DOES YOUR VIDEO-LANGUAGE MODEL ACTUALLY UNDERSTAND THE LANGUAGE INPUT?

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

Paper under double-blind review

Abstract

Driven by the wave of Large Language Models (LLMs), Video-Language Models (VLMs) have become a significant yet challenging technology to bridge the gap between video and text. Although previous VLM works have made significant progress, almost all of them implicitly assume that all the texts are predefined by the specific template. In real-world applications, such an assumption is impossible to satisfy, since predefining all the texts is extremely time-consuming and laborintensive. Besides, these predefined text inputs are too strict and user-unfriendly, limiting their applications. It is observed that given a video input, texts with similar semantics lead to various performances. To this end, in this paper, we propose a novel text-augmented VLM method to improve video-text fusion by text rewriting. Specifically, we first generate various text samples from the original ones based on the pre-trained LLM to target specific text components. A multi-level contrastive learning module is designed to mine the coarse-grained language information. Moreover, we also propose an attribute-based text reasoning strategy to learn fine-grained textual semantics. Extensive experiments on many videolanguage tasks show that the proposed method can serve as the plug-and-play module to effectively improve the performance of state-of-the-art VLM works.

028 1 INTRODUCTION

Due to remarkable success, Large Vision-Language Models (LVLMs) have attracted more and more attention (Tian et al., 2024; Fan et al., 2024; Kim et al., 2024). LVLMs require cooperation from both computer vision and natural language processing for precise semantic alignment and have a wide range of applications such as video summarization and video question answering. Benefiting from the strong knowledge integration ability in large language models (LLMs), LVLMs show superior performances in solving complex image-language tasks by utilizing appropriate human-instructed prompts (Hakim et al., 2023; Duan et al., 2024; Jung et al., 2024). Since the real-world videos contain much temporal information, LVLMs still have difficulty to handle the real-world videos. Besides, the sentence text is the most important input that accompanies the video due to its humanfriendly and descriptive nature.

Current video-language models contain three main popular yet challenging tasks: video question 040 answering (VideoQA) (Gao et al., 2023; Yu et al., 2024), video sentence grounding (VSG) (Zhang 041 et al., 2023b; Qi et al., 2024) and video-text retrieval (VTR) (Zhu et al., 2023; Zhang et al., 2023a). 042 Video QA is a significant multi-modal task where a model is given a video along with a natural 043 language question about the video content, and it must generate or select the correct answer. The task 044 requires the model to understand the visual cues in the video, as well as the language of the question, 045 to provide relevant and accurate responses. Given a language text and an untrimmed video, VSG 046 aims at retrieving the start and end timestamps of the target video moment, semantically according 047 to a sentence text. Given a language text, VTR targets to retrieve relevant videos from a large video 048 database, which can be either text-based (text-to-video retrieval) or video-based (video-to-video retrieval). The significant goal of VTR is to find videos that best match the given input by analyzing visual content, actions, and sometimes audio cues. Their performance in these downstream tasks 051 depends on their capability to extract video features and align them with text features. Some methods only perform data augmentations on the video input to improve the model robustness during training 052 in every epoch. In contrast, existing methods only utilize predefined texts without any augmentation. In real-world applications, these sentence texts with similar textual semantics might be inputted

055

057

058

060

061

062

063

064

065

066

067

068

069



Figure 1: (a-c) Example of the VLM tasks (VSG, VideoQA and VTR), where our method can serve as a plugand-play module for previous VLM models to enhance their efficiency. (d) Pipeline of our proposed method.

071 with different structure/vocabulary variations from various users. As shown in Figure 1(b), the text ("Person pours water into a glass") shares the same semantics as the text ("Water is poured into 072 a glass by person"). However, previous methods yield dissimilar grounding results in Figure 1. 073 The main reason is that these methods cannot utilize their weak text encoder to learn discriminative 074 textual representations, which illustrates the significance of handling the text variations. Therefore, it 075 is important to ensure that the designed VLM-based model is robust enough to deal with various texts 076 with different templates. However, existing language augmentation approaches are not sufficiently 077 effective to integrate the multi-modal inputs. Some methods target to replace or mask some words in a sentence, which only brings limited influence in diversifying the text structure/vocabulary. It 079 is comparatively weaker than video augmentations. The target language augmentation approach should effectively rewrite sentence texts while reserving the core textual semantics. The approach is 081 urgently required for model training to achieve the best results.

082 In this paper, we propose a simple yet highly effective framework to improve the robustness and 083 performance of VLMs. Specifically, we leverage the large language models (LLM) to generate 084 multiple variants of each text in the video-text pairs. To obtain different text types, we generate a 085 small set of variation-origin by two strategies: LLM-based datasets and existing original datasets. After obtaining the variation-origin pairs, we utilize them as examples to prompt LLM model for 087 diversifying all the texts in the VSG datasets. Different from previous sentence augmentation works 880 that only change some words to preserve sentence structures, LLM has a strong language processing ability to generate rich variations for diverse text inputs due to their extensive training datasets and 089 emergent properties. Based on the above sentence augmentation, each video corresponds to diverse 090 texts. Moreover, we introduce GPT (Radford, 2018) to generate various hard negative texts from the 091 original (anchor) texts by changing different sentence parts. In particular, we utilize precise prompt 092 engineering to modify specific parts of the sentence with the rest parts unchanged. Also, we generate 093 positive samples that lie relatively far from the anchor in the embedding space. To further understand 094 the latent textual semantics, we design an attribute-based text reasoning strategy for fine-grained text mining. To analyze the relative significance of each sentence part, we incorporate these generated 096 samples by a weighted contrastive loss function. With these diverse texts, we target to train VLM models with augmentation from the text perspective. For convenience, we randomly choose a text 098 augmentation from multiple diverse texts.

099 Our main contributions are summarized as follows: 1) We make the first attempt to explore the 100 effect of the template-free text for the robust VLM task, where we localize the target activity by a 101 user-friendly text with any form instead of a predefined text. Also, we propose a novel framework 102 that utilizes the LLMs to generate positive and negative texts, where each negative text is used to 103 highlight a sentence component. 2) To obtain diverse positive and negative texts, we augment the 104 text by both word-level and structure-level for rewriting texts. To selectively integrate these gen-105 erated texts, we design two modules (generating multi-level texts module and attribute-based text reasoning module) to understand the text input from different granularities. Besides, a weighted 106 contrastive loss is introduced to integrate these sentence components by assigning adaptive weights 107 to these components. 3) For three downstream tasks (video sentence grounding, video question an-



Figure 2: Illustration of our proposed framework. Best viewed in color.

swering and video-text retrieval), we conduct experiments on many popular yet challenging datasets.Extensive results show that our proposed model outperforms existing approaches by a large margin.Moreover, our method can serve as a plug-and-play module for state-of-the-art VLM methods.

2 RELATED WORKS

130 Large vision-language models. The breakthrough of LLMs in language-oriented tasks (Ma et al., 131 2024; Du et al., 2024) and the emergence of GPT-4 have prompted researchers to explore the poten-132 tial of LLMs in assisting with a range of tasks across multi-modal scenarios (Carolan et al., 2024; 133 Yin et al., 2024). This has led to the development of a new field, namely large vision-language 134 models (LVLMs). A variety of strategies and models have been proposed to address the discrepancy 135 between text and other modalities. Some works employ learnable texts to extract visual information 136 and generate language using LLMs conditioned on the visual features. Models including GPT-40, 137 MiniGPT-4 and LLaVA learn simple projection layers to align the visual features from visual en-138 coders with text embeddings for LLMs. Additionally, parameter-efficient fine-tuning is adopted 139 by introducing lightweight trainable adapters into models. Several benchmarks have verified that LVLMs demonstrate satisfactory performance on visual perception and comprehension. 140

Although these methods have achieved promising results, all of them heavily rely on correctly aligned multi-modal datasets. Therefore, it is highly expected to develop a VLM model that is robust to different texts with similar semantics, which has not been studied as far as we know. Thus, *we make the first attempt to reveal the text understanding problem in VLM task and propose to eliminate the negative impact of the different texts with any template.* More details in Section A.2.

146 147

108

109

110

111 112

113 114

115

116

117

118

119

121

122

123

125

126

127 128

129

3 Methodology

148

We elaborate on the proposed method, which strengthens the text encoder to obtain consistent representations for various semantically similar texts in real-world multi-modal datasets (*e.g.*, Charades-STA (Sigurdsson et al., 2016)), multiple semantics-similar texts often share a video moment with the target activity. For example, "Person opens the door" and "The door is opened by person" have similar semantics. Since the text template is fixed, it is still challenging to diversify the text input. Thus, we design a text augmentation module to generate semantically similar texts. The overall framework is shown in Figure 2.

Problem statement. Due to the strong language processing ability of of LLMs, we utilize LLM to generate various texts by replacing different components for simulating the practical labeling process in the format-free setting. We denote $\{Q_1, \ldots, Q_M\}$ as the textual input set in the VLM task, where *M* denotes the total number of sentences. Previous VLM methods (Yu et al., 2023; Wang et al., 2022b; 2024) cannot well handle these texts with similar semantics since they do not fully understand the textual information in the sentence. To address the existing models' limitations in correlating major sentence parts with suitable video representations, we present a novel plug-andplay method for generating negative and positive samples targeting specific sentence parts. These
 samples facilitate improved perception of specific parts of the sentence, eventually enhancing the
 understanding of video-language correlation. We use the generated samples as auxiliary samples
 alongside the original training samples by employing a novel weighted contrastive loss. The proposed approach is application-agnostic and can be adopted successfully in the multi-modal task.

- 167
- 168 169

170

171

172

173

174

175

189 190

3.1 GENERATING MULTI-LEVEL TEXTS FOR COARSE-GRAINED LANGUAGE ALIGNMENT

By treating original texts as anchors, we target to leverage LLMs for generating positive and negative texts to fully understand the texts, where we regard the generated text sharing similar semantics as a positive text; otherwise, the generated text is negative. For the *m*-th text, we denote the generated positive text as P_m and the negative text as N_m . To obtain diverse generated texts, we adopt three text rewriting approaches: human-based rewriting, chat robot-based rewriting and open-source LLM-based rewriting in Section A.1. For convenience, we term the pair of generated text and

original text as variation-origin text pairs.

176 Fighth text us variation origin text pairs. 177 Real-world multi-modal datasets include multiple video-text pairs (V,Q), where V denotes the 178 video and Q denotes one of corresponding texts. These texts differ in two levels: word level and 179 sentence structure level. Therefore, we have two types of text rewriting by a two-step rewriting pro-180 cess: word-level rewriting and structure-level rewriting. For convenience, we take the positive text 181 augmentation as an example.

Word-level augmentation. In the first step, we directly rewrite the original text by changing some words. With pre-trained LLMs, we can rewrite all the texts by the following prompt: $P'_{i,m} \leftarrow$ LLM(Q_m , "Rewrite the text 'text' concisely by changing the *i*-th word while keeping the meaning"), where text is substituted with the given text, and we utilize the underlined text to prompt our model for producing morphologically diverse text expressions.

To evaluate the significance of q_i on the semantics of Q_m , we can evaluate the semantic change before and after removing this word:

$$S_1(q_i, Q_m, c) = 1 - \cos(c \cup Q_m, c \cup Q_m \setminus \{q_i\}), \tag{1}$$

where c denotes the prompt and $cos(\cdot, \cdot)$ is the cosine similarity function. In real-world applications, larger $S_1(q_i, Q_m, c)$ denotes that removing q_i will lead to significant semantic changing, indicating that q_i is more relevant.

