

Edit As You Wish: Video Caption Editing with Multi-grained User Control

Linli Yao
School of Computer
Science,
Peking University
Beijing, China
linliyao@stu.pku.edu.cn

Yuanmeng Zhang
Alibaba Group
Beijing, China
zhangyuanmeng.zym@alibaba-
inc.com

Ziheng Wang
School of Information,
Renmin University of
China
Beijing, China
zihengwang@ruc.edu.cn

Xinglin Hou
Alibaba Group
Beijing, China
xinglin.hxl@alibaba-
inc.com

Tiezheng Ge
Alibaba Group
Beijing, China
tiezheng.gtz@alibaba-
inc.com

Yuning Jiang
Alibaba Group
Beijing, China
mengzhu.jyn@alibaba-
inc.com

Xu Sun
School of Computer
Science,
Peking University
Beijing, China
xusun@pku.edu.cn

Qin Jin*
School of Information,
Renmin University of
China
Beijing, China
qjin@ruc.edu.cn

Abstract

Automatically narrating videos in natural language complying with user requests, i.e. Controllable Video Captioning task, can help people manage massive videos with desired intentions. However, existing works suffer from two shortcomings: 1) the control signal is single-grained which can not satisfy diverse user intentions; 2) the video description is generated in a single round which can not be further edited to meet dynamic needs. In this paper, we propose a novel Video Caption Editing (VCE) task to automatically revise an existing video description guided by multi-grained user requests. Inspired by human writing-revision habits, we design the user command as a pivotal triplet $\{operation, position, attribute\}$ to cover diverse user needs from coarse-grained to fine-grained. To facilitate the VCE task, we *automatically* construct an open-domain benchmark dataset named VATEX-EDIT and *manually* collect an e-commerce dataset called EMMAD-EDIT. We further propose a specialized small-scale model (i.e., OPA) compared with two generalist Large Multi-modal Models to perform an exhaustive analysis of the novel task. For evaluation, we adopt comprehensive metrics considering caption fluency, command-caption consistency, and video-caption alignment. Experiments reveal the task challenges of fine-grained multi-modal semantics understanding and processing. Our datasets, codes, and evaluation tools are available at <https://github.com/yaolinli/VCE>.

CCS Concepts

• Computing methodologies → Computer vision.

*Qin Jin is the corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '24, October 28–November 1, 2024, Melbourne, VIC, Australia.

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0686-8/24/10

<https://doi.org/10.1145/3664647.3680724>

Keywords

Video Captioning, Caption Editing, Controllable Generation

ACM Reference Format:

Linli Yao, Yuanmeng Zhang, Ziheng Wang, Xinglin Hou, Tiezheng Ge, Yuning Jiang, Xu Sun, and Qin Jin. 2024. Edit As You Wish: Video Caption Editing with Multi-grained User Control. In *Proceedings of the 32nd ACM International Conference on Multimedia (MM '24)*, October 28–November 1, 2024, Melbourne, VIC, Australia. ACM, New York, NY, USA, 18 pages. <https://doi.org/10.1145/3664647.3680724>

1 Introduction

The proliferation of videos on the Internet heralds the era of video-dominated media. Video captioning, i.e. automatically describing videos using natural language, has been a prevalent task to assist people in comprehending and managing massive videos. However, conventional video captioning systems [51, 67] tend to generate intention-agnostic descriptions, ignoring the various demands of different users. Therefore, a new task branch, namely controllable video captioning [4, 7, 23, 69], has been proposed to integrate user intention as a control signal to guide the description generation.

Although controllable video captioning has great potential in practical applications, existing works have two non-negligible drawbacks. First, they all employ fixed control signals that can only express *single-grained controls*, such as Part-of-Speech (POS) [53] for structure control, or specified object tags [23] for semantic control. These single-grained controls can not satisfy flexible and diverse user demands. Second, they are *single-round controls* that generate a video description once which can not be further revised. Whereas iteratively revising sentences until the ideal texts is a natural process for humans [9]. Imagine a real-world scenario, such as E-commerce product promotion, where sellers upload product videos with descriptions to attract customers. There is a good chance that automated video descriptions fail to highlight the seller's preferences. As a result, sellers need to further improve the descriptions by themselves, which is time-consuming and labor-intensive, especially when facing massive long videos.

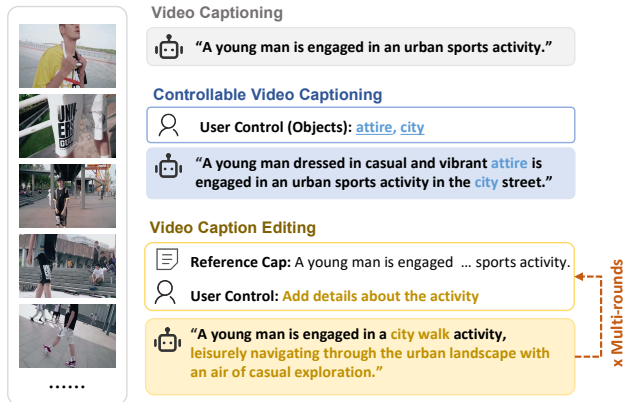


Figure 1: Comparisons between our proposed Video Caption Editing (VCE) task with conventional video captioning and controllable video captioning.

We propose a novel **Video Caption Editing (VCE)** task to automate the video description editing process. The task aims to edit an existing video description conditioned on user commands and video content. As depicted in Figure 1, the inputs of VCE task consist of a video, a reference description, and a user control. It outputs an edited video description based on the user command as a control signal. The reference caption can be initialized using the last edited output sentence which can thus enable *multi-round modification*. The VCE task can facilitate personalized video description generation by fulfilling miscellaneous demands from different users or dynamic demands from the same user.

In the VCE task, how to define the user command to cover multi-grained requests is crucial. Inspired by the human writing-revision habits, we unify the user edit commands into a triplet format $\{operation, position, attribute\}$ (depicted in Figure 2) for three advantages. Firstly, it condenses the core elements in an editing operation. Secondly, it can accommodate two prevalent front-end interface signals including natural language and editing trajectories from tablet computers. Finally, the different combinations of three elements in the triplet can cover *multi-grained* user commands from coarse-grained control (e.g. sentence length change) to fine-grained control (e.g. insert new details in the specific position), as illustrated in Table 1.

We collect two novel benchmark datasets named VATEX-EDIT and EMMAD-EDIT to support the exploration of the VCE task. The VATEX-EDIT dataset is automatically constructed from a large-scale video-text dataset VATEX [54] in the open domain. Meanwhile, to close the gap between research advances and real-life applications, we manually collect an e-commerce editing dataset called EMMAD-EDIT, which is more challenging from two aspects: 1) longer videos (average 27.1 seconds) and longer captions (average ~ 100 words); 2) external domain knowledge needed to generate product-oriented video descriptions. Based on the two benchmark datasets, we propose a specialized model, namely OPA, that converts the command triplet into a textual token sequence to alleviate heterogeneity among multi-grained commands. We demonstrate the feasibility of utilizing a unified framework to handle seven types of user commands. Moreover, we adopt two generalist Large

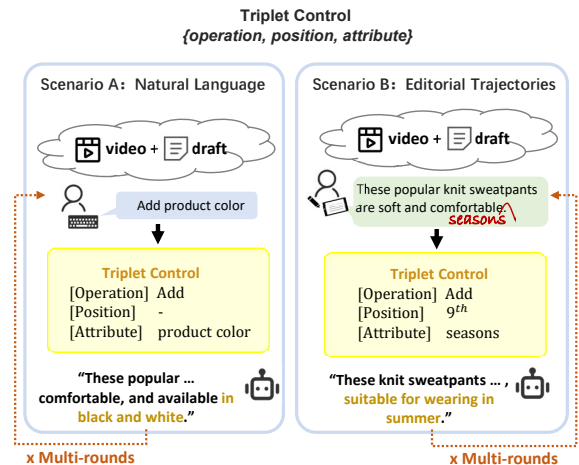


Figure 2: The triplet control designed in the VCE task can pivot two prevalent interaction signals including natural language (Scenario A) and editing trajectories (Scenario B).

Multimodal Models (LMMs) as a comparison to gain an in-depth understanding of the characteristics and challenges of the VCE task.

The main contributions of this paper are four-fold. 1) To the best of our knowledge, we are the first to propose the VCE task to achieve *multi-round* editing and design the user command as a triplet format to express *multi-grained* user requests. 2) We build two benchmark datasets from different domains, including the general domain (VATEX-EDIT) and commercial domain (EMMAD-EDIT), to facilitate the investigation of the VCE task. 3) We develop an evaluation suite to assess the edited video description based on caption fluency, command-caption consistency, and video-caption alignment. 4) We propose a unified specialist framework OPA and adapt two generalist LMM methods to initially tackle the task, followed by a comprehensive analysis.

2 Related Work

Controllable Video Captioning. Video captioning [1, 26, 34, 37, 41, 43, 48, 51, 67] is a challenging cross-modal task to automatically describe the visual contents of a video in natural languages. In order to satisfy the varied pragmatic interests of different users, controllable video captioning [4, 7, 23, 69] has been a newly prevalent task. It aims to derive video descriptions conditioned on a predefined control signal, e.g. visual object tags. Wang et al. [53] introduce Part-of-Speech (POS) information as guidance and Yuan et al. [69] directly utilize an exemplar sentence. Their goal is to generate descriptions with desired syntactic structures. Meanwhile, other works aim to control sentence semantics. Chen et al. [4] proposes a topic-guided model to generate topic-oriented descriptions. Liu et al. [23] focus on producing object-oriented sentences controlled by multiple user-interested objects. However, the above endeavors all generate a sentence once and can't be edited dynamically. Besides, their designed control signals are single-grained which can not cover flexible user intentions. Instead, we define a novel VCE task that can revise a description in multiple rounds covering multi-grained user demands.

Image Caption Editing. Conventional image captioning [18, 21, 22, 46, 52, 57, 60, 62, 64] generates the description for images from

Table 1: Editing commands via different combination of elements $\{operation, position, attribute\}$. It covers seven multi-grained demands from coarse-grained controls (e.g. expand description) to fine-grained controls (e.g. add specified *attributes* at specified *positions*). The atomic operations consist of *add* and *delete*. The multi-grained commands with a reference caption (e.g. “A group of girls is on the field playing a game.”) are unified as a control token sequence (Section 5.1) to guide the model.

Command	Notation	Demand	Unified Input Control
opera. pos attr			
✓	$\langle add, -, - \rangle$	expand description	[ADD] A group of girls is playing a game.
✓ ✓	$\langle add, pos, - \rangle$	expand description at specified <i>positions</i>	[ADD] A group of girls is [MASK] playing a game.
✓	$\langle add, -, attr \rangle$	add specified <i>attributes</i> in description	[ADD] field, hockey; A group of girls is playing a game.
✓ ✓ ✓	$\langle add, pos, attr \rangle$	add specified <i>attributes</i> at specified <i>positions</i>	[ADD] field, hockey; A group of girls is [MASK] playing a game.
✓	$\langle del, -, - \rangle$	shorten description	[DEL] A group of girls is on the field playing a game.
✓ ✓	$\langle del, pos, - \rangle$	shorten description at specified <i>positions</i>	[DEL] A group of girls is on (the filed) [MASK] playing a game.
✓	$\langle del, -, attr \rangle$	delete specified <i>attributes</i> from description	[DEL] field, group; A group of girls is on the field playing a game.

scratch which may lead to factual mistakes. Sammani and Elsayed [39] firstly define the image caption editing task that modifies an existing caption (a.k.a. reference caption) conditioned on the image content to obtain more accurate descriptions. Sammani and Melas-Kyriazi [40] further propose a novel EditNet framework to achieve interactive and adaptive edits. Yuan et al. [68] design an adaptive text-denoising network to alleviate the semantic gap between input images and reference sentences. The above works all edit image captions implicitly. Wang et al. [55] propose the explicit image caption editing task to make the modification process more explainable and efficient. In summary, all these works can only correct the wrong content in the reference captions and ignore specific edit intentions of different users. In this paper, we integrate multi-grained user commands into the video description editing process.

