

---

# TALC: Time-Aligned Captions for Multi-Scene Text-to-Video Generation

---

Hritik Bansal<sup>1</sup>

Yonatan Bitton<sup>2†</sup>

Michal Yarom<sup>2†</sup>

Idan Szpektor<sup>2\*</sup>

Aditya Grover<sup>1\*</sup>

Kai-Wei Chang<sup>1\*</sup>

<sup>1</sup>University of California Los Angeles

<sup>2</sup>Google Research

## Abstract

Recent advances in diffusion-based generative modeling have led to the development of text-to-video (T2V) models that can generate high-quality videos conditioned on a text prompt. Most of these T2V models often produce single-scene video clips that depict an entity performing a particular action (e.g., ‘a red panda climbing a tree’). However, it is pertinent to generate multi-scene videos since they are ubiquitous in the real-world (e.g., ‘a red panda climbing a tree’ followed by ‘the red panda sleeps on the top of the tree’). To generate multi-scene videos from the pretrained T2V model, we introduce **Time-Aligned Captions (TALC)** framework. Specifically, we enhance the text-conditioning mechanism in the T2V architecture to recognize the temporal alignment between the video scenes and scene descriptions. For instance, we condition the visual features of the earlier and later scenes of the generated video with the representations of the first scene description (e.g., ‘a red panda climbing a tree’) and second scene description (e.g., ‘the red panda sleeps on the top of the tree’), respectively. As a result, we show that the T2V model can generate multi-scene videos that adhere to the multi-scene text descriptions and be visually consistent (e.g., entity and background). Further, we finetune the pretrained T2V model with multi-scene video-text data using the TALC framework. We show that the TALC-finetuned model outperforms the baseline methods by 15.5 points in the overall score, which averages visual consistency and text adherence using human evaluation. The project website is <https://talcmst2v.github.io/>.

## 1 Introduction

The ability to generate videos that simulate the physical world has been a long-standing goal of artificial intelligence [1, 2, 3, 4]. In this regard, text-to-video (T2V) models have seen rapid advancements by pretraining on internet-scale datasets of images, videos, and texts [5, 6]. Previous works [7, 8, 9, 10, 11, 12] primarily focus on training conditional denoising diffusion probabilistic models [13] on paired video-text data [14, 15]. After training, these models allow for video generation by sampling from the trained diffusion model, conditioned on a text prompt. However, most of the open-models such as ModelScope[10] VideoCrafter [16, 17], OpenSora [18] are trained with single-scene video-text dataset [14, 19], which is widely available and easy to acquire. However, real-world scenarios often require the generation of multi-scene videos from multi-scene descriptions (e.g., *Scene1*: ‘A koala is napping on a tree.’ *Scene2*: ‘The koala eats leaves on the tree.’). In such cases, the generated video should accurately depict the events in their temporal order (e.g., *Scene2*

---

<sup>†</sup> Equal Contribution. \* Equal Advising. Contact [hbansal@ucla.edu](mailto:hbansal@ucla.edu), [yonatanbitton1@gmail.com](mailto:yonatanbitton1@gmail.com).

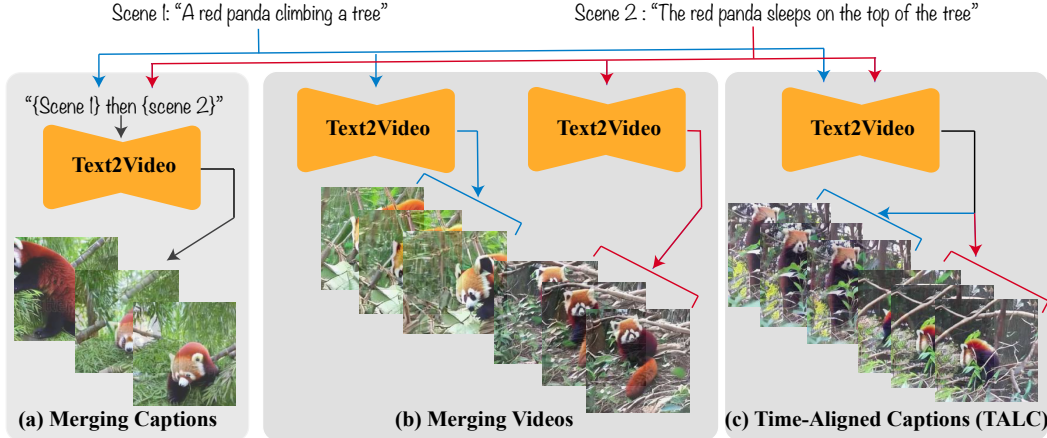


Figure 1: **Multi-scene video generation methods.** (a) Generating a video by merging scene 1 and scene 2 descriptions. (b) The resulting video is composed of the video generated by the description of scene 1 and the video generated by the description of scene 2. (c) In our method (TALC) the generated video is conditioned on the description of scene 1 for the first half of the video frames and on the description of scene 2 for the later video frames.

follows *Scene1*) while maintaining visual consistency, meaning that backgrounds and entities should remain consistent across scenes. While high-performance text-to-video models such as Sora [4] might be able to generate multi-scene videos, we point out that they are closed-source models trained with massive compute resources and lack sufficient details on the model design, training protocol, and datasets. In this work, we present a complementary approach and tackle the challenge of effectively leveraging the capabilities of base T2V models for multi-scene video generation.

The multi-scene text-to-video generation differs from long video synthesis where the goal is to either interpolate (few frames to many frames) [8] or create continuing patterns of the single event in the generated video [11]. Prior works [20, 9] use a transformers [21, 22] to generate video frames for a given scene autoregressively. However, it is hard for their model to generate multiple scenes reliably as the context length increases with history of text descriptions and visual tokens [23] of the previous generated videos (e.g., generating *Scene 4* conditioned on the *Scene1, 2, 3* videos and descriptions). Other works [24] utilize a latent diffusion model [25] to generate video frames autoregressively by conditioning on the entire history of generated videos and scene descriptions. However, the approach is (a) slow due to repeated sampling, (b) generates only one frame per scene description, and (c) shown to work with only limited cartoon characters [26, 27] instead of wide range of visual concepts in the real-world. In this work, our goal is to generate multi-scene videos in the end-to-end manner, using a diffusion text-to-video generative model that is capable of producing content for a wide range of visual entities and actions.

As shown in Figure 1(a), the naive approach to generating a multi-scene video for the scene descriptions  $(T'_1, T'_2)$  would condition the T2V generative model on the merged descriptions. In this setup, the diffusion model processes the entire scene description together, and lacks any information regarding the expected temporal order of events in the generated videos. As a result, we find that this approach leads to poor text-video alignment. As shown in Figure 1(b), an alternative approach generates videos for the individual text descriptions independently and concatenates them in the raw input space along the temporal dimension. While this approach achieves good alignment between the scene description and the scene-specific video segment, the resulting video lacks visual consistency in terms of entity and background appearances.

Prior work [28, 29] generates multi-scene videos by utilizing knowledge of the entity, background, and their movements from large language models [30]. However, these videos are generated independently for each scene before being merged. Moreover, these methods do not offer a way to learn from real-world multi-scene video-text data. To remedy these challenges, we propose TALC (**T**ime-**A**ligned **C**aptions), a simple and effective framework to generate consistent and faithful multi-scene videos. As shown in Figure 1(c), our approach conditions the T2V generative model with the knowledge of the temporal alignment between the parts of the multi-scene video and multi-scene descriptions.



Figure 2: **Examples of multi-scene video generation baselines.** (a) Generating video on the merged descriptions, leads to a poor text-video alignment. (b) Generating videos for the individual text descriptions and concatenate them temporally, leads to a lack of background consistency. (c) Our approach (TALC) enhances the scene-level text-video alignment and maintains background consistency.

Specifically, TALC conditions the visual representations of earlier video frames on the embeddings of the earlier scene description, and likewise, it conditions the representations of later video frames on the embeddings of the later scene description in the temporal dimension. Additionally, the temporal modules in the T2V diffusion architecture allows information sharing between video frames (the first half and the second half) to maintain visual consistency. Thus, TALC enhances the scene-level text-video alignment while providing all the scene descriptions to the diffusion model at once. Further, our TALC framework can enhance the multi-scene text-to-video generation capabilities with real-world multi-scene data (§3.3).

