

SEER: LANGUAGE INSTRUCTED VIDEO PREDICTION WITH LATENT DIFFUSION MODELS

Xianfan Gu³ Chuan Wen^{1,2,3} Weirui Ye^{1,2,3} Jiaming Song⁴ Yang Gao^{1,2,3,*}

¹ IIS, Tsinghua University ² Shanghai Artificial Intelligence Laboratory

³ Shanghai Qi Zhi Institute ⁴ NVIDIA

guxf@sqz.ac.cn {cwen20,ywr20,gaoyangiiis}@mails.tsinghua.edu.cn
jiamings@nvidia.com

ABSTRACT

Imagining the future trajectory is the key for robots to make sound planning and successfully reach their goals. Therefore, text-conditioned video prediction (TVP) is an essential task to facilitate general robot policy learning. To tackle this task and empower robots with the ability to foresee the future, we propose a simple and computation-efficient model, named **Seer**, by inflating the pretrained text-to-image (T2I) stable diffusion models along the temporal axis. We enhance the U-Net and language conditioning model by incorporating computation-efficient spatial-temporal attention. Furthermore, we introduce a novel Frame Sequential Text Decomposer module that dissects a sentence’s global instruction into temporally aligned sub-instructions, ensuring precise integration into each frame of generation. Our framework allows us to effectively leverage the extensive prior knowledge embedded in pretrained T2I models across the frames. With the adaptable-designed architecture, Seer makes it possible to generate high-fidelity, coherent, and instruction-aligned video frames by fine-tuning a few layers on a small amount of data. The experimental results on Something Something V2 (SSv2), Bridgedata and EpicKitchens-100 datasets demonstrate our superior video prediction performance with around 480-GPU hours versus CogVideo with over 12,480-GPU hours: achieving the 31% FVD improvement compared to the current SOTA model on SSv2 and 83.7% average preference in the human evaluation. Our project is available at <https://seervideodiffusion.github.io/>

1 INTRODUCTION

Text-conditioned Video Prediction (TVP), a task that generates future video frames conditioned on a few frames and language instructions, is crucial for diverse downstream tasks requiring instruction alignment and temporal consistency from an initial environment. For instance, in video editing, TVP empowers the generation of diverse temporal movements from an input video clip, guided by a range of language instructions. Importantly, TVP can selectively extend video segments while preserving temporal consistency from the input video. TVP also plays an important role in the scenario of robot learning. TVP samples coherent future frames with aligned motion trajectories based on the initial state of a robot, providing task-level visual guidance for long-horizon planning. This overcomes the challenge for a robot to align abstract language instructions with long-horizon operations. Consequently, learning a TVP model is a fundamental task to achieve both temporal consistency in the transition distribution of the initial state and alignment between task-level language instructions and video motion, facilitating the development of video foundation models.

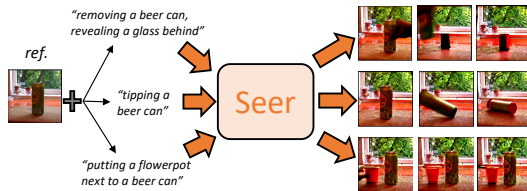


Figure 1: Seer is an efficient video diffusion model that uses natural language instructions and reference frames (*ref.*) to predict multiple variations of future frames.

*Corresponding Author.

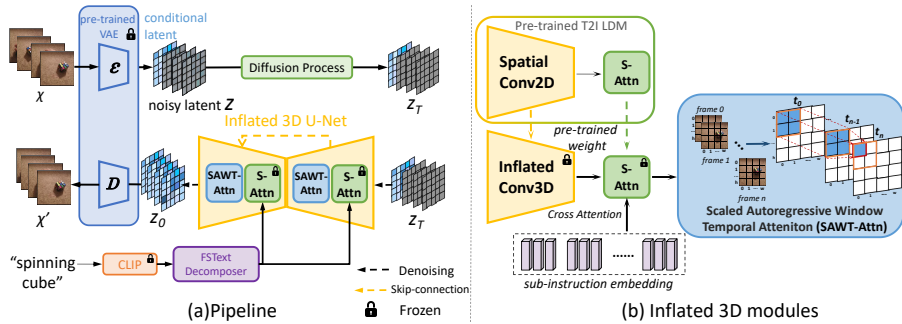


Figure 2: (a) Seer’s pipeline includes an Inflated 3D U-Net for diffusion and a Frame Sequential Text Transformer for text conditioning. (b) Our Inflated 3D U-Net expands the pre-trained 2D Conv kernel to 3D kernels and connects with the Scaled Autoregressive Window Temporal Attention layer.

Despite its potential benefits, text-to-video prediction (TVP) is a challenging task because it requires a deep understanding of the initial frames, the natural language instruction, and the grounding between language and images, while predicting based upon all the information above. In contrast to traditional text-to-video generation tasks (Wu et al., 2021; 2022a; Hong et al., 2023; Ho et al., 2022b; Singer et al., 2022), which do not explicitly condition on initial frames, TVP requires a model to synthesize predictions based on the given initial frames and textual instructions. Merely generating a few prototypical videos corresponding to the input text is no longer a viable solution in TVP; the task necessitates a more detailed comprehension of temporal movement. Besides, the existing text-to-video generation task usually aims to generate short horizon video clips with text specifying the general content, such as “a person is skiing”, while our aim in the TVP task is to use the text as a task descriptor, such as “tipping a beer can”, as shown in Figure 1.

Specifically, there are mainly three problems limiting the performance of the TVP task: **1) Requirement for large-scale labeled text-video datasets and expensive computational cost:** learning to capture the correspondence between two different modalities is non-trivial and needs large amounts of supervised text-video pairs and excessive computation overhead for training. **2) Low fidelity of generated frames:** the frames generated by the models are usually blurry and cannot clearly display the background and objects specified in the reference frames. **3) Lack of fine-grained instruction for each frame in the task-level videos:** the goals specified by text instructions are usually in the task level, making it difficult to understand the progress and generate the corresponding frame in each timestep only conditioned on a global text embedding. To address these issues, we propose **Seer**: a TVP method capable of generating task-level videos according to the text guidance with high data and training efficiency.

Motivated by the recent progress on generative models (Rombach et al., 2022; Ramesh et al., 2022), we propose to leverage text-to-image (T2I) latent diffusion models (Rombach et al., 2022) for the TVP tasks. T2I models are pretrained with billions of text-image pairs crawled from the internet (Schuhmann et al., 2021). They have acquired rich prior knowledge and thus are able to generate high-quality images corresponding to the text descriptions. Therefore, inheriting such prior knowledge by inflating a T2I model along the temporal axis and fine-tuning it with a small text-video dataset is an appealing solution for TVP tasks, which relieves the requirement for extensive labeled data and computational overhead, i.e., Problem 1.

Since the T2I models contain two modalities: image and language, we propose to inflate these two parts to generate high-quality video frames and fine-grained text instruction embeddings for each timestep respectively. For the visual model, we extend the 2D latent diffusion model (Rombach et al., 2022) to data and computation-efficient 3D network to model spatial dependencies and the temporal dynamics simultaneously, which is called Inflated 3D U-Net. By taking advantage of joint modeling of spatial and temporal dimensions, as well as autoregressive generation, we successfully synthesize coherent and high-fidelity frames, which alleviates Problem 2. As for the language module, in contrast to existing approaches (Zhou et al., 2022; Wu et al., 2022b) that encode one text embedding for the whole video with a text encoder, we propose to decompose the single text instruction into fine-grained guidance embeddings for each time step. We achieve the automatic decomposition by a **Frame Sequential Text (FSText) Decomposer** based on the causal attention mechanism. By temporally splitting the instruction into different phases, Seer improves the guidance embeddings for each frame and thus enables task-level video generation (Problem 3).

We conduct extensive experiments on Something-Something V2 (Goyal et al., 2017), Bridge Data (Ebert et al., 2021) and Epic-Kitchens-100 (Damen et al., 2021) datasets. We outperform all the

baselines, such as MCVD (Voleti et al., 2022), TATS (Ge et al., 2022) and *Tune-A-Video* (TAV) (Wu et al., 2022b), and achieve state-of-the-art performance in terms of FVD and KVD. Especially we improve FVD from 163 to 113 on Something-Something V2, compared to the SOTA TAV. Compare to over 480 hours with 13×8 A100 GPUs in CogVideo (Hong et al., 2023), the experiments show the high efficiency of our method: 120 hours with 4 RTX 3090 GPUs. The ablation studies illustrate the effectiveness of our computation-efficient video-inflated model with the newly proposed FSText Decomposer. Furthermore, our method supports video manipulation by modifying the text instructions, and we demonstrate our superior generation quality through a human evaluation study, showing more than 70% preference over TAV and around 90% preference over TATS and MCVD.

2 RELATED WORK

2.1 TEXT-TO-IMAGE GENERATION

Since Scott Reed et al. (Reed et al., 2016) firstly set up the T2I generation task and proposed a GAN-based method, this multi-modal generation task has attracted the attention of the computer vision community. DALL-E (Ramesh et al., 2021) makes a breakthrough by modeling the T2I generation task as a sequence-to-sequence translation task with a VQ-VAE (Van Den Oord et al., 2017) and Transformer (Vaswani et al., 2017). Since then, many variants have been proposed with an improved image tokenizer (Yu et al., 2022a), hierarchical Transformers (Ding et al., 2022) or domain-specific knowledge (Gafni et al., 2022). With the recent progress of Denoising Diffusion Probabilistic Models (DDPM) (Ho et al., 2020), the diffusion models have been widely used for T2I generation tasks (Nichol et al., 2022; Saharia et al., 2022; Ramesh et al., 2022). Specifically, GLIDE (Nichol et al., 2022) proposes classifier-free guidance for T2I diffusion models to improve image quality. For a better alignment between text and image, DALL-E 2 (Ramesh et al., 2022) proposed to denoise CLIP (Radford et al., 2021) image embedding conditioned on CLIP text embedding, which integrated high-level semantic information. To reduce the computation cost of the denoising process in pixel space, Latent Diffusion Model (LDM) employs VAE (Kingma & Welling, 2014) to operate in the latent space. Seer takes advantage of the prior language-vision knowledge of pretrained LDM and inflates it along the time axis.

2.2 TEXT-TO-VIDEO GENERATION

In contrast to the huge success of Text-to-Image (T2I) generation, Text-to-Video (T2V) generation is still underexplored due to the limitation of the large text-video data annotation and computing resources. Inspired by the various variants of T2I generation, recent T2V studies have attempted to explore compatible variants for video generation modeling. GODIVA (Wu et al., 2021) first proposes a VQ-VAE based auto-regressive model with three-dimensional sparse attention for T2V generation. NÚWA (Wu et al., 2022a) further improves it by designing a 3D encoder with 3D nearby attention and achieves competitive performance on multi-task generation. Unlike the single frame-rate T2V approaches trained from scratch on large-scale text-video datasets, CogVideo (Hong et al., 2023) proposes a multi-frame-rate hierarchical model for T2V generation. This approach leverages the pre-trained module of T2I CogView-2 (Ding et al., 2022).

Motivated by the remarkable progress of T2I diffusion models (Ramesh et al., 2022; Saharia et al., 2022; Rombach et al., 2022), Make-A-Video (Singer et al., 2022), MagicVideo (Zhou et al., 2022), Tune-A-Video (Wu et al., 2022b) and Imagen Video (Ho et al., 2022a) transfer the 2D diffusion models to 3D models by incorporating temporal modules in T2V generation. In contrast to Imagen Video, all other three methods utilize the prior knowledge of T2I pre-trained model. Similarly, we use the pre-trained weight of the 2D T2I diffusion model in our 3D T2V model. Varying from the aforementioned methods, our method Seer utilizes autoregressive attention on both spatial and temporal spaces to generate high-fidelity and coherent video frames. And Seer is able to handle the task-level video prediction by decomposing the language condition into fine-grained sub-instruction.

3 PRELIMINARIES

Denoising Diffusion Probabilistic Models with classifier-free guidance: Diffusion models are probabilistic models that approximate the data distribution by iteratively adding noise and denoising

through a forward/reverse Gaussian Diffusion Process (Ho et al., 2020; Song et al., 2021). The forward process applies noise at each time step $t \in 0, \dots, T$ to the data distribution \mathbf{x}_0 , creating a noisy sample \mathbf{x}_t where $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}$ ($\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$), and $\bar{\alpha}_t$ is the accumulation of the noise schedule $\alpha_{0:T}$ defined by $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$. To denoise images, the diffusion process uses a reparameterized variant of Gaussian noise prediction $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$ targeting Gaussian noise $\boldsymbol{\epsilon}$. The reverse process $p(\mathbf{x}_{t-1}|\mathbf{x}_t)$ of the Markov Chain generates new samples from Gaussian noise, which is approximated by Bayes’ theorem as $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$, where \mathbf{x}_0 is derived from the forward process as $\mathbf{x}_0 = \frac{1}{\sqrt{\bar{\alpha}_t}}(\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t))$.