Large Multimodal Models. Recent months have witnessed the tremendous success of Large Language Models (LLMs) [5, 30, 31, 33, 49, 50, 70] towards artificial general intelligence. Further, Large Multimodal Models (LMMs) [2, 6, 16, 19, 24, 38, 65, 71] endow LLMs with the visual understanding ability by incorporating vision backbones [17, 36, 47]. Existing LMMs primarily bifurcate into two categories: Image Large Language Models (ImgLLMs) [2, 3, 6, 24, 65] and Video Large Language Models (VidLLMs) [16, 19, 27, 38, 71]. The standard architecture of an LMM comprises a vision backbone for encoding images or videos, a projector [3, 14, 63] to translate visual embeddings into textual semantic space, and an LLM to process all multimodal contexts. As the VCE task involves capturing video semantics, understanding textual user controls, and enabling text editing, it can serve as a new touchstone for LMMs.

3 Video Caption Editing Task

3.1 Task Definition

Given a video V and a reference caption $R = \{r_1, \dots, r_L\}$, the VCE task aims to generate an edited caption $Y = \{y_1, \dots, y_T\}$ according to the user edit command C . The edited caption Y should satisfy the constraints of V , R and C . Given a ground truth caption Y^* , the maximum likelihood estimation (MLE) training objective of VCE

task can be formulated as:

$$\mathcal{L}_{MLE} = -\frac{1}{T} \sum_{t=1}^T \log p(y_t^* | y_{<t}^*, V, R, C) \quad (1)$$

The reference caption can be initialized with the output caption from the last round of editing. It is also possible to start the editing process using a auto-generated sentence or human-written one as the reference. Due to the reference caption input setting, the VCE task can naturally achieve an interactive editing process with successive editing rounds, which is in line with human writing habits [9]. The interactive and multi-round revisions can help produce descriptions with higher user satisfaction.

3.2 User Edit Command

It is not trivial to define flexible edit commands in the VCE task to meet various realistic user needs. We observe that natural language and writing-revision traces are two natural interactive modes. The former can be received from keyboards or speech converters, while the latter conveniently expresses user intentions with the prevalence of tablets and wireless stylus pens. A command representation compatible with the above two signals is important and meaningful. In this paper, we propose a novel command representation in a triplet format $\{operation, position, attribute\}$, where *operations* control the overall description editing, *positions* specify the editing locations, which could affect the syntax of sentences, and *attributes* guide the editing operation to control the semantic contents of descriptions.

We define the atomic edit operations as *add* and *delete*, considering that the *replace* editing can be decomposed into the two atomic operations (i.e. first *delete* then *add*). Meanwhile, *position* and *attribute* in the triplet are optional, therefore, as shown in Table 1, seven specific commands¹ via different combinations of *operation*, *position*, and *attribute* elements in the triplet can cover multi-grained realistic demands from coarse-grained (global) controls to fine-grained (local) controls. The designed triplet command

¹Note that we omit the command “ $\langle del, pos, attr \rangle$, delete attributes at specified positions”, as it can be covered by “ $\langle del, pos, - \rangle$, delete description at specified positions”.

Table 2: Data statistics of VATEX-EDIT and EMMAD-EDIT dataset. # denotes the number. $VTime$ refers to the average duration of videos in seconds. Len_{Ref} denotes the average length of reference captions and Len_{GT} is the average length of ground-truth captions. $Edit Dist$ means the average minimum edit distance between reference captions and edited captions.

Dataset	Vision	#Videos/Images			#Editing instances			VTime	Len_{Ref}	Len_{GT}	Edit Dist	Vocab
		Train	Val	Test	Train	Val	Test					
COCO-EE [55]	Image	52,587	3,055	2,948	97,567	5,628	5,366	-	10.3	9.7	10.9	11,802
Flickr30K-EE [55]	Image	29,783	1,000	1,000	108,238	4,898	4,910	-	7.3	6.2	8.8	19,124
VATEX-EDIT	Video	25,467	2,935	5,867	784,805	91,513	181,638	10.0	14.4	16.0	11.9	21,634
EMMAD-EDIT	Video	16,176	5,418	5,502	47,569	15,914	16,169	27.1	91.3	93.7	17.8	44,725

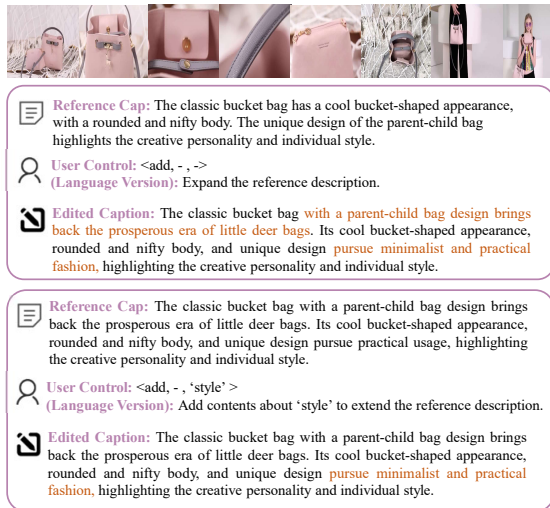


Figure 3: Annotated data instances of the VCE task.

can be flexibly obtained by processing the inputs from front-end interfaces including natural language and writing-revision traces (details in Appendix E). In the following method and experiments sections, we perform video description editing directly based on the triplet command.

4 Data Collection

To facilitate the novel VCE task, we *automatically* construct an open-domain dataset VATEX-EDIT, and *manually* annotate an E-commerce dataset EMMAD-EDIT. Table 2 displays the overall data statistics. Compared to prior image caption editing datasets such as COCO-EE and Flickr30K-EE, our new datasets present several distinct advantages: 1) more challenging with the video input and lengthier captions; 2) more diverse encompassing open-domain and e-commerce data; and 3) larger in scale. Specific annotated data instances are illustrated in Figure 3.

4.1 VATEX-EDIT Construction

It is challenging to construct data samples for the VCE task from scratch, which needs a quadruple (*video*, *command*, *reference caption*, *edited caption*) data, abbreviated as (V, C, R, Y). To mitigate the difficulty, we build the VATEX-EDIT dataset by expanding the widely-used video captioning dataset VATEX [54], which has high-quality caption annotations for each video.

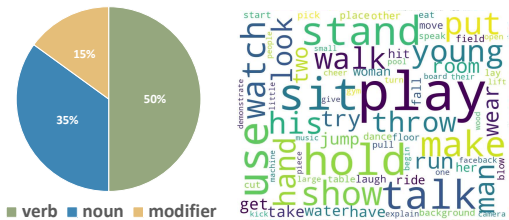


Figure 4: Attribute statistics on the VATEX-EDIT.

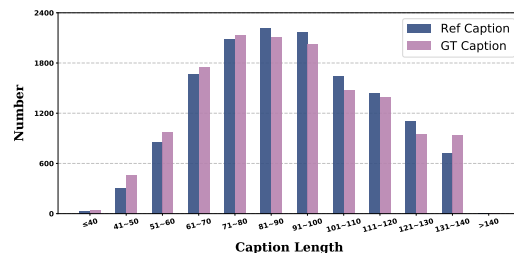


Figure 5: Caption length distributions on EMMAD-EDIT.

We sample an annotated caption of a video as the reference caption, and the next goal is to construct the *command* and the related *edited caption*. In general, we aim to construct related (*command*, *edited caption*) samples including: 1) **coarse-grained length-control commands** referring to the global *add* or *delete* edits that result in length changes, and 2) **finer-grained attribute-related commands** to achieve *add* or *delete* attributes.

Coarse-grained length-control commands. For *add* operation, we directly select a longer caption with remarkable length differences as the edited caption. The *delete* operation does the exact opposite. Considering the *delete* operation can be easily achieved without video content referring, we replace the reference caption with negative captions that have partially misaligned semantics with the video content. Such updated quadruples require models to prioritize removing visually irrelevant content from the reference, which makes the *delete* more challenging.

Fine-grained attribute-related commands. We construct attribute samples in a “degradation” manner, that is firstly detecting attribute words in the reference caption R and then removing them to obtain the edited caption Y . We utilize Spacy Syntactic Dependency² and Semantic Role Labeling [42] to achieve noun, verb or modifier attributes detection and removal in R while maintaining the fluency of $R_{\setminus attr}$ to get Y . After degradation, we can obtain

²<https://spacy.io/>

quadruples for the $\langle del, -, attr \rangle$ command. We exchange the reference caption and the edited caption to obtain the opposite *add* command. The position information can be naturally recorded to support position-related commands.

Statistics and analysis. As shown in Table 2, compared with existing image editing datasets [55], VATEX-EDIT has two salient features: *large-scale* and *diverse*. Figure 4 visualizes the percentages of modifiers, nouns and verbs in the *attribute*. Our automatic construction strategy selects verbs as the dominant attributes because verbs are usually related to temporal visual semantics, which is also one of the core challenges of the video description task. The word cloud of specific attribute words shows the *attribute* diversity.

4.2 EMMAD-EDIT Collection

To satisfy realistic user-demand scenarios, we manually collect a high-quality dataset called EMMAD-EDIT in the E-commerce domain based on a Chinese E-commerce video captioning dataset E-MMAD [72]. The E-MMAD dataset consists of product videos with advertising video descriptions, and additional structure information. Given a product-oriented video V and an original video description R , we recruit crowd workers to annotate three types of edited descriptions as follows.

Simplify original captions to the target length while maintaining sentence fluency and coherence according to the video content. To ensure the challenge of the VCE task, we require that the length of original sentences should be reduced by at least 20%.

Delete specific attributes. It aims to select multiple significant attribute words/phrases from R and remove the attribute-related content to get a new caption Y . The attributes can be nouns, verbs, or modifiers. To ensure the semantic coherence of Y , workers are allowed to modify other parts of R following the “minimal editing” principle.

Delete abstract attributes. We further consider deleting abstract attributes that do not directly appear in R . For example, deleting “Time and Seasons” needs to locate season-related content such as “spring” and “summer”. It is more challenging to edit with abstract attributes and also more down-to-earth since user intentions may be vague.

Statistics and analysis. To ensure annotation quality, extra workers further check the annotated cases. Table 2 shows the specific data statistics. EMMAD-EDIT has two remarkable characteristics, i.e. *long videos* and *long descriptions*. The average video length is 27.1 seconds and the average description length (specified in Figure 5) is around 100 words. We believe the challenging EMMAD-EDIT dataset will promote new technologies for the VCE task.

