

# RACCoon: Remove, Add, and Change Video Content with Auto-Generated Narratives

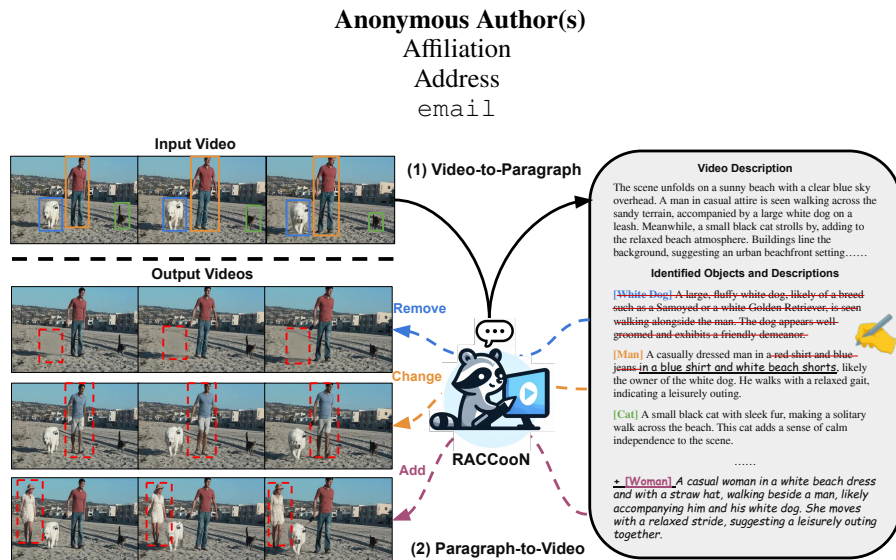


Figure 1: **Overview of RACCoon**, a versatile and user-friendly video-to-paragraph-to-video framework, enables users to remove, add, or change video content via updating auto-generated narratives.

## Abstract

1 This paper proposes **RACCoon**, a versatile and user-friendly **video-to-paragraph-**  
 2 **to-video** generative framework that supports multiple video editing capabilities  
 3 such as removal, addition, and modification, through a unified pipeline. The  
 4 proposed approach stands out from other methods through several significant  
 5 contributions: (1) suggests a multi-granular spatiotemporal pooling strategy to  
 6 generate well-structured video descriptions, capturing both the broad context  
 7 and object details without requiring complex human annotations, simplifying  
 8 precise video content editing based on text for users. (2) RACCoon incorporates  
 9 auto-generated narratives or instructions to enhance the quality and accuracy  
 10 of the generated content. It supports the addition of video objects, inpainting,  
 11 and attribute modification within a unified framework, surpassing existing video  
 12 editing/inpainting benchmarks by demonstrating impressive versatile capabilities  
 13 in video-to-paragraph generation (up to 9.4% $p$   $\uparrow$  absolute improvement in human  
 14 evaluations against the baseline), video content editing (relative 49.7%  $\downarrow$  in FVD).

## 15 1 Introduction

16 Despite advancements in recent video editing models, significant challenges remain in developing a  
 17 versatile and user-friendly framework that facilitates easy video modification for personal use. The  
 18 primary challenges include: 1) the complexity of training a unified framework encompassing multiple  
 19 video editing skills (e.g., remove, add, or change an object). Training a single model to perform  
 20 various editing skills is highly challenging, and recent video editing methods often focus on specific  
 21 tasks, such as background inpainting or attribute editing. 2) the necessity for well-structured textual

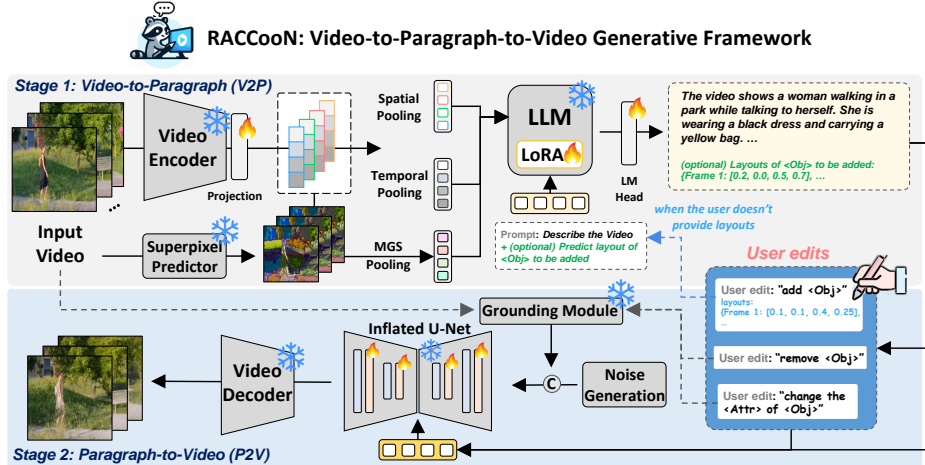


Figure 2: **Illustration of RACCOON framework.** RACCOON generates video descriptions with the three distinct pooled visual tokens, including Multi-Granular Spatiotemporal (MGS) Pooling. Next, users can edit the generated descriptions by adding, removing, or modifying words to create new videos. Note that for adding object tasks, if users do not provide layout information for the objects they want to add, RACCOON can predict the target layout in each frame.