In our experiments, we assess the visual consistency (background and entity consistency) and multi-scene script adherence of the generated videos from Modelscope [10] and Lumiere [6]. Through our automatic and human evaluation, we find that merging scene descriptions leads to high visual consistency but poor text adherence. On the other hand, we observe that merging videos independently achieves the highest text adherence while the visual consistency is compromised. Interestingly, switching to TALC strikes an effective balance between visual consistency and text adherence, outperforming the baseline methods by 11.1 points on the overall score. This score represents the average of visual consistency and text adherence scores, as determined by human evaluation. Furthermore, we construct a multi-scene text-video dataset from real-world videos and fine-tune the T2V generative model using TALC. On our human evaluation, the generated videos from the TALC-finetuned model exhibit higher text adherence than the base model in multi-scene scenarios. Specifically, it outperforms the baseline methods by 15.5 points on the overall score. In summary, our contributions are:

## 2 Preliminaries

In this work, we focus on generating multi-scene videos from scene descriptions using a diffusion-based Text-to-Video (T2V) generative model. The initial step is to equip the generative model with the knowledge of a wide range of visual concepts and actions. This is achieved during the pretraining stage (§2.1). Subsequently, we aim to utilize the base model for multi-scene text-to-video generation task, which we formalize in (§2.3). In §3, we propose our TALC framework and discuss collection of real-world multi-scene text-video data for finetuning the base T2V model.

## 2.1 Diffusion Models for Text-to-Video Generation

Diffusion models [13, 31]  $p_\theta(x)$  are a class of generative models that learn data distribution  $p_{data}(x)$ . Due to their flexible design, we can train their class-conditional versions to learn class-conditional data distributions  $p_{data}(x|y)$  where  $y$  is the conditioning variable, that can take various forms such as labels from a dataset or text description accompanying in a video [32].

We assume a dataset  $\mathcal{S} \subset \mathcal{V} \times \mathcal{T}$  consisting of pairs of  $(V_j, T_j)$  where  $V_j \in \mathbb{R}^{L \times 3 \times H \times W}$  is a raw video consisting of 3 RGB channels,  $L$  frames,  $H$  height,  $W$  width, and  $T_j$  is a text caption. We use  $\mathcal{V}$  and  $\mathcal{T}$  to denote the domain of videos and text, respectively. The aim of T2V generative modeling is to learn the conditional distribution of the videos conditioned on the text  $p_{\mathcal{S}}(V_j|T_j)$ . In this work, we consider diffusion-based generative models that learn the data distribution via iterative denoising of the input video  $z_j \in \mathbb{R}^{L \times C \times H' \times W'}$ . Here,  $z_j$  can either represent the input video in the raw pixel space  $V_j$  [6] or it can represent the latent representation of the video  $z_j = \mathcal{E}(V_j)$  for the latent diffusion models [25] where  $\mathcal{E}$  is an encoder network such as VAE [33].

Given  $z_j$ , diffused variable  $z_{\tau,j} = \alpha_\tau z_j + \beta_\tau \epsilon$  are constructed where  $\epsilon \sim \mathcal{N}(0, I)$  where  $\alpha_\tau$  and  $\beta_\tau$  are sampled from the noise scheduler  $p_\tau$  [34] which define the noise levels the model is trained on. Finally, we train a denoiser network  $f_\theta$  [35, 36] that inputs the diffused variable  $z_\tau$  and embeddings of the text caption to predict the target vector  $y$  where  $y$  can be the original noise  $\epsilon$ , which minimizes the denoising score matching objective [13]:

$$\mathbb{E}_{(V_j, T_j) \in \mathcal{S}, \tau \sim p_\tau, \epsilon \sim \mathcal{N}(0, I)} [\|\epsilon - f_\theta(\tau, z_{\tau,j}, h_j)\|_2^2] \quad (1)$$

where  $h_j = \mathcal{H}(T_j) \in \mathbb{R}^d$  is the embedding of the text caption  $T_j$  where  $\mathcal{H}$  is the text embedding model [37] and  $d$  is the dimension size.

## 2.2 Text Conditioning Mechanism

To ensure the effective textual controllability of video generation, the structure of the denoiser networks is equipped with a cross-attention mechanism [10, 8]. Specifically, it conditions the visual content  $z_\tau \in \mathbb{R}^{L \times C \times H' \times W'}$  on the text. To do so, we first *repeat* the text embeddings of the text caption  $r_j = R(h_j) \in \mathbb{R}^{L \times d}$  where  $R$  is a function that repeats the input text embedding  $h_j$  for  $L$  times in the temporal dimension. Intuitively, the repeat operation represents that the  $L$  frames of the video  $z_j$  are semantically aligned with the textual description  $T_j$  or its text embedding  $r_j$ . In §3, we will manipulate this operation to make the model architecture aware of the video-text alignment in the multi-scene scenario.

These repeated text embeddings  $r_j$  are inputs to the spatial attention block as the key and value in the multi-head attention block. The cross-attention enables the intermediate visual features to capture the semantic information that facilitates an alignment between the language and vision embeddings. Formally,

$$z'_{\tau,j} = CA_{f_\theta}(Q = z_{\tau,j}; K = r_j; V = r_j) \quad (2)$$

where  $CA_{f_\theta}$  is the cross attention mechanism with  $Q, K, V$  as the query, key, and value, respectively, in the spatial blocks of the denoiser network. Additionally,  $z'_{\tau,j}$  is the intermediate representation that is informed with the visual and textual content of the data. In addition to the spatial blocks, the denoiser network also consists temporal blocks that aggregate features across video frames which are useful for maintaining visual consistency in the generated video.

## 2.3 Multi-Scene Text-to-Video Generation

In many real-world scenarios, such as movies, stories, and instructional videos [38], a video may depict multiple transitions with the same or changing entities, as well as multiple actions or events. In addition, the different video segments often share contextual information such as the background or location. These videos are considered multi-scene videos. In this work, we aim to generate multi-scene video  $X = \{x_1, x_2, \dots, x_n\}$  from multi-scene descriptions  $Y = \{y_1, y_2, \dots, y_n\}$  where  $n$  are the number of sentences and each sentence  $y_j$  is a scene description for scene  $j$ . Additionally, the index  $j$  also defines the temporal order of events in the multi-scene script i.e., we want the events

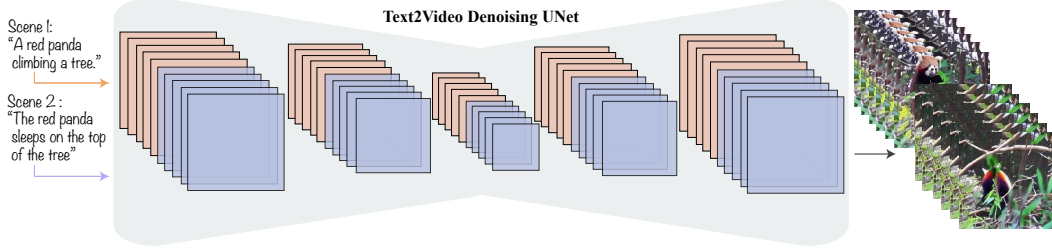


Figure 3: **The architecture of Time-Aligned Captions (TALC).** During the generation process of the video, the initial half of the video frames are conditioned on the embeddings of the description of scene 1 ( $r_{y_1}$ ), while the subsequent video frames are conditioned on the embeddings of the description of scene 2 ( $r_{y_2}$ ).

described in the scene  $j$  to be depicted earlier than the events described in the scene  $k$  where  $k > j$ . Further, we want the parts of the entire generated video  $X$ , given by  $x_j$ , to have high video-text semantic alignment with the corresponding scene description  $y_j$ , also referred to as *text adherence*.

For instance, consider a two-scene description  $Y = \{ \text{'A red panda climbs on a bamboo forest.'}, \text{'The red panda sleeps peacefully in the treetop.'} \}$ . Here, we need the T2V generative model to synthesize the appearance of the red panda (an entity) that remains consistent throughout the generated video, also referred to as *entity consistency*. In addition, we will expect that the context of the multi-scene video of a forest (a background) to remain consistent, also referred to as *background consistency*.

### 3 Method

#### 3.1 TALC: Time-Aligned Captions for Multi-Scene T2V Generation

Most of the existing T2V generative models [10, 16, 6] are trained with large-scale short video-text datasets (10 seconds - 30 seconds) such as WebVid-10M [14]. Here, each instance of the dataset consists of a video and a human-written video description. These videos either lack the depiction of multiple events, or the video descriptions do not cover the broad set of events in the video, instead focusing on the major event shown. As a result, the pretrained T2V generative models only synthesize single video scenes depicting individual events.