194 Structure-level augmentation. Since different users tend to utilize various text structures for video 195 grounding, we need to augment the text structure for more diverse texts. Similarly, we utilize the 196 following prompt for structure-level rewriting: $P_{i,m} \leftarrow \text{LLM}(P'_{i,m}, \text{``Rewrite the text `text' concisely}$ 197 by changing the text structure while keeping the meaning'').

Given a sentence Q_m , we define the sentence-level relevance of Q_m^i as the probability-weighted semantic similarity with other sentences:

$$S_2(Q_m^i, Q_m^j, c) = \sum_{j=1, j \neq i} \cos(Q_m^i, Q_m^j) p(Q_m^j | c),$$
(2)

where $p(Q_m^j, c)$ denotes the generative probability that provides more confidence to Q_m^j , and higher $p(Q_m^j, c)$ makes Q_m^j more acceptable. An intuitive observation is that if a sentence is semantically consistent with other sentences, the sentence is more convincing and more representative.

Similar to positive text augmentation in equation 1 and equation 2, we generate negative texts by changing their words and sentence structure. Thus, based on the multi-level language rewriting, we can conduct coarse-grained language alignment.

209 210 211

201 202

3.2 ATTRIBUTE-BASED TEXT REASONING FOR FINE-GRAINED LANGUAGE ALIGNMENT

In fact, Section 3.1 only considers the semantics of the sentence itself, ignoring the latent information
of the sentence. For example, "a person is driving a car" contains two significant objects: "person"
and "car". "person" corresponds to the following attributes: a head, two eyes, two arms, etc, while
the attributes of "car" include: four wheels, a steering wheel, etc. There attributes will assist VLMs to understand videos and texts for bridging the visual and textual gap.

216 Attribute generation. For some semantically similar sentences, they always have similar attributes. 217 Therefore, we generate the attributes for all positive and negative texts. Although embedding at-218 tributes can help us to understand the sentence, current VLM models cannot fully understand the 219 latent semantics. For example, "a person is driving a car" and "a car is running on the road" have 220 similar semantics. Therefore, rather than directly using the original sentence, we design a model with high confidence in visual attributes. Two intuitions are considered in this model: 1) different 221 from the original sentence, aligning explicitly with visual attributes can push the deigned model to 222 mine the inherent semantics in the given sentence. 2) visual attributes contain more fine-grained features, which can provide more details for cross-modal reasoning. 224

Firstly, we utilize video and text encoders to extract the video and text features. Since our framework is plug-and-play, it does not depend on specific feature encoders. For the fair comparison, we adopt the same video and text encoders with compared methods. For the text Q with J words, we denote word-level text feature as $f^W = \{f_j^w\}_{j=1}^J \in \mathbb{R}^{J \times d}$ and the sentence-level text feature as $f^q \in \mathbb{R}^d$, where d is feature dimension. Similarly, we denote the extracted video features as $f^V = \{f_i^v\}_{i=1}^{N_v} \in \mathbb{R}^{N_v \times d}$, where N_v is the frame number.

231 Attribute sampling. We 232 find that some generated at-233 tributes have a stronger se-234 mantic correlation with vi-235 sual features than others, and 236 some attributes have less significance, which will lead 237 to high computational cost. 238 Therefore, removing some 239



Figure 3: Our attribute selection module.

low significance can not only decrease the computational cost but also improve the model gener-240 alization. As shown in Figure 3, we address the problem by selecting effective attributes from an 241 attribute pool. Two main criteria are utilized during the attribute selection: Firstly, we prioritize at-242 tributes that are both representative and non-redundant. Secondly, we seek attributes with the highest 243 semantic relevance to the images when compared to other attributes. Finally, we use the following 244 steps for attributes: 1) For the attributes a_m associated with sentence Q, we partition them into N_c 245 clusters based on their feature similarity. This clustering strategy aims to ensure that each cluster 246 represents a distinct aspect, e.g., color or shape, in the descriptions. 2) In each cluster, we rank the 247 attributes by assessing their similarity to visual features, and select the one with the highest rele-248 vance. By the above strategy, the following attributes will be filtered out: non-visual attributes and incorrect visual attributes that are semantics-unrelated to the videos. To obtain the optimal attributes, 249 we introduce the following attribute selection strategies: 250

251 Semantic-based selection. Firstly, we want to make the sentence text has the similar semantics 252 with its generated attributes. Since the NLI model can mine the relationship between texts and the 253 attributes by inferring the logical entailment, we introduce an NLI-based binary filter (f_{nli}) as a 254 critic, and discard the pairs which do not achieve the entailment score over the threshold γ_1 :

$$O_1(x,y) = \mathbb{1}\{f_{nli}(x \Rightarrow y) \ge \gamma_1\},\$$

257 where x denotes the input, and y means the output.

Format-based selection. When we rewrite the given sentence, we need to make the format of the given sentence various, and preserve its original meaning. Thus, we want to filter the originvariance pair to learn the format-free dissimilarity. Especially, two metrics are used to evaluate the dissimilarity: 1) the token overlap between different sentences and 2) their syntactic difference. For the first, we filter the pairs with a higher Rouge-L (Lin, 2004) than a threshold γ_3 . As for the syntactic difference, we first parse the constituency tree of the origin and variety, and then filter the pairs based on their tree edit distance:

255 256

$$O_2(x,y) = \mathbb{1}\{D_t(x,y) \ge \gamma_2 \land f_{rou}(x,y) \le \gamma_3\},\$$

where $D_t(\cdot, \cdot)$ denotes the tree edit distance. In equation 3.2, the two dimensions of dissimilarity complement each other. On the one hand, $f_{rou}(\cdot, \cdot)$ promotes lexical divergence in each pair. On the other hand, $D_t(\cdot, \cdot)$ can be used to preempt "hacking" the word-overlap metric by simply switching a few words in the source sentence with corresponding synonyms. **Diversity-based selection.** For sentence rewriting, we need a diverse range of generated sentences since the diversity of attributes can directly affect the robustness of the trained model. Therefore, we introduce a critic O_3 for the diversity. We define two pairs (x_1, y_1) and (x_2, y_2) to be duplicates when one pair entails another, either on the input side $(x_1 \Rightarrow x_2)$ or on the output side $(y_1 \Rightarrow y_2)$. In the diversity filter, we first cluster all entailing pairs, and then discard all but one with the largest entailment score. Thus, we can utilize the graph traversal for the diversity filter.

Based on the above critics, we can filter the attribute candidate pool \mathcal{A} into an updated pool \mathcal{U} :

$$\mathcal{U} = \{(x, y) | (x, y) \in \mathcal{A}, O_1 \land O_2 \land O_3(x, y) = 1\}.$$

278 279 280

281

296 297 298

306

307

313 314

315

316 317 318

319

276

277

3.3 WEIGHTED SENTENCE INCORPORATION FOR CROSS-MODAL FUSION

In fact, different words (e.g., noun, verb, and adjective) have distinct significance in text understand-282 ing. For instance, some adjectives are more important for video grounding in some cases, while 283 some verbs are more significant for distinguishing different target moments. Previous VLM meth-284 ods treat all sentence components equally, which might limit these methods to fully understand the 285 entire sentence. For example, if there is no adjective in the anchor text, the negative text with adjec-286 tives cannot contribute to our model since the adjective is not discriminative for the text. Thus, we 287 aim to analyze the relative significance of each word to adaptively integrate different words, where we adaptively predict the salience of sentence components for each anchor text. Without any super-289 vision, we can obtain the significance score which means which word is more significant for text 290 understanding. By the module, we can find an optimal integration strategy of sentence components, 291 which makes VLM selectively understand different sentence components for a given text.

Incorporating generated sentences. Based on these positive and hard negative samples, we can encourage the designed VLM models to distinguish the difference between different words in each sentence part. For supervising the VLM model to understand the text input, we introduce a contrastive loss based on three types of text input:

$$\mathcal{L}_{cl}^{i} = -\log \frac{\beta \cdot \exp[1/\tau \cdot \cos\left(f^{V}, g_{i}^{n}\right)]}{(1-\beta) \cdot \exp[1/\tau \cdot \cos\left(f^{V}, g_{i}^{p_{i}}\right)] + \beta \cdot \exp[1/\tau \cdot \cos\left(f^{V}, g_{i}^{n}\right)]},\tag{3}$$

where $\beta \in (0,1)$ is a parameter; g_i denotes the *i*-th text; $g_i^{n_{i,j}}$ and $g_i^{p_i}$ denote the negative text and the positive text, respectively; τ denotes the temperature parameter. By equation 3, we can enhance the effectiveness of the designed model by these generated auxiliary texts.

Weighted contrastive loss. Since the visual features have higher computational complexity, we generate the positive and negative texts only by the original text (*i.e.*, anchor text) without considering the video input. Since different words contribute variously to sentence understanding, we target to find the most discriminative word for better text understanding by the following loss:

$$\mathcal{L}_{CL}^{i} = \max(\mathcal{L}_{cl}^{i,1}, \mathcal{L}_{cl}^{i,2}, \dots, \mathcal{L}_{cl}^{i,C}).$$
(4)

For C contrastive losses $(\mathcal{L}_{cl}^{i,1},\ldots,\mathcal{L}_{cl}^{i,C})$, each contrastive loss computed by equation 3 corresponds to a specific negative text, where the corresponding sentence component is changed. In equation 4, the maximum of these decomposed losses corresponds to the sentence component that is most clearly identified. Considering the significance score in equation 1 and equation 2, we can obtain the finally weighted contrastive loss as follows:

$$\mathcal{L}_{weighted} = \sum_{i,j,m,c} S_1(q_i, Q_m, c) \cdot S_2(Q_m^i, Q_m^j, c) \cdot \mathcal{L}_{CL}^i.$$
⁽⁵⁾

Since our method is plug-and-play, we borrow the cross-modal fusion module from an open-source works into our framework, which is the base version of our method.

4 EXERIMENTS

Datasets. For a fair comparison, we utilize the following datasets for evaluation. 1) For the task, we utilize three datasets: ActivityNet Captions (Caba Heilbron et al., 2015), and Charades-STA (Sigurdsson et al., 2016) and TACoS (Regneri et al., 2013). 2) For the VTR task, we adopt two datasets: MSRVTT (Xu et al., 2016) and LSMDC. 3) For the VideoQA task, we use two datasets: NExT-QA (Xiao et al., 2021) and STAR (Wu et al., 2021). More details are placed in Section B.1.