22 prompts that accurately describe videos and can be edited to support diverse video editing skills.  
 23 The quality of prompts critically influences the models’ capabilities and the quality of their outputs.  
 24 Generating detailed prompts is time-consuming and costly, and the quality varies depending on the  
 25 expertise of the annotators. Although Multimodal Large Language Models (MLLMs) have been  
 26 explored for automatically describing videos, they often overlook critical details in complex scenes.  
 27 This oversight compromises the development of a seamless pipeline, hindering both user convenience  
 28 & the effectiveness of video generative models.

29 To tackle these limitations, we introduce **RACCOON**, a novel **video-to-paragraph-to-video (V2P2V)**  
 30 generative framework that facilitates diverse video editing capabilities based on auto-generated  
 31 narratives, as illustrated in Fig. 1. RACCOON allows for the seamless removal and modification  
 32 of subject attributes, as well as the addition of new objects to videos **without requiring densely**  
 33 **annotated video prompts or extensive user planning**. Our framework operates in two main stages:  
 34 *video-to-paragraph (V2P)* and *paragraph-to-video (P2V)* (Please see Fig. 2). In the V2P stage, we  
 35 introduce a new video descriptive framework built on a pre-trained Video-LLM backbone. We find  
 36 that existing Video-LLMs effectively capture holistic video features, yet often overlook detailed  
 37 cues that are critical for accurate video editing, as users may be interested in altering these missing  
 38 contexts. To address this, we propose a novel multi-granular video perception strategy that leverages  
 39 superpixels to capture diverse and informative localized contexts throughout a video. Next, in the P2V  
 40 stage, to integrate multiple editing capabilities into a single model, we fine-tuned a video inpainting  
 41 model that can paint video objects accurately with detailed text, object masks, and condition video.  
 42 Then, by utilizing user-modified prompts from generated descriptions in the V2P stage, our video  
 43 diffusion model can accurately *paint* corresponding video regions, ensuring that textual updates from  
 44 prompts are reflected in various editing tasks. Moreover, to better support our model training, we  
 45 have collected the **Video Paragraph with Localized Mask (VPLM)** dataset—a collection of over 7.2K  
 46 high-quality video-paragraph descriptions and 5.5k detailed object descriptions with corresponding  
 47 masks, annotated from the publicly available ROVIDataset using GPT-4V.

48 We emphasize that RACCOON enhances the quality and versatility of video editing by leveraging  
 49 detailed, automatically generated textual prompts that minimize ambiguity and refine the scope of  
 50 generation. We validate the extensive capabilities of the RACCOON framework in both V2P generation,  
 51 text-based video content editing, and video generation on ActivityNetYouCook2UCF101DAVIS and  
 52 our proposed VPLM datasets. On the V2P side, RACCOON outperforms several strong video  
 53 captioning baselines, particularly improving by average **+9.1%p** on VPLM and up to **+9.4%p** on  
 54 YouCook2 compared to PG-VL, based on both automatic metrics and human evaluation, as shown  
 55 in Tabs. 1 and 2. On the P2V side, RACCOON surpasses previous strong video editing/inpainting  
 56 baselines over three subtasks of video content editing (remove, add, and change video objects) over 9  
 57 metrics. Detailed results and visualizations are in Tab. 3 and Fig. 3.