Classifier-free guidance (Ho & Salimans, 2022) is introduced for conditional diffusion models to generate images without requiring an extra image classifier. A conditional model with a parameterized reverse process $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c})$ uses a conditional identifier \mathbf{c} through $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, \mathbf{c})$. To predict an unconditional score, the conditional identifier is replaced with a null token \emptyset and denoted as $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, \mathbf{c} = \emptyset)$. Classifier-free guidance can then be approximated as a linear combination of conditional and unconditional predictions:

$$\bar{\boldsymbol{\epsilon}}_\theta(\mathbf{x}_t, t, \mathbf{c}) = (1 + w)\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, \mathbf{c}) - w\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, \mathbf{c} = \emptyset), \quad (1)$$

where w is the guidance scale. Text-video and text-image-based diffusion models (Rombach et al., 2022; Saharia et al., 2022; Nichol et al., 2022; Ho et al., 2022b; Singer et al., 2022) use DDPM with classifier-free guidance. This diffusion method can be adapted to various tasks with flexibility.

Latent Diffusion Models: Compared with image diffusion, video diffusion has significantly higher computation costs because it needs to process multiple frames. Recent works have explored the computation-efficient version of diffusion modeling, such as latent diffusion model (LDM) (Rombach et al., 2022). LDM proposes the VAE-based latent diffusion, including a KL-regularized autoencoder for encoding/decoding latent representation $\boldsymbol{\varepsilon}(\mathbf{x})$, and a diffusion model to operate on the latent space \mathbf{z}_t . For the conditional generation, LDM introduces a domain-specific encoder $\boldsymbol{\tau}_\theta$ to the projection of condition \mathbf{y} for various modality generations. Thus, the objective of LDM is:

$$L_{\text{LDM}} = \mathbb{E}_{t, \boldsymbol{\varepsilon}(\mathbf{x}), \mathbf{y}, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{z}_t, t, \boldsymbol{\tau}_\theta(\mathbf{y}))\|^2 \right] \quad (2)$$

4 METHODOLOGY

In this paper, we aim to explore an efficient diffusion method to predict coherent video frames guided by language instructions, which requires learning to parse natural language, understand the scene, and ground the language and scene together. However, it is challenging to directly apply conventional video diffusion models for TVP due to the following problems: (1) The limited labeled text-video data and computational resources. (2) Low fidelity of frame generation. (3) Lack of fine-grained instruction for each frame in the task-level videos.

4.1 OVERVIEW OF SEER

Motivated by the robust generative capabilities of text-to-image (T2I) diffusion models, we leverage the prior knowledge implied in pretrained T2I models by inflating the 2D U-Net (Rombach et al., 2022) and incorporating temporally consistent layers. However, the inflated video diffusion model guided solely by coarse global language instruction tends to generate irrelevant T2I outcomes and fails to maintain temporal coherency between video frames. To address this limitation and provide precise and controllable guidance for our inflated model, we introduce a novel temporal decomposition component for language instruction, this component decomposes global instruction as temporally aligned sub-instruction for delicate task-level guidance, which significantly enhances the fidelity and coherency of predicted video.

Our Seer method comprises two main components: the video diffusion and the language conditioning modules. We propose to enhance these two components to facilitate high-fidelity frame synthesis and the temporal alignment of text instructions, respectively. Specifically, as shown in Figure 2 (a), we utilize two pathways to implement the conditional diffusion process guided by reference frames and language: **1)** We incorporate the spatial-temporal module discussed in Section 4.2 into the Inflated 3D U-Net. This integration enables the propagation of contextual information from reference frames to future frames within the spatial-temporal space, allowing for coherent motion prediction based on the reference frames. **2)** To plan continuous motion with fine-grained language guidance, we introduce a Frame Sequential Text (FSText) Decomposer in Section 4.3. This module transforms

global language instructions into multi-timestep sub-instructions that are synchronized with video. Subsequently, we inject these frame-wise subinstruction tokens into the intermediate latent space of the video frames at each time step. With this design, we merely train the spatial-temporal layers and FSText module from scratch while freezing the remaining pretrained modules within our 3D inflated U-Net. These two modules are jointly trained by the diffusion objective, where f_θ is our FSText decomposer, τ is the frozen CLIP text encoder, and \mathbf{y} is the input text:

$$L_{\text{diffusion}} = \mathbb{E}_{t, \epsilon(\mathbf{x}), \mathbf{y}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\|\epsilon - \epsilon_\theta(\mathbf{z}_t, t, f_\theta(\tau(\mathbf{y})))\|^2 \right], \quad (3)$$

4.2 DATA & COMPUTATION-EFFICIENT 3D NETWORK

To design a computation-efficient visual backbone for our video diffusion model, we refer to some relevant works on lifting 2D to 3D video modeling (Carreira & Zisserman, 2017) and efficient attention computation (Liu et al., 2021; 2022). In general, we leverage the latent diffusion model (LDMs) pretrained on T2I tasks to build a text-video model. Our inflated 3D U-Net consists of two principal components as illustrated in Figure 2 (b): **1)** The 3D spatial layers, where we draw inspiration from I3D (Carreira & Zisserman, 2017) and enhance the 2D convolution kernel from (3×3) to a 3D counterpart $(1 \times 3 \times 3)$ with an added video frames axis from the pre-trained 2D modules, consisting of a series of 2D ResNet blocks and Spatial Attention Blocks.

2) The temporal layers, play a crucial role in our visual backbone for propagating contextual information from the reference frame’s image prior across the temporal sequence. We investigated various temporal attention and incorporated them into our 3D U-Net architecture. Our empirical observations indicate that bi-directional temporal attention tends to disregard guidance from reference frames, and both bi-directional and directed temporal attention struggle to capture dependencies among spatial regions, as discussed in Section 5.5. To address these limitations while reducing complexity, we employ an efficient approach that builds upon the concept of window attention (Liu et al., 2021) in 3D space: the implementation of local window attention in an autoregressive manner across spatial-temporal dimensions. As illustrated in Figure 3, we establish fixed local windows for each spatial region with a window size of $m \times m$ relative to the global frame sequence n . Within this framework, we compute self-attention using a causal mask, considering both local spatial and global temporal dimensions within the 3D space. This effectively constrains pixel propagation from the future temporal-spatial sequence.

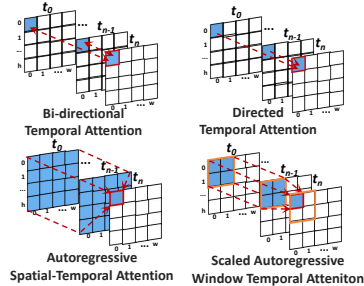


Figure 3: Variants of temporal attention, only the blue tokens attend to the current token in the red box. Red dashed arrows indicate the direction of attention. And the orange boxes indicate the local window region (2×2 window in this case).

Finally, We maintain the acquired knowledge from the 2D modules by freezing all pretrained weights and exclusively training the spatiotemporal attention layers during fine-tuning. Overall, through a combination of frozen pre-trained spatial layers and lightweight spatiotemporal layers, our inflated 3D U-Net not only retains crucial knowledge but also enhances fine-tuning efficiency.

4.3 FRAME SEQUENTIAL TEXT DECOMPOSER

For the language conditioning module, since our 3D inflated U-Net is built upon a pretrained text-to-image model, we noticed that using a text-to-image prior alongside a global instruction tends to provide strong semantic guidance, which can override the scene in reference frames, deviating from the intended guidance for prediction based on the existing scenes. To address the above limitation and better capture long-term dependencies from both text and reference frames, we introduce the Frame Sequential Text (FSText) Decomposer. This novel approach decomposes the global instruction into fine-grained sub-instructions, aligning with each frame. We further explore the interpretability of sub-instruction embeddings in Section 5.6. To derive a sequence of temporally aligned sub-instruction embeddings from the global instruction generated by the CLIP text encoder Radford et al. (2021), we employ a transformer-based temporal network designed to fulfill three essential properties for meaningful sub-instructions: **1)** Contextual aggregation, which ensures that the inner tokens of each sub-instruction aggregate contextual information within the

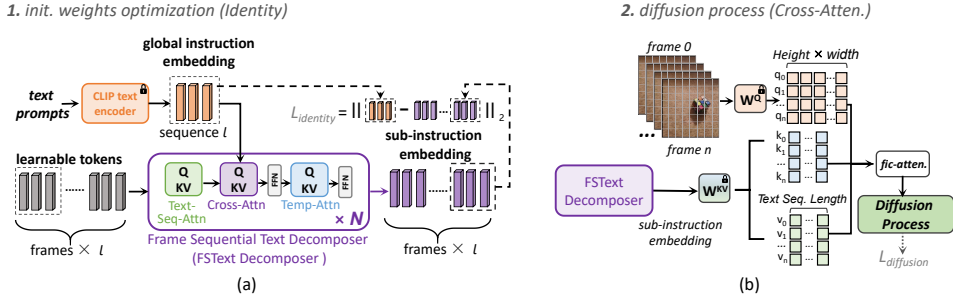


Figure 4: Frame Sequential Text Decomposer is shown in (a). We start by initializing the weight of the network to project identity vectors from CLIP text tokens. We then optimize the generated text tokens via the diffusion process (b), where frame-individual cross-attention is denoted by “fic-attn.”

sentence. **2)** Semantic inheritance, the semantic information of these sub-instructions is inherited directly from the global instruction **3)** Temporal consistency ensures alignment between the sub-instructions and the time sequence, thereby facilitating the generation of temporally consistent video. Based on these properties, our network consists of three key components: **a)** To achieve the property of contextual aggregation, we employ Text-Sequential Attention, akin to BERT, a bidirectional self-attention layer (Devlin et al., 2019) to capture global dependencies among different positions within text sentences. **b)** To ensure semantic inheritance, we use Cross-Attention, responsible for projecting the global instruction’s textual sequence onto the inner tokens of each sub-instruction, this component ensures that all sub-instructions contain essential global instruction signals for guiding video frame generation. **c)** To maintain temporal consistency, we adopt temporal Attention, a directed attention layer to capture temporal dependencies along the frame axis, which enhances temporal consistency among the generated sub-instructions throughout the video.

Specifically, as shown in Figure 5, we start with a global CLIP text embedding, denoted as (l, C) , where l signifies the text sentence length and C is the channel size, we initialize learnable tokens with shape (n, l, C) where n denotes the number of frames. The tokens are fed into the text sequential attention layer to perform self-attention along the l axis. Subsequently, the cross-attention layer employs these learnable tokens as queries and the global text embedding as keys and values, resulting in a one-to-multiple projection from the global text into n time steps. This yields (n, l, C) tokens for n frame containing task instruction information. Finally, the temporal attention layer conducts directed attention along the n axis for each token in the textual sequence, transforming the macro-instruction progress into frame-specific guidance.

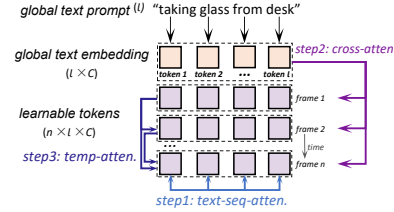


Figure 5: The FSText attention of sub-instruction tokens.

After getting n sub-instruction embeddings corresponding to each frame, the next step is to inject this guidance into the diffusion process, which is commonly completed by a cross-attention layer. As shown in Figure 4 (b), different from the existing works (Zhou et al., 2022; Wu et al., 2022b) that calculate the cross-attention between the global instruction embedding and n frames. In our cross-attention layer, where cross-attention is calculated separately between visual latent vectors and sub-instruction embeddings for each frame, and the results from all frames are then concatenated, an attention mechanism we refer to as frame-individual cross-attention (*fic-attn*).

Initialization We find initialization is critical to FSText decomposer. Especially, the random initialization fails to approximate the distribution of text embeddings in the pretrained T2I model and results in poor performance. To guarantee the sub-instruction embeddings become a close approximation of the CLIP text embedding, we employ an initialization strategy by enforcing the FSText decomposer to be an identity function (Note that this initialization step is completed before the diffusion process. We ablate this design in Section 5.5). It can be achieved by this objective:

$$L_{\text{identity}} = \|f_{\theta}(\tau(\mathbf{y})) - \tau(\mathbf{y})\|^2 \quad (4)$$

5 EXPERIMENTS

In this section, we evaluate Seer on the text-conditioned video prediction task. We compare against various recent methods and conduct ablation studies on the techniques presented in Section 4.