5 Methodology

In this section, we begin by introducing how to transform the triplet control into a unified textual sequence. Subsequently, we explore three approaches for the VCE task to facilitate a comprehensive comparison. We propose the **Operation-Position-Attribute (OPA)** model as a small-scale specialist. Additionally, we utilize an Image Large Language Model (ImgLLM) pipeline, and an end-to-end Video Large Language Model (VidLLM) to observe the performance of large multimodal models. Lastly, we develop an evaluation protocol for the novel task.

5.1 Input Format Design

We first integrate the seven specific edit commands introduced in Table 1 into a unified format to achieve multi-grained control.

The main challenge is the heterogeneity among the three elements of command, including *operation*, *position*, and *attribute*. On the one hand, *operations* and *attributes* change the textual semantics while *positions* mainly influence sentence syntax. On the other hand, *attributes* are specific textual words while *positions* are absolute position indexes.

To tackle the above challenges, we unify the input format as a textual token sequence. As shown in Table 1, we define two special tokens, [ADD] and [DEL], to represent different *add* or *delete* operations. *Attribute* words are naturally text tokens. For *position*, we put special tokens [MASK] in the reference caption to indicate the absolute position indexes. For example, a positioned reference caption “A group of girls is [MASK] playing a game” guides the model to generate new details between words “is” and “playing”. Finally, we concatenate the operation token, the attribute words, and a positioned reference caption as a control sequence to guide the model for description generation. Table 1 visualizes the input control sequences under seven specific commands when the reference caption is “A group of girls is playing a game”.

5.2 OPA: A Small-Scale Model as the Specialist

We construct a small-scale encoder-decoder Transformer architecture, i.e. multi-modal BART [13], to achieve the video description editing task under the guidance of processed control sequences. We utilize the pre-trained BART weights and endow it with the multi-modal ability to understand video content. The specific architecture is depicted in Appendix A.

Input Representation. Given a video V , a reference caption $R = \{r_1, \dots, r_L\}$, and a triplet command C , we first extract frame-level visual features and map them to the same dimension as textual embedding. We denote the input attributes as $A = \{a_1, \dots, a_M\}$ and the indicated position index as $l \in [1, L]$. Taking the most fine-grained command “*add* specified *attributes* at specified *positions*” as an example, the concatenated control sequence \tilde{C} for the command is defined as $\{[ADD], A, \tilde{R}\}$. Using special tokens to separate each part, it is formulated as:

$$\tilde{C} = \{[opera] [ADD] [/opera] [attr] A [/attr] [ref] \tilde{R} [/ref]\} \quad (2)$$

where the positioned reference caption \tilde{R} is formulated as:

$$\tilde{R} = \{r_1, \dots, r_{l-1}, [MASK], r_{l+1}, \dots, r_L\} \quad (3)$$

Finally, we input the visual features $\{V_1, \dots, V_N\}$ and the textual control sequence embedding $W_{\tilde{C}} = \{W_{[o]}, W_{[ADD]}, \dots, W_{[r]}\}$ to the Transformer encoder. If the *position* is empty in the command, we input the original reference caption R . When the *attribute* is empty in the command, we set A as an empty set.

Leverage Pre-trained Knowledge. The overall training objective as formulated in Section 3.1 is to generate an edited description conditioned on the video features and the control sequence. It is worth noting that we keep the [MASK] token consistent with the same token in the *Text Infilling* pre-trained task of BART to leverage the intrinsic pre-training textual ability.

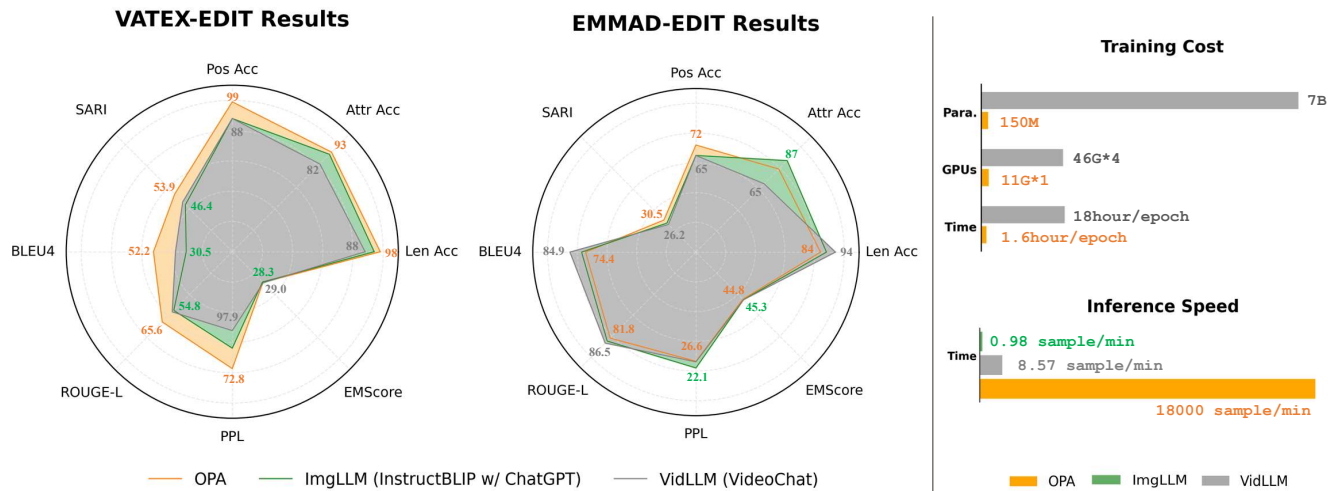


Figure 6: Overall performance of the small-scale specialist model (i.e., OPA) and large-scale generalist models (i.e., ImgLLM pipeline and end-to-end VidLLM) on the VATEX-EDIT and EMMAD-EDIT dataset. We utilize InstructBLIP [6] w/ ChatGPT [30] as the ImgLLM pipeline without training. Meanwhile, we conduct instruction tuning on the VideoChat-7B [16] as the end-to-end VidLLM method. The training cost (parameters, used GPUs, and training time) and inference speed via a single GPU are on the right. It is noted that we take the negative of PPL↓ and rescale it on the radar coordinate for better visualization.

5.3 LMMs as Contrastive Generalists

Large multimodal models integrate the advantages of visual understanding and remarkable natural language processing abilities (e.g., text editing) from LLMs. It is significant to probe their performance on the VCE task. Consequently, we explore two typical branches of LMMs including an ImgLLM pipeline and an end-to-end VidLLM.

ImgLLM Pipeline. We utilize the InstructBLIP [6] as the ImgLLM. Nevertheless, ImgLLM can only handle images or a single video frame as visual input. To adapt the ImgLLM to the VCE task, we combine InstructBLIP with ChatGPT [30]. In this way, InstructBLIP transforms visual semantics at the frame level into textual context, while ChatGPT consolidates all textual task context and achieves caption editing. Specifically, we extract frames from a given video and utilize InstructBLIP to produce detailed visual descriptions for each frame. The frame descriptions with the VCE task definition, task guidelines, and in-context demonstrations [8] of the relative command type will be combined as the final prompt to the ChatGPT. Instruction details are provided in Appendix A.

End-to-end VidLLM. As VidLLM can handle video tasks directly, we employ the VideoChat [16] as an end-to-end LMM solution. Specifically, we reformat the VATEX-EDIT and EMMAD-EDIT datasets into question-answer chat samples (refer to Appendix A for specifics) and conduct further instruct-tuning on VideoChat-7B using two datasets respectively.

5.4 Evaluation Suite

How to evaluate the novel VCE task is another noteworthy challenge. Conventional video captioning tasks adopt widely-used captioning metrics such as BLEU4, METEOR, and CIDEr. However, these reference-based metrics only measure the consistency between generated captions and ground-truth annotations, which

are insufficient. In this paper, we evaluate the VCE task from three aspects: 1) fluency, 2) controllability, and 3) text-video alignment.

Fluency. Following the previous work [55], we adopt widely-used BLEU4 [35] and ROUGE-L [20] metrics to measure the overall generation quality. We also use the **Perplexity (PPL)** [11] metric that reflects the grammatical correctness and semantic meaningfulness.

Controllability. Measuring whether an edited caption strictly follows control signals is important for the VCE task. Inspired by previous work [9], we first utilize the SARI [58] metric to measure the overall edit quality, i.e. the consistency between expected-to-edit and actually-edited spans. Moreover, we design three breakdown metrics namely **Length Accuracy**, **Attribute Accuracy** and **Position Accuracy** to measure whether the edited caption satisfies the $\{operation, position, attribute\}$ triplet control. Concretely, Len-Acc reflects the length change accuracy. Attr-Acc checks the appearance of commanded attribute words. Pos-Acc evaluates whether the model inserts/removes content in the specified positions.

Text-Video alignment. The VCE task inherently requires the alignment between edited descriptions and the given video. We use **EMScore** [44] to calculate the semantic similarity between edited captions and videos. It focuses on both coarse-grained similarity (video-sentence) and fine-grained similarity (frame-word).

6 Experiments

6.1 Implementation Details

We implement the small-scale OPA model based on Huggingface Transformers library [56]. The default setting is initialized by the BART_{base}. We get video frames using fps=1. For the VATEX-EDIT dataset in *English*, we adopt BLIP [15] ViT-B/16 to extract frame

Table 3: Ablation study of the OPA model on the VATEX-EDIT dataset. *Multimodal BART* is the backbone of OPA framework. *Pure Transformer* is the same model without pre-trained BART parameters. *Vision Align* means the vision-text alignment.

	Model	Controllability				Fluency			Vision Align
		Len-Acc	Attr-Acc	Pos-Acc	SARI	BLEU4	ROUGE-L	PPL↓	EMScore
1	Multimodal BART	-	-	-	49.7	48.0	62.0	73.9	28.7
2	Multimodal BART _{Opera}	97	-	-	52.1	49.6	63.2	70.6	28.7
3	Multimodal BART _{Opera+Attr}	97	93	-	53.8	52.3	65.7	72.7	28.7
4	OPA	98	93	99	53.9	52.2	65.6	72.8	28.7
5	Pure Transformer _{Opera+Pos+Attr}	98	82	97	52.6	50.5	64.4	77.2	28.7
6	3 Single-grained Models	96	69	99	53.3	51.5	64.6	72.8	28.6

Table 4: Overall and breakdown performances on the EMMAD-EDIT dataset.

	Command	Controllability				Fluency			Vision Align
		Len-Acc	Attr-Acc	Pos-Acc	SARI	BLEU4	ROUGE-L	PPL↓	EMScore
1	$\langle add, -, - \rangle$	57	-	-	26.2	62.1	73.9	23.7	44.7
2	$\langle add, pos, - \rangle$	75	-	59	27.0	83.6	90.3	25.1	44.9
3	$\langle add, -, attr \rangle$	80	74	-	31.7	84.7	89.9	26.6	45.2
4	$\langle add, pos, attr \rangle$	92	70	85	32.9	88.1	93.3	25.5	44.8
5	$\langle del, -, - \rangle$	100	-	-	33.5	66.8	73.8	30.5	44.6
6	$\langle del, pos, - \rangle$	99	-	-	30.7	83.9	90.6	28.8	44.6
7	$\langle del, -, attr \rangle$	100	93	-	33.6	75.2	83.6	28.9	44.7
8	Overall	84	79	72	30.5	74.4	81.8	26.6	44.8

features. The max frame sequence N is set to 20. For the EMMAD-EDIT dataset in *Chinese*, we initialize our model with the Chinese BART. We adopt CN-CLIP [59] ViT-B-16 to extract video frame-level features. The max frame sequence N is set to 30 and the max decoding length is set to 150. For training, we use AdamW [25] with a learning rate of $1e-5$ and optimize for 20 epochs with a batch size of 20. During inference, we set the beam size of generation as 5. The ImgLLM pipeline utilizes the identical frame number N as the OPA model. This pipeline doesn't involve any training. We choose one in-context learning sample for every command type integrated into the ChatGPT prompt. For VideoChat model, we set frame number N as 8 to fit its default setting. We conduct further instruct-tuning on the official 7B checkpoints with batch size 64.