Method		Without	text augme	entation			With t	ext augmen	tation	
Method	R@11	` R@5↑	R@10↑	MdR↓	MnR↓	R@1↑	R@5↑	R@10↑	MdR↓	MnR↓
			Text-to	o-video reti	rieval					
CLIP-ViT-B/32	16.0	72.9	8 2 2	2.0	14.2	40.1	(0.2	765	4.0	10.0
A-Pool (Gorti et al., 20	40.9	72.8	82.2 83 5	2.0	14.5	40.1	08.2	/0.5 91 /	4.0	18.9
CLIP-ViP (Xue et al. 2	023) 50.1	74.8	84.6	1.0	-	42.3	69.4	77.8	3.0	16.5
+Ours	51.7	75.3	85.9	1.0	11.8	50.4	73.6	84.5	1.0	12.4
T-MASS (Wang et al., 2	024) 50.2	75.3	85.1	1.0	11.9	42.1	68.9	79.2	2.0	15.8
+Ours	52.3	77.9	87.6	1.0	10.9	51.4	70.2	81.3	1.0	11.7
CLIP-ViT-B/16										
X-Pool (Gorti et al., 20	48.2	73.7	82.6	2.0	12.7	39.7	68.5	78.4	4.0	16.5
+Ours	50.7	76.2	85.2	1.0	12.4	48.9	75.3	84.0	1.0	13.8
CLIP-ViP (Xue et al., 2	023) 54.2	77.2	84.8	1.0	-	51.2	73.9	80.4	2.0	14.8
+Ours	50.8	79.4	85.9	1.0	10.5	53.0	77.8	84.2	1.0	12.5
1-MASS (wang et al., 2	52.7 54 Q	82.6	85.0 86.8	1.0	10.5	49.2 53.4	70.3 81.0	86.2	2.0	11.5
Tours	54.9	02.0	Video	to text retr	ioval	55.4	01.0	00.2	1.0	11.5
CLIP-ViT-B/32			video		icvai					
X-Pool (Gorti et al., 20	()22) 44.4	73.3	84.0	2.0	9.0	41.2	68.5	80.4	3.0	13.8
+Ours	45.8	76.4	87.3	1.0	7.5	42.7	74.5	86.0	2.0	8.3
UATVR (Fang et al., 20	023) 46.9	73.8	83.8	2.0	8.6	43.0	67.9	78.3	3.0	11.7
+Ours	49.7	75.6	86.4	1.0	7.3	47.8	74.0	83.9	2.0	7.8
T-MASS (Wang et al., 2	(024) 47.7	78.0	86.3	2.0	8.0	42.9	73.5	82.6	3.0	13.9
+Ours	51.5	79.9	89.8	1.0	6.4	49.5	78.1	87.5	1.0	8.2
CLIP-ViT-B/16		72.0	04.1	2.0	0.4	12.0	72.0	01.6	2.0	10.5
A-Pool (Gorti et al., 20	(122) 46.4 50.2	73.9	84.1 86 3	2.0	8.4 6.1	42.8	72.0	81.0 84.2	3.0	10.5
UATVR (Fanglet al. 20)23) 48.1	76.3	85.4	2.0	8.0	41.6	73.0	81.9	3.0	10.6
+Ours	50.9	77.4	90.5	2.0	6.8	48.9	76.3	87.9	2.0	7.6
T-MASS (Wang et al., 2	024) 50.9	80.2	88.0	1.0	7.4	48.3	75.6	84.9	2.0	8.9
+Ours	53.7	84.2	91.5	1.0	3.4	50.8	82.7	90.4	1.0	4.9
Fable 2: VideoQA	performance	e comparis	son on N	ExT-QA	dataset	, where	the value	e means t	he accu	racy of
providing the right a	iswer.	•								•
Method		# F	Witl	hout text a	ugmenta	tion	Wi	th text aug	mentatio	n
		# Frames	Tempor	al Caus	al Des	cription	Tempora	ıl Causa	1 Desc	ription
All-in-One (Wang e	t al., 2023)	32	48.6	48.0)	63.2	40.2	37.9	5	3.8
+Ours		32	50.1	51.9)	64.7	48.6	50.2	6	1.3
Just Ask (Yang et a	l., 2021a)	20	51.4	49.6	<u>,</u>	63.1	42.7	40.1	5	4.0

3	Ę	5	2
3	Ę	5	ļ
3	Ē	5	6

+Ours

MIST (Gao et al., 2023)

+Ours

HiTeA (Ye et al., 2022)

+Ours

InternVideo (Wang et al., 2022a)

+Ours

BLIP-2 (Li et al., 2023b)

+Ours

SeViLA (Yu et al., 2023)

+Ours

357 358

359

360 361 362

Table 3: Comparison Results on STAR VideoQA benchmark.

54.3

56.6

60.3

58.3

62.8

58.5

62.5

67.2

70.1

67.7

72.4

20

32

32

16

16

8

8

4

4

4

4

52.9

54.6

56.9

62.4

65.7

62.5

66.3

70.3

72.9

72.1

74.9

67.8

66.9

69.8

75.6

77.3

75.8

76.4

79.8

80.4

82.2

85.3

50.9

51.9

57.2

52.2

60.4

52.9

61.8

64.0

69.2

64.0

70.5

49.3

48.2

55.4

57.6

63.9

57.4

59.7

61.9

70.1

66.8

72.7

62.7

55.3

67.9

59.3

74.9

70.3

74.5

72.3

78.4

76.9

83.9

Method (Frames Number)		Without text	augmentation	1	With text augmentation			
Method (Frances Number)	Interaction	Sequence	Prediction	Feasibility	Interaction	Sequence	Prediction	Feasibility
All-in-One (Wang et al., 2023) (32)	47.5	50.8	47.7	44.0	42.9	48.5	44.0	40.2
+Ours (32)	48.3	51.9	49.6	45.7	47.9	51.3	48.7	44.3
MIST (Gao et al., 2023) (32)	55.5	54.2	54.2	44.4	50.7	51.4	50.2	38.4
+Ours (32)	58.6	59.5	58.4	47.0	57.0	56.3	57.2	45.8
InternVideo (Wang et al., 2022a) (8)	62.7	65.6	54.9	51.9	55.6	61.0	50.3	47.2
+Ours (8)	63.8	67.7	58.9	55.2	61.8	64.9	57.4	54.3
SeViLA (Yu et al., 2023) (4)	63.7	70.4	63.1	62.4	58.7	62.2	57.9	57.8
+Ours (4)	66.7	72.9	66.4	65.3	65.3	68.9	64. 2	63. 4
BLIP-2 (Li et al., 2023b) (4)	65.4	69.0	59.7	54.2	60.9	66.3	54.3	50.1
+Ours (4)	67.8	72.5	61.4	56.8	66.2	71.6	58.7	55.3

370 371

372 **Evaluation metrics.** For the VTR task, we utilize Recall at rank $\{1, 5, 10\}$ (R@1, R@5, and 373 R@10), Median Rank (MdR), and Mean Rank (MnR) for evaluating the retrieval performance. 374 For the VSG task, we evaluate the grounding performance by "R@n, IoU=m", which means the 375 percentage of queries having at least one result whose Intersection over Union (IoU) with ground truth is larger than m. In our experiments, we use $n \in \{1,5\}$ for all datasets, $m \in \{0.5, 0.7\}$ for 376 ActivityNet Captions and Charades-STA, $m \in \{0.3, 0.5\}$ for TACoS. As for the VideoQA task, 377 we introduce the following metrics: temporal, causal, description, interaction, sequence, prediction

	means v	veakiy-supe	rvised.						
		W	/ithout text	augmentatio	on		With text at	gmentation	l
Method	Туре	R@1,	R@1,	R@5,	R@5,	R@1,	R@1,	R@5,	R@5,
		IoU=0.3	IoU=0.5	IoU=0.3	IoU=0.5	IoU=0.3	IoU=0.5	IoU=0.3	IoU=0.5
	ActivityNet Captions								
2D-TAN	I FS	59.45	44.51	85.53	77.13	48.32	29.38	71.36	62.30
+Ours	FS	60.46	45.29	87.94	77.43	51.86	32.64	72.98	63.75
MMN	FS	65.05	48.59	87.25	79.50	55.30	31.76	74.88	71.52
+Ours	FS	66.05	49.31	89.75	81.27	58.76	33.08	75.33	73.59
G2L	FS	-	51.68	-	81.32	55.75	33.01	75.25	70.89
+Ours	FS	66.34	54.26	91.77	84.29	60.90	46.86	84.39	80.62
VCA	WS	50.45	31.00	71.79	53.83	31.74	25.37	46.98	42.76
+Ours	WS	51.72	33.19	72.85	55.11	32.99	28.56	48.31	44.07
WSTAN	WS	52.45	30.01	79.38	63.42	33.72	25.74	49.30	45.88
+Ours	WS	53.10	31.56	80.24	65.77	35.20	27.99	51.84	48.69
CNM	WS	55.68	33.33	-	-	35.72	28.95	50.06	48.72
+Ours	WS	56.11	34.08	81.09	67.34	39.56	31.77	52.88	51.99
	•			Chara	ides-STA				
2D-TAN	I FS	39.81	23.25	79.33	52.15	20.18	11.35	47.05	33.82
+Ours	FS	40.27	24.95	82.96	53.28	23.99	14.75	49.22	34.18
MMN	FS	47.31	27.28	83.74	58.41	25.33	18.80	45.97	35.08
+Ours	FS	49.07	29.32	85.06	60.13	26.87	22.48	46.03	37.85
G2L	FS	47.91	28.42	84.80	59.33	26.54	19.85	48.06	36.70
+Ours	FS	55.77	32.97	91.38	60.39	34.85	27.96	74.28	46.70
VCA	WS	38.13	19.57	78.75	37.75	17.87	12.39	45.70	22.13
+Ours	WS	40.95	20.31	80.42	39.26	18.63	15.72	46.17	23.88
WSTAN	WS	29.35	12.28	76.13	41.53	8.15	5.43	35.27	11.86
+Ours	WS	30.24	14.06	77.35	42.99	10.77	6.92	37.40	13.88
CNM	WS	35.15	14.95	-	-	14.34	9.65	43.88	18.79
+Ours	WS	35 72	16 33	76 52	43 18	16.83	12.05	45 60	21 64

9	Table 4: VSG performance comparison under official train/test splits, where "FS" denotes "fully-supervised"
)	and "WS" means "weakly-supervised".



Table 5: Figure: Performance comparison with state-of-the-art methods on the TACoS for the VSG task, where left figure compares the effectiveness (R@5, IoU=0.5) and the efficiency (QPS), right figure shows that our method can serve as a plug-and-play module to enhance their efficiency. Table: Efficiency comparison for VSG on TACoS without text augmentation. "Aug time" denotes the time of generating multi-level texts.

Table 6: Main ablation study on the VSG task with G2L as the base model, where we remove each key individual component to investigate its effectiveness.

		ActivityNe	et Captions		Charades-STA			
Model	R@1	R@1	R@5	R@5	R@1	R@1	R@5	R@5
	IoU=0.3	IoU=0.5	IoU=0.3	IoU=0.5	IoU=0.5	IoU=0.7	IoU=0.5	IoU=0.7
Ours(a)	53.77	40.28	76.94	72.25	28.51	20.34	67.85	38.71
Ours(b)	55.35	42.03	79.50	74.91	30.88	23.92	70.66	41.58
Ours(c)	57.63	43.86	81.34	77.99	32.50	24.03	71.76	42.92
Ours(full)	60.90	46.86	84.39	80.62	34.85	27.96	74.28	46.70

and feasibility. In these metrics, lower MdR and MnR denotes better performance. For the metrics, higher value means better performance. Bold denotes the best performance.

Implementation details. For video encoding, we utilize the 112×112 pixels shape of every frame of videos. As for the text encoder, we feed the texts to a pre-trained txt encoder to embed word-level features. The dimensions d of video and text tokens are 512. We set $\gamma_1 = 0.8, \gamma_2 = 0.6, \gamma_3 = 0.7$ and $\mu = 0.6$ in our experiments to achieve the best performance. We train our model for 100 epochs with an Adam optimizer with the learning rate 3×10^{-4} .



Figure 4: Training performance of each ablation module with text augmentation on the ActivityNet Captions dataset (left, VSG), the NexT-QA dataset (middle, VideoQA) and the MSR-VTT dataset (right, VTR).

4.1 PERFORMANCE COMPARISON

Following previous open-source methods, we directly cite the corresponding results from compared methods. In this paper, we treat our as the plug-and-play module for state-of-the-art VLM models to improve their performance.