We introduce TALC, a novel and effective framework to generate multi-scene videos from diffusion T2V generative models based on the scene descriptions. Our approach focuses on the role of text conditioning mechanism that is widely used in the modern T2V generative models (§2.2). Specifically, we take inspiration from the fact that the parts of the generated video  $x_j$  should depict the events described in the scene description  $y_j$ . To achieve this, we ensure that the representations for the part of the generated video aggregates language features from the scene description  $y_j$ .

Consider that we want to generate a multi-scene video  $X \in \mathbb{R}^{L \times 3 \times H \times W}$  from the scene descriptions  $y_j \in Y$ , using a T2V generative model  $f_\theta$ . Furthermore, we assume that individual video segments  $x_j$  are allocated  $L/n$  frames within the entire video  $X$ . Let  $z_X = [z_{x_1}; z_{x_2}; \dots; z_{x_n}] \in \mathbb{R}^{L \times C \times H' \times W'}$  represent the representation for the entire video  $X$ , and  $z_{x_j} \in \mathbb{R}^{(L/n) \times C \times H' \times W'}$  for the  $j^{\text{th}}$  part of the video that are concatenated in the temporal dimension. In addition, consider  $r_Y = \{r_{y_1}, \dots, r_{y_n}\}$  be the set of text embeddings for the multi-scene description  $Y$  and  $y_j$  be an individual scene description. In the TALC framework, the Eq. 2 is changed to:

$$z'_{\tau, x_j} = CA_{f_\theta}(Q = z_{\tau, x_j}, K = r_{y_j}, V = r_{y_j}) \quad (3)$$

$$z'_{\tau, X} = [z'_{x_1}; z'_{x_2}; \dots; z'_{x_n}] \quad (4)$$

Here,  $\tau$  represents the timestamp in the diffusion modeling setup, which is applied during training as well as inference. We illustrate the framework in Figure 3. While TALC aims to equip the generative model with the ability to depict all the events in the multi-scene descriptions, the visual consistency

is ensured by the temporal modules (attentions and convolution blocks) in the denoiser network. By design, our approach can be applied to the pretrained T2V model during inference.

### 3.2 Baselines

Here, we describe the baseline methods that could be used to generate videos for the multi-scene descriptions from a given diffusion text-to-video generative model.

#### 3.2.1 Merging Captions

In this setup, we create a single caption by merging all the multi-scene descriptions. Specifically, the multi-scene descriptions  $Y = \{y_1, y_2, \dots, y_n\}$  can be written as a single prompt  $\mathcal{P} = y_1.$ Then,  $y_2.$ ... Then,  $y_n.$  For instance, the two-scene description  $Y = \{‘A red panda climbs on a bamboo forest.’, ‘The red panda sleeps peacefully in the treetop.’\}$  will change to  $\mathcal{P} = ‘A red panda climbs on a bamboo forest. Then, the red panda sleeps peacefully in the treetop.’$  Subsequently, we generate a video from the T2V model  $f_\theta$  by conditioning it on  $\mathcal{P}$ . While this approach mentions the temporal sequence of the events in a single prompt, the T2V model does not understand the temporal boundaries between the two events. Specifically, the Eq. 2 suggests that the visual features for all the video frames will aggregate information from the entire multi-scene description, at once, without any knowledge about the alignment between the scene description and its expected appearance in the generated video.

#### 3.2.2 Merging Videos

In this setup, we generate videos for each scene description individually and merge them in the raw input space. Formally, the individual scene description  $y_i$  conditions the T2V model  $f_\theta$  to generate the parts of the multi-video  $x_i$ . Finally, we stitch the individual videos together to synthesize the entire video  $X = x_1, x_2, \dots, x_n$ . In this process, the parts of the multi-scene video closely adhere to the scene descriptions, leading to high text fidelity. However, since the generated videos do not have access to all the multi-scene descriptions (e.g., the video for Scene 2 is not informed about Scene 1), the visual consistency across the entire video is quite poor.

### 3.3 Multi-Scene Video-Text Data Generation

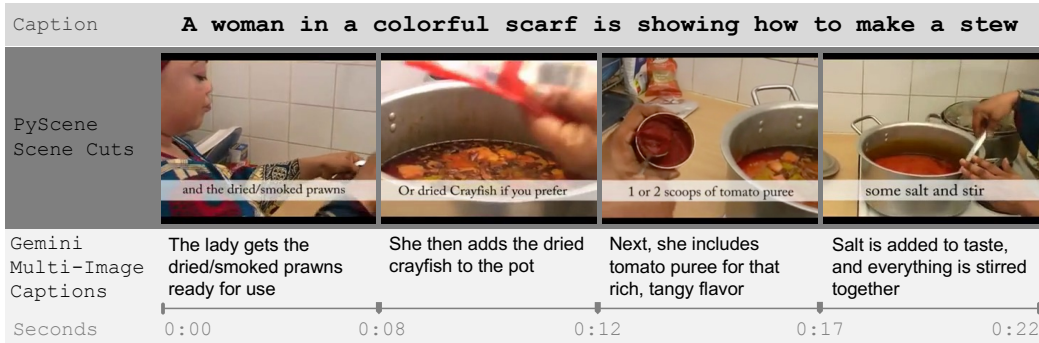


Figure 4: **Our approach for generating time-aligned video captions.** The process begins with PyScene cuts identifying the boundaries of distinct scenes within a video. Keyframes are then selected from the median of each scene. These frames are processed collectively through the Gemini model to produce multi-image captions that maintain narrative continuity by contextualizing each scene within the video’s overall sequence.

While our approach generates better multi-scene videos, the text adherence capabilities of the pretrained T2V generative model are limited. This is due to the lack of multi-scene video-text data during its pretraining. Unlike single video-text datasets, the multi-scene video-text datasets are not widely available and are hard to curate for model training. This is attributed to the fact that high-quality caption generation requires a lot of human labor which is time-consuming and expensive. Prior work such as ActivityNet [39] has curated human captions for specific video scenes

depicting useful actions in long videos. However, the video scenes are either overlapping or have a large temporal gap between them that will be harmful for natural and smooth variations between the generated multi-scene videos. Hence, the absence of high-quality captions for continuous video scenes in the dataset makes unsuitable for T2V generative training.

To this end, we aim to create a real-world multi-scene video-text dataset to allow further training of the pretrained T2V models. Specifically, we leverage the capability of the multimodal foundation model, Gemini-Pro-Vision [40], to generate high-quality synthetic data for enhanced video-text training [41]. Formally, we start with a video-text dataset  $\mathcal{M} = \mathcal{A} \times \mathcal{B}$  consisting of pairs of  $(A_i, B_i)$  where  $A_i$  is a raw video and  $B_i$  is the corresponding video description from the dataset. Subsequently, we utilize PySceneDetect library <sup>1</sup> to generate continuous video scenes from  $A_i = \{A_{i,1}, A_{i,2}, \dots, A_{i,m}\}$  where  $m$  is the number of scene cuts in the video. A similar approach was used in a prior work [12] to detect scene changes in the video data. Then, we sample the middle video frame  $F_{i,j}$  as a representative of the semantic content in the video scene  $A_{i,j}$ . Finally, we input all the video frames  $F_i = \{F_{i,1}, \dots, F_{i,m}\}$  for a single video  $A_i$  and the entire video caption  $B_i$  to a large multimodal model [40]. Specifically, the model is prompted to generate high-quality captions for each of the frames  $F_{i,j}$  such they form a coherent narrative guided by the common caption  $B_i$ . We provide the prompt provided to the multimodal model in Appendix §A. In Figure 4 we provide an instance for the multi-scene video-text data generation.

**Datasets.** To construct a multi-scene video-text dataset, we utilize existing dataset that include natural (real) videos and associated high-quality human-written captions that summarize the entire video. Specifically, we choose MSR-VTT [42] and VaTeX [43]. Most of the videos in MSR-VTT are 10-30 seconds long while VaTeX consists 10 seconds long videos. In addition, each video in MSR-VTT and VaTeX consists 20 captions and 10 captions, respectively, out of which one is selected at random for multi-scene data generation. As described above, a single video is cut into multiple video segments using Pyscene library. In our experiments, we retain the first four video segments and discard any additional segments if the library generates more than four. Since the accuracy of the multi-scene captioning and the computational demands during finetuning are influenced by the number of scenes, we opt to limit the scene count to four for our experiments. However, future work could employ similar methodologies to scale the number of scenes, given more computing power and advanced multi-scene captioning models. We provide the data statistics for the final multi-scene data in Appendix §G.