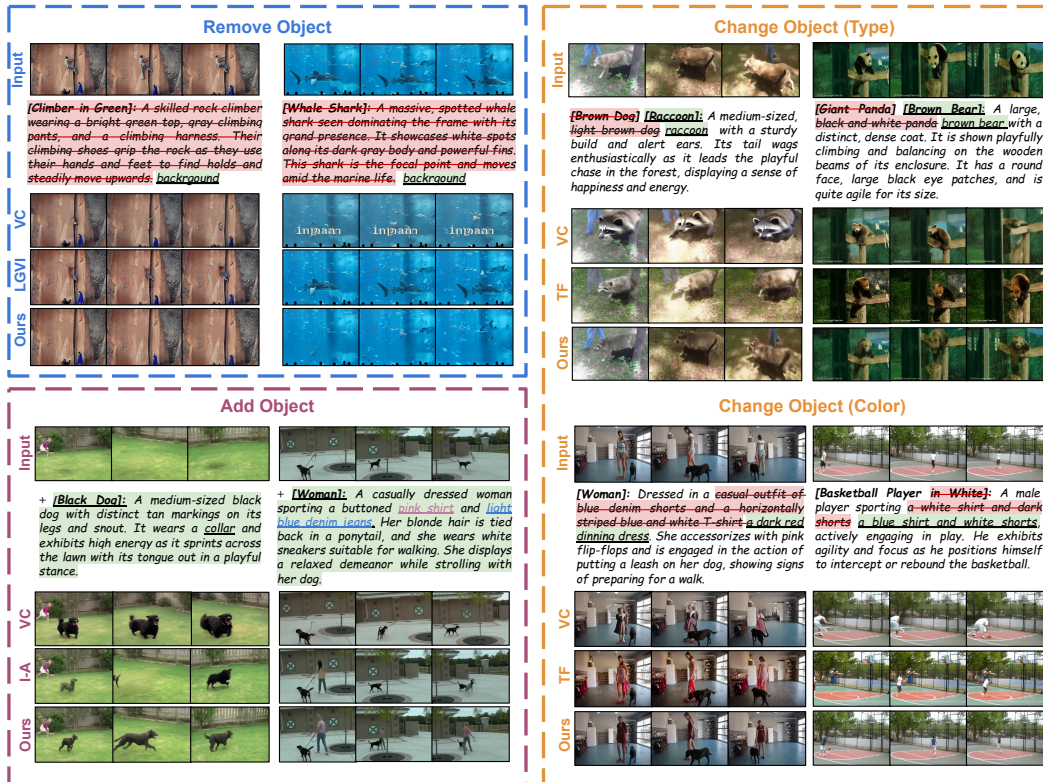


Figure 3: **Qualitative Comparison between RACCooN and other baselines.** Baseline names are abbreviated: **VC**: VideoComposer**I-A**: Inpainting Anything**TF**: TokenFlowWe underlined visual details in our caption. Best viewed in color.

Table 1: **Results of Single Object Prediction** on VPLM test set. Metrics are abbreviated:**S**: *SPICE*, **B**: *BLEU-4*, **C**: *CIDEr*.

Methods	S	B	C	IoU	FVD	CLIP
<i>open-source MLLMs</i>						
LLaVA	17.4	27.5	18.5	-	-	-
Video-Chat	18.2	25.3	19.1	-	-	-
PG-VL	18.2	27.4	14.6	-	-	-
<i>proprietary MLLMs</i>						
Gemini 1.5 Pro	19.2	23.5	11.0	0.115	371.63	0.978
GPT-4o	20.6	28.0	37.4	0.179	447.67	0.977
<b>RACCooN</b>	<b>23.1</b>	<b>31.0</b>	<u>33.5</u>	<b>0.218</b>	<u>432.42</u>	<b>0.983</b>

Table 2: **Results of Human Evaluation** on YouCook2. We measure the quality of the description through four metrics: Logic Fluency (Logic), Language Fluency (Lang.), Video Summary (Summ.), and Video Details (Details). We report the normalized score  $s \in [0, 100]$ .

Methods	Logic	Lang.	Summ.	Details	Avg.
Ground Truth	66.7	42.2	41.7	<b>72.2</b>	55.7
PG-VL	77.2	81.1	69.4	62.8	72.6
<b>RACCooN</b>	<b>80.6</b>	<b>85.0</b>	<b>72.2</b>	<b>72.2</b>	<b>77.5</b>

Table 3: **Results of Video Content Editing on three sub-tasks** on VPLM test. We gray out models that conduct the DDIM inversion process and have a different focus on our inpainting-based model.

Model	Change Object			Remove Object			Add Object		
	CLIP-T $\uparrow$	CLIP-F $\uparrow$	Qedit $\uparrow$	FVD $\downarrow$	SSIM $\uparrow$	PSNR $\uparrow$	FVD $\downarrow$	SSIM $\uparrow$	PSNR $\uparrow$
<i>Inversion-based Models</i>									
FateZero	25.18	94.47	1.01	1037.05	47.35	15.16	1474.80	47.65	15.45
TokenFlow	29.25	96.23	1.31	1317.29	47.06	15.83	1373.20	49.95	15.95
<i>Inpainting-Based Models</i>									
Inpaint Anything	24.86	92.01	1.01	383.81	82.33	27.69	712.59	77.75	22.41
LGVI	23.82	<b>95.33</b>	1.04	915.24	56.16	19.14	1445.43	47.93	16.09
VideoComposer	27.61	94.18	<b>1.25</b>	827.04	47.34	17.55	1151.90	48.01	15.76
<b>RACCooN</b>	<b>27.85</b>	<u>94.78</u>	<u>1.15</u>	<b>162.03</b>	<b>84.38</b>	<b>30.34</b>	<b>415.82</b>	<b>77.81</b>	<b>23.38</b>