5.1 DATASETS

We conduct experiments on three text-video datasets: Something Something-V2 (SSv2) (Goyal et al., 2017), which contains videos of human daily behaviors with language instructions, BridgeData (Ebert et al., 2021) that is rendered by a photo-realistic kitchen simulator with text prompts, and EpicKitchens-100 (Damen et al., 2021) (Epic100), which collects human daily activities in the kitchen in egocentric vision with multi-language narrations. For SSv2, we follow (Soomro et al., 2012) to evaluate the first 2048 samples during evaluation to save testing time. For BridgeData, we split the dataset into an 80% training set and 20% validation set for evaluation. To reduce complexity, we downsample each video clip to 12 frames for SSv2 and Epic100, and 16 frames for BridgeData during training/evaluation. Besides, We provide the zero-shot evaluation on EGO4D (Grauman et al., 2022) dataset in section D.2. Moreover, we also included an additional evaluation on the UCF-101 dataset (Soomro et al., 2012) in Appendix B.2.

5.2 IMPLEMENTATION DETAILS

We use the pre-trained weights of Stable Diffusion-v1.5 (Rombach et al., 2022) to initialize the VAE, ResNet Blocks and Spatial-Cross Attention layers of the 3D U-Net. We freeze both the pre-trained VAE, and only fine-tune the SAWT-Atten layers in our 3D U-Net. To fine-tune the FSText Decomposer, we initialized it as the identity function of the CLIP text embedding, as described in Section 4.3. We train the models with an image resolution of 256×256 on Something Something-V2 for 200k training steps, EpicKitchens-100 and BridgeData for 80k training steps. In the evaluation stage, we speed up the sampling process with the fast sampler DDIM (Song et al., 2020) and conditional guidance of 7.5 for 30 timesteps. See more details in Appendix C.

5.3 EVALUATION SETTINGS

Baselines. We compare Seer with seven publicly released baselines for video generation (1) conditional video diffusion methods: *Tune-A-Video* (Wu et al., 2022b), *Masked Conditional Video Diffusion* (MCVD)(Voleti et al., 2022), *Video Probabilistic Diffusion Models* (PVDm)(Yu et al., 2023) and *VideoFusion*(Luo et al., 2023); (2) autoregressive-based transformer method: *Time-Agnostic VQGAN and Time-Sensitive Transformer* (TATS)(Ge et al., 2022) and *Make It Move* (MAGE)(Hu et al., 2022b); and (3) CNN-based encoder-decoder: *SimVP*(Gao et al., 2022).

Machine Evaluation. We evaluate the text-conditioned video prediction of several baseline methods on Something Something-V2 (SSv2) (with 2 reference frames), Bridgedata (with 1 reference frame) and Epic-Kitchens-100 (Epic100) (with 1 reference frame). Additionally, we conduct several ablation studies of our proposed modules on SSv2. We report the Fréchet Video Distance (FVD) and Kernel Video Distance (KVD) metrics in our evaluation. FVD and KVD are calculated with the Kinetics-400 pre-trained I3D model (Carreira & Zisserman, 2017). We evaluate FVD and KVD on 2,048 SSv2 samples, 5,558 Bridgedata samples and 9,342 Epic100 samples in the validation sets. For FVD metrics, we follow the evaluation code of VideoGPT (Yan et al., 2021). We further evaluate the class-conditioned video prediction of our method on the UCF-101 dataset and present the comparison results in Appendix B.2.

Human Evaluation. Besides evaluating the models on the standard validation sets, we also manually modify the text prompts to provide richer testing results, called text-conditioned video manipulation. Because of the absence of ground-truth frames, we conducted a human evaluation of text-conditioned video manipulation (TVM) using 99 video clips from the validation set of SSv2. We manually modified partial

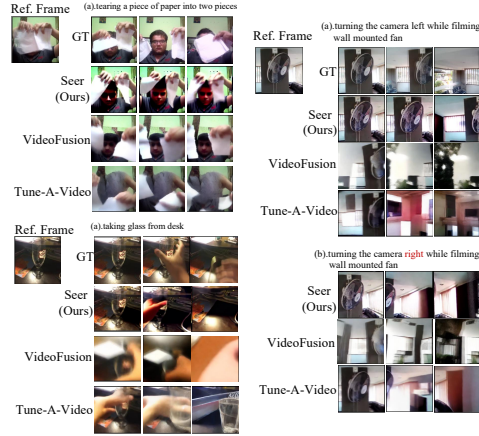


Figure 6: Visualization of Text-conditioned Video Prediction with the original (a) and manually modified (b) text prompts on SSv2.

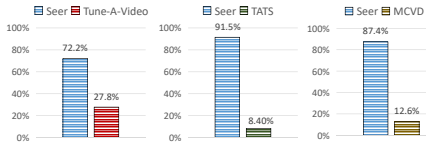


Figure 7: Human evaluation results. Preference percentage for TVM task on SSv2.

Table 1: **Text-conditioned video prediction (TVP) results on Something-Something V2 (SSv2), Bridgedata (Bridge), and Epic-Kitchens-100 (Epic100).** We report both FVD and KVD metrics.

Method	Pre.-weight	Text	Resolution	SSv2		Bridge		Epic100	
				FVD↓	KVD↓	FVD↓	KVD↓	FVD↓	KVD↓
TATS (Ge et al., 2022)	video	No	128 × 128	428.1	2177	1253	6213	920.0	5065
MCVD (Voleti et al., 2022)	No	No	256 × 256	1407	3.80	1427	2.50	4804	5.17
SimVP (Gao et al., 2022)	No	No	64 × 64	537.2	0.61	681.6	0.73	1991	1.34
MAGE (Hu et al., 2022b)	video	Yes	128 × 128	1201.8	1.64	2605	3.19	1358	1.61
PVDM (Yu et al., 2023)	No	No	256 × 256	502.4	61.08	490.4	122.4	482.3	104.8
VideoFusion (Luo et al., 2023)	txt-video	Yes	256 × 256	163.2	0.20	501.2	1.45	349.9	1.79
Tune-A-Video (Wu et al., 2022b)	txt-img	Yes	256 × 256	291.4	0.91	515.7	2.01	365.0	1.98
Seer (Ours)	txt-img	Yes	256 × 256	112.9	0.12	246.3	0.55	271.4	1.40

text prompts and generated 99 predicted videos for each method. Then, we invited 54 anonymous evaluators to rate the quality of the prediction, with a higher priority placed on the semantic contents in the videos and an intermediate priority placed on the fidelity of the video frames. We report the overall preference choices among the 99 video clips. More details are introduced in Appendix E

5.4 MAIN RESULTS

Quantitative Results. Table 1 presents the text-conditioned video prediction results on Something-Something-V2 (SSv2), BridgeData and Epic-kitchens-100 (Epic100). Seer achieves the best performance among all baselines, with the lowest Fréchet Video Distance (FVD) of 112.9 and Kinematic Distance (KVD) of 0.12 in SSv2, the lowest FVD of 246.3 and KVD of 0.55 in BridgeData, and the lowest FVD of 271.4 in Epic100. Notably, Seer, MAGE, VideoFusion and Tune-A-Video all incorporate text conditioning, and the results highlight Seer’s superior text-video alignment performance.

The results of the human evaluation in the text-conditioned video manipulation experiment are shown in Figure 7. Our proposed Seer outperforms the other baselines in terms of both semantic content and fidelity of video, with a preference rate of at least 72.2% in comparison. This indicates that Seer is effective in generating high-quality video clips that are faithful to the input text prompts.

Qualitative Results. Figure 6 compares the text-conditioned video prediction and manipulation performance of Seer, VideoFusion and Tune-A-Video on Something-Something-V2 (SSv2). Seer performs better in handling the temporal dynamics of the video and achieving more precise text-video alignment in video manipulation. For instance, consider the task of “turning the camera left” Seer seamlessly generates coherent movement while preserving the background during the camera view adjustment. In contrast, both VideoFusion and Tune-A-Video exhibit semantic movement but fail to maintain temporal consistency in the video background. Additionally, Seer can imagine hidden objects by utilizing its text-to-image diffusion prior. This flexibility allows Seer to effectively address occlusion in video prediction. In the “tearing paper” sample, Seer accurately predicts that a man is hidden behind the paper and generates coherent frames including the man’s face. Figure 8 compares Seer, VideoFusion and Tune-A-Video’s TVP performance on Bridgedata, illustrating that Seer achieves better text-video alignment of instructed behavior and target objects in future frames, and predicts a more coherent video with higher fidelity.

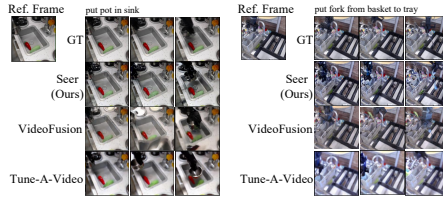


Figure 8: Visualization of Text-conditioned Video Prediction on Bridgedata.

5.5 ABLATION STUDY

In this section, we evaluate the effect of different components of our method in the TVP task on the SSv2 dataset. We also evaluate the zero-shot TVP task of different models on the EGO4D dataset.

Temporal Attention. As shown in Table 2 studies the effectiveness of different types of temporal attention. Our scaled autoregressive window temporal attention (win-auto.) outperforms autoregressive spatial-temporal attention (autoreg.), bi-directional temporal attention (bi-direct.) and directed temporal attention (directed.), resulting in the lowest FVD and KVD scores. We also find that “directed.” further improves video prediction performance compared to “bi-direct.” because it utilizes the inductive bias of sequential generation.

Table 2: **Ablation study of temporal attention**

temp. attn.	FVD↓	KVD↓
<i>bi-direct.</i>	258.2	0.56
<i>directed.</i>	222.3	0.40
<i>autoreg.</i>	200.1	0.30
<i>win-auto.</i> (Ours)	112.9	0.12

FSText Decomposer. Table 3 compares different weight initialization strategies of FSText decomposer. The results show that using identity initialization described in Section 4.3 yields higher prediction quality compared with random initialization. This finding demonstrates that identity initialization is necessary for the temporal-text projection of FSText decomposer. See additional ablation results in Appendix B.5.

Fine-tune Setting. We compare various fine-tuning settings of 3D Inflated U-Net in Table 4. Our default setting involves fine-tuning both FSText decomposer (FSText.) and scaled autoregressive window temporal attention (SAWT-Attn.) layers (*temp-attn.*), while freezing the remaining modules in 3D U-Net. For the “*temp-attn.*” setting, we only finetune the SAWT-Attn. layers and freeze all other components. In the “*cross+temp-attn.*” setting, we jointly update the parameters of spatial-cross attn. and SAWT-Attn. layers. We observe that our default setting achieves the highest quality of video prediction among all these settings. Based on our default setting, further fine-tuning “*cross+temp-attn.*” causes the performance of Seer to drop a lot. These results suggest that the optimization of the FSText decomposer is strongly guided by the frozen conditional diffusion prior.

Table 3: **Init. weight ablation results of FSText**

init. weight	FVD↓	KVD↓
<i>random</i>	367.9	0.75
<i>identity(Ours)</i>	112.9	0.12

Table 4: **Ablation study of Fine-tune settings**

fine-tune	FSText.	FVD↓	KVD↓
<i>temp-attn.</i>		328.2	1.26
<i>cross+temp-attn.</i>		249.9	0.73
<i>temp-attn.(Ours)</i>	✓	112.9	0.12
<i>cross+temp-attn.</i>	✓	1807	5.12

5.6 VISUAL ANALYSIS OF INSTRUCTION EMBEDDING

To assess the impact of sub-instruction guidance on frame generation and unveil the implicit semantic information contained within a single sub-instruction at specific time steps, we compare text-conditioned video prediction (TVP) results conditioned on default frame-wise sub-instructions with those conditioned on the constant clone of a sub-instruction from the third or twelfth frame along the video axis, as depicted in Figure 9. Unlike the default frame-wise sub-instruction guidance, the sub-instruction clone from the twelfth frame tends to produce a transition from the reference frame to the video’s termination state without temporal coherence. In contrast, the sub-instruction clone from the third frame tends to guide the motion until an intermediate state of the video, indicating that temporal sub-instructions provide proximate semantic guidance for motion at each time step. These findings underscore that sub-instructions align with the temporal sequence of the global instruction, offering fine-grained guidance across multiple steps. See more results in the Appendix D.3.

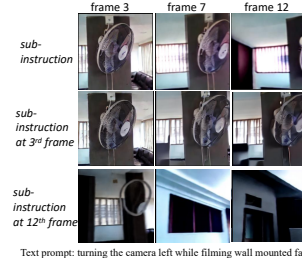


Figure 9: TVP results conditioned on frame-wise sub-instruction and constant clone of sub-instruction.