## 4 Evaluation

In this section, we describe the evaluation scheme for videos generated from multi-scene text descriptions. First, we describe the evaluation metrics that we aim to assess in this work (§4.1). Then, we generate multi-scene descriptions for a diverse set of tasks (§4.2). Finally, we present the details for automatic and human evaluation of the generated videos (§4.3).

### 4.1 Metrics

The ability to assess the quality of the generated multi-scene videos is a challenging task itself. As humans, we can judge the multi-scene videos across diverse perceptual dimensions [44] that the existing automatic methods often fails to capture [45]. Following [28], we focus on the visual consistency of the generated video, text adherence capabilities of the T2V models, and video quality of the video. Here, we present the metrics with the aspects that they intend to assess in the generated video for multi-scene text description.

**Visual Consistency.** This metric aims to assess the (entity or background) consistency between the frames of the multi-scene videos. Here, the **entity consistency** aims to test whether the entities in the multi-scene video are consistent across the video frames. For instance, the appearance of an animal should not change without a change described in the text description. In addition, the **background consistency** aims to test whether the background of the multi-scene video remains consistent across the video frames. For instance, the room should not change without a change description in the text.

<sup>1</sup><https://github.com/Breakthrough/PySceneDetect>

**Text Adherence.** This metric aims to test whether the generated video adheres to the multi-scene text description. For instance, the events and actions described in the text script should be presented in the video accurately, and in the correct temporal order.

In our experiments, we compute the visual consistency and text adherence with the automatic and human evaluators. Further, we compute the overall score, which is the average of the visual consistency and text adherence scores. In addition, we also assess the visual quality of the generated videos using human evaluation to understand whether the video contains any flimsy frames, shaky images, or undesirable artifacts (Table 1).

## 4.2 Task Prompts

Here, we curate a set of task prompts for diverse scenarios, aiming to holistically assess the quality of the generated videos.

**Single character in multiple visual contexts (S1).** In this scenario, we instruct an LLM, GPT-4, to create a coherent script consisting of four scenes. Each scene features a specific animal character performing diverse activities in every scene. This task assesses the capability of the T2V model to generate consistent appearance of the entity and its background while adhering to the different actions (or events) described in the multi-scene text script. For instance, a generated script could be ‘Scene 1: A red panda is climbing a tree. Scene 2: The red panda eats the leaves on the tree. Scene 3: The red panda lies down on the branch of the tree. Scene 4: The red panda sleeps on the branch’. In total, we generate 100 prompts in this scenario.

**Different characters in a specific visual context (S2).** In this scenario, we instruct a language model, GPT-4, to create a coherent script consisting of four scenes. Each scene features different animal characters engaging in the same activity in every scene [20]. This task assesses the capability of the T2V model to generate consistent appearance of the background while adhering to the appearance of the different characters in the multi-scene text script. For instance, a generated script could be ‘Scene 1: A cat leaps onto countertop. Scene 2: A dog leaps onto the same countertop. Scene 3: A rabbit leaps onto the same countertop. Scene 4: A raccoon leaps onto the same countertop’. In total, we generate 100 prompts in this scenario.

**Multi-scene captions from real videos (S3).** Here, we aim to assess the ability of the model to generate multi-scene videos for open-ended prompts that are derived from real-world videos. This task also assesses the ability of the T2V model to generate consistent appearances of the entity and its background while adhering to multi-scene descriptions. Specifically, we use our multi-scene video-text data generation pipeline (§3.3) to create such prompts for the real videos from the test splits of the video-text datasets. For example, a multi-scene text script could be ‘Scene 1: A beauty vlogger introduces her skincare routine. Scene 2: She applies a serum to her face, smoothing it in’. We present a sample of the various task prompts in the Appendix §B. In total, we generate 100 prompts in this scenario.

## 4.3 Evaluator

In this work, we devise an automatic evaluation framework and perform human evaluation to assess the quality of the multi-scene generated videos.

**Automatic Evaluation.** Here, we utilize the capability of a large multimodal model, GPT-4-Vision [46], to reason over multiple image sequences. First, we sample four video frames, uniformly, from each scene in the generated video (e.g., 8 videos frames for two-scene video). Then, we prompt the multimodal model with the temporal sequence of video frames from different scenes and the multi-scene text description. Specifically, we instruct the multimodal model to decide the quality of the generated video across various metrics including entity consistency, background consistency, and text adherence. For each metric, the multimodal model assigns one of three possible response  $\{yes = 1, partial = 0.5, no = 0\}$ . For instance, *yes* for the entity consistency metric implies that the video frames sampled from the generated video have consistent appearance of the entity described in the multi-scene script. In this work, we do not utilize any existing video-text alignment models [47, 41] for evaluating text adherence as they are trained on single-scene video-text datasets. We present the automatic evaluation prompt in Appendix §C.



**Human Evaluation.** We also conduct a human evaluation to assess the multi-scene generated videos along the dimensions of visual consistency, text adherence, and visual quality. Specifically, we ask the annotators from Amazon Mechanical Turk (AMT) to choose one of three options for each metric  $\{yes, partial, no\}$ , similar to the automatic evaluation. In addition, we choose the annotators that pass a preliminary qualification exam. We present the screenshot of the UI in Appendix §D.

#### 4.4 Evaluation Setup

Since merging captions (§3.2) and TALC (§3.1) methods input the entire multi-scene text description at once, the quality of the video generated by these methods is influenced by the number of scenes described in the text script. Hence, we calculate the performance of the baselines and TALC by averaging the scores assigned to videos generated for two, three, and four scenes. Additionally, we report on visual consistency by averaging the performance across the entity and background consistency metrics. Here, the entity consistency scores are calculated for the task prompts S1 and S3 (since S2 aims to change the characters across scenes), and the background consistency and text adherence scores are computed for all the task prompts. We also evaluate the impact of TALC-based finetuning on the single scene generation in Appendix §I.

## 5 Experiments

### 5.1 Text-to-Video Generative Models

In this work, we utilize ModelScope [10] and Lumiere [6] T2V models for multi-scene video generation. Here, ModelScope is an open-source T2V model with 1.7 billion parameters including the video encoder, text encoder, and denoising U-net network. Specifically, it is trained to generate 16 video frames on the mix of WebVid [14] video-text dataset and LAION [48] image-text dataset. We perform most of our experiments on ModelScope due to its easy-of-access and adoption in prior works [28]. In addition, we also include Lumiere-T2V, a model that leverages space-time U-Net denoising networks to generate high-quality videos. In this work, we include early experiments with Lumiere to showcase the flexibility of the TALC approach for multi-scene video generation.

**Base model with TALC.** As described in §3.1, our approach modifies the traditional text-conditioning mechanism to be aware of the alignment between text descriptions and individual video scenes. By design, the TALC framework can be applied to the base T2V model during inference, without any multi-scene finetuning. Thus, we compare the performance of the multi-scene videos generated from ModelScope and Lumiere T2V base models under three settings: merging captions, merging videos, and TALC. In this setting, we generate 16 frames per scene from ModelScope and 80 frames per scene from Lumiere. We provide more details on the inference in Appendix §F.

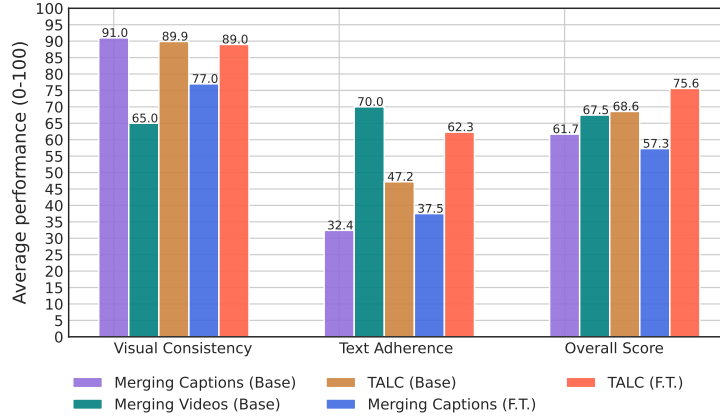
**Finetuning with TALC.** Since the base model is pretrained with single-scene data, we aim to show the usefulness of TALC framework when we have access to the multi-scene video-text data. To this end, we finetune ModelScope on the multi-scene video-text data (§3.3) with TALC framework. As a pertinent baseline, we also finetune the ModelScope without TALC framework by naively merging the scene-specific captions in the raw text space. In this setting, we finetune the T2V model with 8 frames per scene and the maximum number of scenes in an instance is set to 4. We provide further details on the finetuning setup in Appendix §H. The inference settings are identical to the prior method of generating videos from the base model without finetuning.

