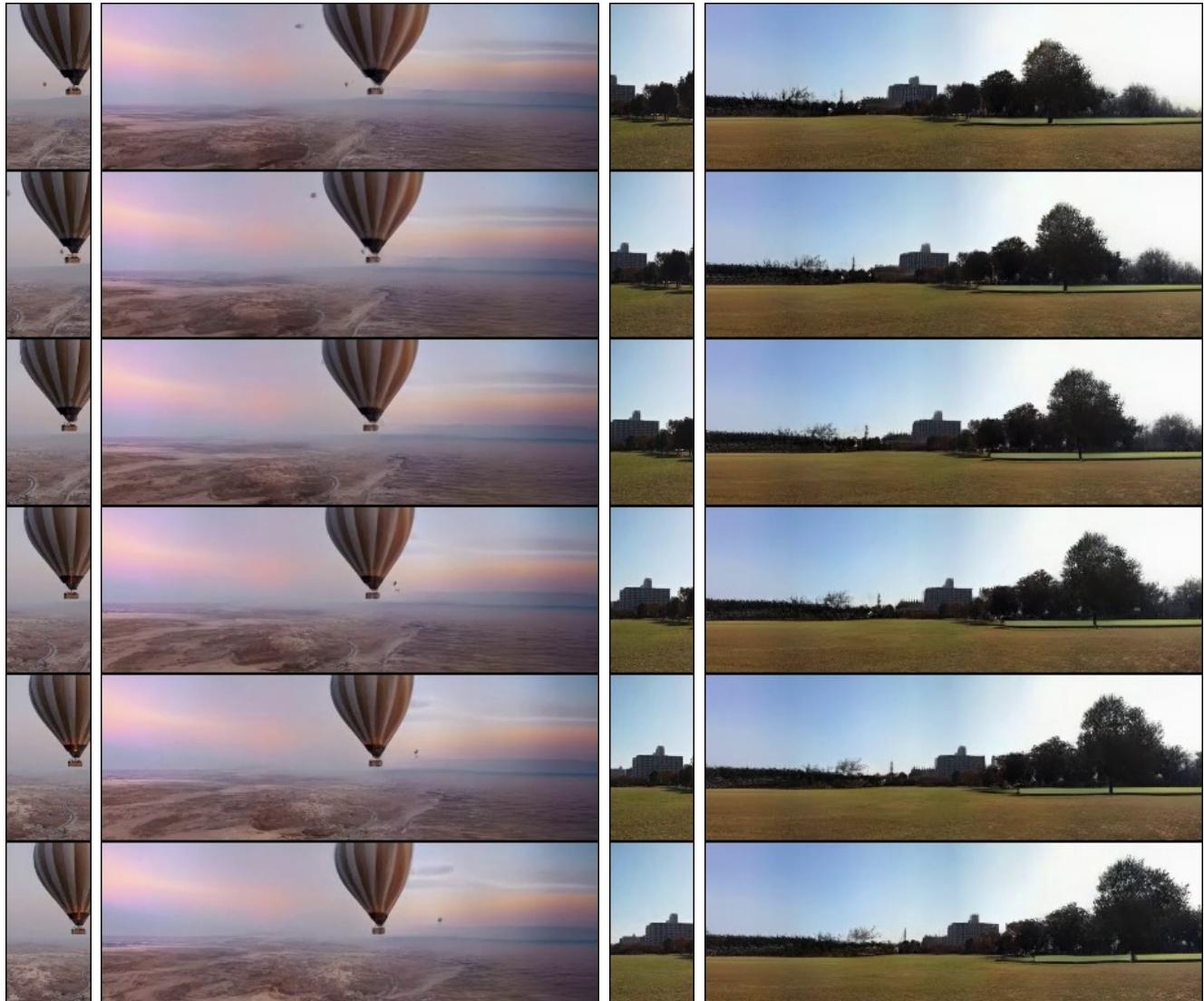


Figure 11. **Ultrawide outpainting results.** Given a vertical slice of 64×128 , MAGVIT expands it into a panorama video of 384×128 by doing vertical outpainting for 5 times on each side.