5.7 ZERO-SHOT EVALUATION

To assess the model’s generalizability on an unseen dataset, we conducted a comparative analysis between our proposed Seer model and the current state-of-the-art baseline VideoFusion. Both models were fine-tuned on the Something-V2 dataset and evaluated on the EGO4D dataset, with the results summarized in Table 5. Notably, Seer method outperformed VideoFusion, as evidenced by its superior performance in terms of FVD and KVD metrics. Drawing insights from these results, we can conclude that our default setup demonstrates strong generalizability while effectively adapting to the EGO4D dataset with higher-quality video generation.

Table 5: **Zero-shot text-conditioned video prediction on EGO4D**

temp. attn.	FVD↓	KVD↓
<i>VideoFusion</i>	618.6	1.85
<i>Seer(Ours)</i>	301.7	0.55

6 CONCLUSION

In this paper, we propose Seer, a sample and computation-efficient model, for the challenging text-conditioned video prediction (TVP) task. We design a data and computation-efficient video network with Frame Sequential Text (FSText) Decomposer to inflate the pretrained text-to-image (T2I) stable diffusion models along the temporal axis. With the rich prior knowledge contained in pretrained T2I models and the well-designed architecture, Seer successfully generates high-quality videos by only fine-tuning the SAWT-Attn and FSText Decomposer, which significantly reduces the data and computation costs. The experiments illustrate our superior performance over all the recent models.

7 REPRODUCIBILITY STATEMENT

The main implementations of our proposed method are in Section 5.2 and 5.1. In addition, the settings of the experiments and hyper-parameters we choose are in Appendix C. And the implementation details are in Appendix B.1.

ACKNOWLEDGEMENT

This work is supported by the Ministry of Science and Technology of the People’s Republic of China, the 2030 Innovation Megaprojects “Program on New Generation Artificial Intelligence” (Grant No. 2021AAA0150000). This work is also supported by the National Key R&D Program of China (2022ZD0161700).

REFERENCES

- Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6299–6308, 2017.
- Yilun Dai, Mengjiao Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Josh Tenenbaum, Dale Schuurmans, and Pieter Abbeel. Learning universal policies via text-guided video generation. *arXiv preprint arXiv:2302.00111*, 2023.
- Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. The epic-kitchens dataset: Collection, challenges and baselines. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 43(11):4125–4141, 2021. doi: 10.1109/TPAMI.2020.2991965.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- Ming Ding, Wendi Zheng, Wenyi Hong, and Jie Tang. Cogview2: Faster and better text-to-image generation via hierarchical transformers. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=GkDbQb6qu_r.
- Frederik Ebert, Yanlai Yang, Karl Schmeckpeper, Bernadette Bucher, Georgios Georgakis, Kostas Daniilidis, Chelsea Finn, and Sergey Levine. Bridge data: Boosting generalization of robotic skills with cross-domain datasets. *arXiv preprint arXiv:2109.13396*, 2021.
- Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XV*, pp. 89–106. Springer, 2022.
- Zhangyang Gao, Cheng Tan, Lirong Wu, and Stan Z. Li. Simvp: Simpler yet better video prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3170–3180, June 2022.
- 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 *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XVII*, pp. 102–118. Springer, 2022.

- 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 *Proceedings of the IEEE international conference on computer vision*, pp. 5842–5850, 2017.
- Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, Miguel Martin, Tushar Nagarajan, Ilija Radosavovic, Santhosh Kumar Ramakrishnan, Fiona Ryan, Jayant Sharma, Michael Wray, Mengmeng Xu, Eric Zhongcong Xu, Chen Zhao, Siddhant Bansal, Dhruv Batra, Vincent Cartillier, Sean Crane, Tien Do, Morrie Doulaty, Akshay Erapalli, Christoph Feichtenhofer, Adriano Fragomeni, Qichen Fu, Abrham Gebreselasie, Cristina González, James Hillis, Xuhua Huang, Yifei Huang, Wenqi Jia, Weslie Khoo, Jáchym Kolář, Satwik Kottur, Anurag Kumar, Federico Landini, Chao Li, Yanghao Li, Zhenqiang Li, Karttikeya Mangalam, Raghava Modhugu, Jonathan Munro, Tullie Murrell, Takumi Nishiyasu, Will Price, Paola Ruiz, Merey Ramazanova, Leda Sari, Kiran Somasundaram, Audrey Southerland, Yusuke Sugano, Ruijie Tao, Minh Vo, Yuchen Wang, Xindi Wu, Takuma Yagi, Ziwei Zhao, Yunyi Zhu, Pablo Arbeláez, David Crandall, Dima Damen, Giovanni Maria Farinella, Christian Fuegen, Bernard Ghanem, Vamsi Krishna Ithapu, C. V. Jawahar, Hanbyul Joo, Kris Kitani, Haizhou Li, Richard Newcombe, Aude Oliva, Hyun Soo Park, James M. Rehg, Yoichi Sato, Jianbo Shi, Mike Zheng Shou, Antonio Torralba, Lorenzo Torresani, Mingfei Yan, and Jitendra Malik. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 18995–19012, June 2022.
- Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffusion models for high-fidelity video generation with arbitrary lengths. *arXiv preprint arXiv:2211.13221*, 2022.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- 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*, 2022a.
- Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *arXiv preprint arXiv:2204.03458*, 2022b.
- Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pre-training for text-to-video generation via transformers. In *International Conference on Learning Representations*, 2023.
- Tobias Höppe, Arash Mehrjou, Stefan Bauer, Didrik Nielsen, and Andrea Dittadi. Diffusion models for video prediction and infilling. *arXiv preprint arXiv:2206.07696*, 2022.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022a. URL <https://openreview.net/forum?id=nZeVKeeFYf9>.
- Yaosi Hu, Chong Luo, and Zhenzhong Chen. Make it move: Controllable image-to-video generation with text descriptions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022b.
- Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 5(2):3019–3026, 2020.
- Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 1725–1732, 2014.

- Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
- Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 3202–3211, 2022.
- 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10209–10218, 2023.
- Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *International Conference on Machine Learning*, pp. 16784–16804. PMLR, 2022.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pp. 8821–8831. PMLR, 2021.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- Scott Reed, Zeynep Akata, Xinchun Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. In *International conference on machine learning*, pp. 1060–1069. PMLR, 2016.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022.
- Masaki Saito, Shunta Saito, Masanori Koyama, and Sosuke Kobayashi. Train sparsely, generate densely: Memory-efficient unsupervised training of high-resolution temporal gan. *International Journal of Computer Vision*, 128(10-11):2586–2606, 2020.
- Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport: What and where pathways for robotic manipulation. In *Conference on Robot Learning*, pp. 894–906. PMLR, 2022.
- 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 preprint arXiv:2209.14792*, 2022.

- Ivan Skorokhodov, Sergey Tulyakov, and Mohamed Elhoseiny. Stylegan-v: A continuous video generator with the price, image quality and perks of stylegan2. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3626–3636, 2022.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021.
- Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pp. 4489–4497, 2015.
- Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. Mocogan: Decomposing motion and content for video generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1526–1535, 2018.
- Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.
- 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.
- Vikram Voleti, Alexia Jolicoeur-Martineau, and Christopher Pal. Masked conditional video diffusion for prediction, generation, and interpolation. *arXiv preprint arXiv:2205.09853*, 2022.
- Chenfei Wu, Lun Huang, Qianxi Zhang, Binyang Li, Lei Ji, Fan Yang, Guillermo Sapiro, and Nan Duan. Godiva: Generating open-domain videos from natural descriptions. *arXiv preprint arXiv:2104.14806*, 2021.
- 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 *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XVI*, pp. 720–736. Springer, 2022a.
- Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Weixian Lei, Yuchao Gu, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. *arXiv preprint arXiv:2212.11565*, 2022b.
- Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. Videogpt: Video generation using vq-vae and transformers. *arXiv preprint arXiv:2104.10157*, 2021.
- 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 preprint arXiv:2206.10789*, 2022a.
- Lijun Yu, Yong Cheng, Kihyuk Sohn, José Lezama, Han Zhang, Huiwen Chang, Alexander G Hauptmann, Ming-Hsuan Yang, Yuan Hao, Irfan Essa, et al. Magvit: Masked generative video transformer. *arXiv preprint arXiv:2212.05199*, 2022b.
- 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. *arXiv preprint arXiv:2202.10571*, 2022c.
- Sihyun Yu, Kihyuk Sohn, Subin Kim, and Jinwoo Shin. Video probabilistic diffusion models in projected latent space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.

Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on robot learning*, pp. 1094–1100. PMLR, 2020.

Daquan Zhou, Weimin Wang, Hanshu Yan, Weiwei Lv, Yizhe Zhu, and Jiashi Feng. Magicvideo: Efficient video generation with latent diffusion models. *arXiv preprint arXiv:2211.11018*, 2022.

A APPENDIX

B ADDITIONAL EXPERIMENTAL RESULTS

B.1 IMPLEMENTATION DETAILS OF BASELINES

We compare three baselines in our paper. For Tune-A-Video, to ensure a fair comparison, we use the pre-trained weight of Stable Diffusion-v1.5¹ (same as our model) to initialize the UNet and we fine-tune the model with an image resolution of 256×256 on the training sets of Something Something-V2 (SSv2), Bridgedata and Epic-Kitchens-100 for 200k training steps. Similarly, we further fine-tune pre-trained VideoFusion on Something Something-V2, Epic-Kitchens-100, and BridgeData for 200k training steps. For MCVD, SimVP and PVDM, we train the model with an image resolution of 256×256 , 64×64 , 256×256 respectively on the training sets of SSv2, Epic-Kitchens-100, and Bridgedata for 300k training steps. For TATS and MAGE, we fine-tune the pre-trained UCF-101 model with an image resolution of 128×128 on the training sets of SSv2, Epic-Kitchens-100, and Bridgedata for 300k training steps.

B.2 EVALUATION DETAILS AND RESULTS OF UCF-101

Most prior text-conditioned video generation methods (Hong et al., 2023; Ho et al., 2022b; Singer et al., 2022; Zhou et al., 2022) evaluate their performance on the UCF-101 (Soomro et al., 2012) benchmark. However, since our proposed method, Seer, is designed for text-conditioned video prediction (TVP) on task-level video datasets, the UCF-101 benchmark, which evaluates class-conditioned video prediction on random short-horizon video clips, is not an ideal evaluation benchmark for TVP. Nonetheless, in order to fairly compare these baselines, we still evaluate the class-conditioned video prediction performance of Seer on UCF-101.

Settings Specifically, we fine-tune our model with a video resolution of $16 \times 256 \times 256$ on UCF-101. Following the evaluation protocols of (Hong et al., 2023), Seer predicts the videos conditioned on 5 reference frames during fine-tuning and inference stage. We report FVD and Inception score (IS) metrics on the UCF-101 dataset (Soomro et al., 2012). The IS is calculated by a C3D model (Tran et al., 2015) that is pre-trained on the Sports-1M dataset (Karpathy et al., 2014) and fine-tuned on UCF101. We follow the evaluation code of TGAN-v2 (Saito et al., 2020) to calculate IS metric. Following (Hong et al., 2023; Ho et al., 2022b; Singer et al., 2022), we evaluate the FVD metric with 2,048 samples and IS metric with 100k samples in the validation set of UCF-101.

Results Table 6 presents the class-conditioned video prediction results on UCF-101, demonstrating that Seer outperforms CogVideo (Hong et al., 2023) and MagicVideo (Zhou et al., 2022), but falls short of Make-A-Video (Singer et al., 2022). Make-A-Video employs unlabelled video pre-training on temporal layers and achieves the best performance among all other methods. While Make-A-Video shows superior performance on FVD and IS, Seer has the potential to further improve its generation performance by addressing the following two limitations. First, Seer has not been pre-trained on video data. Second, Seer obtains latent vectors via a pre-trained 2D VAE, which has not been fine-tuned on UCF-101 and limits the video generation quality of Seer (with 259.4 FVD and 68.16 IS reconstruction quality). However, as we focus on the text-conditioned video prediction task, addressing the above limitations on UCF-101 is out of the scope of this paper.

B.3 EVALUATION RESULTS OF SAMPLING STEPS

Sampling steps: To assess the influence of sampling steps on the quality of video predictions across varying sequence lengths, we conducted an evaluation on both 12-frame and 16-frame video predictions using a series of DDIM sampling steps (10, 20, 30, 40, 50, 60 DDIM steps). All generated outputs were sampled utilizing a 12-frame SSv2 fine-tuned model. The comparative results are presented in Figure 10. Notably, the 16-frame curve exhibits a more rapid decline from DDIM steps 10 to 20 compared to the 12-frame curve. As both curves progress beyond DDIM step 30, they tend to stabilize, showing marginal gains. These findings collectively underscore that increasing

¹<https://github.com/CompVis/stable-diffusion>

Table 6: **Class-conditioned video prediction performance on UCF-101** we evaluate the Seer on the UCF-101 with 16-frames-long videos. Ex.data indicates that the model has been pre-trained or fine-tuned on extra datasets.