In this section, we present the results for the baselines and TALC framework averaged over a diverse task prompts and multiple scenes using automatic evaluation (§5.2) and human evaluation (§5.3). Finally, we provide qualitative examples for the multi-scene generated videos to showcase the usefulness of our approach (§5.4).

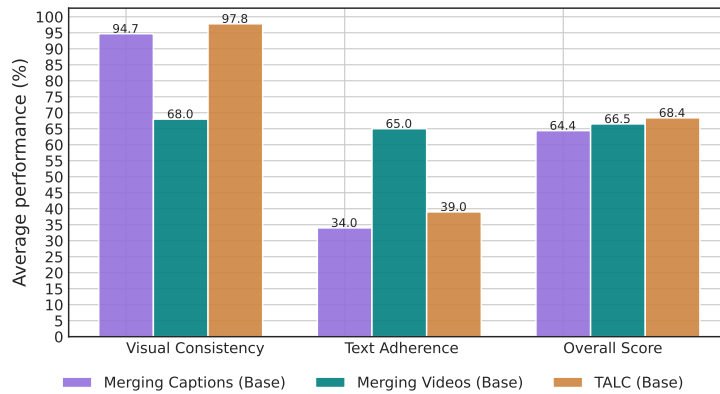
### 5.2 Automatic Evaluation

We compare the performance of the baselines (e.g., merging captions and merging videos) with the TALC framework for ModelScope and Lumiere using the automatic evaluation in Figure 5.

**TALC outperforms the baselines without any finetuning.** In Figure 5(a), we find that the overall score, average of visual consistency and text adherence, of the multi-scene videos generated using



(a) Performance on ModelScope T2V model.



(b) Performance on Lumiere T2V model.

Figure 5: **Automatic evaluation results for (a) ModelScope and (b) Lumiere.** In (a), we observe that TALC-finetuned ModelScope model achieves the highest overall score, that is the average of the visual consistency and text adherence scores. In (b), we find that TALC framework with the Lumiere base model outperforms merging captions and merging videos on the overall scores. We report the average performance across the diverse multi-scene prompts and the number of generated scenes.

the base ModelScope with TALC (68.6 points), outperforms the overall score achieved by the videos generated using merging captions (61.7 points) and merging videos (67.5 points) with the base ModelScope. Specifically, we observe that the visual consistency of the generated video is high for merging captions (91 points) and TALC (89.9 points) while it is low for merging videos (65 points). This indicates that merging videos independently for the individual scene descriptions does not preserve the background and entity appearances across the different frames. In addition, we observe that the text adherence using TALC outperforms merging captions by 14.8 points, while the text adherence is the highest with a score of 70 points using merging videos. This can be attributed to the design of the merging videos baseline where individual video scenes adhere to the scene-specific descriptions well. Hence, merging videos independently approach can be viewed as an upper bound on the text adherence metric.

In Figure 5(b), we observe similar trends for the Lumiere T2V generative model. Specifically, we find that the overall score for TALC outperforms merging captions and merging videos by 4 points and 2 points, respectively. In addition, we observe that merging captions and TALC achieve a high visual consistency score while merging videos independently has poor visual consistency. Further, we find that TALC outperforms merging captions by 5 points on text adherence, while merging videos achieves the highest text adherence 65 points. This highlights that the model more easily generates

multi-scene videos that adhere to individual text scripts, whereas adherence to the text diminishes when the model is given descriptions of multiple scenes all at once.

**Finetuning with TALC achieves the best performance.** Earlier, we evaluated the usefulness of the TALC framework with the base model. However, the base models are trained with the single-scene video-text data that might limit their capability for multi-scene video generation. To alleviate this issue, we finetune ModelScope T2V model on the multi-scene video-text data (§3.3). Specifically, we finetune the model using the merging captions method and TALC framework, independently.

In Figure 5(a), we find that finetuning with TALC achieves the highest overall score of 75.6 points in comparison to all the baselines. Specifically, we observe that the visual consistency does not change much with finetuning using the TALC method (89.9 points vs 89 points). Interestingly, we observe that finetuning with merging captions reduces the visual consistency by a large margin of 14 points. This can be attributed to the lack of knowledge about the natural alignment between video scenes and individual scene descriptions, which gets lost during the merging of captions.

Additionally, we find that the text adherence of the TALC-finetuned model is 15.1 points more than the text adherence of the TALC-base model. Similarly, we find that the text adherence of the merging captions-finetuned model is 5.1 points more than the text adherence of the merging captions-base model. This highlights that finetuning a T2V model with multi-scene video-text data helps the most with enhancing its text adherence capability.

**Fine-grained Results.** To perform fine-grained analysis of the performance, we assess the visual consistency and text adherence scores for the baselines and TALC framework across diverse task prompts and number of scenes on ModelScope. We present their results in Appendix §E. In our analysis, we find that finetuning with TALC achieves the highest overall score over the baselines across all the scenarios. In addition, we notice that the highest performance is achieved in the scenario that consist of the different entities in a specific visual context. Further, we observe that the performance of the all the methods reduces when the task prompts get more complex i.e., multi-scene captions from real videos. In addition, we observe that finetuning with TALC achieves the highest overall score over the baselines across all the number of scenes. Specifically, we observe that the performance of the merging captions and TALC framework reduces as the number of scenes being generated increases. Overall, we show that the TALC strikes a good balance between visual consistency and text adherence to generate high-quality multi-scene videos.

### 5.3 Human Evaluation

**TALC achieves the best performance in human evaluation.** We compare the performance of the baselines and TALC framework for ModelScope using human evaluation in Figure 6. We find that TALC-finetuned model outperforms the merging captions and merging video methods with the base model by 12 points and 15.5 points, respectively, on the overall score. In addition, we find that using TALC framework in the base model outperforms the merging captions and merging video methods with the base model by 7.6 points and 11.1 points, respectively, on the overall score. Further, we observe that the merging captions with the base model achieves the highest visual consistency score of 96.5 points while it is the lowest for merging videos generated from the base model. In addition, we find that the text adherence of the TALC-finetuned and TALC-base model is better than merging captions-finetuned and merging captions-base model, respectively. Our results highlight at the benefit of including the inductive bias of temporal alignment between the video scenes and their scene descriptions for multi-scene video generation.

Table 1: **Human evaluation results on the visual quality of the generated videos from ModelScope.** We observe that the visual quality of the generated videos are close to each other for the base model. However, finetuning the model with merging captions reduces the video quality by a large margin while TALC-finetuned model retains the video quality.

Method	Quality
Merging Captions (Base)	80.5
Merging Videos (Base)	86.5
TALC (Base)	84.5
Merging Captions (F.T.)	63.4
TALC (F.T.)	83.3

**Visual quality of the generated videos.** We compare the visual quality of the generated videos using human evaluation in Table 1. We find that the visual quality of videos generated from the base

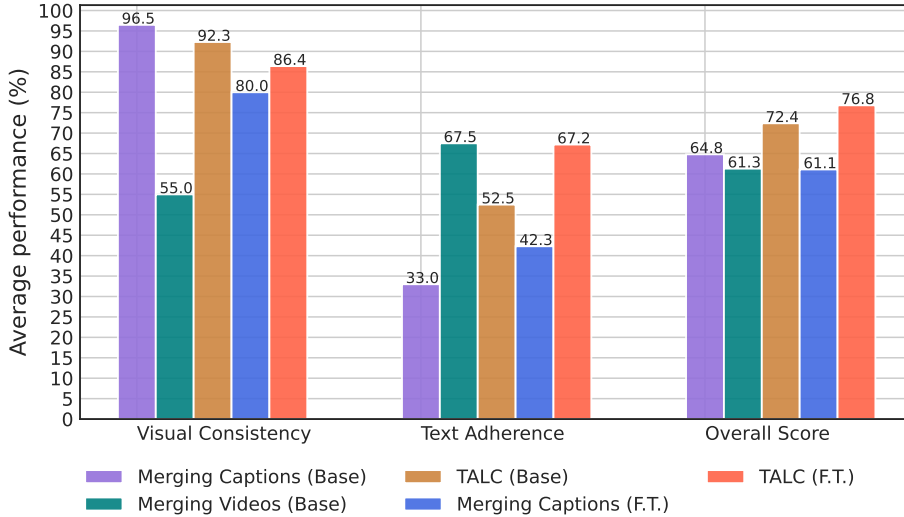


Figure 6: **Human evaluation results for ModelScope model.** We observe that the base model using the TALC framework outperforms the merging captions and merging videos baselines on the overall score. In addition, TALC-finetuned model enhances the text adherence and achieves the highest overall score. We report the average performance across the diverse multi-scene prompts and the number of generated scenes.

model ranges from 80.5 – 86.5 using the baselines and TALC framework. However, we observe that the visual quality of generated videos is quite poor for the model finetuned with merging captions with a score of 63.4 points. This highlights that finetuning a T2V model with multi-scene video-text data by naively merging the scene-specific descriptions in the raw text space leads to undesirable artifacts in the generated video. Finally, we find that the TALC-finetuned model (83.3) achieves a video quality score similar to that of the TALC-base model (84.5), indicating that our finetuning data preserves the visual quality observed during the model’s pretraining. While our work is centered around multi-scene evaluation, we also perform single-scene evaluation in Appendix §I.