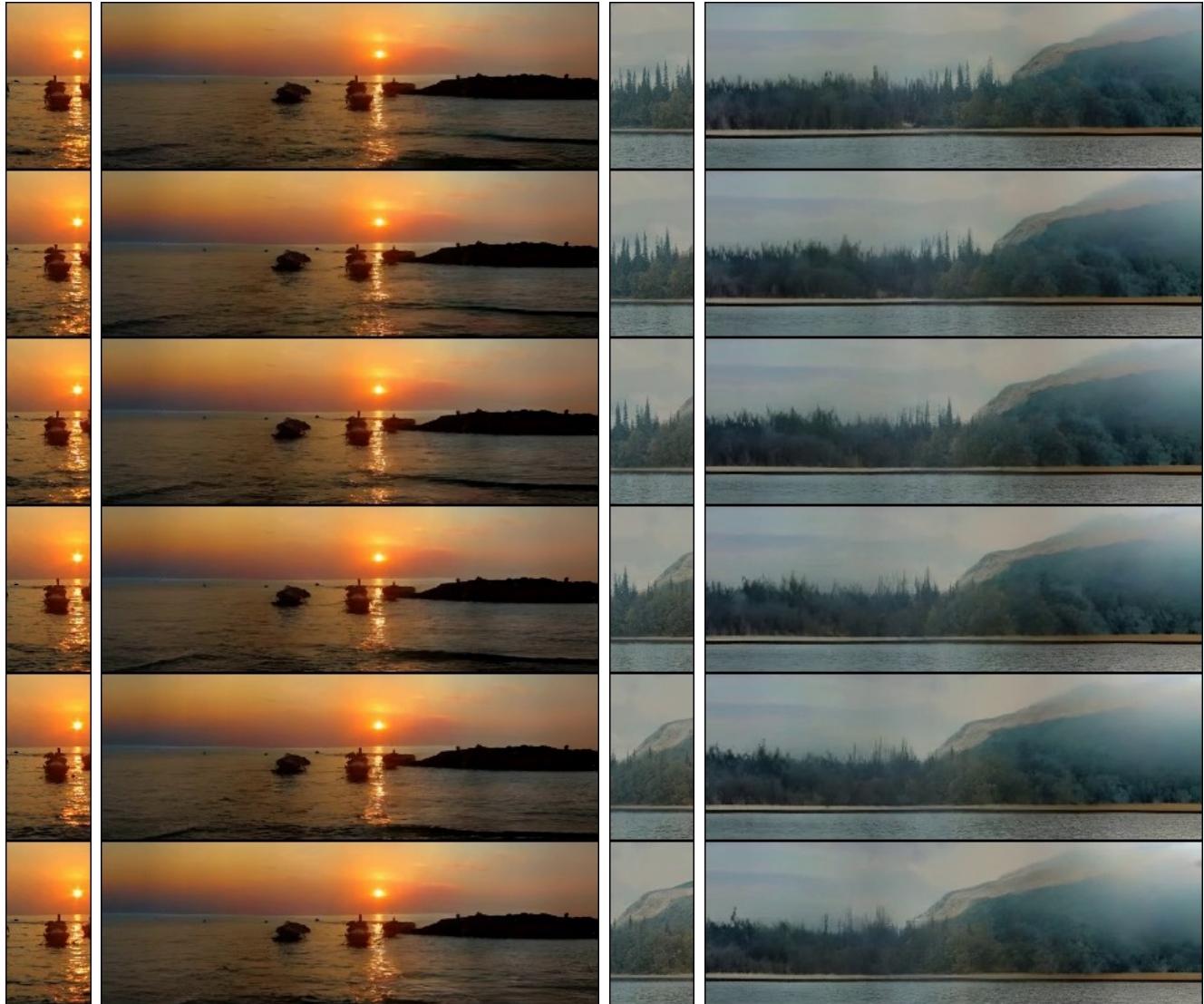


Figure 12. **Ultrawide outpainting results.** Given a vertical slice of 64×128 , MAGVIT expands it into a panorama video of 384×128 by doing vertical outpainting for 5 times on each side.