Performance comparison on the VTR task. VTR is a challenging multi-modal task, which re-449 quires the designed model can effectively bridge the gap between videos and texts. In this paper, 450 we consider two subtask: text-to-video retrieval and video-to-text retrieval. Table 1 illustrates the 451 effectiveness of our model as the plug-and-play module for previous VTR methods. We can find 452 that when using augmented text, all the compared methods suffer performance degradation. The 453 core reason is that previous VTR methods pay less attention to the language input, and ignore much 454 language information in the sentence query. By using our model as the plug-and-play module, pre-455 vious method can obtain significant performance improvement since our proposed model can fully 456 mine latent language semantics.

Performance comparison on the VideoQA task. Similar to the VTR task, we conduct performance comparison VideoQA performance comparison. The experimental results are summarized in Table 2 and Table 3, where the performance of previous methods was unsatisfactory. The key reason is that previous methods have difficulty in understanding the rewritten question. Different from them, we can explore more deep and fine-grained language information by attribute-based text reasoning.

462 Performance comparison on the VSG task. We conduct VSG performance comparison on all three 463 datasets with official train/test splits under both fully-supervised (Gao et al., 2017; Li et al., 2023a; 464 Liu et al., 2018; Li et al., 2023c; Yuan et al., 2019a; Zhang et al., 2019b; 2020b; Zeng et al., 2020; 465 Gao & Xu, 2021; Zhang et al., 2021; Gao et al., 2021; Wang et al., 2022b) and weakly-supervised 466 setting (Chen et al., 2022; Yang et al., 2021b; Zhang et al., 2020c; Wang et al., 2021b;a; Zheng et al., 467 2022). Table 4 and 5 reports the quantitative comparison results. Obviously, our proposed model can 468 help state-of-the-art VSG methods for performance improvement over all metrics on three datasets, 469 which demonstrates the superiority of our proposed model. It is mainly because our model can fully understand the query knowledge by the text augmentation process. 470

Efficiency comparison. We evaluate the efficiency of our proposed model, by fairly comparing its running time and model size in the inference phase with existing open-source methods for the VSG task on TACoS. As shown in Table 5, it can be observed that we achieve much faster processing speeds with relatively fewer learnable parameters.

475 476

477

441

442 443

444 445

446

447

448

4.2 Ablation Study and Analysis

478 Main ablation studies. To demonstrate the effectiveness of each component in our model, we 479 conduct ablation studies regarding the components (*i.e.*, Augmenting texts by Pre-trained LLMs, 480 Generating Positive and Negative texts, Significance Estimation for Sentence Component Integra-481 tion and Cross-modal Fusion) in Table 6. In particular, we remove each key individual module to 482 investigate its contribution. For convenience, we design four ablation models: 1) Ours(a). We re-483 move the "Augmenting texts by Pre-trained LLMs" module while keeping the other three modules. 2) Ours(b). We remove the "Generating Positive and Negative texts" module while keeping the 484 other three modules. 3) Ours(c). We remove the "Significance Estimation for Sentence Component 485 Integration" module while keeping the other three modules. Besides, we use our full model as the

487	Table 7: Ablatic	on study on different word types for the text-to	p-video task on DiDeMo (Anne Hendricks et al	1.,
488	2017) and VATE	X (Wang et al., 2019), where T-MASS (Wang	g et al., 2024) is the base model with CLIP-Vi	Г-
100	B/32 as backbon	e.		
403		DiDeMo	VATEX	

	Mathad	DiDeMo			VATEX						
	Methou	R@1↑	R@5↑	R@10↑	$MdR\downarrow$	MnR↓	R@1↑	R@5↑	R@10↑	$MdR\downarrow$	$MnR\downarrow$
	w/o Verb	46.0	74.1	82.7	2.0	14.3	61.2	93.2	94.0	2.0	2.7
	w/o Noun	43.1	71.8	82.3	2.0	15.1	61.3	91.0	95.6	2.0	3.3
	w/o Subject	48.6	77.1	84.4	2.0	13.4	65.2	92.7	95.9	1.0	3.0
-	Full model	52.7	79.3	88.6	1.0	10.4	64.9	93.7	98.2	1.0	1.4

baseline: Ours(full). As shown in Table 6, all four modules contribute a lot to the final performances on all three datasets, demonstrating their effectiveness under the VSG task.

496 Training process of different ablation 497 models. Following (Lin et al., 2020b), we 498 analyze the training process and retrieval 499 performance of different ablation models 500 in Figure 4. We can obtain the follow-501 ing representative observations: (i) Dur-502 ing training, Our(full) outperforms other 503 ablation models, which further demon-504 strates the effectiveness of each module. (ii) Our(full) converges faster than abla-505 tion models, showing that our full model 506 is more efficient. For instance, Our(full) 507 converges within 70 epochs, while Our(c) 508 converges after 80 epochs. Thus, our 509 full model can process these challenging 510 datasets more efficiently. 511

Analysis on different word types. Dur-512 ing generating negative texts, we change 513 a part of the sentence to help us under-514 stand the whole sentence. As shown in 515 Table 4.2, we analyze the effect of dif-516 ferent word types. Among three word 517 types (noun, verb, and subject), the noun 518 is the most significant. It is because the 519 noun can help our model localize the ob-520 ject in the given video for text-to-video 521 task. Besides, the noun can be used to 522 generate more semantic-rich attributes for 523 fine-grained language alignment. On the contrary, the subject brings minimal per-524 formance improvement. 525

526 Visualization Figure 5 depicts the grounding visualizations. Our model can significantly improve the state-of-the-art

(a) Visualization for the VSG task on Charades-STA.

Q1: *Kids are in a classroom finger painting.* **O2:** *Children are painting in room.*

Ground truth (GT	') 1 st video (🗸)	
T-MASS(Q1)	1 st video (🖌)	
+Ours(Q1)	1 st video (🗸)	
T-MASS(Q2)	2^{nd} video (×)	
+Ours(Q2)	1 st video (🗸)	

(b) Visualization for the VTR task on MSR-VTT.

Q1: How did the woman in yellow support the boy in blue at the start?

Ground truth (GT)	Hold baby up.							
MIST(Q1)	Hold baby up. (🗸)							
+Ours(Q1)	Hold baby up. (🗸)							
MIST(Q2)	Pull his finger. (×)							
+Ours(Q2)	Hold him up. (🗸)							
(c) Visualization for the VideoQA task on NexT-QA.								

Figure 5: Visualization results.

529 VLM methods for different tasks. This is because our model can fully understand the textual input 530 by attribute-based text reasoning.

531 532

533

486

495

5 CONCLUSION

In this paper, we rethink the LLM task from the user-friendly language input. We observe that many VLMs cannot fully understand the language texts. Given some texts with similar semantics and a video, these VLMs output various results. Thus, we design a plug-and-play framework to improve the generation ability of previous methods on various text templates. Extensive experiments on many challenging datasets show that our framework can serve as the plug-and-play module for state-of-the-art VLM works to improve their performance on various video-language tasks. In our future work, we will extend our model into more multi-modal tasks.

540 REFERENCES

547

562

588

589

590

Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell.
 Localizing moments in video with natural language. In *ICCV*, 2017.

- Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet:
 A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 961–970, 2015.
- Kilian Carolan, Laura Fennelly, and Alan F Smeaton. A review of multi-modal large language and vision models. *arXiv preprint arXiv:2404.01322*, 2024.
- Jiaming Chen, Weixin Luo, Wei Zhang, and Lin Ma. Explore inter-contrast between videos via
 composition for weakly supervised temporal sentence grounding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 267–275, 2022.
- Jingyuan Chen, Xinpeng Chen, Lin Ma, Zequn Jie, and Tat-Seng Chua. Temporally grounding natural sentence in video. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 162–171, 2018.
- Long Chen, Chujie Lu, Siliang Tang, Jun Xiao, Dong Zhang, Chilie Tan, and Xiaolin Li. Rethink ing the bottom-up framework for query-based video localization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.
- Zhenfang Chen, Lin Ma, Wenhan Luo, and Kwan-Yee Kenneth Wong. Weakly-supervised spatiotemporally grounding natural sentence in video. In *ACL*, 2019.
- Zhuoyun Du, Chen Qian, Wei Liu, Zihao Xie, Yifei Wang, Yufan Dang, Weize Chen, and Cheng
 Yang. Multi-agent software development through cross-team collaboration. *arXiv preprint arXiv:2406.08979*, 2024.
- Jinhao Duan, Hao Cheng, Shiqi Wang, Alex Zavalny, Chenan Wang, Renjing Xu, Bhavya Kailkhura, and Kaidi Xu. Shifting attention to relevance: Towards the predictive uncertainty quantification of free-form large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 5050–5063, 2024.
- Xuguang Duan, Wen-bing Huang, Chuang Gan, Jingdong Wang, Wenwu Zhu, and Junzhou Huang.
 Weakly supervised dense event captioning in videos. In *NeurIPS*, 2018.
- Lijie Fan, Dilip Krishnan, Phillip Isola, Dina Katabi, and Yonglong Tian. Improving clip training
 with language rewrites. *Advances in Neural Information Processing Systems*, 36, 2024.
- ⁵⁷⁵ Bo Fang, Chang Liu, Yu Zhou, Min Yang, Yuxin Song, Fu Li, Weiping Wang, Xiangyang Ji, Wanli Ouyang, et al. Uatvr: Uncertainty-adaptive text-video retrieval. In *ICCV*, 2023.
- Difei Gao, Luowei Zhou, Lei Ji, Linchao Zhu, Yi Yang, and Mike Zheng Shou. Mist: Multi-modal iterative spatial-temporal transformer for long-form video question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14773–14783, 2023.
- Jialin Gao, Xin Sun, Mengmeng Xu, Xi Zhou, and Bernard Ghanem. Relation-aware video reading
 comprehension for temporal language grounding. In *EMNLP*, pp. 3978–3988, 2021.
- Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. Tall: Temporal activity localization via language query. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pp. 5267–5275, 2017.
 - Junyu Gao and Changsheng Xu. Fast video moment retrieval. In ICCV, pp. 1523–1532, 2021.
 - Satya Krishna Gorti, Noël Vouitsis, Junwei Ma, Keyvan Golestan, Maksims Volkovs, Animesh Garg, and Guangwei Yu. X-pool: Cross-modal language-video attention for text-video retrieval. In *CVPR*, 2022.
- Zaber Ibn Abdul Hakim, Najibul Haque Sarker, Rahul Pratap Singh, Bishmoy Paul, Ali Dabouei,
 and Min Xu. Leveraging generative language models for weakly supervised sentence component
 analysis in video-language joint learning. *arXiv preprint arXiv:2312.06699*, 2023.

- Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet:
 A large-scale video benchmark for human activity understanding. In 2015 IEEE conference on computer vision and pattern recognition (CVPR), pp. 961–970. IEEE, 2015.
- Sarah Ibrahimi, Xiaohang Sun, Pichao Wang, Amanmeet Garg, Ashutosh Sanan, and Mohamed
 Omar. Audio-enhanced text-to-video retrieval using text-conditioned feature alignment. In *ICCV*, 2023.
- Peng Jin, Hao Li, Zesen Cheng, Kehan Li, Xiangyang Ji, Chang Liu, Li Yuan, and Jie Chen. Diffusionret: Generative text-video retrieval with diffusion model. In *ICCV*, 2023.
- Jaehun Jung, Peter West, Liwei Jiang, Faeze Brahman, Ximing Lu, Jillian Fisher, Taylor Sorensen, and Yejin Choi. Impossible distillation for paraphrasing and summarization: How to make high-quality lemonade out of small, low-quality model. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 4439–4454, 2024.
- Hyunjae Kim, Seunghyun Yoon, Trung Bui, Handong Zhao, Quan Tran, Franck Dernoncourt,
 and Jaewoo Kang. Fine-tuning clip text encoders with two-step paraphrasing. *arXiv preprint arXiv:2402.15120*, 2024.
- Hongxiang Li, Meng Cao, Xuxin Cheng, Yaowei Li, Zhihong Zhu, and Yuexian Zou. G2l: Semantically aligned and uniform video grounding via geodesic and game theory. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 12032–12042, 2023a.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023b.
- Pandeng Li, Chen-Wei Xie, Hongtao Xie, Liming Zhao, Lei Zhang, Yun Zheng, Deli Zhao, and Yongdong Zhang. Momentdiff: Generative video moment retrieval from random to real. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023c.
- 623 Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Annual Meeting of the* 624 Association for Computational Linguistics, 2004.
- ⁶²⁵
 ⁶²⁶
 ⁶²⁷
 ⁶²⁷
 ⁶¹⁷
 ⁶²⁸
 ⁶²⁷
 ⁶²⁹
 ⁶²⁷
 ⁶²⁹
 ⁶²⁷
 ⁶²⁷
 ⁶²⁸
 ⁶²⁷
 ⁶²⁹
 ⁶²⁷
 ⁶²⁹
 ⁶²⁷
 ⁶²⁹
 ⁶²⁷
 ⁶²⁹
 ⁶²⁹
 ⁶²⁹
 ⁶²⁹
 ⁶²⁹
 ⁶²¹
 ⁶²¹
 ⁶²²
 ⁶²⁵
 ⁶²⁵
 ⁶²⁶
 ⁶²⁶
 ⁶²⁷
 ⁶²⁷
 ⁶²⁸
 ⁶²⁷
 ⁶²⁹
 ⁶²⁹
 ⁶²⁹
 ⁶²⁹
 ⁶²⁹
 ⁶²⁹
 ⁶²⁰
 ⁶²¹
 ⁶²¹
 ⁶²²
 ⁶²²
 ⁶²⁵
 ⁶²⁵
 ⁶²⁶
 ⁶²⁶
 ⁶²⁷
 ⁶²⁷
 ⁶²⁸
 ⁶²⁸
 ⁶²⁹
 ⁶²⁹
- Zhijie Lin, Zhou Zhao, Zhu Zhang, Qi Wang, and Huasheng Liu. Weakly-supervised video moment retrieval via semantic completion network. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 11539–11546, 2020b.
- ⁶³¹ Daizong Liu, Xiaoye Qu, Jianfeng Dong, Pan Zhou, Yu Cheng, Wei Wei, Zichuan Xu, and Yulai
 ⁶³² Xie. Context-aware biaffine localizing network for temporal sentence grounding. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021.
- Meng Liu, Xiang Wang, Liqiang Nie, Xiangnan He, Baoquan Chen, and Tat-Seng Chua. Attentive
 moment retrieval in videos. In *Proceedings of the 41nd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pp. 15–24, 2018.
- Yuqi Liu, Pengfei Xiong, Luhui Xu, Shengming Cao, and Qin Jin. Ts2-net: Token shift and selection transformer for text-video retrieval. In *ECCV*, 2022.
- Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li. Clip4clip: An
 empirical study of clip for end to end video clip retrieval and captioning. *Neurocomputing*, 508: 293–304, 2022.
- Yueen Ma, Zixing Song, Yuzheng Zhuang, Jianye Hao, and Irwin King. A survey on visionlanguage-action models for embodied ai. *arXiv preprint arXiv:2405.14093*, 2024.
- 647 Niluthpol Chowdhury Mithun, Sujoy Paul, and Amit K. Roy-Chowdhury. Weakly supervised video moment retrieval from text queries. In *CVPR*, 2019.

648 Guoshun Nan, Rui Qiao, Yao Xiao, Jun Liu, Sicong Leng, Hao Zhang, and Wei Lu. Interven-649 tional video grounding with dual contrastive learning. In Proceedings of the IEEE Conference on 650 Computer Vision and Pattern Recognition (CVPR), 2021. 651 Zhaobo Qi, Yibo Yuan, Xiaowen Ruan, Shuhui Wang, Weigang Zhang, and Qingming Huang. Col-652 laborative debias strategy for temporal sentence grounding in video. IEEE Transactions on Cir-653 cuits and Systems for Video Technology, 2024. 654 655 Alec Radford. Improving language understanding by generative pre-training. 2018. 656 657 Michaela Regneri, Marcus Rohrbach, Dominikus Wetzel, Stefan Thater, Bernt Schiele, and Manfred Pinkal. Grounding action descriptions in videos. Transactions of the Association for Computa-658 tional Linguistics, 1:25-36, 2013. 659 660 Anna Rohrbach, Marcus Rohrbach, Niket Tandon, and Bernt Schiele. A dataset for movie descrip-661 tion. In CVPR, 2015. 662 663 Gunnar A Sigurdsson, Gül Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav Gupta. 664 Hollywood in homes: Crowdsourcing data collection for activity understanding. In ECCV, 2016. 665 Xinyu Tian, Shu Zou, Zhaoyuan Yang, and Jing Zhang. Argue: Attribute-guided prompt tuning for 666 vision-language models. In Proceedings of the IEEE/CVF Conference on Computer Vision and 667 Pattern Recognition, pp. 28578–28587, 2024. 668 669 Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. Journal of machine 670 learning research, 9(11), 2008. 671 Jiamian Wang, Guohao Sun, Pichao Wang, Dongfang Liu, Sohail Dianat, Majid Rabbani, Raghuveer 672 Rao, and Zhiqiang Tao. Text is mass: Modeling as stochastic embedding for text-video retrieval. 673 In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 674 16551-16560, 2024. 675 676 Jinpeng Wang, Yixiao Ge, Rui Yan, Yuying Ge, Kevin Qinghong Lin, Satoshi Tsutsui, Xudong Lin, 677 Guanyu Cai, Jianping Wu, Ying Shan, et al. All in one: Exploring unified video-language pre-678 training. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6598-6608, 2023. 679 680 Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. Vatex: A 681 large-scale, high-quality multilingual dataset for video-and-language research. In ICCV, 2019. 682 683 Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan 684 Xu, Yi Liu, Zun Wang, et al. Internvideo: General video foundation models via generative and 685 discriminative learning. arXiv preprint arXiv:2212.03191, 2022a. 686 Yuechen Wang, Jiajun Deng, Wengang Zhou, and Houqiang Li. Weakly supervised temporal adja-687 cent network for language grounding. TMM, 2021a. 688 689 Zheng Wang, Jingjing Chen, and Yu-Gang Jiang. Visual co-occurrence alignment learning for 690 weakly-supervised video moment retrieval. In Proceedings of the 29th ACM International Con-691 *ference on Multimedia*, pp. 1459–1468, 2021b. 692 Zhenzhi Wang, Limin Wang, Tao Wu, Tianhao Li, and Gangshan Wu. Negative sample matters: A 693 renaissance of metric learning for temporal grounding. In Proceedings of the AAAI Conference 694 on Artificial Intelligence, 2022b. 695 696 Bo Wu, Shoubin Yu, Zhenfang Chen, Joshua B Tenenbaum, and Chuang Gan. Star: A benchmark 697 for situated reasoning in real-world videos. In Thirty-fifth Conference on Neural Information 698 Processing Systems Datasets and Benchmarks Track (Round 2), 2021. 699 Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. Next-qa: Next phase of question-700 answering to explaining temporal actions. In Proceedings of the IEEE/CVF conference on com-701 puter vision and pattern recognition, pp. 9777–9786, 2021.

- 702 Huijuan Xu, Kun He, Bryan A Plummer, Leonid Sigal, Stan Sclaroff, and Kate Saenko. Multilevel 703 language and vision integration for text-to-clip retrieval. In Proceedings of the AAAI Conference 704 on Artificial Intelligence, volume 33, pp. 9062–9069, 2019. 705 Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging 706 video and language. In CVPR, 2016. 707 708 Hongwei Xue, Yuchong Sun, Bei Liu, Jianlong Fu, Ruihua Song, Houqiang Li, and Jiebo Luo. 709 Clip-vip: Adapting pre-trained image-text model to video-language representation alignment. In 710 ICLR, 2023. 711 Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. Just ask: Learning to 712 answer questions from millions of narrated videos. In Proceedings of the IEEE/CVF international 713 conference on computer vision, pp. 1686-1697, 2021a. 714 715 Wenfei Yang, Tianzhu Zhang, Yongdong Zhang, and Feng Wu. Local correspondence network for 716 weakly supervised temporal sentence grounding. *IEEE Transactions on Image Processing*, 30: 717 3252-3262, 2021b. 718 Qinghao Ye, Guohai Xu, Ming Yan, Haiyang Xu, Qi Qian, Ji Zhang, and Fei Huang. Hitea: Hierar-719 chical temporal-aware video-language pre-training. arXiv preprint arXiv:2212.14546, 2022. 720 721 Zhenfei Yin, Jiong Wang, Jianjian Cao, Zhelun Shi, Dingning Liu, Mukai Li, Xiaoshui Huang, 722 Zhiyong Wang, Lu Sheng, Lei Bai, et al. Lamm: Language-assisted multi-modal instruction-723 tuning dataset, framework, and benchmark. Advances in Neural Information Processing Systems, 724 36, 2024. 725 Shoubin Yu, Jaemin Cho, Prateek Yadav, and Mohit Bansal. Self-chained image-language model 726 for video localization and question answering. arXiv preprint arXiv:2305.06988, 2023. 727 728 Shoubin Yu, Jaemin Cho, Prateek Yadav, and Mohit Bansal. Self-chained image-language model for 729 video localization and question answering. Advances in Neural Information Processing Systems, 730 36, 2024. 731 732 Yitian Yuan, Lin Ma, Jingwen Wang, Wei Liu, and Wenwu Zhu. Semantic conditioned dynamic modulation for temporal sentence grounding in videos. In Advances in Neural Information Pro-733 cessing Systems (NIPS), pp. 534–544, 2019a. 734 735 Yitian Yuan, Tao Mei, and Wenwu Zhu. To find where you talk: Temporal sentence localization 736 in video with attention based location regression. In Proceedings of the AAAI Conference on 737 Artificial Intelligence, volume 33, pp. 9159–9166, 2019b. 738 739 Runhao Zeng, Haoming Xu, Wenbing Huang, Peihao Chen, Mingkui Tan, and Chuang Gan. Dense regression network for video grounding. In Proceedings of the IEEE Conference on Computer 740
- Da Zhang, Xiyang Dai, Xin Wang, Yuan-Fang Wang, and Larry S Davis. Man: Moment alignment network for natural language moment retrieval via iterative graph adjustment. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1247–1257, 2019a.

Vision and Pattern Recognition (CVPR), pp. 10287–10296, 2020.

- Gengyuan Zhang, Jisen Ren, Jindong Gu, and Volker Tresp. Multi-event video-text retrieval. In
 Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 22113–22123, 2023a.
- Hao Zhang, Aixin Sun, Wei Jing, and Joey Tianyi Zhou. Span-based localizing network for natural language video localization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 6543–6554, 2020a.
- Hao Zhang, Aixin Sun, Wei Jing, and Joey Tianyi Zhou. Temporal sentence grounding in videos: A
 survey and future directions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(8):10443–10465, 2023b.