#### 5.4 Qualitative Analysis

We provide qualitative examples of generating multi-scene videos using ModelScope (fine-tuned with TALC) and Lumiere (base model with TALC) for diverse scenarios in Figure 12. Our analysis reveals that both ModelScope and Lumiere are capable of producing multi-scene videos that exhibit high text adherence and visual consistency.

Considering the case of the same animal engaging in multiple actions (referred to as "one character multiple contexts"). The videos generated by ModelScope successfully maintained the same animal while varying the background and action between the scenes. Conversely, the videos generated by Lumiere displayed the same animal performing different actions with minimal background alterations. We believe that this distinction is attributed to ModelScope’s fine-tuning with TALC. Considering different animals within a particular visual setting (referred to as "multiple-characters same context"), both ModelScope and Lumiere demonstrated impressive abilities in preserving the consistency of the background across the videos and adhering closely to the provided text. During our analysis, we noticed that the multi-scene captions derived from real videos (referred to as "open-ended captions") exhibited a substantial number of changes between the various scenes. In this scenario, Lumiere, when employed without fine-tuning, displayed challenges in adhering to the text, while ModelScope achieved a higher degree of text adherence but was also prone to visual artifacts.

### 6 Related Work

**Text-to-Video Generative Modeling.** The field of text-to-video (T2V) synthesis has significantly evolved from its inception with models like VGAN [2] and MoCoGAN [49], leveraging the foun-

dational technologies of GANs [50] and VAEs [51] to produce concise, single-scene videos. The narrative depth was further expanded through transformer-based architectures such as CogVideo [52] and VideoGPT [53], enhancing the complexity of video content yet remaining within the confines of single scenes. The advent of diffusion models, exemplified by Imagen Video [54], marked a notable advancement in T2V synthesis. Despite these strides, the challenge of creating multi-scene videos that reflect the complexity of the physical world [1, 2, 3] remains. Our work, TALC, extends the capabilities of T2V models to multi-scene storytelling, filling a crucial gap in the synthesis landscape.

**Image-to-Video Animation.** The exploration of multi-scene video generation, innovative methods such as Lumiere [6] and Make-a-Video [55] have employed a two-step process, transforming text to images and then animating these images into videos. While these approaches have advanced visual quality, they often fall short in weaving seamless multi-scene narratives. This limitation is echoed in the work of Emu Video [8], which underscores the difficulty of achieving narrative coherence across multiple scenes. TALC focuses on direct generation of multi-scene narratives from textual prompts aiming for a narrative flow and visual consistency across scenes.

**Multi-Scene Video Generation.** The pursuit of multi-scene T2V synthesis has been furthered by recent innovations like Phenaki [20] and Stable Video Diffusion [12], which have explored new frontiers in video generation from textual prompts and the scaling of latent diffusion models, respectively. Additionally, Dreamix [56] and Pix2Video [57] have broadened the scope of diffusion models, applying them to video editing and animation. Despite these advancements, the task of generating videos that convey coherent narratives across multiple scenes remains formidable, highlighted by recent works such as VideoPoet [9], ModelScope [10] and Make-A-Scene [58]. TALC tackles this task and offers a framework produces videos spanning multiple scenes. We also introduce nuanced evaluation approach. This approach integrates both automated assessments and human evaluations to rigorously gauge the quality and narrative coherence of the generated content, evaluating text adherence, object consistency and background consistency, contributing to the ongoing refinement of T2V synthesis.

## 7 Conclusion

We introduced TALC, a simple and effective method for improving the text-to-video (T2V) models for multi-scene generation. Specifically, it incorporates the knowledge of the natural alignment between the video segments and the scene-specific descriptions. Further, we show that TALC-finetuned T2V model achieve high visual consistency and text adherence while the baselines suffer from one or both of the metrics. Given its design, our framework can be easily adapted into any diffusion-based T2V model. An important future direction will be to scale the amount of multi-scene video-text data and deploy TALC framework during pretraining of the T2V models.

## 8 Acknowledgement

We would like to thank Ashima Suvarna for providing feedback on the draft. Hritik Bansal is supported in part by AFOSR MURI grant FA9550-22-1-0380.

## References

- [1] Dinesh Acharya, Zhiwu Huang, Danda Pani Paudel, and Luc Van Gool. Towards high resolution video generation with progressive growing of sliced wasserstein gans. *arXiv preprint arXiv:1810.02419*, 2018.
- [2] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. *Advances in neural information processing systems*, 29, 2016.
- [3] Jake Bruce, Michael Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative interactive environments. *arXiv preprint arXiv:2402.15391*, 2024.
- [4] OpenAI. Video generation models as world simulators. <https://openai.com/research/video-generation-models-as-world-simulators>, 2024.

- [5] T Brooks, B Peebles, C Homes, W DePue, Y Guo, L Jing, D Schnurr, J Taylor, T Luhman, E Luhman, et al. Video generation models as world simulators, 2024.
- [6] Omer Bar-Tal, Hila Chefer, Omer Tov, Charles Herrmann, Roni Paiss, Shiran Zada, Ariel Ephrat, Junhwa Hur, Yuanzhen Li, Tomer Michaeli, et al. Lumiere: A space-time diffusion model for video generation. *arXiv preprint arXiv:2401.12945*, 2024.
- [7] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
- [8] Rohit Girdhar, Mannat Singh, Andrew Brown, Quentin Duval, Samaneh Azadi, Sai Saketh Rambhatla, Akbar Shah, Xi Yin, Devi Parikh, and Ishan Misra. Emu video: Factorizing text-to-video generation by explicit image conditioning. *arXiv preprint arXiv:2311.10709*, 2023.
- [9] Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, Jonathan Huang, Rachel Hornung, Hartwig Adam, Hassan Akbari, Yair Alon, Vighnesh Birodkar, et al. Videopoet: A large language model for zero-shot video generation. *arXiv preprint arXiv:2312.14125*, 2023.
- [10] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023.
- [11] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22563–22575, 2023.
- [12] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [14] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *IEEE International Conference on Computer Vision*, 2021.
- [15] Hongwei Xue, Tiankai Hang, Yanhong Zeng, Yuchong Sun, Bei Liu, Huan Yang, Jianlong Fu, and Baining Guo. Advancing high-resolution video-language representation with large-scale video transcriptions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5036–5045, 2022.
- [16] Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, et al. Videocrafter1: Open diffusion models for high-quality video generation. *arXiv preprint arXiv:2310.19512*, 2023.
- [17] Haoxin Chen, Yong Zhang, Xiaodong Cun, Menghan Xia, Xintao Wang, Chao Weng, and Ying Shan. Videocrafter2: Overcoming data limitations for high-quality video diffusion models. *arXiv preprint arXiv:2401.09047*, 2024.
- [18] Zangwei Zheng, Xiangyu Peng, and Yang You. Open-sora: Democratizing efficient video production for all, March 2024. URL <https://github.com/hpcaitech/Open-Sora>.
- [19] Wenjing Wang, Huan Yang, Zixi Tuo, Huiguo He, Junchen Zhu, Jianlong Fu, and Jiaying Liu. Videofactory: Swap attention in spatiotemporal diffusions for text-to-video generation. *arXiv preprint arXiv:2305.10874*, 2023.
- [20] 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 preprint arXiv:2210.02399*, 2022.