References

- [1] Sami Abu-El-Haija, Nisarg Kothari, Joonseok Lee, Paul Natsev, George Toderici, Balakrishnan Varadarajan, and Sudheendra Vijayanarasimhan. Youtube-8m: A large-scale video classification benchmark. *arXiv preprint arXiv:1609.08675*, 2016. [17](#)
- [2] Dinesh Acharya, Zhiwu Huang, Danda Pani Paudel, and Luc Van Gool. Towards high resolution video generation with progressive growing of sliced wasserstein gans. *arXiv:1810.02419*, 2018. [16](#)
- [3] Adel Ahmadyan, Liangkai Zhang, Artiom Ablavatski, Jianing Wei, and Matthias Grundmann. Objectron: A large scale dataset of object-centric videos in the wild with pose annotations. In *CVPR*, 2021. [12, 17, 26](#)
- [4] Mohammad Babaeizadeh, Mohammad Taghi Saffar, Suraj Nair, Sergey Levine, Chelsea Finn, and Dumitru Erhan. Fitvid: Overfitting in pixel-level video prediction. *arXiv:2106.13195*, 2021. [15, 16](#)
- [5] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liang, Qiang Xu, Anush Krishnan, Yu Pan, Giacarlo Baldan, and Oscar Beijbom. nuScenes: A multi-modal dataset for autonomous driving. In *CVPR*, 2020. [12, 17, 26](#)
- [6] Joao Carreira, Eric Noland, Andras Banki-Horvath, Chloe Hillier, and Andrew Zisserman. A short note about Kinetics-600. *arXiv:1808.01340*, 2018. [12, 16](#)
- [7] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the Kinetics dataset. In *CVPR*, 2017. [15](#)
- [8] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. MaskGIT: Masked generative image transformer. In *CVPR*, 2022. [17, 19](#)
- [9] Aidan Clark, Jeff Donahue, and Karen Simonyan. Adversarial video generation on complex datasets. *arXiv:1907.06571*, 2019. [15, 16](#)
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2019. [12](#)
- [11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020. [12](#)
- [12] Frederik Ebert, Chelsea Finn, Alex X Lee, and Sergey Levine. Self-supervised visual planning with temporal skip connections. In *CoRL*, 2017. [12, 16](#)
- [13] Songwei Ge, Thomas Hayes, Harry Yang, Xi Yin, Guan Pang, David Jacobs, Jia-Bin Huang, and Devi Parikh. Long video generation with time-agnostic VQGAN and time-sensitive transformer. In *ECCV*, 2022. [12, 13, 15, 16, 17, 18, 19, 22](#)
- [14] Raghad Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The “something something” video database for learning and evaluating visual common sense. In *ICCV*, 2017. [12, 17, 25](#)
- [15] Agrim Gupta, Stephen Tian, Yunzhi Zhang, Jiajun Wu, Roberto Martín-Martín, and Li Fei-Fei. MaskViT: Masked visual pre-training for video prediction. *arXiv:2206.11894*, 2022. [16](#)
- [16] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv:1606.08415*, 2016. [12](#)
- [17] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. In *ICLR Workshops*, 2022. [16, 17](#)
- [18] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. CogVideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv:2205.15868*, 2022. [16](#)
- [19] Tobias Höppe, Arash Mehrjou, Stefan Bauer, Didrik Nielsen, and Andrea Dittadi. Diffusion models for video prediction and infilling. *arXiv:2206.07696*, 2022. [16, 18, 23, 24](#)
- [20] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. [12](#)
- [21] Emmanuel Kahembwe and Subramanian Ramamoorthy. Lower dimensional kernels for video discriminators. *Neural Networks*, 132:506–520, 2020. [16](#)
- [22] Guillaume Le Moing, Jean Ponce, and Cordelia Schmid. CCVS: Context-aware controllable video synthesis. In *NeurIPS*, 2021. [15, 16, 18, 22](#)