Method	Ex.data	Cond.	Resolution	FVD↓	IS↑
MoCoGAN-HD (Tulyakov et al., 2018)	No	Class.	256 × 256	700±24	33.95±0.25
VideoGPT (Yan et al., 2021)	No	No	128 × 128	-	24.69±0.30
RaMViD (Höppe et al., 2022)	No	No	128 × 128	-	21.71±0.21
StyleGAN-V (Skorokhodov et al., 2022)	No	No	128 × 128	-	23.94±0.73
DIGAN (Yu et al., 2022c)	No	No	-	577±22	32.70±0.35
TGANv2 (Saito et al., 2020)	No	Class.	128 × 128	1431.0	26.60±0.47
VDM (Ho et al., 2022b)	No	No	64 × 64	-	57.80±1.3
TATS-base (Ge et al., 2022)	No	Class.	128 × 128	278±11	79.28±0.38
MCVD (Voleti et al., 2022)	No	No	64 × 64	1143.0	-
LVDM (He et al., 2022)	No	No	256 × 256	372±11	27±1
MAGVIT-B (Yu et al., 2022b)	No	Class.	128 × 128	159±2	83.55±0.14
VideoFusion (Luo et al., 2023)	txt-video	Class.	128 × 128	173	80.03
CogVideo (Hong et al., 2023)	txt-img & txt-video	Class.	160 × 160	626	50.46
Make-A-Video (Singer et al., 2022)	txt-img & video	Class.	256 × 256	81.25	82.55
MagicVideo (Zhou et al., 2022)	txt-img & txt-video	Class.	-	699	-
Seer(Ours)	txt-img	Class.	256 × 256	260.7	57.74
pre-trained VAE*	-	-	256 × 256	259.4	68.16

*we evaluate the reconstruction quality of pre-trained 2D VAE in this table, the pre-trained 2D VAE is initialized with the pre-trained weight from Stable Diffusion-v1.5 without extra fine-tuning.

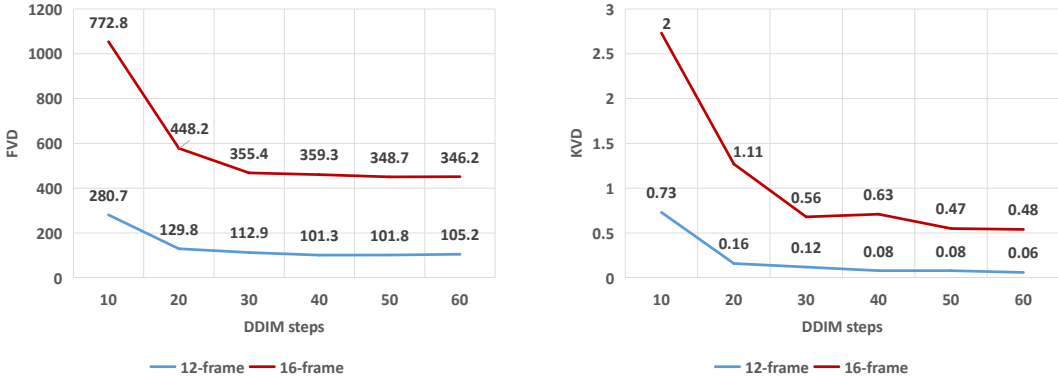


Figure 10: Evaluation results of sampling 12-frame video and 16-frame video using 12-frame Seer model with DDIM sampling steps ranging from 10 to 60 on the Something-Nothing V2 dataset.

DDIM sampling steps notably enhances video quality for longer sequences (10 to 30 DDIM steps). However, the quality improvements of longer videos diminish as DDIM steps exceed 30.

Video length: Furthermore, we present the qualitative outcomes of the 16-frame video in Figure 12. A comparison with the results of the 12-frame video in Figure 23 reveals that both the 16-frame and 12-frame videos can effectively capture task-described motion, as demonstrated in the example of "moving pen." However, the 16-frame video exhibits a gradual loss of appearance information from the reference frame and a decline in temporal consistency along the temporal axis with the increasing sequence length. Given that 16-frame videos are beyond the anticipated sequence length for a 12-frame model, enhancing the generation quality of the 16-frame video could be achieved through the training of a dedicated 16-frame Seer model.

Methods comparison: In addition, we explore the impact of fast sampling on generation results during evaluation by comparing Seer with Tune-A-Video. We apply a series of DDIM sampling steps (10, 20, 30, 40, 50 DDIM steps), as shown in Figure 11. Seer consistently outperformed Tune-A-Video in terms of both FVD and KVD, with improvements observed from 20 DDIM steps to 50 DDIM steps. Particularly noteworthy is Seer's advantage in video quality (280.7 FVD and 0.73 KVD) compared to Tune-A-Video (419.3 FVD and 1.5 KVD) when using only 10 DDIM steps, demonstrating Seer's ability to sample high-fidelity videos efficiently with minimal denoising steps.

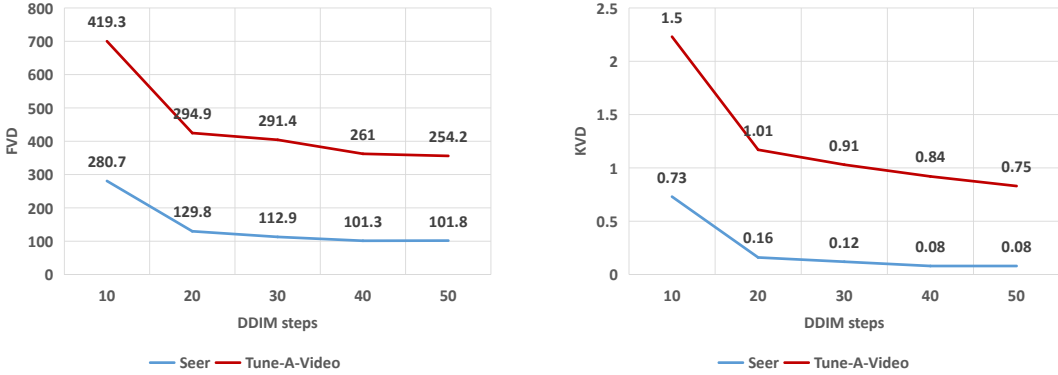


Figure 11: Evaluation results of Seer and Tune-A-Video with DDIM sampling steps ranging from 10 to 50 on the Something-Something V2 dataset.

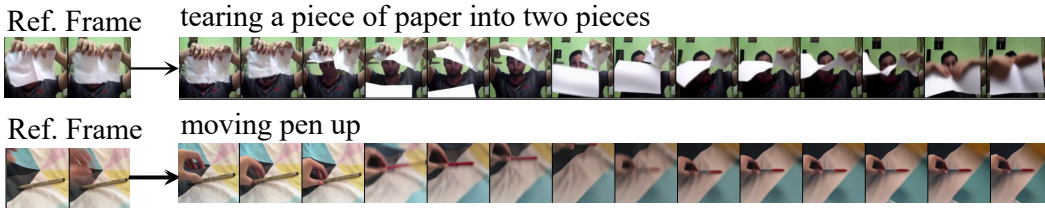


Figure 12: Visualization of 16-frame text-conditioned video prediction sampled by 12-frame Seer on SSv2.

B.4 COMPUTATION EFFICIENCY OF SEER

To enhance the computational efficiency of Seer, both the inclusion of frozen 2D layers and the thoughtful design of temporal layers contribute significantly. A comprehensive evaluation of the computational cost for different temporal layer types, conducted on a single 24GB NVIDIA 3090 GPU, is presented in Table 7. Remarkably, in comparison to plain autoregressive spatial-temporal attention (autoreg.), both bi-directional/directed temporal attention (bi-direct./directed.) and SAWT-Atten (win-auto.) substantially reduce computational overhead. Furthermore, SAWT-Atten demonstrates superior generation quality compared to bi-directional and directed temporal attentions, as evidenced by the ablation results in Section 5.5 of this paper. In addition, a compar-

Table 7: **Training time of the variants of temporal attention on a 16-frame video**

temp. attn.	(sec./iter.)
<i>bi-direct.</i>	2.35
<i>directed.</i>	2.35
<i>autoreg.</i>	5.50
<i>win-auto.</i> (Ours)	2.40

Table 8: **Training time (time.) and GPU memory (Mem.) consumption of the models (16-frame)**

model	2D. frozen	time. (sec./iter.)	Mem. (GB)
<i>Seer (Ours)</i>	Yes	0.75	24.7
<i>Seer</i>	No	0.96	39.2
<i>VideoFusion.</i>	No	1.07	45.0

Table 9: **Training time (time.) and GPU memory (Mem.) consumption of the models (16-frame) ($\geq 90\%$ GPU memory usage)**

model	2D. frozen	time. (sec./iter.)	Mem. (GB)
<i>Seer (Ours)</i>	Yes	3.10	72.9
<i>Seer</i>	No	6.89	75.8
<i>VideoFusion.</i>	No	7.63	78.7

ison involving our ablated setting and the baseline method VideoFusion (Luo et al., 2023), which shares a network design based on the Stable Diffusion U-Net, is presented. In our proposed Seer setting, the 2D spatial layers of the U-Net are frozen during fine-tuning, while the ablated setting maintains all 2D spatial layers trainable. GPU memory consumption and training time for different models were assessed running with a single 80GB NVIDIA A800 GPU. These results are available in Table 8, and highlight the proposed Seer setting with frozen 2D layers exhibiting approximately half the GPU memory consumption (24.7GB) compared to Seer without frozen 2D layers (39.2GB) and VideoFusion (45.0GB). To ensure a fair comparison of training speed among various settings on the single A800 GPU, we conducted additional assessments of models with GPU memory usage exceeding 90%, as detailed in Table 9. The results in the table reveal a notable advantage in com-

Table 10: **Layer depth in FSText. (SSv2).**

num. layers.	FVD↓	KVD↓
2	238.6	0.51
4	229.7	0.23
8(Ours)	112.9	0.12

Table 11: **Components in FSText Decomposer (SSv2).** Our settings are marked in gray

Temp.	Seq.	FVD↓	KVD↓
✓		125.8	0.13
	✓	127.7	0.14
✓	✓	112.9	0.12

Table 12: **Ablation study of temporal attention (BridgeData)**

temp. attn.	FVD↓	KVD↓
<i>bi-direct.</i>	284.5	0.71
<i>directed.</i>	258.0	0.64
<i>autoreg.</i>	261.5	0.83
<i>win-auto.(Ours)</i>	246.3	0.55

Table 13: **Ablation study of Fine-tune settings (BridgeData)**

fine-tune	FSText.	FVD↓	KVD↓
<i>temp-attn.</i>		410.7	0.97
<i>cross+temp-attn.</i>		319.9	1.01
<i>temp-attn.(Ours)</i>	✓	246.3	0.55
<i>cross+temp-attn.</i>	✓	2058.4	9.43

putation efficiency for the proposed Seer setting with frozen 2D layers, exhibiting reduced training time (3.10 sec/iter) compared to Seer without frozen 2D layers (6.89 sec/iter) and VideoFusion (7.63 sec/iter). These findings firmly establish the enhanced computational efficiency of Seer relative to baseline models.

B.5 ADDITIONAL ABLATION RESULTS

FSText layer depth In this section, we additionally investigate the impact of FSText Decomposer’s layer depth in Table 10. Our default setting (8-layer FSText Decomposer) outperforms shallower models (2-layer and 4-layer) in terms of FVD. Though the 4-layer model shows a marginal advantage over the 8-layer model in terms of KVD, our experiments indicate that the 8-layer FSText Decomposer shows a remarkable advantage on FVD metrics and exhibits robustness in text-video alignment. Therefore, we adopt the 8-layer FSText Decomposer as the default setting for Seer.

FSText componentets To evaluate the impact of individual attention layers within the FSText decomposer, we conducted an ablation study on the FSText component, as presented in Table 8. In this study, we ablate the temporal attention layer, labeled as ”Temp.,” and the text-sequential attention layer, denoted as ”Seq.,” within the FSText network. To maintain consistency in model size across different settings, we replaced the ablated component with a Cross-Attention layer. The results, shown in Table 11, highlight the superiority of our proposed setting, which integrates both text-sequential-attention layers and temporal attention layers. Our proposed setting outperforms the other two settings, underscoring the significant attributes of Text-Sequential-Attention layers and Temporal Attention layers to capture text contextual information and model temporal dependencies, collectively enhancing the overall performance of the FSText decomposer.