- [21] 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.
- [22] Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.
- [23] 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10459–10469, 2023.
- [24] Tanzila Rahman, Hsin-Ying Lee, Jian Ren, Sergey Tulyakov, Shweta Mahajan, and Leonid Sigal. Make-a-story: Visual memory conditioned consistent story generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2493–2502, 2023.
- [25] 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*, pages 10684–10695, 2022.
- [26] Yitong Li, Zhe Gan, Yelong Shen, Jingjing Liu, Yu Cheng, Yuexin Wu, Lawrence Carin, David Carlson, and Jianfeng Gao. Storygan: A sequential conditional gan for story visualization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6329–6338, 2019.
- [27] Shengming Yin, Chenfei Wu, Huan Yang, Jianfeng Wang, Xiaodong Wang, Minheng Ni, Zhengyuan Yang, Linjie Li, Shuguang Liu, Fan Yang, et al. Nuwa-xl: Diffusion over diffusion for extremely long video generation. *arXiv preprint arXiv:2303.12346*, 2023.
- [28] Han Lin, Abhay Zala, Jaemin Cho, and Mohit Bansal. Videodirectorgpt: Consistent multi-scene video generation via llm-guided planning. *arXiv preprint arXiv:2309.15091*, 2023.
- [29] Fuchen Long, Zhaofan Qiu, Ting Yao, and Tao Mei. Videodrafter: Content-consistent multi-scene video generation with llm. *arXiv preprint arXiv:2401.01256*, 2024.
- [30] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [31] Calvin Luo. Understanding diffusion models: A unified perspective. *arXiv preprint arXiv:2208.11970*, 2022.
- [32] 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*, pages 8821–8831. Pmlr, 2021.
- [33] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [34] Ting Chen. On the importance of noise scheduling for diffusion models. *arXiv preprint arXiv:2301.10972*, 2023.
- [35] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015.
- [36] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.

- [37] 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*, pages 8748–8763. PMLR, 2021.
- [38] Abhay Zala, Jaemin Cho, Satwik Kottur, Xilun Chen, Barlas Oguz, Yashar Mehdad, and Mohit Bansal. Hierarchical video-moment retrieval and step-captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23056–23065, 2023.
- [39] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 961–970, 2015.
- [40] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [41] Hritik Bansal, Yonatan Bitton, Idan Szepktor, Kai-Wei Chang, and Aditya Grover. Videocon: Robust video-language alignment via contrast captions. *arXiv preprint arXiv:2311.10111*, 2023.
- [42] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5288–5296, 2016.
- [43] Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. Vatex: A large-scale, high-quality multilingual dataset for video-and-language research. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4581–4591, 2019.
- [44] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. *arXiv preprint arXiv:2311.17982*, 2023.
- [45] Emanuele Bugliarello, H Hernan Moraldo, Ruben Villegas, Mohammad Babaeizadeh, Mohammad Taghi Saffar, Han Zhang, Dumitru Erhan, Vittorio Ferrari, Pieter-Jan Kindermans, and Paul Voigtlaender. Storybench: A multifaceted benchmark for continuous story visualization. *Advances in Neural Information Processing Systems*, 36, 2024.
- [46] OpenAI. Gpt-4v(ision) system card, 2023b. <https://openai.com/research/gpt-4v-system-card>, 2023.
- [47] Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metze, Luke Zettlemoyer, and Christoph Feichtenhofer. Videoclip: Contrastive pre-training for zero-shot video-text understanding. *arXiv preprint arXiv:2109.14084*, 2021.
- [48] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022.
- [49] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. Mocogan: Decomposing motion and content for video generation. In *CVPR*, pages 1526–1535, 2018.
- [50] Antonia Creswell, Tom White, Vincent Dumoulin, Kai Arulkumaran, Biswa Sengupta, and Anil A Bharath. Generative adversarial networks: An overview. *IEEE signal processing magazine*, 35(1):53–65, 2018.
- [51] Carl Doersch. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*, 2016.
- [52] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022.



- [53] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. Videogpt: Video generation using vq-vae and transformers. *arXiv preprint arXiv:2104.10157*, 2021.
- [54] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
- [55] 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. In *ICLR*, 2023.
- [56] Eyal Molad, Eliahu Horwitz, Dani Valevski, Alex Rav Acha, Yossi Matias, Yael Pritch, Yaniv Leviathan, and Yedid Hoshen. Dreamix: Video diffusion models are general video editors. *arXiv preprint arXiv:2302.01329*, 2023.
- [57] Duygu Ceylan, Chun-Hao P Huang, and Niloy J Mitra. Pix2video: Video editing using image diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 23206–23217, 2023.
- [58] 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 *European Conference on Computer Vision*, pages 89–106. Springer, 2022.
- [59] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

## A Prompt for Multi-Scene Caption Generation

We present the prompt used to generate multi-scene captions using large multimodal model, Gemini-Pro-Vision, in Figure 7. In particular, we utilize Gemini-Pro-Vision since it can reason over multiple image sequences. Specifically, we provide the multimodal model with a video frame from each of the segmented videos and the single caption for the entire video present in the original video-text datasets.

---

Your task is to create captions for a series of images, each taken from different video scenes. For every image shown, craft a 7-10 word caption (7-10 words) that precisely describes what's visible, while also linking these captions into a fluid, engaging story. A common caption will be given to help guide your narrative, ensuring a smooth transition between scenes for a cohesive story flow. Remember to not hallucinate in your responses.

Common Caption: {caption}

---

Figure 7: Prompt to generate multi-scene caption using large multimodal models.

## B Task Prompts

### B.1 Single character in multiple visual contexts

We present the prompt used to generate multi-scene text descriptions for single character in multiple visual contexts from GPT-4 in Figure 8.

---

Create four concise continuous movie scenes (7-10 words) focusing on a specific real-world character. The scenes should form a cohesive narrative.

Guidelines:

Choice of Character: Select a real-world animal as the focal point of your scenes.

Scene Description: Clearly describe the setting, actions, and any notable elements in each scene.

Connection: Ensure that the scenes are logically connected, with the second scene following on from the first.

Brevity and Precision: Keep descriptions short yet vividly detailed.

Example:

Character: polar bear

Scene 1: A polar bear navigates through a icy landscape.

Scene 2: The polar bear hunts seals near a crack in the ice.

Scene 3: The polar bear feasts on the seal.

Scene 4: The polar bear curls up for a nap.

Now it's your turn.

---

Figure 8: GPT-4 Prompt to generate multi-scene prompts for single character in multiple visual contexts.

### B.2 Different characters in a specific visual context

We present the prompt used to generate multi-scene text descriptions for different characters in a specific visual context from GPT-4 in Figure 9.

## C Automatic Evaluation Prompt

We present the prompt used to perform automatic evaluation of the multi-scene generated videos using large multimodal model, GPT-4-Vision, in Figure 10. We utilize GPT-4-Vision for automatic evaluation since it can reason over multiple images. Specifically, we provide the multimodal model with four video frames for each scene in the generated video. The model has to provide its judgments based on the entity consistency, background consistency, and text adherence of the video frames.

---

Create four concise scene descriptions (7-10 words) where different characters perform identical action/events.

Choice of Characters: Select four real-world animals as the focal point of the individual scenes.

Background Consistency: Ensure that the background is consistent in both the scenes.

Brevity and Precision: Keep descriptions short yet vividly detailed.

Example:

Characters: teddy bear, panda, grizzly bear, polar bear

Scene 1: A teddy bear swims under water.

Scene 2: A panda swims under the same water.

Scene 3: A panda swims under the same water.

Scene 4: A panda swims under the same water.

Now it's your turn.

---

Figure 9: GPT-4 Prompt to generate multi-scene prompts for different characters in a specific visual context.

---

You are a capable video evaluator. You will be shown a text script with two-scene descriptions where the events/actions . Video generating AI models receive this text script as input and asked to generate relevant videos. You will be provided with eight video frames from the generated video. Your task is to answer the following questions for the generated video.

1. Entity Consistency: Throughout the video, are entities consistent? (e.g., clothes do not change without a change described in the text script)

2. Background Consistency: Throughout the video, is the background consistent? (e.g., the room does not change described in the text script)

3. Text Adherence: Does the video adhere to the script? (e.g., are events/actions described in the script shown in the video accurately and in the correct temporal order)

Respond with NO, PARTIALLY, and YES for each category at the end. Do not provide any additional explanations.