6.2 Compare Specialist and Generalist Models

Figure 6 shows the overall performance of three baselines on the VATEX-EDIT and EMMAD-EDIT datasets. Interestingly, overall performances are divergent across the two datasets. On the large-scale open-domain VATEX-EDIT dataset, the small-scale specialist OPA model with only 150M parameters outperforms the LMM approaches. It suggests that with sufficient training instances (784,805 samples in VATEX-EDIT), a small specialized model has the potential to perform more effectively and efficiently.

On the e-commerce EMMAD-EDIT dataset, the LMM methods achieve higher scores across most metrics. EMMAD-EDIT is more challenging because it requires domain knowledge, such as unseen product attributes and advertising description style, and has limited

training data (refer to Table 2). Results show that ImgLLM with InstructBLIP and ChatGPT achieve highest Attr-Acc (87%) even without training. We argue that on this domain, generalist methods are more promising to leverage their intrinsic knowledge to edit product-related descriptions.

Despite performance advantages, the training cost and inference speed must be taken into account due to the booming number of videos. As illustrated in Figure 6 (right), compared to the VidLLM and ImgLLM, the OPA model demonstrates significant benefits over both VidLLM and ImgLLM in terms of lower training costs and faster inference speeds. Designing a model that balances performance with speed and cost represents a crucial trade-off.

Although both specialist and generalist models offer unique advantages, there remains considerable scope for further enhancements to develop an effective editing system, especially on the EMMAD-EDIT dataset. The *Controllability* (Pos-Acc 72%, Attr-Acc 87% and Len-Acc 94%) on the EMMAD-EDIT dataset is insufficient. Moreover, the alignment between video and caption, as indicated by the EMScore, requires significant improvement. In conclusion, the key of the VCE task lies in the combination of fine-grained video and user control understanding and precise text editing capabilities.

6.3 Further Task Analysis

We conduct further ablation studies on the small-scale OPA model to delve into a detailed analysis of the VCE task.

Increase control signals. We analyze the editing performance under different control signals in Table 3. The proposed OPA framework achieves high *controllability* accuracy (Len-Acc 98%, Attr-Acc

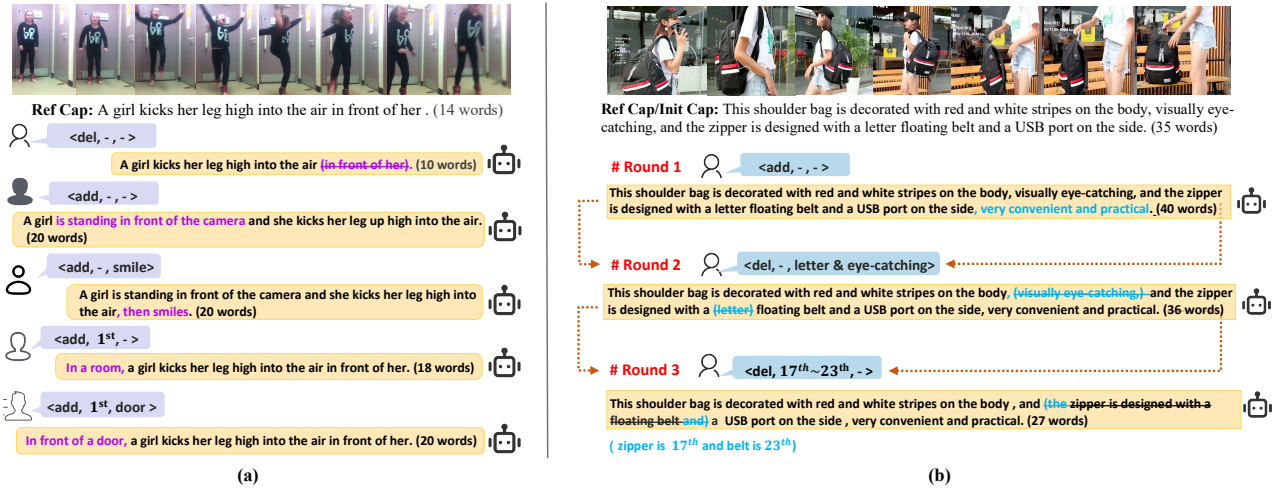


Figure 7: Visualization of (a) multi-grained command editing and (b) successive multi-round editing using the OPA model.

Table 5: The effects of visual modality on the VATEX-EDIT.

Command	Video	EMScore	SARI	BLEU4
Overall	✗	28.1	53.5	51.7
	✓	28.7 (+0.6)	53.9 (+0.4)	52.2 (+0.5)
$\langle del, -, - \rangle$	✗	27.0	39.4	12.1
	✓	28.5 (+1.5)	40.4 (+1.0)	13.5 (+1.4)
$\langle add, -, - \rangle$	✗	28.3	41.9	10.6
	✓	29.0 (+0.7)	42.0 (+0.1)	11.1 (+0.5)

93%, and Pos-Acc 99%) while maintaining sentence quality. The first block (lines 1-4) shows the controllability accuracy and caption quality when progressively inputting more control signals into the model. With the increasing aspects of control signals, there is no decline in sentence fluency and text-vision alignment. It indicates that our model can edit the reference caption with reasonable syntactic and semantic changes under multi-aspect guidance. Line 5 shows the result of Pure Transformer trained from scratch. Without BART pre-trained parameters, the overall controllability and fluency metrics decrease (SARI from 53.9 to 52.6, BLEU4 from 52.2 to 50.5), which verifies the benefits of textual pre-training knowledge. **Unified framework vs separate models.** To satisfy different-granularity commands, we compare the performances of training a unified OPA model vs. training multiple separate models in Table 3. In the *3 Single-grained Models* setting (line 6), we train three models respectively to deal with three control granularities, i.e. $\{operation\}$, $\{operation, attribute\}$, and $\{operation, position, attribute\}$. The OPA model reaches a remarkably higher score on the Attr-Acc (93% vs 69%) with better SARI, BLEU4, and ROUGE-L. It demonstrates that our unified input design can alleviate the confusion and heterogeneity of multi-grained commands.

Difficulty level of various commands. Table 4 (lines 1-7) displays the performances of different commands, indicating their respective difficulty levels. Generally, the *add* operation proves to be more challenging than *delete*, primarily because it requires the constraint of video content. For the *add* operations, the finer the command

Table 6: Mean score (rated 1-5) of the human evaluation on the two datasets. *Trans.* is short for Pure Transformer.

Dataset	Model	Control.	Fluency	Vision Align
EMMAD-EDIT	Trans.	2.94	3.02	3.27
	OPA	3.98	3.85	3.76
	GT	4.67	4.37	4.22
VATEX-EDIT	Trans.	4.18	4.23	3.76
	OPA	4.36	4.43	3.93
	GT	4.48	4.34	4.41

granularity (lines 1-4), the higher the *controllable* and *fluency* scores. It reveals that when provided with more detailed control signals, the model can generate desired captions more easily.

Effects of vision modality. We compare the model performance with and without video input in Table 5. Adding vision modality brings overall metric improvements since it provides visual semantics to guide edited video description generation. For $\langle del, -, - \rangle$ command, we especially construct challenging samples in which reference captions have misalignments with videos (Section 4.1). With the visual semantics, our model prioritizes removing the video-misalignment contents and achieving a higher EMScore (from 27.0 to 28.5). Similarly, $\langle add, -, - \rangle$ command requires enriching the original caption referring to the video content.

6.4 Quantitative Results

Multi-grained editing controls. Provided with various commands, the OPA model can output different edited descriptions to satisfy multi-grained user requests. As Figure 7 (a) shows, our OPA model successfully generates different desired descriptions given the same video, the same reference caption but different commands from coarse-grained (e.g. $\langle add, -, - \rangle$) to fine-grained (e.g. $\langle add, 1^{st}, door \rangle$).

Successive editing controls. The OPA model also supports interactive editing with successive controls in the VCE task, depicted in Figure 7 (b). The edited description can serve as the reference

caption in the next round to make further editing to satisfy dynamic user demands.

Human Evaluation. We further adopt human evaluation to assess the quality of edited video descriptions. We recruit 20 evaluators to score the generated descriptions. We randomly sample 200 test cases from VATEX-EDIT and 350 cases from EMMAD-EDIT respectively. During the evaluation, we randomly order the edited captions generated from *Pure Transformer baseline*, *OPA*, and *groundtruths (GT)*. The evaluators are asked to rate each description from three aspects on a scale of 1 to 5 points. Table 6 shows the OPA model exceeds the controllable Transformer baseline in three aspects, especially the *controllability*.

7 Conclusion

We propose a novel multi-modal task named Video Caption Editing (VCE), which aims to automatically edit video descriptions under the guidance of multi-grained user commands. To satisfy diverse and varied user demands, we design the user control signal as a $\{operation, position, attribute\}$ triplet to flexibly cover both coarse-grained and fine-grained controls. We collect two datasets named VATEX-EDIT and EMMAD-EDIT from different domains. We further employ comprehensive metrics to assess fluency, controllability, and vision-text alignment. Finally, we introduce a small specialized model called OPA, an ImgLLM pipeline, and an end-to-end VidLLM to dive into the task challenges and provide good starting points.

Limitations and Future Work. This paper primarily introduces appropriate baseline solutions for the VCE task, aiming to provide a thorough analysis. Nonetheless, it falls short of designing architectural innovations, leaving ample room for exploration in the future. Further insights into significant future directions are discussed in Appendix F.

Acknowledgments

This work was partially supported by the National Natural Science Foundation of China (No. 62072462), the Fundamental Research Funds for the Central Universities, and the National Natural Science Foundation of China (No. 62176002).