- Songyang Zhang, Houwen Peng, Jianlong Fu, and Jiebo Luo. Learning 2d temporal adjacent networks for moment localization with natural language. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020b.
- Yuanhan Zhang, Kaiyang Zhou, and Ziwei Liu. What makes good examples for visual in-context learning? *Advances in Neural Information Processing Systems*, 36, 2024.
- Zhu Zhang, Zhijie Lin, Zhou Zhao, and Zhenxin Xiao. Cross-modal interaction networks for query based moment retrieval in videos. In *Proceedings of the 42nd International ACM SIGIR Confer- ence on Research and Development in Information Retrieval (SIGIR)*, pp. 655–664, 2019b.
- 765
 766
 767
 768
 768
 769
 769
 769
 760
 760
 761
 762
 763
 764
 765
 765
 765
 766
 766
 767
 768
 768
 768
 768
 769
 769
 769
 760
 760
 760
 760
 760
 760
 761
 761
 762
 763
 764
 765
 766
 766
 767
 768
 768
 768
 768
 769
 769
 769
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
 760
- Zongmeng Zhang, Xianjing Han, Xuemeng Song, Yan Yan, and Liqiang Nie. Multi-modal inter action graph convolutional network for temporal language localization in videos. *IEEE TIP*, 30:
 8265–8277, 2021.
- Shuai Zhao, Linchao Zhu, Xiaohan Wang, and Yi Yang. Centerclip: Token clustering for efficient text-video retrieval. In *ACM SIGIR*, 2022.
- Minghang Zheng, Yanjie Huang, Qingchao Chen, Yuxin Peng, and Yang Liu. Weakly supervised temporal sentence grounding with gaussian-based contrastive proposal learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15555–15564, 2022.
 - Cunjuan Zhu, Qi Jia, Wei Chen, Yanming Guo, and Yu Liu. Deep learning for video-text retrieval: a review. *International Journal of Multimedia Information Retrieval*, 12(1):3, 2023.
 - A APPEN
 - A APPENDIX
- 785 A.1 AUGMENTING TEXTS BY PRE-TRAINED LLMS

Considering the strong natural language processing of pre-trained LLMs, we will generate various text augmentations by pre-trained LLMs. Inspired by the effectiveness of LLMs, we rewrite the texts in VLM datasets to generate variation-origin pairs based on the following approaches.

790 791

779

780

781 782 783

784

786

A.1.1 AUGMENTATION BY HUMAN.

792 To make our model more user-friendly, we randomly invite ten persons from different countries to 793 rewrite some texts. For convenience, we randomly choose some video-text pairs from these VLM 794 datasets. For obtaining diverse text variations, we encourage rewriters to rewrite these texts based 795 on the target video segments. The human-based rewriting approach will enhance the creation of text 796 augmentation. Finally, we can obtain the variation-origin text pairs that include the original texts 797 and the corresponding human-rewritten version.

Augmentation by chat robots. Recently, LLMs-based chat robots (*e.g.*, ChatGPT or Bard) have achieved impressive performance in natural language processing. Thus, we target to rewrite the texts by chat robots. Firstly, we randomly choose some texts from the VSG datasets. Then, we utilize the web portals of chat robots to generate target texts by providing prompts. Some examples of the chat robot rewriting is shown in our supplementary material. With the powerful language processing ability of these chat robots, we can rewrite texts by utilizing different templates and vocabularies. The rewriting approach can preserve most textual semantics corresponding to the target segment.

Augmentation by open-source LLMs. Since it will lead to significant financial and time costs if
 we generate all the augmented texts by these closed-source chat robots (*e.g.*, ChatGPT and Bard),
 we utilize an open-source LLMs, LLaMA, to generate the positive texts. Due to the strong gener alization ability, LLaMA can be directly used to rewrite all the texts in the video-text datasets. For
 better generation ability, we leverage the LLaMA-7B model for text generation to guarantee that the
 generated texts are diverse and semantically relevant to the original texts.

Based on the above approaches, we can obtain three types of variation-origin text pairs: Bard-based,
ChatGPT-based, and LLaMA-based. Then, we treat them as inputs for the In-Context Learning
strategy (Zhang et al., 2024). In each approach, we randomly choose some texts from the video-text
datasets for target text generation, which will generate some variation-origin text pairs. These pairs
contain comprehensive and diverse training samples as the input of our framework.

816 A.2 MORE DETAILS FOR RELATED WORKS
817

818 Fully-supervised VSG. VSG is a new task introduced recently (Gao et al., 2017; Anne Hendricks 819 et al., 2017). Most previous algorithms (Anne Hendricks et al., 2017; Gao et al., 2017; Chen et al., 2018; Zhang et al., 2019b; Yuan et al., 2019a; Zhang et al., 2020b; Liu et al., 2021) have been 820 proposed within the propose-and-rank framework, which first generates moment candidates and 821 then utilizes multimodal matching to retrieve the most relevant candidate for a text. Some of them 822 (Anne Hendricks et al., 2017; Gao et al., 2017) take multiple sliding windows as candidates. To 823 improve the quality of the candidates, (Zhang et al., 2019b; Yuan et al., 2019a) pre-cut the video 824 on each frame by multiple pre-defined temporal scales, and directly integrate sentence information 825 with fine-grained video clips for scoring. For instance, Xu et al. (Xu et al., 2019) introduce a 826 multi-level model to integrate visual and textual features earlier and further re-generate texts as 827 an auxiliary task. Zhang et al. (Zhang et al., 2019a) model relations among candidate moments 828 produced from a convolutional neural network with the guidance of the text information. Although 829 these methods achieve great performance, they are severely limited by the heavy computation on 830 proposal matching/ranking, and sensitive to the quality of pre-defined proposals. Recently, many methods (Chen et al., 2020; Yuan et al., 2019b; Zeng et al., 2020; Zhang et al., 2020a; Nan et al., 831 2021) propose to utilize the boundary-regression framework. Specifically, they directly predict two 832 probabilities at each frame by leveraging cross-modal interactions between video and text, which 833 indicate whether this frame is a start/end frame of the ground truth video moment. 834

835 Weakly-supervised VSG. Despite the decent progress on the grounding performance, fully-836 supervised methods severely rely on the numerous annotations, which are significantly laborintensive and time-consuming to obtain. To alleviate this dense reliance to a certain extent, some 837 weakly-supervised VSG methods are proposed (Mithun et al., 2019; Chen et al., 2019; Lin et al., 838 2020a). For a weakly supervised VSG task, (Duan et al., 2018) decomposes it into two sub-tasks: 839 event captioning and text localization. (Duan et al., 2018) first assumes that each caption describes 840 only one temporal moment, and then designs a cycle network to train the model. As the pioneering 841 work for weakly-supervised VSG, (Mithun et al., 2019) learns a joint representation between the 842 video and the text by proposing Text-Guided-Attention network and utilizing an attention weight. 843 To improve the exploration and exploitation, (Lin et al., 2020a) chooses the top-K proposals and 844 measures the semantic similarity between the video and the text for localization. By proposing a 845 Semantic Completion Network, (Lin et al., 2020a) treats the masked text as input and predicts the 846 masked words from the video features.

847 848

815

B MORE EXPERIMENTS

849 850 851

B.1 MORE DATASET DETAILS

For convenience, we only utilize three datasets on the VSG task as example. We can utilize similar process for the other tasks.

1) ActivityNet Captions contains 19,209 videos from ActivityNet (Heilbron et al., 2015) with
71,953 textual descriptions. The videos are of diverse and open contents with an average length
of 2 minutes, and the annotated segments are significantly various in length, ranging from several
seconds to over 3 minutes. Following the public experimental setting (Zhang et al., 2019b), we use
37,417, 17,505, and 17,031 segment-text pairs for training, validation, and testing. For each dataset
split (training, validation, and testing), we generate 10,000 positive texts by each rewriting approach
(ChatGPT, Bard, and LLaMA) for each split.

2) Charades-STA (Gao et al., 2017) is an extension of the Charades dataset (Sigurdsson et al., 2016)
with temporal annotations. It contains 9,848 videos with an average length of 30 seconds and mainly focuses on daily indoor activities. There are 12,408 and 3,720 segment-text pairs in the training and

66	tranificst spins.						
67	·	Mathad	Tuna	R@1,	R@1,	R@5,	R@5,
68		Method	Type	IoU=0.3	IoU=0.5	IoU=0.3	IoU=0.5
60		CTRL	FS	10.25	8.92	27.85	16.96
70		ACRN	FS	11.70	10.38	28.04	18.30
70		CMIN	FS	12.87	9.96	27.99	17.43
/1		SCDM	FS	14.35	10.96	28.33	18.82
72		DRN	FS	16.39	11.77	30.85	20.99
73		2D-TAN	FS	16.98	13.79	33.02	21.13
74		MMN	FS	19.90	15.32	34.77	23.85
75		FVSG	FS	22.30	17.85	35.23	23.07
76		RaNet	FS	22.95	18.86	37.44	22.91
77		G2L	FS	23.48	19.87	37.52	24.88
78		MomentDiff	FS	26.76	18.73	38.49	25.84
79		G2L+Ours	-	37.92	27.48	46.73	35.95

Table 8: Performance comparison on TACoS dataset (with text augmentation) under official train/test splits.

Table 9: Effectiveness comparison on ActivityNet Captions dataset (without text augmentation) under official train/test splits.

1					
Mathad	Tuna	R@1,	R@1,	R@5,	R@5,
Method	Type	IoU=0.3	IoU=0.5	IoU=0.3	IoU=0.5
CTRL	FS	-	29.01	-	59.17
2D-TAN	FS	59.45	44.51	85.53	77.13
DRN	FS	-	45.45	-	77.97
RaNet	FS	-	45.59	-	75.93
MIGCN	FS	-	48.02	-	78.02
MMN	FS	65.05	48.59	87.25	79.50
G2L	FS	-	51.68	-	81.32
ICVC	WS	46.62	29.52	80.92	66.61
LCNet	WS	48.49	26.33	82.51	62.66
VCA	WS	50.45	31.00	71.79	53.83
WSTAN	WS	52.45	30.01	79.38	63.42
CNM	WS	55.68	33.33	-	-
G2L+Ours	-	66.34	54.26	91.77	84.29

testing sets, respectively. Similar to the ActivityNet Captions dataset, we generate 2,000 positive
 texts by each rewriting approach (ChatGPT, Bard, and LLaMA) for each dataset split.

3) TACoS (Regneri et al., 2013) consists of 127 videos with an average length of 7 minutes. These videos are selected from the MPII Cooking Composite Activities video corpus, which contains activities of cooking scenarios, thus lacking diversity. Following the standard split of (Gao et al., 2017), we use 10,146, 4,589, and 4,083 segment-text pairs for training, validation, and testing, respectively. Similar to the ActivityNet Captions dataset, we generate 4,000 positive texts by each rewriting approach (ChatGPT, Bard, and LLaMA) for each dataset split.

Similarly, we conduct the similar augmentation operation in other datasets.

911 Effect of negative samples. To evaluate the effect of the generated negative samples, we conduct
912 the ablation study for the VideoQA task on the NExT-QA dataset. Table 19 shows the results.
913 Obviously, we can utilize the generated negative samples for performance improvement in terms of
914 all the metrics.

Effect of attribute. Similarly, we analyze the significance of the attributes in our proposed framework by performing the ablation study in Table 20. In this table, with the attributes, our model
achieves the significant performance improvement. It is because the attributes can help our model
fully understand the language input by mining latent fine-grained language semantics.

9	1	8
9	1	9

933

Table 10: Text-to-video comparisons on DiDeMo (Anne Hendricks et al., 2017) and VATEX (Wang et al., 2019). Bold denotes the best performance. "–": result is unavailable.