Additional ablation on BridgeData To validate the robustness and consistency of Seer ablation results across different datasets, we conducted additional experiments on the BridgeData dataset, extending our analysis from the Something Something-V2 (SSv2) dataset. The corresponding ablation studies on Seer fine-tuning settings are presented in Table 13, while the temporal layer settings are detailed in Table 12. These results mirror the ablation outcomes reported in Table 4 and Table 2 for the SSv2 dataset. Notably, the consistent improvement observed in both Seer fine-tuning and temporal layer ablation across different datasets, as demonstrated in Table 13 and Table 12 on the BridgeData dataset, demonstrates the robustness of the Seer component design.

Qualitative results of fine-tuning ablation We conduct a qualitative analysis of various fine-tune settings. We provide additional visualizations of Fine-tune Setting ablation in Section 5.5 of the main paper. Figure 13 shows the results of different settings. Among these settings, our default setting ”*temp+FSText*” stands out as it preserves a higher-level temporal consistency in video prediction starting from reference frames and also delivers superior text-based video motion compared to the other fine-tune settings.

B.6 EVALUATION OF POLICY LEARNING ON SEER

To investigate whether Seer can help policy learning, we choose the UniPi (Dai et al., 2023) as our baseline in the simulated robotics environment Meta-World (Yu et al., 2020), which gener-

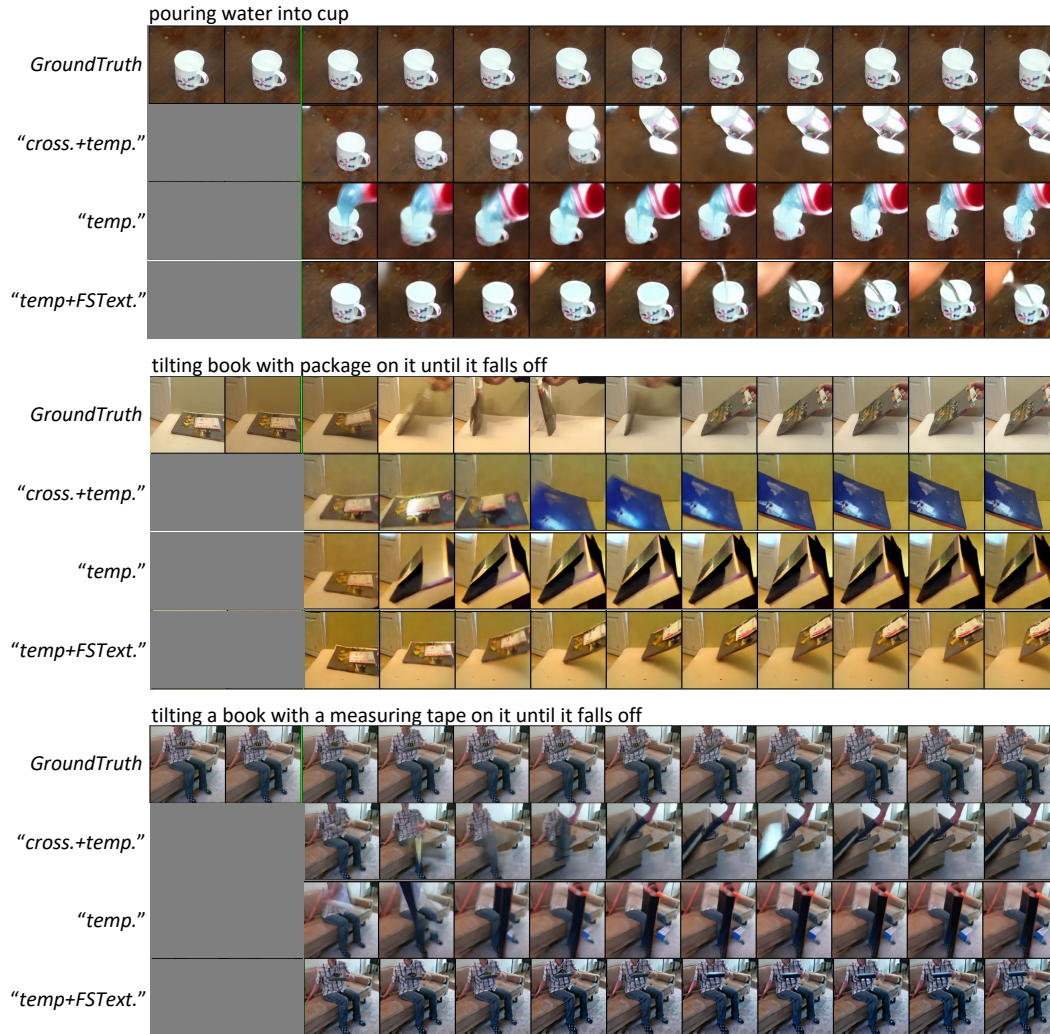


Figure 13: Additional qualitative results of fine-tuning ablation. *"temp+FSText."* is our default setting.

ates videos from the initial state and infers actions from the adjacent frames via a pre-trained inverse-dynamics model. Specifically, we distill a policy model from the videos generated by the Seer and the labeled actions from the pretrained inverse-dynamics model. We choose 3 tasks and use 10 in-domain videos for each task to fine-tune the Seer. And we make comparisons between **a) Policy a:** Distilled policy from the dataset generated by the Fine-tuned Seer model (1000 generated videos for each task). **b) Policy b:** Distilled policy from the 10 in-domain videos for each task. The results are shown in Table 14, where we find that compared to utilizing the given 10 in-domain videos to generate policy, fine-tuning the Seer with them and generating more videos can be better because it can acquire more scalable data, which is of comparable quality. Therefore, the videos generated by Seer can help policy learning in simulated robotics tasks in a way. The visualization of Seer within a robot simulation environment is presented in Section D.5

Table 14: Success rate of distilled policy

tasks	policy a	policy b
<i>button-press-topdown-v2</i>	0.45	0.4
<i>drawer-close-v2</i>	0.1	0.0
<i>drawer-open-v2</i>	0.05	0.0

C IMPLEMENTATION DETAILS

C.1 FINE-TUNING AND SAMPLING

In this section, we list the hyperparameters, fine-tuning details, sampling details, and hardware information of our model in Table 15.

C.2 ARCHITECTURE INFORMATION

In this section, we list the hyperparameters of 3D U-Net in Table 16 and hyperparameters of FSText Decomposer in Table 17.

Table 15: Hyperparameters and details of Fine Tuning/Inference

param.	value
optim.	AdamW
Adam- β_1	0.9
Adam- β_2	0.99
Adam- ϵ	$1e^{-8}$
weight decay	$1e^{-2}$
lr	$1.024e^{-4}$
end lr	0.0
lr sche.	cosine
noise sche.	cosine
train batch size	1/GPU
grad. acc.	2
warmup steps	10k
resolution	256×256
train. steps	200k
train. hardware	4 RTX 3090
val. batch size	2/GPU
sampler	DDIM
sampling steps	30
guidance scale	7.5

Table 16: Hyperparameters of 3D U-Net

hyperparam.	value
input/output channels	4
Base channels	320
Channel multipliers	1,2,4,4
3D Downsample blocks	4
3D Upsample blocks	4
Number of layers (per block)	2
Modules of layer	3D ResnetBlock Spatial-cross Atten. SAWT-Atten. Down./Up. 3D ResnetBlock
Dimension of atten. heads	8
activation function	SiLU
Dimension of cross-atten.	768

Table 17: Hyperparameters of FSText Decomposer

hyperparam.	value
learnable tokens channels	768
output channels	768
Base channels	768
Number of layers	8
Modules of layer	Seq-cross Atten. Feedforward Directed temporal Atten. Feedforward
Number of atten. heads	8
Dimension of cross-atten.	768

D VISUALIZATION

D.1 ADDITIONAL QUALITATIVE RESULTS

We provide additional visualization on Something-Something v2 (SSv2) of our text-conditioned video prediction in Figure 22, and text-conditioned video prediction/manipulation results in Figure 23. Additionally, we provide the visualization on BridgeData of text-conditioned video prediction in Figure 24 and text-conditioned video prediction/manipulation in Figure 25. We also provide the visualization results of text-conditioned video prediction on Epic-Kitchens-100 in Figure 26.

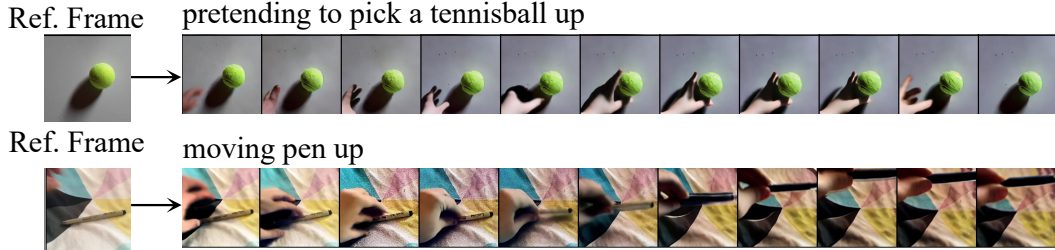


Figure 14: Additional visualization results of 12-frame Text-conditioned video prediction (reference frame=1) on SSv2 dataset.

D.2 GENERALIZABILITY OF SEER ON DOWNSTREAM TASKS

To thoroughly assess Seer’s adaptability on downstream tasks such as video manipulation on BridgeData, we conducted an investigation into diverse fine-tuning strategies on this dataset. Our default setting proceeded to fine-tune only the temporal block of 3D U-Net, freezing the FSText decomposer on a mixed dataset including both the pre-training video dataset and the task-specific downstream dataset. Additionally, we compare our default setting with LORA modules (Hu et al., 2022a), which are integrated into each temporal layer of the UNet architecture and throughout the entire FSText decomposer. As depicted in Figure 15, both our proposed setting and the LORA setting have effectively identified the previously unseen ”black mug” in the training set of BridgeData. Besides, our proposed setting maintains a higher background consistency from the reference frame (frame 0) and generates future frames with superior fidelity when compared to the LORA setting. We also provide additional visualization in Appendix D.4.

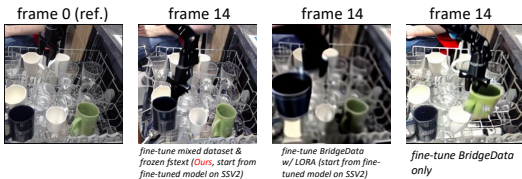


Figure 15: The visualization of comparison with different fine-tune settings on BridgeData (given text instructions ”pick up black mug”, where ”black” text prompt is unseen in BridgeData). Compare to other two settings, our default setting successfully recognizes ”black mug” and generates RGB frame with higher fidelity

D.3 ADDITIONAL VISUALIZATION OF SUB-INSTRUCTION EMBEDDING

We present additional visualizations comparing frame-specific sub-instructions to those constant clones from the third or twelfth frame along the temporal axis (Figure 16). These visualizations reveal that each frame’s sub-instruction represents the motion at its corresponding time step. For example, in the case of “pushing iphone adapter from left to right” the default sub-instruction guides the entire instructed behavior with continuous temporal motion, while the sub-instruction embedding from the third frame directs the motion to a halfway point in the case of pushing the iphone adapter to right. In contrast, utilizing the sub-instruction embedding from the twelfth frame leads to frames that are disconnected from the reference frame’s scene, resulting in the absence of transition between frames. This observation demonstrates that sub-instructions closely follow the temporal sequence of the global instruction, providing precise guidance across multiple steps.

D.4 ADDITIONAL VISUALIZATION OF VIDEO MANIPULATION WITH UNSEEN OBJECTS AND ZERO-SHOT VIDEO MANIPULATION

We also include visualizations of video manipulation with unseen objects in BridgeData and zero-shot generation on EGO4D. These experiments involved fine-tuning the SSv2 video model on a mixed dataset consisting of SSv2 and BridgeData. The fine-tuning process lasted for 800k iterations, employing a learning rate of $4.096e^{-5}$. Subsequently, the model’s performance was evaluated on Bridgedata (as shown in Figure 17(a)) and the EGO4D dataset (illustrated on the right side of Figure 17(b)), where we sampled 15 frames with one reference frame. In the evaluation on BridgeData, Seer demonstrated the ability to recognize previously unseen objects such as “red plate” and “cabbage” and successfully performed a “pick up” motion, leveraging prior knowledge learned from SSv2. In certain scenarios, such as picking up the plate, where Seer has not encountered the specific action during training, the generated video depicted some limitations in the interaction between the robotic arm and the plate. Regarding the evaluation on EGO4D, although Seer had not been fine-tuned on this dataset, it exhibited the capability to accurately identify objects in the EGO4D environment, such as “laptop” and “book,” and execute actions based on prior knowledge acquired from observing human activities. However, Seer still faced challenges in predicting future frames based on the understanding of the scene in the reference frame. For instance, in the case of “closing a book,” Seer tended to generate a hand outside the camera view instead of manipulating the object with the main body, such as the hand holding the book, within the scene.