References

- Nayyer Aafaq, Naveed Akhtar, Wei Liu, Syed Zulqarnain Gilani, and Ajmal Mian. 2019. Spatio-temporal dynamics and semantic attribute enriched visual encoding for video captioning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 12487–12496.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966* (2023).
- Junbum Cha, Wooyoung Kang, Jonghwan Mun, and Byungseok Roh. 2023. Honeybee: Locality-enhanced Projector for Multimodal LLM. *arXiv preprint arXiv:2312.06742* (2023).
- Shizhe Chen, Qin Jin, Jia Chen, and Alexander G. Hauptmann. 2019. Generating Video Descriptions With Latent Topic Guidance. *IEEE Transactions on Multimedia* 21, 9 (2019), 2407–2418. <https://doi.org/10.1109/TMM.2019.2896515>
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. PaLM: Scaling Language Modeling with Pathways. *CoRR abs/2204.02311* (2022).
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Albert Li, Pascale Fung, and Steven C. H. Hoi. 2023. InstructBLIP: Towards General-purpose Vision-Language Models with Instruction Tuning. *ArXiv abs/2305.06500* (2023). <https://api.semanticscholar.org/CorpusID:258615266>
- Aditya Deshpande, Jyoti Aneja, Liwei Wang, Alexander Schwing, and David A Forsyth. 2019. Diverse and controllable image captioning with part-of-speech guidance. *CVPR* (2019).
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2022. A survey for in-context learning. *arXiv preprint arXiv:2301.00234* (2022).
- Wanyu Du, Vipul Raheja, Dhruv Kumar, Zae Myung Kim, Melissa Lopez, and Dongyeop Kang. 2022. Understanding Iterative Revision from Human-Written Text. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Dublin, Ireland, 3573–3590. <https://aclanthology.org/2022.acl-long.250>
- Zeyao Du. 2019. GPT2-Chinese: Tools for training GPT2 model in Chinese language. <https://github.com/Morizeyao/GPT2-Chinese>.
- Fred Jelinek, Robert L Mercer, Lalit R Bahl, and James K Baker. 1977. Perplexity—a measure of the difficulty of speech recognition tasks. *The Journal of the Acoustical Society of America* 62, S1 (1977), S63–S63.
- Peng Jin, Ryuichi Takanobu, Caiwan Zhang, Xiaochun Cao, and Li Yuan. 2023. Chat-UniVi: Unified Visual Representation Empowers Large Language Models with Image and Video Understanding. *arXiv preprint arXiv:2311.08046* (2023).
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 7871–7880. <https://doi.org/10.18653/v1/2020.acl-main.703>
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models. *arXiv:2301.12597* [cs.CV]
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*. PMLR, 12888–12900.
- KunChang Li, Yanan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. 2023. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355* (2023).
- Kunchang Li, Yali Wang, Yanan He, Yizhuo Li, Yi Wang, Limin Wang, and Yu Qiao. 2022. UniFormerV2: Spatiotemporal Learning by Arming Image ViTs with Video UniFormer. *arXiv preprint arXiv:2211.09552* (2022).
- Xiangyang Li, Xinhang Song, Luis Herranz, Yaohui Zhu, and Shuqiang Jiang. 2016. Image captioning with both object and scene information. In *Proceedings of the 24th ACM international conference on Multimedia*. 1107–1110.
- Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. 2023. Video-LLaVA: Learning Unified Visual Representation by Alignment Before Projection. *arXiv preprint arXiv:2311.10122* (2023).
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*. 74–81.
- Chang Liu, Changhu Wang, Fuchun Sun, and Yong Rui. 2016. Image2Text: a multimodal image captioner. In *Proceedings of the 24th ACM international conference on Multimedia*. 746–748.
- Daqing Liu, Zheng-Jun Zha, Hanwang Zhang, Yongdong Zhang, and Feng Wu. 2018. Context-aware visual policy network for sequence-level image captioning. In *Proceedings of the 26th ACM international conference on Multimedia*. 1416–1424.
- Fenglin Liu, Xuancheng Ren, Xian Wu, Bang Yang, Shen Ge, and Xu Sun. 2021. O2NA: An Object-Oriented Non-Autoregressive Approach for Controllable Video Captioning. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Association for Computational Linguistics, Online, 281–292. <https://doi.org/10.18653/v1/2021.findings-acl.24>
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. *arXiv preprint arXiv:2304.08485* (2023).
- Ilya Loshchilov and Frank Hutter. 2017. Decoupled Weight Decay Regularization. In *International Conference on Learning Representations*.
- Huaisiao Luo, Lei Ji, Botian Shi, Haoyang Huang, Nan Duan, Tianrui Li, Jason Li, Taroon Bharti, and Ming Zhou. 2020. UniVL: A Unified Video and Language Pre-Training Model for Multimodal Understanding and Generation. *arXiv preprint arXiv:2002.06353* (2020).
- Ruipu Luo, Ziwang Zhao, Min Yang, Junwei Dong, Minghui Qiu, Pengcheng Lu, Tao Wang, and Zhongyu Wei. 2023. Valley: Video Assistant with Large Language model Enhanced ability. *arXiv preprint arXiv:2306.07207* (2023).

- [28] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. 2023. Video-ChatGPT: Towards Detailed Video Understanding via Large Vision and Language Models. *ArXiv abs/2306.05424* (2023). <https://api.semanticscholar.org/CorpusID:259108333>
- [29] Meinard Müller. 2007. Dynamic time warping. *Information retrieval for music and motion* (2007), 69–84.
- [30] OpenAI. 2022. OpenAI: introducing ChatGPT. <https://openai.com/blog/chatgpt>
- [31] OpenAI. 2023. GPT-4 Technical Report. *arXiv:2303.08774 [cs.CL]*
- [32] OpenAI. 2023. GPT-4V(ision) System Card. https://cdn.openai.com/papers/GPTV_System_Card.pdf.
- [33] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. *CoRR abs/2203.02155* (2022).
- [34] Boxiao Pan, Haoye Cai, De-An Huang, Kuan-Hui Lee, Adrien Gaidon, Ehsan Adeli, and Juan Carlos Niebles. 2020. Spatio-temporal graph for video captioning with knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10870–10879.
- [35] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*. 311–318.
- [36] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *ICML*.
- [37] Vasili Ramanishka, Abir Das, Dong Huk Park, Subhashini Venugopalan, Lisa Anne Hendricks, Marcus Rohrbach, and Kate Saenko. 2016. Multimodal video description. In *Proceedings of the 24th ACM international conference on Multimedia*. 1092–1096.
- [38] Shuhuai Ren, Linli Yao, Shicheng Li, Xu Sun, and Lu Hou. 2023. TimeChat: A Time-sensitive Multimodal Large Language Model for Long Video Understanding. *ArXiv abs/2312.02051* (2023).
- [39] Fawaz Sammani and Mahmoud Elsayed. 2019. Look and Modify: Modification Networks for Image Captioning.
- [40] Fawaz Sammani and Luke Melas-Kyriazi. 2020. Show, Edit and Tell: A Framework for Editing Image Captions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [41] Rakshith Shetty and Jorma Laaksonen. 2016. Frame-and segment-level features and candidate pool evaluation for video caption generation. In *Proceedings of the 24th ACM international conference on Multimedia*. 1073–1076.
- [42] Peng Shi and Jimmy Lin. 2019. Simple bert models for relation extraction and semantic role labeling. *arXiv preprint arXiv:1904.05255* (2019).
- [43] Xiangxi Shi, Jianfei Cai, Shaifu Joty, and Jiuxiang Gu. 2019. Watch it twice: Video captioning with a refocused video encoder. In *Proceedings of the 27th ACM international conference on multimedia*. 818–826.
- [44] Yaya Shi, Xu Yang, Haiyang Xu, Chunfeng Yuan, Bing Li, Weiming Hu, and Zheng-Jun Zha. 2022. EMScore: Evaluating Video Captioning via Coarse-Grained and Fine-Grained Embedding Matching. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18–24, 2022*.
- [45] Yin Shukang, Fu Chaoyou, Zhao Sirui, Xu Tong, Wang Hao, Sui Dianbo, Shen Yunhang, Li Ke, Sun Xing, and Chen Enhong. 2023. Woodpecker: Hallucination Correction for Multimodal Large Language Models. *arXiv preprint arXiv:2310.16045* (2023).
- [46] Matteo Stefanini, Marcella Cornia, Lorenzo Baraldi, Silvia Cascianelli, Giuseppe Fiameni, and Rita Cucchiara. 2022. From show to tell: a survey on deep learning-based image captioning. *IEEE transactions on pattern analysis and machine intelligence* 45, 1 (2022), 539–559.
- [47] Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. 2023. Eva-clip: Improved training techniques for clip at scale. *arXiv preprint arXiv:2303.15389* (2023).
- [48] Mingkang Tang, Zhanyu Wang, Zhenhua Liu, Fengyun Rao, Dian Li, and Xiu Li. 2021. Clip4caption: Clip for video caption. In *Proceedings of the 29th ACM International Conference on Multimedia*. 4858–4862.
- [49] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and Efficient Foundation Language Models. *CoRR* (2023).
- [50] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288* (2023).
- [51] Subhashini Venugopalan, Marcus Rohrbach, Jeffrey Donahue, Raymond Mooney, Trevor Darrell, and Kate Saenko. 2015. Sequence to Sequence – Video to Text. In *2015 IEEE International Conference on Computer Vision (ICCV)*. 4534–4542. <https://doi.org/10.1109/ICCV.2015.515>
- [52] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3156–3164.
- [53] Bairui Wang, Lin Ma, Wei Zhang, Wenhao Jiang, Jingwen Wang, and Wei Liu. 2019. Controllable Video Captioning with POS Sequence Guidance Based on Gated Fusion Network. *arXiv preprint arXiv:1908.10072* (2019).
- [54] Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. 2019. Vatec: A large-scale, high-quality multilingual dataset for video-and-language research. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4581–4591.
- [55] Zhen Wang, Long Chen, Wenbo Ma, Guangxing Han, Yulei Niu, Jian Shao, and Jun Xiao. 2022. Explicit Image Caption Editing. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVI*. 113–129.
- [56] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*. 38–45.
- [57] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*. PMLR, 2048–2057.
- [58] Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, and Chris Callison-Burch. 2016. Optimizing Statistical Machine Translation for Text Simplification. *Transactions of the Association for Computational Linguistics* 4 (2016), 401–415. https://doi.org/10.1162/tacl_a_00107
- [59] An Yang, Junshu Pan, Junyang Lin, Rui Men, Yichang Zhang, Jingren Zhou, and Chang Zhou. 2022. Chinese CLIP: Contrastive Vision-Language Pretraining in Chinese. *arXiv preprint arXiv:2211.01335* (2022).
- [60] Dingyi Yang, Hongyu Chen, Xinglin Hou, Tiezheng Ge, Yuning Jiang, and Qin Jin. 2023. Visual captioning at will: Describing images and videos guided by a few stylized sentences. In *Proceedings of the 31st ACM International Conference on Multimedia*. 5705–5715.
- [61] Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023. The Dawn of LMMs: Preliminary Explorations with GPT-4V(ision). *arXiv:2309.17421 [cs.CV]*
- [62] Linli Yao, Weijing Chen, and Qin Jin. 2023. CapEnrich: Enriching Caption Semantics for Web Images via Cross-modal Pre-trained Knowledge. In *TheWebConf*.
- [63] Linli Yao, Lei Li, Shuhuai Ren, Lean Wang, Yuanxin Liu, Xu Sun, and Lu Hou. 2024. DeCo: Decoupling Token Compression from Semantic Abstraction in Multimodal Large Language Models. *arXiv preprint arXiv:2405.20985* (2024).
- [64] Linli Yao, Weiying Wang, and Qin Jin. 2022. Image difference captioning with pre-training and contrastive learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 36. 3108–3116.
- [65] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. 2023. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178* (2023).
- [66] Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Anwen Hu, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. 2023. mPLUG-Owl2: Revolutionizing Multi-modal Large Language Model with Modality Collaboration. *arXiv:2311.04257 [cs.CL]*
- [67] Youngjae Yu, Hyungjin Ko, Jongwook Choi, and Gunhee Kim. 2017. End-to-End Concept Word Detection for Video Captioning, Retrieval, and Question Answering. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 3261–3269. <https://doi.org/10.1109/CVPR.2017.347>
- [68] Mengqi Yuan, Bing-Kun Bao, Zhiyi Tan, and Changsheng Xu. 2023. Adaptive Text Denoising Network for Image Caption Editing. *ACM Trans. Multimedia Comput. Commun. Appl.* 19, 1s, Article 41 (feb 2023), 18 pages. <https://doi.org/10.1145/3532627>
- [69] Yitian Yuan, Lin Ma, Jingwen Wang, and Wenwu Zhu. 2020. Controllable Video Captioning with an Exemplar Sentence. In *The 28th ACM International Conference on Multimedia (MM '20)*.
- [70] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma, Yufei Xue, Jidong Zhai, Wenguang Chen, Peng Zhang, Yuxiao Dong, and Jie Tang. 2022. GLM-130B: An Open Bilingual Pre-trained Model. *abs/2210.02414* (2022).
- [71] Hang Zhang, Xin Li, and Lidong Bing. 2023. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858* (2023).
- [72] Zhipeng Zhang, Xinglin Hou, Kai Niu, Zhongzhen Huang, Tiezheng Ge, Yuning Jiang, Qi Wu, and Peng Wang. 2022. Attract me to Buy: Advertisement Copywriting Generation with Multimodal Multi-structured Information. *arXiv preprint arXiv:2205.03534* (2022).