Method		DiE	DeMo Retrie	eval			VA	TEX Retrie	eval
Wethod	R@1↑	R@5↑	R@10↑	$MdR\downarrow$	MnR↓	R@1↑	R@5↑	R@10↑	$MdR\downarrow$
CLIP-ViT-B/32									
X-Pool (Gorti et al., 2022)	44.6	73.2	82.0	2.0	15.4	60.0	90.0	95.0	1.0
+Ours	47.8	75.5	84.6	1.0	14.1	61.2	93.4	98.9	1.0
UATVR (Fang et al., 2023)	43.1	71.8	82.3	2.0	15.1	61.3	91.0	95.6	1.0
+Ours	45.5	72.9	84.8	1.0	13.4	65.2	92.7	98.9	1.0
T-MASS (Wang et al., 2024)	50.9	77.2	85.3	1.0	12.1	63.0	92.3	96.4	1.0
+Ours	52.7	79.3	88.6	1.0	10.4	64.9	93.7	98.2	1.0

Table 11: Performance comparison on TACoS dataset (without text augmentation) under official train/test splits.

1	Mathad	Tuna	R@1,	R@1,	R@5,	R@5,
	Method	Type	IoU=0.3	IoU=0.5	IoU=0.3	IoU=0.5
	CTRL	FS	18.32	13.30	36.69	25.42
	ACRN	FS	19.52	14.62	34.97	24.88
	CMIN	FS	24.64	18.05	38.46	27.02
	SCDM	FS	26.11	21.17	40.16	32.18
	DRN	FS	-	23.17	-	33.36
	2D-TAN	FS	37.29	25.32	57.81	45.04
	MMN	FS	39.24	26.17	62.03	47.39
	FVSG	FS	41.48	29.12	64.53	50.00
	G2L	FS	42.74	30.95	65.83	49.86
	RaNet	FS	43.34	33.54	67.33	55.09
	MomentDiff	FS	44.78	33.68	-	-
	G2L+Ours	-	53.56	40.31	71.62	62.06

946 947 948

949 Influence of corse-grained language alignment. To show the importance of our corse-grained
950 language alignment strategy, we conduct the ablation study for the VideoQA task on the NExT-QA
951 dataset in Table 21. Table 21 illustrates the effectiveness of the corse-grained language alignment
952 strategy.

Influence of fine-grained language alignment. To evaluate the effectiveness of the fine-grained language alignment module, we conduct the ablation study for the VideoQA task on the NExT-QA dataset in Table 22. Obviously, our fine-grained language alignment module can effectively improve the performance over all metrics.

957 958 959 960 Analysis on the hyper-parameters. Moreover, we investigate the robustness of the proposed model to different hyper-parameters in Figure 6. We find that we can obtain the best performance when $\gamma_1 = 0.4, \gamma_2 = 0.85, \gamma_3 = 0.25, \beta = 0.6, \tau = 0.3.$

Feature visualization. To investigate the feature distributions of the sentences during language alignment, we randomly choose some origin-variation pairs, and show the t-SNE Van der Maaten & Hinton (2008) visualizations of "before language alignment" and "after language alignment" in Figure 7. We can find that there is a large distribution gap between the origin and the variation of "before language alignment".

- 966
- 967
- 968
- 969
- 970
- 971

~	 _

Table 12: Text-to-video comparisons on MSRVTT (Xu et al., 2016) and LSMDC (Rohrbach et al., 2015). Bold denotes the best performance.

Mathad			MSRVTT					LSMDC		
Method	R@1↑	R@5↑	R@10↑	$MdR\downarrow$	MnR \downarrow	R@1↑	R@5↑	R@10↑	$MdR\downarrow$	MnR↓
CLIP-ViT-B/32										
X-Pool (Gorti et al., 2022)	46.9	72.8	82.2	2.0	14.3	25.2	43.7	53.5	8.0	53.2
+Ours	47.8	74.9	83.5	1.0	12.3	26.7	45.8	55.5	7.0	50.9
DiffusionRet (Jin et al., 2023)	49.0	75.2	82.7	2.0	12.1	24.4	43.1	54.3	8.0	40.7
+Ours	51.7	77.9	84.5	1.0	11.8	25.7	45.2	55.8	7.0	38.5
TEFAL (Ibrahimi et al., 2023)	49.4	75.9	83.9	2.0	12.0	26.8	46.1	56.5	7.0	44.4
+Ours	51.9	77.4	83.5	1.0	11.8	28.4	46.9	58.2	6.0	42.1
CLIP-ViP (Xue et al., 2023)	50.1	74.8	84.6	1.0	-	25.6	45.3	54.4	8.0	-
+Ours	51.7	75.3	85.9	1.0	11.8	26.8	47.6	58.5	6.0	42.3
T-MASS (Wang et al., 2024)	50.2	75.3	85.1	1.0	11.9	28.9	48.2	57.6	6.0	43.3
+Ours	52.3	77.9	87.6	1.0	10.9	30.8	50.4	59.1	5.0	40.8
CLIP-ViT-B/16										
X-Pool (Gorti et al., 2022)	48.2	73.7	82.6	2.0	12.7	26.1	46.8	56.7	7.0	47.3
+Ours	50.7	76.2	85.2	1.0	12.4	26.9	49.5	57.4	6.0	45.0
CLIP-ViP (Xue et al., 2023)	54.2	77.2	84.8	1.0	-	29.4	50.6	59.0	5.0	-
+Ours	56.8	79.4	85.9	82.8	1.0	32.6	51.8	60.9	3.0	38.7
T-MASS (Wang et al., 2024)	52.7	77.1	85.6	1.0	10.5	30.3	52.2	61.3	5.0	40.1
+Ours	54.9	82.6	86.8	1.0	10.2	35.1	55.7	64.8	3.0	38.9

Table 13: Example of LLM-based text generation, where "Origin" denotes the given original text and "Variation" denotes the generated texts. Although the only difference in the sentence structure between original text and negative text is "door" and "closet", they have different semantics. In contrast, positive text that are more distinct from the original text has similar semantics with original text.

Origin	Variation	Туре	Semantically similar?
Person opens	Person opens the closet.	Negative	×
	Door is opened by person.	Positive	~



Figure 6: Parameter analysis on the MSR-TT dataset for the text-to-video retrieval task with X-Pool as the base model.

1027
1028
1029Table 14: Performance comparison on Charades-STA dataset (without text augmentation) under
official train/test splits.1029R@1R@1R@5R@5

	Mathad	Type	R@1,	R@1,	R@5,	R@5,
	Method	туре	IoU=0.5	IoU=0.7	IoU=0.5	IoU=0.7
	VSA-RNN	FS	10.50	4.32	48.43	20.21
	VSA-STV	FS	16.91	5.81	53.89	23.58
	CTRL	FS	23.62	8.89	58.92	29.52
	2D-TAN	FS	39.81	23.25	79.33	52.15
	RaNet	FS	43.87	26.83	86.67	54.22
	DRN	FS	45.40	26.40	88.01	55.38
	MMN	FS	47.31	27.28	83.74	58.41
	G2L	FS	47.91	28.42	84.80	59.33
	MomentDiff	FS	53.79	30.18	-	-
	IVG-DCL	FS	50.24	32.88	-	-
· · · · · · · · · · · · · · · · · · ·	SCN	WS	23.58	9.97	71.80	38.87
	CTF	WS	27.30	12.90	-	-
	WSTAN	WS	29.35	12.28	76.13	41.53
	ICVC	WS	31.02	16.53	77.53	41.91
	MARN	WS	31.94	14.18	70.00	37.40
	CCL	WS	33.21	15.68	73.50	41.87
	CRM	WS	34.76	16.37	-	-
	CNM	WS	35.15	14.95	-	-
	VCA	WS	38.13	19.57	78.75	37.75
	LCNet	WS	39.19	18.17	80.56	45.24
	G2L+Ours	-	55.77	32.97	91.38	60.39

Table 15: Video-to-text comparisons on MSRVTT without text augmentation.

CLIP-ViT-B/32 CLIP4Clip (Luc et al. 2022)	42.7			•	
CLID4Clip (Luc at al. 2022)	427				
CLIF4CIIP (Luo et al., 2022)	/	70.9	80.6	2.0	11.6
+Ours	44.2	73.8	84.3	2.0	10.4
CenterCLIP (Zhao et al., 2022)	42.8	71.7	82.2	2.0	10.9
+Ours	44.5	73.0	84.1	1.0	9.7
X-Pool (Gorti et al., 2022)	44.4	73.3	84.0	2.0	9.0
+Ours	45.8	76.4	87.3	1.0	7.5
TS2-Net (Liu et al., 2022)	45.3	74.1	83.7	2.0	9.2
+Ours	48.6	77.5	85.7	2.0	7.9
DiffusionRet (Jin et al., 2023)	47.7	73.8	84.5	2.0	8.8
+Ours	51.0	75.9	87.4	2.0	6.9
UATVR (Fang et al., 2023)	46.9	73.8	83.8	2.0	8.6
+Ours	49.7	75.6	86.4	1.0	7.3
T-MASS (Wang et al., 2024)	47.7	78.0	86.3	2.0	8.0
+Ours	51.5	79.9	89.8	1.0	6.4
CLIP-ViT-B/16					
X-Pool (Gorti et al., 2022)	46.4	73.9	84.1	2.0	8.4
+Ours	50.2	77.4	86.3	2.0	6.1
TS2-Net (Liu et al., 2022)	46.6	75.9	84.9	2.0	8.9
+Ours	48.8	78.3	86.1	1.0	7.6
CenterCLIP (Zhao et al., 2022)	47.7	75.0	83.3	2.0	10.2
+Ours	49.8	78.0	86.4	2.0	6.5
UATVR (Fang et al., 2023)	48.1	76.3	85.4	2.0	8.0
+Ours	50.9	77.4	90.5	2.0	6.8
T-MASS (Wang et al., 2024)	50.9	80.2	88.0	1.0	7.4
+Ours	53.7	84.2	91.5	1.0	3.4

1089
1090
1091
1092
1093
1094
1095
1096

Table 16: Ablation study on different word types on the VSG task.ActivityNet Captions

			1	
Model	R@1	R@1	R@5	R@5
	IoU=0.3	IoU=0.5	IoU=0.3	IoU=0.5
w/o Noun	58.87	44.23	83.90	78.04
w/o Verb	59.15	44.92	83.71	78.95
w/o Subject	59.36	45.78	84.02	79.88
Full	60.90	46.86	84.39	80.62
	Cha	rades-STA		
Modal	R@1	R@1	R@5	R@5
Widdei	IoU=0.5	IoU=0.7	IoU=0.5	IoU=0.7
w/o Noun	31.14	25.88	74.29	45.28
w/o Verb	31.70	26.25	73.98	45.33
w/o Subject	33.64	27.61	74.13	46.45
Full	34.85	27.96	74.28	46.70

Table 17: Effectiveness comparison with text augmentation under official train/test splits, where "FS" denotes "fully-supervised" and "WS" means "weakly-supervised".