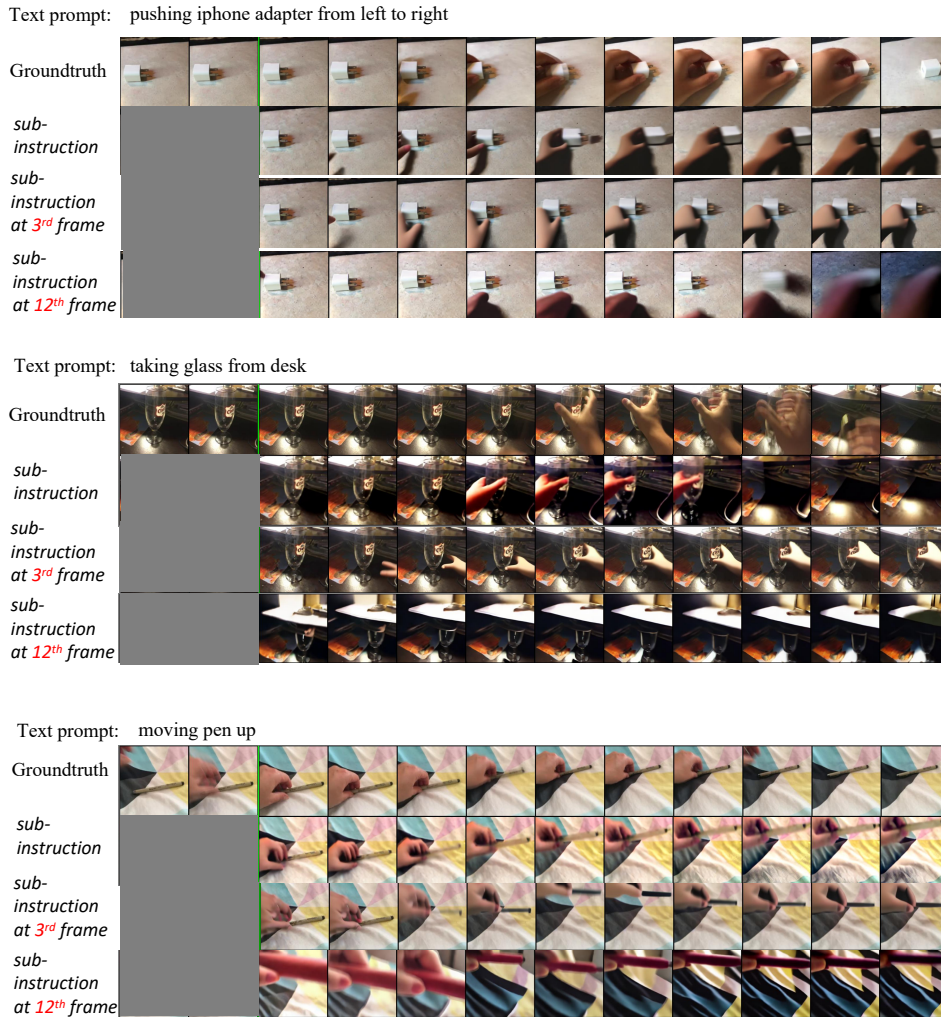


Figure 16: Additional visualization results of Seer’s Text-conditioned video prediction conditioned on frame-wise sub-instruction and duplicate sub-instruction at the third/twelfth frame.

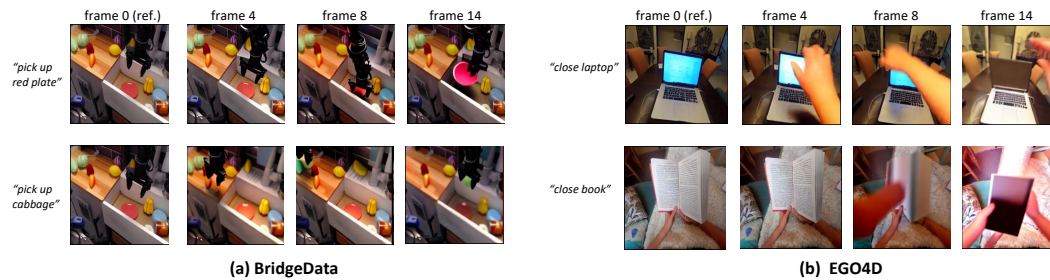


Figure 17: Additional visualization of Seer generalizability evaluation on unseen dataset EGO4D (b) and BridgeData (a), where "cabbage" and "red plate" are unseen objects in the training set of BridgeData.

D.5 VISUALIZATION OF ROBOT SIMULATION ENVIRONMENT

To investigate the applicability of Seer in robot environments, we assess its performance on robot simulation datasets, including Meta-World (Yu et al., 2020), CLIPort (Shridhar et al., 2022), and RLbench (James et al., 2020). Beginning with the SSv2-finetuned Seer model, we further fine-tune it on a mixed dataset comprising Something Something-V2 (SSv2), Bridgedata, 236 MetaWorld video clips, 785 CLIPort video clips, and over 2000 video clips from RLbench for 80,000 steps, using a learning rate of 4e-5. Throughout the fine-tuning process, all 2D layers of the 3D inflated U-Net and the FSText decomposer remain frozen. Figure 18 presents visualizations of Seer on these robot simulation datasets. In these visualizations, Seer successfully generates coherent videos with continuous motion trajectories aligned with language instructions. In the CLIPort video clips, Seer accurately detects small target objects in complex scenes and places them in the correct positions as instructed by the task descriptors. Moreover, in the RLbench videos, Seer demonstrates the ability to perform multi-step tasks such as "stacking a pyramid with the boxes." These observations highlight Seer’s adaptability to a multi-task environment, encompassing human video, robot manipulation, and simulation scenarios, while maintaining robust temporal alignment with task descriptors. We conduct an assessment of Seer performance on policy learning, and a further analysis is presented in Section B.6.

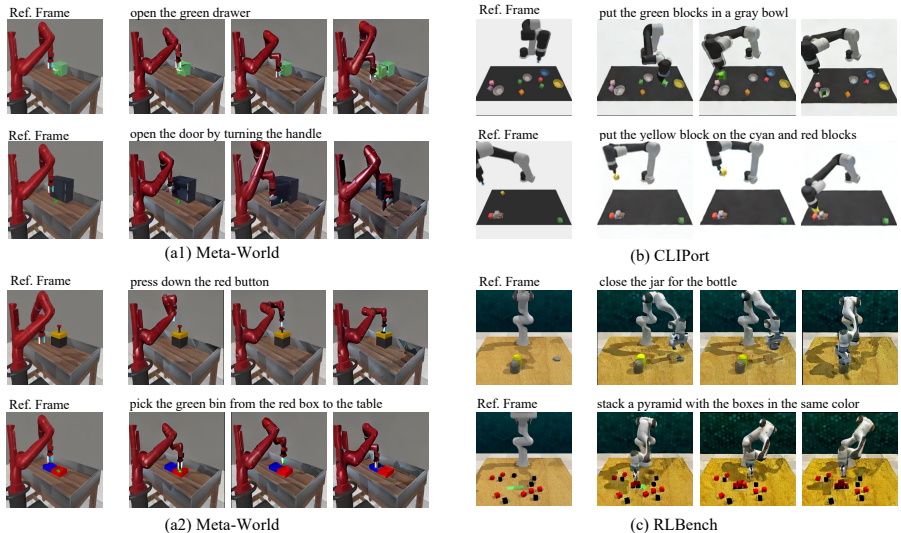


Figure 18: Text-conditioned video prediction of Seer (conditioned on first frame) on Meta-World (a1-2), CLIPort (b), and RLbench (c) datasets

D.6 LONG VIDEO PREDICTION

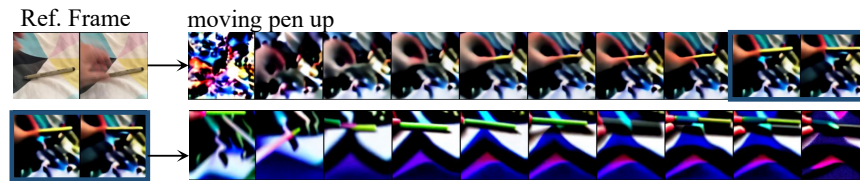
To evaluate Seer’s ability to generate extended video frames during the inference stage, we conducted a qualitative test on long video prediction using a 12-frame SSv2 fine-tuned Seer model. Our video sampling approach involves the sequential generation of video clips, where the first 12 frames are conditioned on 2 reference frames, and the last two frames of the first clip serve as the condition for the subsequent clip. The concatenation of these clips results in a long video. A crucial element for aligning long video frames with language is the presence of a prolonged sequence of sub-instructions. To achieve this, we employ two text conditioning strategies for expanding the sub-instruction embedding along the frame axis, enabling the iterative sampling of longer frames.

The first strategy involves interpolating sub-instruction embedding along the frame axis, while the second strategy entails repeating the sub-instruction embedding from the first video clip to guide the second clip. In Figure 19, we compare these two strategies and observe that direct interpolation (Figure 19 (a)) tends to degrade overall generation quality, introducing unexpected noise. Conversely, utilizing the same sub-instruction (Figure 19 (b)) can maintain coherent motion, persisting until the target object is no longer present or the current motion is in a terminated state. Although this strategy facilitates the generation of coherent movements, it is distinct from a simple upsampling of video

clips along the frame axis. It is noteworthy that all results were obtained using the 12-frame Seer model.

Intuitively, appending an additional upsampler network could potentially enhance Seer’s ability to expand video frames, which causes extra computational costs. Observing the results of the evaluation, we believe that expanding the frame length of generated subinstruction embeddings during the training stage represents a promising direction for enabling Seer to generate longer video frames in multi-steps without incurring additional computational overhead.

(a) Interpolate sub-instruction embedding



(b) Repeat sub-instruction embedding

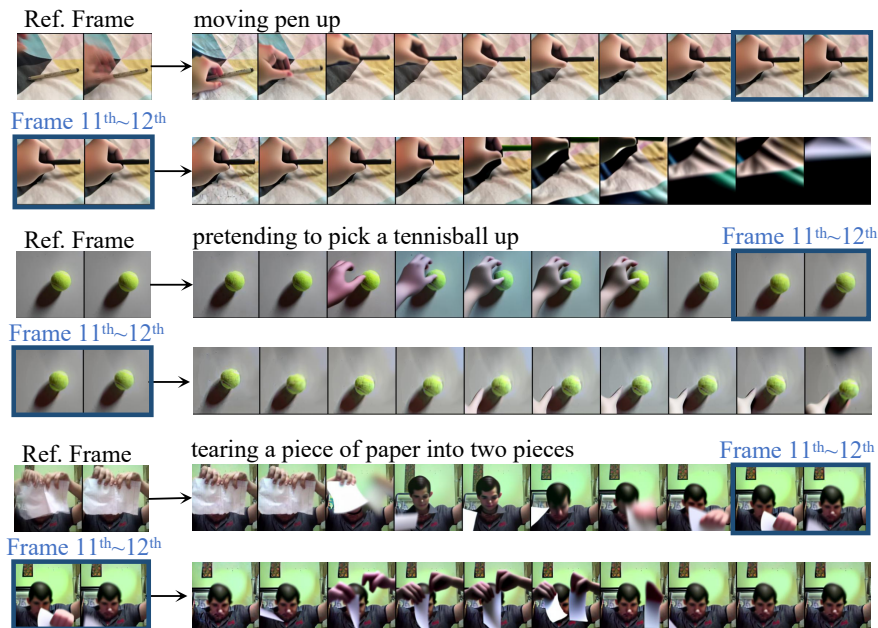


Figure 19: 22-frame video prediction conditioned on 2 reference frames using a 12-frame Seer model on SSv2. (a) Interpolation of 12-frame subinstruction embeddings to a 22-frame sequence. (b) Repetition of 12-frame subinstruction embeddings for the second video clip prediction from the eleventh frame to the twenty-second frame.