In the supplementary material, we first introduce details of three approaches (Section A), then present more construction details about the VATEX-EDIT dataset (Section B) and EMMAD-EDIT dataset (Section C). Moreover, we provide the pseudo-code of the Position Accuracy metric (Section D) and the conversion instructions from interface signals to triplet format (Section E). Finally, based on a good start on the task, dataset and method foundation, we further discuss a variety of interesting aspects worth exploring in the future (Section F) and the related social impact (Section G).

A Model Details

A.1 Architecture of OPA Model

The overall architecture of the specialist model OPA is depicted in Figure 8, which is built on the BART model.

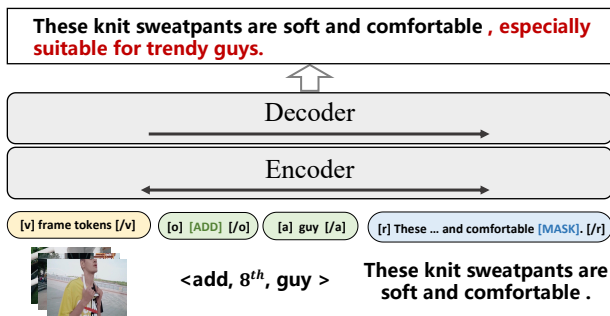


Figure 8: The overall OPA framework. 8^{th} denotes the specified position of the 8^{th} word.

A.2 Prompt of ImgLLM Pipeline

We build a ChatGPT pipeline with the visual expert model, i.e. InstructBLIP, to verify the vision-enhanced LLM performance on the VCE task. Specifically, we extract key frames from a given video

Question:
 "You are an AI visual assistant to tackle a novel task called Video Caption Editing (VCE) Task. The task goal is to automatically revise an existing video description guided by flexible user requests and the original video content. The inputs of VCE task consist of a video, a reference description, and a user command. It outputs an edited video description based on the user command as a control signal.
 1) Reference Description: This is the original sentence that describes a scene or action from the video. The [POS] token within this description indicates where specific changes should be made.
 2) User Command: A sentence with specific guidance to clarify the edit requirements.

Inputs:
 Reference Description: A man with a goatee is outside smoking a pipe on a cloudy day.
 User Command: Shorten the reference description. "

Answer:
 "A man lighting a tobacco pipe and smoking the pipe ."

Figure 9: A question-answer chat example used to conduct instruction tuning on the VidLLM.

Table 7: Results of GPT-4o on the EMMAD-EDIT dataset.

	Control.	Fluency	Text-Vision Align
OPA	66.4	76.5	44.8
VidLLM (VideoChat)	62.6	81.7	45.3
ImgLLM (w/ ChatGPT)	66.9	79.8	45.3
GPT-4o	71.1	81.9	45.3

and utilize InstructBLIP to produce detailed descriptions for each frame. We uniformly select 20 frames for VATEX-EDIT dataset and 30 frames for EMMAD-EDIT dataset, which is exactly the same frame number that our OPA model uses. These frame descriptions with timestamps help the ChatGPT to understand the exhaustive video content.

Designed Prompts. As Figure 10 illustrates, we conduct a well-designed prompt for ChatGPT including the task definition, helpful guidelines, an in-context learning [8] example, the video information, and the new case to be solved. With the input prompt, ChatGPT can output the desired edited caption with the required format.

In-Context Learning. Considering the differences between seven multi-grained commands, we manually select seven high-quality input-output demonstrations involving all command types. The command type of given examples will be aligned with the new case to fulfill better results. For example, if we want ChatGPT to edit a video description under the $\langle add, pos, attr \rangle$ command, we will give a matched example of $\langle add, pos, attr \rangle$ command to help it better solve the task.

A.3 Instruction-tuning Data of VidLLM

We convert the original samples from the VATEX-EDIT and EMMAD-EDIT datasets into an instruction-tuning format to facilitate the training of end-to-end video large language models. Figure 9 illustrates the converted data sample.

A.4 Results of GPT4-o

As Table 7 shows, we report the results of GPT-4o³ in the challenging EMMAD-EDIT dataset as a reference upper bound of the task. Given that the video interface for GPT-4o is not yet available, we extract 8 frames from each video and assess the model using these multiple images as inputs.

B VATEX-EDIT Dataset Details

B.1 Automatic Dataset Construction

We construct quadruples (*video, command, reference caption, edited caption*) according to two types of commands including **coarse-grained length-control commands** and **fine-grained attribute-related commands**. We complement more details about constructing attribute-related commands.

Attribute-related commands. Our goal is to construct related (*command, edited caption*) samples for each (*video, reference caption*) to support attribute related commands in a “degradation” manner. The main challenges lie in two aspects: 1) extracting meaningful

³<https://openai.com/index/hello-gpt-4o/>

You are an AI visual assistant to tackle a novel task called Video Caption Editing (VCE) Task. The task goal is to automatically revise an existing video description guided by flexible user requests and the original video content.

Task Definition:

The inputs of VCE task consist of a video, a reference description, and a user command. It outputs an edited video description based on the user command as a control signal.

Inputs:

1. **Reference Description:** This is the original sentence that describes a scene or action from the video. The [POS] token within this description indicates where specific changes should be made.
2. **User Command:** A sentence with specific guidance to clarify the edit requirements.
3. **Video Information:** Descriptions of extracted frames from the video. These are provided to understand the video context. Note: Sometimes, this information may be marked as 'None', indicating no specific video context is provided.

Guidelines:

- Make sure the output edited description is fluent and coherent.
- If the video context is unclear or not provided, base your edits on common or neutral assumptions about the scenario.
- If object identification from the Video Information is uncertain or seems erroneous, provide a more general description without detailed specifications.
- Ensure that the final edited description feels as if it's generated by an AI visual assistant who is watching the video in real time.

Examples:

1. **Reference Description:** A person is [POS] painting a picture.
2. **User command:** Add contents related to "brush" at [POS].
3. **Video Information:** *The frame at the second 1 shows ... The frame at the second 2 shows ... The frame at the second 3 shows*
4. **Expected Output:** A person is using a fine-tip brush and carefully painting a picture. "

Now I need your help to handle the following Video Caption Editing Task case and output the Edited Description. The output format should be "Expected Output: the edited sentence" without other outputs.

New Inputs:

Reference Description: A [POS] person is eating a cake.

User command: Expand the reference description and add video-related content at [POS].

Video Information:

The frame at the second 1 features a close-up view of a person's hand holding a paintbrush, which they are using to paint on a piece of paper or canvas. The person appears to be focused on their work, possibly creating a painting or illustration. In the background, there is a table with various objects, such as cups and a bottle, suggesting that the person might be working in a studio or creative space. The overall scene showcases the artist's attention to detail and dedication to their craft, as they skillfully use the paintbrush to bring their artistic vision to life.

The frame at the second 2 ...

.....

Figure 10: Designed prompts for the ImgLLM pipeline. It encompasses the task definition, helpful guidelines, an in-context learning [8] example, exhaustive visual descriptions, and the new case to be solved. With the well-designed prompt, ChatGPT can output the desired edited description with the required format. We manually construct seven in-context demonstrations involving each specific command type. For each input prompt, the specific command type of the in-context demonstration is aligned with the new case to help ChatGPT better understand the task.

noun, verb, or modifier attributes in the reference caption R and 2) deleting the attribute-related semantic spans while maintaining the fluency of rest content $R_{\setminus attr}$ to get an attribute-removed caption Y . In detail, we adopt four steps as follows:

- First, we use the Spacy syntactic dependency parser to build a textual dependency tree that contains the Part-of-Speech information and relationships between tokens. We select reasonable branches in the parsed tree as attributes and further prune the branch to ensure the fluency of the rest caption.
- Second, we use a Semantic Role Labeling model [42] to analyze the semantic roles of each span in a sentence. It can help to better judge whether a parsed attribute span in the first step can be deleted or not, especially noun attributes.
- Third, we merge the attributes which only modify one or two tokens to improve the task challenge, that is, the model may be required to edit with multiple attributes in one round.
- Finally, considering the intrinsic error of the parsing model, we adopt a post-processing stage to filter low-quality sentences considering sentence fluency, edited token length, and attribute diversity. The ground-truth sentence fluency is further verified in Section B.2.

When the attribute-related spans are removed in a sentence, the edited positions can be naturally recorded. Through the above steps, we can get high-quality samples for the $\langle del, pos, attr \rangle$ commands. Meanwhile, reversed samples can be obtained for the $\langle add, pos, attr \rangle$ command by exchanging the reference caption and the edited caption. For these two fine-grained commands, the sentence length, structure, and semantics are controlled at the same time.

For commands $\langle add, -, attr \rangle$ and $\langle del, -, attr \rangle$ that omit positions, they mainly control the sentence length and semantics, not structure. We get relevant samples by relaxing the structure constraint based on the above "degradation" manner. In detail, we replace the edited caption by retrieving desired sentences satisfying

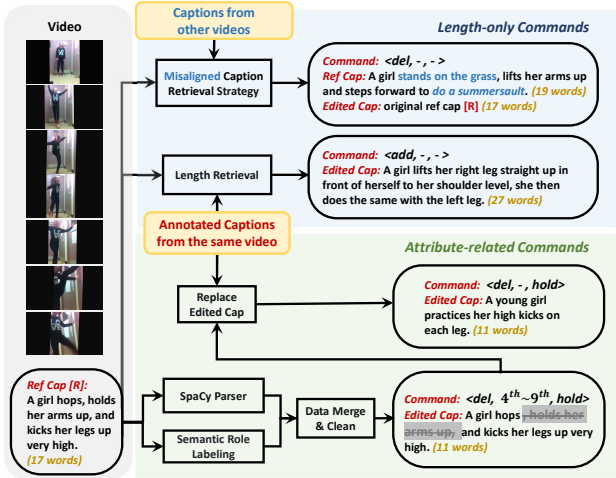


Figure 11: The automatic procedure of VATEX-EDIT dataset construction. The omitted commands, e.g. $\langle add, pos, attr \rangle$ and $\langle del, pos, - \rangle$, can be easily converted from the data samples of shown commands.

both length and semantics (attributes) constraints from original annotated descriptions for the same video.