- [23] Pauline Luc, Aidan Clark, Sander Dieleman, Diego de Las Casas, Yotam Doron, Albin Cassirer, and Karen Simonyan. Transformation-based adversarial video prediction on large-scale data. *arXiv:2003.04035*, 2020. [16](#)
- [24] Charlie Nash, João Carreira, Jacob Walker, Iain Barr, Andrew Jaegle, Mateusz Malinowski, and Peter Battaglia. Transframer: Arbitrary frame prediction with generative models. *arXiv:2203.09494*, 2022. [16](#)
- [25] Ruslan Rakhimov, Denis Volkhonskiy, Alexey Artemov, Dennis Zorin, and Evgeny Burnaev. Latent video transformer. In *VISIGRAPP (5: VISAPP)*, 2021. [16](#)
- [26] Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. In *ICLR Workshops*, 2018. [12](#)
- [27] Masaki Saito, Eiichi Matsumoto, and Shunta Saito. Temporal generative adversarial nets with singular value clipping. In *ICCV*, 2017. [16](#)
- [28] Masaki Saito, Shunta Saito, Masanori Koyama, and Sosuke Kobayashi. Train sparsely, generate densely: Memory-efficient unsupervised training of high-resolution temporal gan. *IJCV*, 128(10):2586–2606, 2020. [15, 16](#)
- [29] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv:2209.14792*, 2022. [16](#)
- [30] Ivan Skorokhodov, Sergey Tulyakov, and Mohamed Elhosseiny. StyleGAN-V: A continuous video generator with the price, image quality and perks of StyleGAN2. In *CVPR*, 2022. [16](#)
- [31] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. *arXiv:1212.0402*, 2012. [12, 16](#)

- [32] Yu Tian, Jian Ren, Menglei Chai, Kyle Olszewski, Xi Peng, Dimitris N Metaxas, and Sergey Tulyakov. A good image generator is what you need for high-resolution video synthesis. In *ICLR*, 2021. 16
- [33] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3D convolutional networks. In *ICCV*, 2015. 15
- [34] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. MoCoGAN: Decomposing motion and content for video generation. In *CVPR*, 2018. 16
- [35] 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:1812.01717*, 2018. 12, 14, 15, 16
- [36] Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang, Mohammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. Phenaki: Variable length video generation from open domain textual description. *arXiv:2210.02399*, 2022. 16
- [37] Vikram Voleti, Alexia Jolicoeur-Martineau, and Christopher Pal. Masked conditional video diffusion for prediction, generation, and interpolation. In *NeurIPS*, 2022. 15, 16
- [38] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. In *NeurIPS*, 2016. 16
- [39] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE TIP*, 13(4):600–612, 2004. 15
- [40] Dirk Weissenborn, Oscar Täckström, and Jakob Uszkoreit. Scaling autoregressive video models. In *ICLR*, 2019. 16
- [41] Chenfei Wu, Jian Liang, Lei Ji, Fan Yang, Yuejian Fang, Dixin Jiang, and Nan Duan. NÜWA: Visual synthesis pre-training for neural visual world creation. In *ECCV*, 2022. 16
- [42] Yuxin Wu and Kaiming He. Group normalization. In *ECCV*, 2018. 12
- [43] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. VideoGPT: Video generation using vq-vae and transformers. *arXiv:2104.10157*, 2021. 16
- [44] Sihyun Yu, Jihoon Tack, Sangwoo Mo, Hyunsu Kim, Junho Kim, Jung-Woo Ha, and Jinwoo Shin. Generating videos with dynamics-aware implicit generative adversarial networks. In *ICLR*, 2022. 16
- [45] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 15