D.7 FAILURE CASE

In this section, we present instances where Seer encounters challenges in handling environmental motion in human-generated videos. The generated videos highlight situations such as "dropping a card in front of a coin" and "book falling like a rock" (refer to Figure 20), where Seer successfully predicts task-descriptive motions like "dropping" and "falling" and correctly identifies text-described objects such as "card" and "book." However, the generated future frames fall short in capturing appearance consistencies, such as the color of the card in the previous frame and the cover of the book in the reference environment. In the scenario of "pouring red wine into a glass," Seer tends to generate a wine glass based on its knowledge of pouring red wine but overlooks the transition distribution from the reference environment.

Notably, in the Epic-Kitchens-100 dataset, where scene transitions are prevalent, Seer exhibits a preference for predicting camera pose movements and generating novel views of the environment, reflecting its imaginative capabilities. However, these outcomes extend beyond the scope of Seer's primary objective, which is to learn human behavior. Consequently, addressing challenges such as filtering out irrelevant background information, including camera pose and object occlusion, and refining Seer's awareness of temporal motion becomes imperative for its adaptation to learning from internet videos.

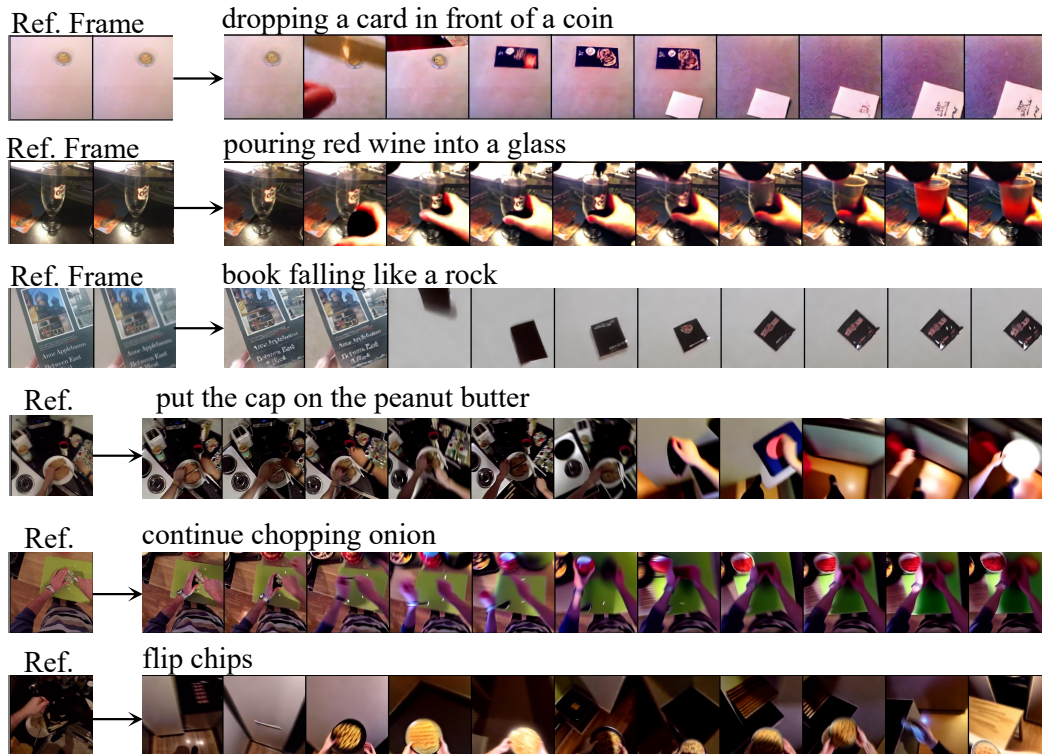


Figure 20: Seer failure cases on Something Something-V2 (row 1,2,3) and Epic-Kitchens-100 (row 4,5,6)

E HUMAN EVALUATION DETAILS

To evaluate the quality of video predictions according to human preferences, we conducted a human evaluation with 99 video clips on the validation set of the Something-Something V2 dataset (SSv2), the evaluation process involved 54 anonymous evaluators. To eliminate biases towards specific baselines, we randomly selected 20 questions for each evaluator. Each single-choice question consisted of a ground-truth video as a reference, a manually modified text instruction, and two video prediction results generated by Seer and another baseline method. The evaluators were required to choose the video clip that is more consistent with the text instruction and has higher fidelity from the two options. To ensure the clarity of the questions, we provided an example to explain the options in each questionnaire. Moreover, we recommended that evaluators prioritize video predictions with strong text-based motions as their first preference and the fidelity of the generated video as their second preference. For reference, Figure 21 provides a screenshot of an example questionnaire.

In total, we collected 342 responses for the Seer vs. TATS comparison, 363 responses for the Seer vs. Tune-A-Video comparison, and 357 responses for the Seer vs. MCVD comparison. And the results in the main paper Figure 7 are calculated based on the collected questionnaires.

Text-conditioned Video Prediction

Please review the text prompt and the reference video and choose the video result that is closer to the text description from the two options presented. If neither option accurately matches the text description, please choose the result that is more similar to the reference video, regardless of the resolution. For example, if the resolution is low but the generated result closely matches the text description, it is recommended to select that option.

Example:

According to the following text description, the video that needs to be generated is: "turning the camera right while filming wall-mounted fan". Please note that the given reference video may not match the language description exactly, but your task is to choose a video that accurately reflects the given language description.

Reference frame



(A)

(B)

Figure 21: Screenshot of a questionnaire example shown to human evaluators.

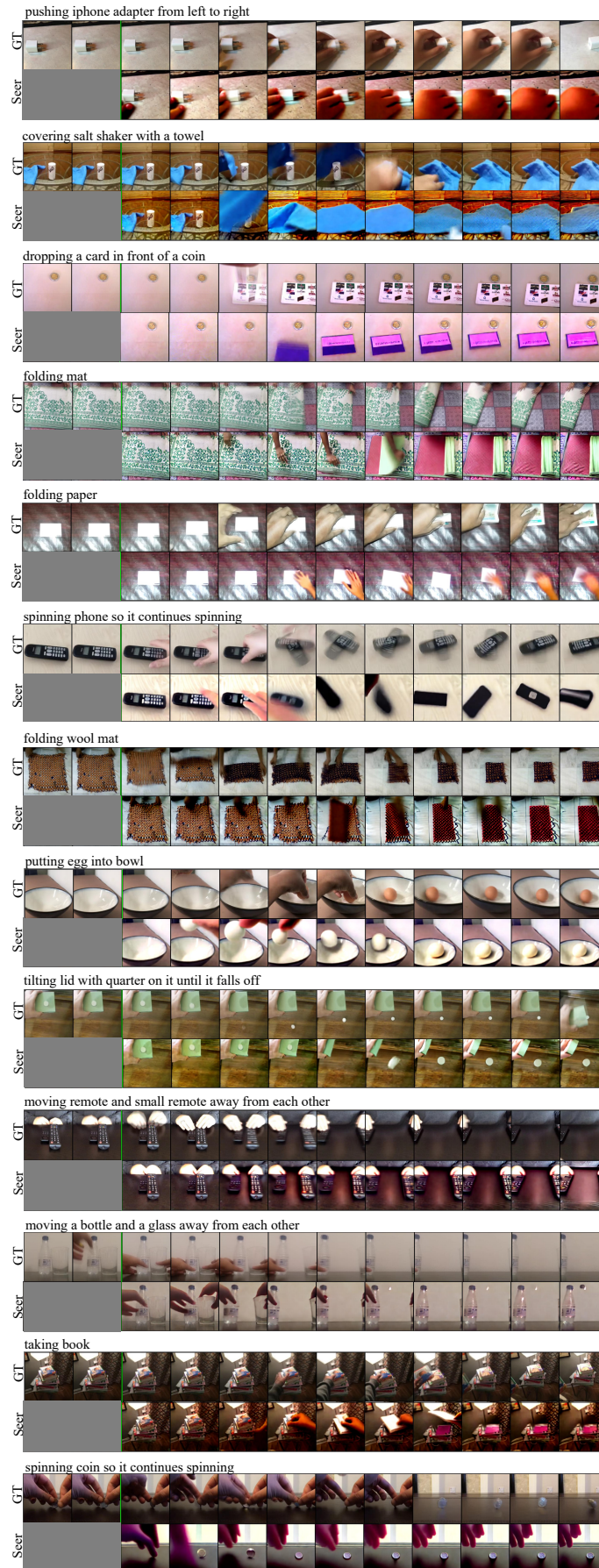


Figure 22: Text-conditioned video prediction of Seer on SSV2.



Figure 23: Text-conditioned video prediction/manipulation of Seer on SSv2, where “pred.” refers to prediction, “mani.” refers to manipulation.



Figure 24: Text-conditioned video prediction of Seer on BridgeData.



Figure 25: Text-conditioned video prediction/manipulation of Seer on BridgeData, where “pred.” refers to prediction, “mani.” refers to manipulation.

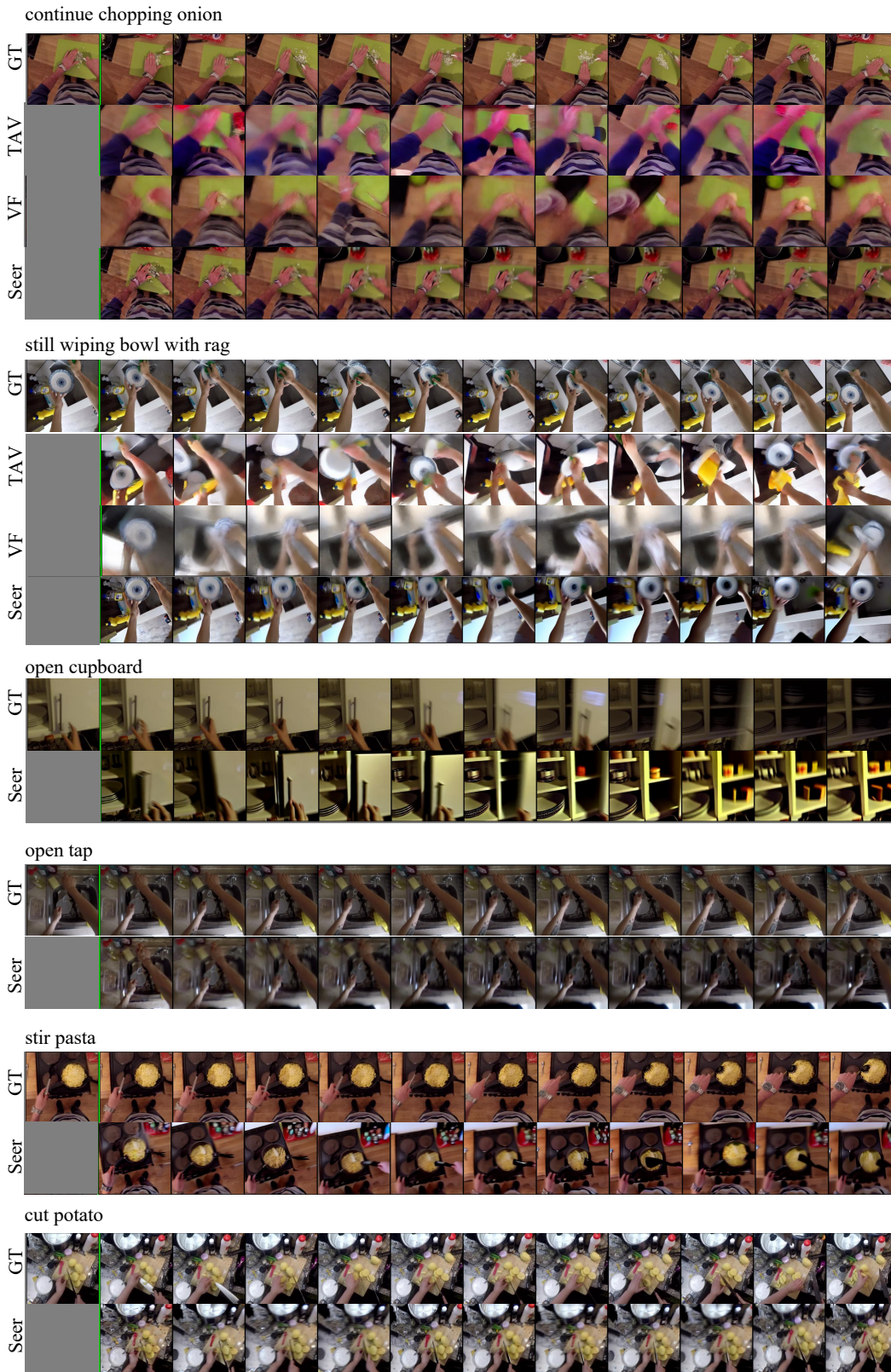


Figure 26: Text-conditioned video prediction (conditioned on first frame) on Epic-Kitchens-100. TAV refers to Tune-A-Video, VF indicates VideoFusion.