B.2 VATEX-EDIT Test set Quality

In the VATEX-EDIT dataset, it is worth noting that only the “ $\langle del, pos, - \rangle$, shorten description at specified positions” command uses the auto-constructed sentences as ground-truth captions in order to control the sentence structure. Meanwhile, the data samples of other commands utilize human-annotated sentences from the original VATEX dataset as ground-truth captions. To assess the quality of the auto-constructed sentences, we evaluate the test set of VATEX-EDIT through human and ChatGPT evaluations. For the ChatGPT evaluation, we design suitable prompts (depicted in Figure 12) to guide ChatGPT to judge whether a sentence is fluent or not. We randomly sample 100 ground-truth sentences on the test set and calculate the fluency rate obtained by ChatGPT. For human evaluation, we recruit 20 crowd workers to rate the fluency score (ranging from 1 to 5) of 200 randomly sampled ground-truth sentences. As Table 8 shows, the automatically constructed ground-truth sentences of the “ $\langle del, pos, - \rangle$ ” command are as fluent as the human-annotated ground-truth sentences of other commands (87% vs 91% and 4.40 vs 4.37), which indicates the high quality of the VATEX-EDIT test set.

C EMMAD-EDIT Dataset Details

We manually collect the high-quality dataset EMMAD-EDIT in the E-commerce domain based on the Chinese E-commerce video captioning dataset E-MMAD. The data sample in E-MMAD dataset consists of a product video, an advertising video description, and additional information including video titles and structure attributes.

It collects 120,984 videos with average duration of 30.4 seconds and the annotated Chinese description length is 67 words on average. The characteristics of long videos and long captions make it

Input Prompt: Forget everything you have been asked before, now you are a fluency assessment machine. Sentence fluency is defined as a sentence without spelling errors, grammatical errors, punctuation errors, etc. You don't have to worry about whether the sentence is redundant or not. You will be given several sentences in English, and you only need to determine whether the sentence is fluent but do not need to care about the specific content of the sentence. For each sentence, you must first state its number, then repeat the unchanged sentence, and further say whether it is fluent or not. If it is fluent, say “Yes”, otherwise say “No”. Please do not modify the original sentence, and do not output anything else.

Figure 12: Input prompts of ChatGPT fluency evaluation.

Table 8: Fluency evaluation of ground-truth sentences on the VATEX-EDIT test set. ChatGPT measures the fluency (YES/NO) rate of 100 sentences. Human measures the fluency rate (ranging from 1 to 5, 5 is the best) of 200 sentences.

	Constructed	Annotated
ChatGPT	87%	91%
Human	4.40	4.37

suitable for building a challenging VCE dataset. During annotation, we select data samples from E-MMAD dataset with relatively long descriptions. In addition to videos and descriptions, we also provide product structure information and video titles for reference.

C.1 EMMAD-EDIT dataset statistics

Table 9 shows breakdown data statistics of the EMMAD-EDIT dataset. It has overall 79,652 editing instances for 16,176 product videos. Diverse unique attributes and vocabulary also indicate data richness. We separate the abstract attribute-related data samples as an extra challenging subset. Considering realistic demands, we utilize all these data samples to construct “ $\langle add, -, attr \rangle$, add attributes in description” command cases. Note that in the Experiments section, we present the EMMAD-EDIT results training without the abstract subset data by default.

C.2 The Impact of Data Volume.

Considering the limited scale of manually collected data in EMMAD-EDIT, we analyze the results under different volume data with 4K, 8K, 12K samples. Figure 13 shows that a growing volume of data consistently increases the controllable scores. Breakdown analysis of multi-grained commands reveals that more challenging commands, e.g. $\langle add, -, - \rangle$, require higher volume of training data samples to get desired performance.

Table 9: Data statistics of EMMAD-EDIT dataset. $VTime$ denotes the average time length (seconds) of videos. Len_{Ref} denotes the average length of reference captions and Len_{GT} is the length of groundtruth captions. $Uni. Attrs$ means the vocabulary of annotated attributes. $Overall_{Abstract}$ is the challenging subset of abstract attribute-related samples.

Command	#Videos			#Editing instances			VTime	Len _{Ref}	Len _{GT}	Edit Dist	Uni. Attrs	Vocab
	Train	Val	Test	Train	Val	Test						
<add, -, - >	7,751	2,599	2,633	7,752	2,599	2,633	26.6	72.2	100.9	29.8	-	29,241
<add, pos, - >	3,221	1,077	1,094	3,221	1,077	1,094	26.3	91.6	99.4	8.7	-	18,640
<add, -, attr >	3,221	1,078	1,094	3,221	1,078	1,094	27.1	93.1	100.6	8.7	2,388	18,501
<add, pos, attr >	3,221	1,077	1,093	3,221	1,077	1,093	26.6	92.0	99.9	8.7	2,907	18,584
<del, -, - >	7,751	2,599	2,633	7,753	2,599	2,633	26.4	100.1	71.4	29.6	-	29,456
<del, pos, - >	3,221	1,078	1,095	3,221	1,078	1,095	27.1	102.1	86.1	17.2	-	18,734
<del, -, attr >	3,221	1,078	1,095	3,221	1,078	1,095	27.2	102.7	86.2	17.8	3,038	18,924
Overall _{Specific}	16,176	5,418	5,502	31,610	10,586	10,737	26.9	91.3	90.4	20.8	6,003	44,725
Overall _{Abstract}	15,955	5,328	5,432	15,959	5,328	5,432	26.8	90.9	100.9	11.4	648	44,347

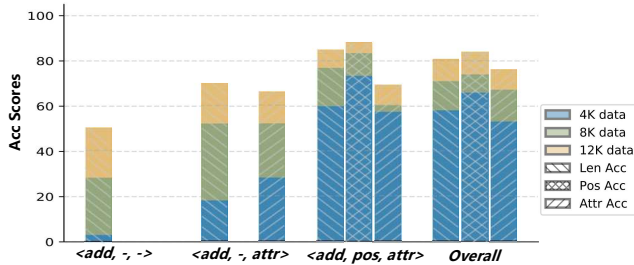


Figure 13: Overall and breakdown performance under different data volumes (4K/8K/12K data samples).

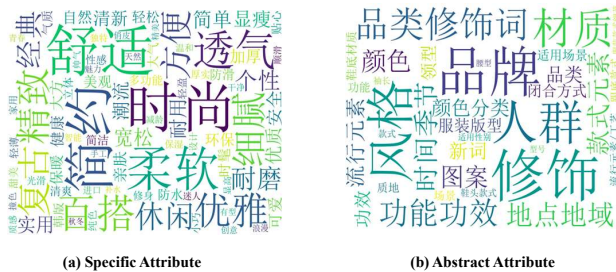


Figure 14: Word clouds of attributes on the EMMAD-EDIT dataset. The left shows the diversity of annotated specific attributes and the right shows the abstract attributes.

C.3 Data Visualization

We manually collect data samples that support the editing of two types of attributes: *specific* and *abstract*. *Specific* attributes directly appear in the reference caption to support straightforward content editing such as “comfortable” and “fashion”. Meantime, *abstract* attributes are more high-level concepts that may consist of multiple specific attributes. For example, “Style”, “Target People”, and “Time and Seasons” are annotated *abstract* attributes in the dataset. The

word clouds (shown in Figure 14) of *specific* and *abstract* attributes show the diversity of the EMMAD-EDIT dataset.

We visualize annotated samples of the EMMAD-EDIT dataset in Figure 15. Besides quadruples (*video*, *command*, *reference caption*, *edited caption*) that support video caption editing, we also provide additional product information such as structured information and the video title. The annotation interface during data construction is shown in Figure 16.

D Metric Details

To support evaluations in Chinese, we utilize a Chinese GPT-2 [10] to calculate Chinese sentence likelihood and then obtain PPL scores. For EMScore, we replace the core vision-language alignment model EN-CLIP [36] with CN-CLIP [59].

D.1 Position Accuracy Design

We propose a novel *Position Accuracy* (*Pos-Acc*) metric to measure whether models insert/remove the content in the specified positions under fine-grained editing control. The challenge of calculating *Pos-Acc* is the misalignment between the reference caption R and the edited caption Y . The two textual sentences have variable lengths and the edited operation may appear in the other positions to maintain the overall fluency. To tackle the above challenge, we propose a novel Dynamic Sequence Aligning (DSA) algorithm to align two variable-length textual sequences based on the absolute positions, inspired by classical Dynamic Time Warping (DTW) [29].

The pseudo-code is presented in Figure 17, which can align two variable-length text sequences in positional indexes, resulting in related spans $\{S_{m_1}, S_{m_2}, \dots, S_{m_K}\}$ in Y aligned to [MASK] tokens $\{m_1, m_2, \dots, m_K\}$ in R . We count the percentage of correct samples that insert/remove new content S_{m_K} in given position m_K .

As Figure 18 illustrates, we visualize two aligned cases in English and Chinese respectively obtained by the DSA algorithm.

E Conversion from Interface Signals to Triplet Format

Though the main focus of this paper lies in exploring the novel VCE task utilizing triplet commands, we also shed light on how to

Video ID: 200563409291.mp4	
	
Video Title	VH女包2020新款潮流单肩包时尚简约小鹿水桶包休闲女士手提斜挎包 <i>VH Women's Bag 2020 New Trend Single-shoulder Bag Fashionable and Simple Little Deer Bucket Bag Casual Ladies' Handheld Crossbody Bag</i>
Structured Info.	品类:斜挎包,单肩包,水桶包,女包;时间季节:2020;新品:新款;风格:时尚,休闲,简约,潮流,欧美时尚;修饰:手提,小鹿;人群:女士;上市时间:2018年春夏;款式:单肩包;里料材质:织物;背包方式:单肩斜挎手提;品牌:VANESSA HOGAN;颜色分类:香草白1,黑色,香草白,婴儿粉;皮革材质:牛皮; <i>Category: Crossbody Bags, Shoulder Bags, Bucket Bags, Women's Bags; Season: 2020; New Product: New Model; Style: Fashionable, Casual, Simple, Trendy, European and American Fashion; Embellishment: Handheld, Little Deer; Target Audience: Women; Launch Date: Spring/Summer 2018; Style: Shoulder Bag; Lining Material: Fabric; Backpack Method: Single Shoulder Crossbody Handheld; Brand: VANESSA HOGAN;</i>
Command	<add, -, ->
Ref Cap	经典水桶包, 外型酷似水桶, 包身圆润又不失俏皮, 别致的字母包设计, 凸现自我的创意个性与个人色彩。 <i>The classic bucket bag has a cool bucket-shaped appearance, with a rounded and nifty body. The unique design of the parent-child bag highlights the creative personality and individual style.</i>
Edited Cap	经典水桶包, 字母包设计, 再现繁盛时代的小鹿包。外型酷似水桶外型, 包身圆润又不失俏皮的造型, 别致的字母包设计, 追求简约实用的时尚, 凸现自我的创意个性与个人色彩。 <i>The classic bucket bag with a parent-child bag design brings back the prosperous era of little deer bags. Its cool bucket-shaped appearance, rounded and nifty body, and unique design pursue minimalist and practical fashion, highlighting the creative personality and individual style.</i>
Command	<del, -, “实用, 俏皮”> <del, -, 'nifty, practical'>
Ref Cap	经典水桶包, 字母包设计, 再现繁盛时代的小鹿包。外型酷似水桶外型, 包身圆润又不失俏皮的造型, 别致的字母包设计, 追求简约实用的时尚, 凸现自我的创意个性与个人色彩。 <i>The classic bucket bag with a parent-child bag design brings back the prosperous era of little deer bags. Its cool bucket-shaped appearance, rounded and nifty body, and unique design pursue minimalist and practical fashion, highlighting the creative personality and individual style.</i>
Edited Cap	经典水桶包, 字母包设计, 再现繁盛时代的小鹿包。外型酷似水桶(外型), 包身圆润(又不失俏皮的造型), 别致的字母包设计, 追求简约(实用的)、时尚, 凸现自我的创意个性与个人色彩。 <i>The classic bucket bag with a parent-child bag design brings back the prosperous era of little deer bags. Its cool bucket-shaped appearance, rounded (and-nifty) body, and unique design pursue minimalist (and-practical) fashion, highlighting the creative personality and individual style.</i>
Command	<add, -, “风格”> <add, -, 'style'>
Ref Cap	经典水桶包, 字母包设计, 再现繁盛时代的小鹿包。外型酷似水桶, 包身圆润又不失俏皮, 别致的字母包设计, 追求实用, 凸现自我的创意个性与个人色彩。 <i>The classic bucket bag with a parent-child bag design brings back the prosperous era of little deer bags. Its cool bucket-shaped appearance, rounded and nifty body, and unique design pursue practical usage, highlighting the creative personality and individual style.</i>
Edited Cap	经典水桶包, 字母包设计, 再现繁盛时代的小鹿包。外型酷似水桶外型, 包身圆润又不失俏皮的造型, 别致的字母包设计, 追求简约实用的时尚, 凸现自我的创意个性与个人色彩。 <i>The classic bucket bag with a parent-child bag design brings back the prosperous era of little deer bags. Its cool bucket-shaped appearance, rounded and nifty body, and unique design pursue minimalist and practical fashion, highlighting the creative personality and individual style.</i>