Two-scene descriptions:

Scene 1: {scene1}

Scene 2: {scene2}

---

Figure 10: Prompt used to perform automatic evaluation of the multi-scene generated videos. We use this prompt when the number of scenes in the task prompt is two.

## D Human Evaluation Screenshot

We present the screenshot for the human evaluation in Figure 11. Specifically, we ask the annotators to judge the visual quality, entity consistency, background consistency, and text adherence of the multi-scene generated videos across diverse task prompts and number of scenes.

## E Fine-grained Analysis

We present the automatic evaluation results across task prompts (§E.1) and number of scenes (§E.2).

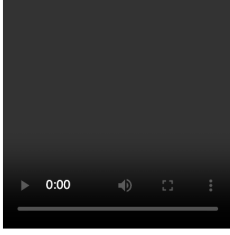
### E.1 Task Prompts

We compare the performance of the baselines and TALC framework across different task prompts in Table 2. We find that the TALC-finetuned model achieves the highest overall score over all the baselines. Specifically, we find that the TALC framework achieves a high visual consistency with scores close to the merging captions baseline. Further, we observe that the TALC framework achieves a higher text adherence in comparison to the merging captions, with or without finetuning, across all the task prompts.

### E.2 Number of Scenes

We compare the performance of the baselines and TALC framework across different number of generated scenes in Table 3. We find that the TALC-finetuned model outperforms all the baselines in this setup. In addition, we find that the visual consistency of the TALC framework does not change much with the number of the scenes. However, we find that the text adherence of the baselines and TALC framework reduces with the number of generated scenes. The text adherence scores of

Answer the following questions based on the multiple-scene descriptions and the candidate generated video.



Multi-Scene Descriptions:  $\$(multi\_scene\_description)$

Does the video exhibit a good **Visual Quality**? (are there any filmy frames, shaky images, and undesirable artifacts in the video?)  
 Yes  Partially  No

Does the video exhibit **Entity Consistency**? (entities are consistent e.g., the shape and features of the animal do not change unless specified)  
 Yes  Partially  No

Does the video exhibit **Background Consistency**? (background is consistent when required e.g., the room does not change without a change described in the scene description)  
 Yes  Partially  No

Does the video exhibit **Text Adherence**?  
 Yes  Partially  No

Submit

Figure 11: Human Annotation Layout

Table 2: **Automatic evaluation results across task prompts.** Here, S1 refers to the single character in multiple visual contexts. S2 refers to the different characters in a specific visual context. S3 refers to the multi-scene captions from real videos. We abbreviate Finetuning as F.T., Visual consistency as V.C., Text adherence as T.A.

Method	S1			S2			S3		
	V.C.	T.A.	Overall	V.C.	T.A.	Overall	V.C.	T.A.	Overall
Merging captions (Base)	95.9	43.8	69.9	99.5	21.5	60.5	81.9	32.0	56.9
Merging videos (Base)	69.4	71.5	70.5	71.2	82.3	76.7	57.9	58.7	58.3
TALC (Base)	93.6	48.3	71.0	99.3	57.3	78.3	81.4	36.1	58.8
Merging captions (F.T.)	94.2	57.1	75.7	94.9	31.2	63.1	50.8	24.3	37.5
TALC (F.T.)	93.3	62.6	<b>77.9</b>	98.9	76.3	<b>87.6</b>	79.7	47.9	<b>63.8</b>

the merging videos does not change with the number of scenes as it generates the videos for the individual scenes independently.

Table 3: **Automatic evaluation results across different number of scenes in the task prompts.** We abbreviate Finetuning as F.T., visual consistency as V.C., and Text Adherence as T.A.

Method	# scenes = 2			# scenes = 3			# scenes = 4		
	V.C.	T.A.	Overall	V.C.	T.A.	Overall	V.C.	T.A.	Overall
Merging captions (Base)	93.2	34.4	63.8	92.5	33.0	62.7	87.4	29.9	58.6
Merging videos (Base)	66.7	69.9	68.3	65.2	71.9	68.6	63.5	70.7	67.1
TALC (Base)	92.6	54.4	73.5	89.8	48.0	68.9	87.3	39.3	63.3
Merging captions (F.T.)	87.7	45.7	66.7	83.2	39.5	61.4	60.0	27.4	43.7
TALC (F.T.)	88.5	66.6	<b>77.5</b>	90.7	64.1	<b>77.4</b>	87.8	56.1	<b>71.9</b>

## F Inference Details

We provide the details for sampling multi-scene videos from the ModelScope and Lumiere T2V models in Table 4 and Table 5, respectively.

## G Multi-Scene Data Statistics

We provide the details for the multi-scene video-text dataset in Table 6.

Table 4: Sampling setup for ModelScope T2V model.

Resolution	256 × 256
Number of video frames per scene	16
Guidance scale	12
Sampling steps	100
Noise scheduler	DPMSolverMultiStepScheduler

Table 5: Sampling setup for Lumiere T2V model.

Resolution	1024 × 1024
Number of video frames per scene	80
Guidance scale	8
Sampling steps	256
Noise scheduler	DPMSolverMultiStepScheduler

## H Finetuning Details

We provide the details for finetuning ModelScope T2V model with TALC framework in Table 7.

## I Single-Scene Video Generation

To ascertain that our model’s new multi-scene generation function does not detract from its single-scene generation performance, we conducted a series of evaluations using the VBench framework [44]. VBench offers a robust analysis of various video generation aspects such as adherence to text prompts, stylistic integrity, semantic coherence, and overall aesthetic quality.

Our analysis, shown in Table 8, establishes a refined baseline: ModelScope (Single-Scene Finetuning), fine-tuned on single-scene video generation data. This process yielded an average score of 0.48, indicating a decrease from the base ModelScope’s average score of 0.63. This suggests that the optimizations in the base model, such as integrating high-quality images, are not fully utilized in single-scene fine-tuning.

Interestingly, fine-tuning the model on multi-scene data (ModelScope - Multi-Scene Finetuning) resulted in improved performance with an average score of 0.59, surpassing the single-scene fine-tuned version. This indicates that multi-scene data enriches the model’s understanding of video content, enhancing both multi-scene and single-scene video generation.

This comparison highlights the importance of data curation and fine-tuning strategies, showing that our approach not only enables complex multi-scene narratives but also improves single-scene video generation.

<sup>2</sup><https://huggingface.co/damo-vilab/text-to-video-ms-1.7b/tree/main>

Table 6: Multi Scene Video-Text Data Statistics

Number of entire videos	7107
% of single scene	27.3%
% of two scenes	25.4%
% of three scenes	31.3%
% of four scenes	16.0%
Number of video scene-caption instances	20177

Table 7: Training details for the TALC-finetuned ModelScope T2V model.

Base Model	ModelScope <sup>2</sup>
Trainable Module	UNet
Frozen Modules	Text Encoder, VAE
Batch size	20
Number of GPUs	5 Nvidia A6000
Resolution	256 × 256
Crop	CenterCrop
Learning Rate Scheduler	Constant
Peak LR	1.00E-05
Warmup steps	1000
Optimizer	Adam (0.9, 0.999, 1e-2, 1e-8) [59]
Max grad norm	1
precision	fp16
Noise scheduler	DDPM
Number of frames per video scene	8
prediction type	epsilon

Table 8: Single-Scene Evaluation Results using VBench, comparing the base model, when it is fine-tuned on multi-scene data, and single-scene data (‘f.t.’ stands for fine-tuned). Our analysis shows ModelScope (Single-Scene Finetuning) as a refined baseline with an average score of 0.48, compared to the base ModelScope’s 0.63. Fine-tuning on multi-scene data (ModelScope - Multi-Scene Finetuning) yields an improved score of 0.59, highlighting the efficacy of multi-scene data in enhancing video generation performance.

Dimension	ModelScope (Base)	ModelScope (Multi-scene f.t.)	ModelScope (Single-Scene f.t.)
Appearance style	0.23	0.24	0.21
Color	0.85	0.83	0.78
Human action	0.96	0.92	0.75
Object class	0.86	0.77	0.42
Overall consistency	0.26	0.26	0.22
Spatial relationship	0.35	0.29	0.14
Subject consistency	0.90	0.83	0.75
Temporal flickering	0.97	0.89	0.86
Temporal style	0.26	0.25	0.21
<b>Average</b>	0.63	0.59	0.48

**One character multiple contexts**



**Multiple character same context**



**Open-ended scenes**



Figure 12: Examples of videos generated by Time-Aligned Captions (TALC)