Figure 15: Data samples of annotated EMMAD-EDIT dataset.

convert the prevalent interface signals, i.e. natural language and handwriting editing trajectories, to the triplet formatted controls.

E.1 From Natural Language to Triplets

To convert natural language signals, we can adopt fuzzy matching to recognize add or delete operations and use text parser tools to get specific semantic roles of a sentence. We can also leverage the ability of LLMs such as ChatGPT and LLaMA-2, which already have outstanding language understanding and summarization capabilities. Specifically, we can design suitable prompts to convert text sentences into triplets to meet the pre-defined requirements.

E.2 From Triplets to Natural Language

The triplet command in our annotated datasets can be easily converted to natural languages to support more diverse scenarios and applications. For example, using natural language to explore the capability of LLMs in the VCE task (A). Specifically, we design a series of templates shown in Table 10 to convert the triplet command to

natural language. We will also release the two benchmark datasets with user commands both in the triplet and natural language format to benefit the community.

E.3 From Handwriting-revision Traces to Triplets

The recent advancement of GPT-4Vision (GPT-4V) has shown its powerful capability of Visual Referring Prompting [61]. GPT-4V can well understand visual pointers (such as circles, arrows or traces) directly drawn on images, therefore, revealing a novel human-model interaction method called “visual referring prompting”. Combined with its accurate OCR capability, GPT-4V can serve as an ideal tool to input the handwriting editing traces as an image and convert it to the triplet control output, as illustrated in Figure 19.

F Future Directions

In this paper, we make the first attempt to propose the novel Video Caption Editing (VCE) task and collect two benchmark datasets



Figure 16: Annotation interface for constructing the EMMAD-EDIT dataset.

```

1 # R is the word sequence of Reference Caption
2 # Y is the word sequence of Edited Caption
3 def DynamicSeqAlign(R, Y):
4     # Initialize DTW distance matrix
5     DTW = {}
6     for i in range(len(R)):
7         DTW[(i, -1)] = float('inf')
8     for j in range(len(Y)):
9         DTW[(-1, j)] = float('inf')
10    DTW[(-1, -1)] = 0
11    # Initialize the aligned path matrix
12    PATH = {}
13    # Dynamic programming to align R and Y
14    for i in range(len(R)):
15        for j in range(len(Y)):
16            # get distance between word i and word j
17            dist = get_dist(R[i], Y[j])
18            # get the min distance from last time step
19            min_dist = min(DTW[(i-1, j)], DTW[(i, j-1)], DTW[(i-1, j-1)])
20            # update DTW matrix for current time step
21            DTW[(i, j)] = dist + min_dist
22            # record the related distance path
23            PATH = update(PATH)
24    # get the minimum distance path aligning R and Y
25    best_path = trace(PATH)
26    # filter repetitive words in Y
27    align_path = filter(best_path)
28    return align_path
29
30 # get distance between two words from R and Y
31 def get_dist(R_word, Y_word):
32     if R_word == Y_word: # same words
33         return 0
34     elif R_word == '[MASK]': # ([MASK], Y_word)
35         return 100
36     else: # unmatched words
37         return float('inf')

```

Figure 17: Pseudo-code for Dynamic Sequence Aligning in a Python-like style.

to support it. We further propose a unified framework OPA for VCE and compare it with a ChatGPT pipeline. Based on the task, dataset and method foundation, we aim to make a good start for the community. There are various interesting aspects worth exploring in the future:

- **A versatile system for video captioning and editing.** The dense annotation of the quadruple (*video, command, reference caption, edited caption*) in our dataset has the potential to support building a versatile system encompassing conventional video captioning, controllable video caption and video caption editing. For initialization, conventional video captioning can produce a generated description for a given video. subsequently, video caption editing can be utilized to update and revise the original description. When omitting the reference caption, the rest annotation (*video, command, edited caption*) can be adjusted to achieve controllable video captioning.
- **A more robust system for poor reference caption.** In the VCE task, the reference caption can be the edited caption from the last round to fulfill multi-round editing. Furthermore, it can also be extended with human-written drafts or machine-generated captions. Under these circumstances, the robustness of the VCE system to low-quality reference captions should be taken into consideration.
- **The abstract-attribute subset of EMMAD-EDIT is under-explored.** The abstract-attribute subset (details in Sec 4.2.) involves abstract attributes that do not directly appear in the reference caption. It requires models to understand and reason the video content at a higher semantic level. In the experiments, we only assess the performances of OPA and

Command: <add, pos, attr>		
EN Case	Ref Cap	“A woman gives a demonstration [MASK-1] to come to her [MASK-2] sessions [MASK-3].”
	Edited Cap	“A woman is giving a demonstration and invitation to come to her gong therapy sessions.”
	Pos Align	([MASK-1], “and invitation”, ✓); ([MASK-2], “gong therapy”, ✓); ([MASK-3], None, ✗)
CN Case	Ref Cap	“此款皮鞋精选 [MASK-1] 头层牛皮，采用擦色 [MASK-2] 打磨工艺；鞋帮选用松紧套脚，穿着上更加的便捷。”
	Edited Cap	“此款高帮皮鞋精选 优质 头层牛皮，采用擦色 手工 打磨工艺；鞋帮选用 的是 松紧套脚，穿着上更加的便捷。”
	Pos Align	([MASK-1], “优质”, ✓); ([MASK-2], “手工”, ✓)

Figure 18: Cases of aligned reference captions and edited captions by the proposed Dynamic Sequence Aligning Algorithm.

ChatGPT pipeline on the easier *specific* subset, whose performances are not yet so desirable, the *abstract* subset therefore remains a challenge for future exploration.

- **Serve as a touchstone for video large language models (VidLLMs).** The recent emergence of VidLLMs such as GPT4-Vision [32], VideoChat [16], and VideoChatGPT [28] has ignited sparks for a generalist video assistant. However, existing research also points out their limitations for long video understanding [12], visual hallucination [45] and multi-modal instruction following capability [66]. The VCE benchmark can serve as a new multi-modal evaluation for assessing both long video understanding and multi-grained text editing abilities.

G Social Impacts

The proposed Video Caption Editing (VCE) task addresses significant challenges in the realm of video content management, offering far-reaching social implications.

Enhanced Accessibility. By allowing for multi-grained user control over video captions, our system can significantly improve accessibility for individuals with disabilities. Users who rely on captions for understanding video content, such as the hearing-impaired, can benefit from more precise and customizable descriptions tailored to their specific needs.

E-commerce and Marketing. The integration of our VCE system into e-commerce platforms can revolutionize how product videos are presented. By enabling dynamic and detailed video descriptions that cater to individual consumer preferences, businesses can enhance user engagement and improve product understanding.

User Empowerment and Creativity. The ability to edit video captions interactively empowers users to tailor content to their specific and creative desires. This democratizes content creation, allowing individuals and small creators to produce high-quality, customized video content without requiring extensive technical expertise.

Table 10: Defined templates that conveniently convert triplet commands to natural language format. The '}' represents the placeholder of specific attributes.

Command	Conversion Template
<del, -, attr >	Delete contents about '}' from the reference description.
<add, pos, attr>	Add contents about '}' at [POS].
<add, -, attr>	Add contents about '}' to expand the reference description.
<add, pos, - >	Add video-related contents at [POS].
<del, pos, - >	Contents have been deleted at [POS], please make the rest sentence fluent and coherent.
<del, -, - >	Shorten the reference description.
<del, -, attr>	Expand the reference description.

You are a writing-revision trace conversion assistant. Given an image with red hand-written edited traces from users, please help me to summarize the user editing intention to a triplet (operation, position, attribute). In the triplet, operations mean the overall description length editing (add or delete), positions specify the edited locations which can be indicated by a [POS] token as a placeholder, and attributes mean the the specific semantic contents change of original descriptions. I will give you several cases to understand the conversion from edited traces to a summarized edited triplet.

Guidelines:

- The different combinations of operation, position, and attribute elements in the editing triplet can cover multi-grained realistic demands from coarse-grained controls (add/delete some contents) to fine-grained controls (add/delete specified attribute-related contents at specified positions). The *operation* and *attribute* elements in the triplet can be omitted and denoted as "None".
- In some complex editing scenarios that can not convert to a single editing triplet, you can split the complex editing into firstly delete and then add operations and output multiple related triplets as Example 5 shows.

Example 1:

Operation: delete
Position: These [POS] knit sweatpants are soft and comfortable.
Attribute: popular

Example 2:

Operation: add
Position: These popular knit sweatpants are soft and comfortable [POS].
Attribute: None

Example 3:

Operation: add
Position: None
Attribute: targeted customer

Example 4:

Operation: add
Position: These popular knit sweatpants are soft and comfortable, suitable for [POS] wearing in summer.
Attribute: man

Example 5:

Operation: delete
Position: men are climbing [POS] with harnesses and safety helmets [POS] boots.
Attribute: frozen waterfalls; ice

Operation: add
Position: men are climbing with harnesses [POS] and safety helmets boots.
Attribute: environment

Now, based on the above examples and task definitions, please help me convert the revision traces in the uploaded image to a triplet. The output triplet format should be the same as the output in the examples.

New Input:

GPT4-V Output:
Operation: add
Position: a worker uses [POS] a high speed grinder to work on a large piece of metal.
Attribute: protective gear

New Input:

GPT4-V Output:
Operation: delete
Position: an overhead view of a golf course [POS] and an older man examines a younger golfer perform a stroke .
Attribute: is shown

Operation: add
Position: an overhead view of a golf course and an older man examines a younger golfer perform a stroke [POS].
Attribute: sky

Figure 19: Designed prompts and output instances of GPT-4V to convert the human-writing editing trace as an image to the triplet format control.