1107						j				
1108				ActivityNe	et Captions			Charad	es-STA	
1100	Method	Туре	R@1,	R@1,	R@5,	R@5,	R@1,	R@1,	R@5,	R@5,
1110			IoU=0.3	IoU=0.5	IoU=0.3	IoU=0.5	IoU=0.3	IoU=0.5	IoU=0.3	IoU=0.5
1110	CTRL	FS	20.27	19.40	47.82	40.78	10.32	3.54	36.98	13.40
1111	+Ours	FS	22.85	21.73	48.66	43.12	12.88	4.19	37.46	16.72
1112	2D-TAN	FS	48.32	29.38	71.36	62.30	20.18	11.35	47.05	33.82
1113	+Ours	FS	51.86	32.64	72.98	63.75	23.99	14.75	49.22	34.18
1114	DRN	FS	48.94	30.26	69.34	64.79	23.51	13.76	47.35	34.10
1115	+Ours	FS	50.75	32.92	73.86	67.48	26.44	15.83	49.07	36.52
1116	RaNet	FS	52.93	31.42	73.80	65.73	22.86	15.73	46.21	30.49
1117	+Ours	FS	53.88	33.74	75.96	68.31	23.11	16.97	48.05	33.21
1118	MMN	FS	55.30	31.76	74.88	71.52	25.33	18.80	45.97	35.08
1119	+Ours	FS	58.76	33.08	75.33	73.59	26.87	22.48	46.03	37.85
1120	G2L	FS	55.75	33.01	75.25	70.89	26.54	19.85	48.06	36.70
1121	+Ours	FS	60.90	46.86	84.39	80.62	34.85	27.96	74.28	46.70
1122	ICVC	WS	21.88	18.59	42.76	36.82	9.27	6.89	37.55	13.98
1123	+Ours	WS	22.36	20.17	44.19	38.77	8.64	6.89	39.57	16.63
1124	LCNet	WS	30.15	22.08	45.80	39.25	20.48	16.61	42.32	20.89
1125	+Ours	WS	33.94	21.73	46.35	40.57	23.72	17.59	43.80	23.77
1126	VCA	WS	31.74	25.37	46.98	42.76	17.87	12.39	45.70	22.13
1107	+Ours	WS	32.99	28.56	48.31	44.07	18.63	15.72	46.17	23.88
1127	WSTAN	WS	33.72	25.74	49.30	45.88	8.15	5.43	35.27	11.86
1128	+Ours	WS	35.20	27.99	51.84	48.69	10.77	6.92	37.40	13.88
1129	CNM	WS	35.72	28.95	50.06	48.72	14.34	9.65	43.88	18.79
1130	+Ours	WS	39.56	31.77	52.88	51.99	16.83	12.05	45.60	21.64
1131		1								

Table 18: Ablation study on different augmentation approaches.ActivityNet CaptionsModelR@1R@1R@5R@5ModelIoU=0.3IoU=0.5IoU=0.3IoU=0.5w/o ChatGPT58.7344.2082.1978.55w/o Bard59.6744.9583.0379.36w/o LLaMA59.8845.4083.2579.57Full60.9046.8684.3980.62Charades-STAModelR@1R@1R@5R@5W/o ChatGPT33.0426.9373.5045.27w/o Bard33.9227.4073.8345.62w/o Bard33.9227.4073.8345.62w/o LLaMA34.0827.4673.5545.39Full34.8527.9674.2846.70	- - - - -	R@5 IoU=0.5 78.55 79.36 79.57 80.62 R@5 IoU=0.7	mentation ap 18 R@5 IoU=0.3 82.19 83.03 83.25 84.39 R@5	different aug yNet Caption R@1 IoU=0.5 44.20 44.95 45.40 46.86	tion study on Activit R@1 IoU=0.3 58.73 59.67 59.88 60.90	Table 18: Ablat Model w/o ChatGPT w/o Bard w/o LLaMA
Table 18: Ablation study on different augmentation approaches.ActivityNet CaptionsModelR@1R@1R@5R@5ModelIoU=0.3IoU=0.5IoU=0.3IoU=0.5w/o ChatGPT58.7344.2082.1978.55w/o Bard59.6744.9583.0379.36w/o LaMA59.8845.4083.2579.57Full60.9046.8684.3980.62Charades-STAModelIoU=0.5IoU=0.7W/o ChatGPT33.0426.9373.5045.27w/o Bard33.9227.4073.8345.62w/o LLaMA34.0827.4673.5545.39Full34.8527.9674.2846.70	- - - - -	R@5 IoU=0.5 78.55 79.36 79.57 80.62 R@5 IoU=0.7	mentation ap 18 R@5 IoU=0.3 82.19 83.03 83.25 84.39 R@5	different aug yNet Caption R@1 IoU=0.5 44.20 44.95 45.40 46.86	tion study on Activit R@1 IoU=0.3 58.73 59.67 59.88 60.90	Table 18: Ablat Model w/o ChatGPT w/o Bard w/o LLaMA
ActivityNet CaptionsModelR@1R@1R@5R@5ioU=0.3ioU=0.5ioU=0.3ioU=0.5w/o ChatGPT58.7344.2082.1978.55w/o Bard59.6744.9583.0379.36w/o LLaMA59.8845.4083.2579.57Full60.9046.8684.3980.62Charades-STAModelR@1R@1R@5R@5ioU=0.5ioU=0.7ioU=0.5ioU=0.7w/o ChatGPT33.0426.9373.5045.27w/o Bard33.9227.4073.8345.62w/o LLaMA34.8527.9674.2846.70	-	R@5 IoU=0.5 78.55 79.36 79.57 80.62 R@5 IoU=0.7	R@5 IoU=0.3 82.19 83.03 83.25 84.39	yNet Caption R@1 IoU=0.5 44.20 44.95 45.40 46.86	Activit R@1 IoU=0.3 58.73 59.67 59.88 60.90	Model w/o ChatGPT w/o Bard w/o LLaMA
Model R@1 R@1 R@5 R@5 Wo ChatGPT 58.73 44.20 82.19 78.55 w/o Bard 59.67 44.95 83.03 79.36 w/o LLaMA 59.88 45.40 83.25 79.57 Full 60.90 46.86 84.39 80.62 Charades-STA Model R@1 R@1 R@5 R@5 Model IoU=0.5 IoU=0.7 IoU=0.5 IoU=0.7 w/o ChatGPT 33.04 26.93 73.50 45.27 w/o Bard 33.92 27.40 73.83 45.62 w/o LLaMA 34.08 27.46 73.55 45.39 Full 34.85 27.96 74.28 46.70	- - - -	R@5 IoU=0.5 78.55 79.36 79.57 80.62 R@5 IoU=0.7	R@5 IoU=0.3 82.19 83.03 83.25 84.39	R@1 IoU=0.5 44.20 44.95 45.40 46.86	R@1 IoU=0.3 58.73 59.67 59.88 60.90	Model w/o ChatGPT w/o Bard w/o LLaMA
Model IAUE 0.3 IAUE 0.5 IAUE 0.3 IAUE 0.5 w/o ChatGPT 58.73 44.20 82.19 78.55 w/o Bard 59.67 44.95 83.03 79.36 w/o LLaMA 59.88 45.40 83.25 79.57 Full 60.90 46.86 84.39 80.62 Charades-STA Model R@1 R@1 R@5 R@5 W/o ChatGPT 33.04 26.93 73.50 45.27 w/o Bard 33.92 27.40 73.83 45.62 W/o LLaMA 34.08 27.46 73.55 45.39 Full 34.85 27.96 74.28 46.70	-	IoU=0.5 78.55 79.36 79.57 80.62 R@5 IoU=0.7	IoU=0.3 82.19 83.03 83.25 84.39	IoU=0.5 44.20 44.95 45.40 46.86	IoU=0.3 58.73 59.67 59.88 60.90	Model w/o ChatGPT w/o Bard w/o LLaMA
Image: Note of the indext	-	78.55 79.36 79.57 80.62 R@5 IoU=0.7	82.19 83.03 83.25 84.39	44.20 44.95 45.40 46.86	58.73 59.67 59.88 60.90	w/o ChatGPT w/o Bard w/o LLaMA
Instruction Solution	-	79.36 79.57 80.62 R@5 IoU=0.7	83.03 83.25 84.39	44.95 45.40 46.86	59.67 59.88 60.90	w/o Bard w/o LLaMA
w/o LLaMA 59.88 45.40 83.25 79.57 Full 60.90 46.86 84.39 80.62 Charades-STA Model R@1 R@1 R@5 R@5 IoU=0.5 IoU=0.7 IoU=0.5 IoU=0.7 w/o ChatGPT 33.04 26.93 73.50 45.27 w/o Bard 33.92 27.40 73.83 45.62 w/o LLaMA 34.08 27.46 73.55 45.39 Full 34.85 27.96 74.28 46.70	-	79.57 80.62 R@5 IoU=0.7	83.25 84.39	45.40 46.86	59.88 60.90	w/o LLaMA
Full60.9046.8684.3980.62Charades-STAModelR@1R@1R@5R@5IoU=0.5IoU=0.7IoU=0.5IoU=0.7w/o ChatGPT33.0426.9373.5045.27w/o Bard33.9227.4073.8345.62w/o LLaMA34.0827.4673.5545.39Full34.8527.9674.2846.70	-	80.62 R@5 IoU=0.7	84.39	46.86	60.90	
TunOutputCharades-STAModelR@1R@1R@5R@5IoU=0.5IoU=0.7IoU=0.5IoU=0.7w/o ChatGPT33.0426.9373.5045.27w/o Bard33.9227.4073.8345.62w/o LLaMA34.0827.4673.5545.39Full34.8527.9674.2846.70	-	R@5 IoU=0.7	D @5	1.000		Full
Characes-STAModelR@1R@1R@5R@5IoU=0.5IoU=0.7IoU=0.5IoU=0.7w/o ChatGPT33.0426.9373.5045.27w/o Bard33.9227.4073.8345.62w/o LLaMA34.0827.4673.5545.39Full34.8527.9674.2846.70Table 19: Ablation study on for the negative samples on the NExT-QA VideoQA dataset with 2 as our base model.ModelTemporal↑Causal↑Descripting 1.5W/o negative samples61.363.471.5w/o negative samples69.270.178.4	-	R@5 IoU=0.7	DQ5	rodoc VIA	Che	
ModelK@1K@1K@3K@3IoU=0.5IoU=0.7IoU=0.5IoU=0.7w/o ChatGPT33.0426.9373.5045.27w/o Bard33.9227.4073.8345.62w/o LLaMA34.0827.4673.5545.39Full34.8527.9674.2846.70	-	IoU=0.7		D@1		
Table 10: -0.3 $100-0.7$ $100-0.7$ $100-0.7$ w/o ChatGPT 33.04 26.93 73.50 45.27 w/o Bard 33.92 27.40 73.83 45.62 w/o LLaMA 34.08 27.46 73.55 45.39 Full 34.85 27.96 74.28 46.70 Table 19: Ablation study on for the negative samples on the NExT-QA VideoQA dataset with $2 \frac{as our base model}{2 as our base model}$. $Model$ Temporal \uparrow Causal \uparrow Descripting the dataset with $2 \frac{as our base model}{2 \frac{Model}{2 \frac{100}{2 \frac$	-	100=0.7				Model
w/o ChatGPT 33.04 20.93 73.30 43.27 w/o Bard 33.92 27.40 73.83 45.62 w/o LLaMA 34.08 27.46 73.55 45.39 Full 34.85 27.96 74.28 46.70 Table 19: Ablation study on for the negative samples on the NExT-QA VideoQA dataset with 2 as our base model.ModelTemporal \uparrow Causal \uparrow Descripting the colspan="4">Descripting the colspan="4">Causal \uparrow W/o negative samples 61.3 63.4 71.5 w/o negative samples 69.2 70.1 78.4		45.27	73 50	26.03	33.04	w/o ChotGPT
w/o LlaMA 33.92 21.40 73.83 43.02 w/o LLaMA 34.08 27.46 73.55 45.39 Full 34.85 27.96 74.28 46.70 Table 19: Ablation study on for the negative samples on the NExT-QA VideoQA dataset with 2 as our base model.ModelTemporal \uparrow Causal \uparrow Descripting the samples of the sample of the samples of the sample of the samp		45.27	73.30	20.93	33.04	w/o Bard
W/0 LLaWA34.0821.4013.5343.53Full34.8527.9674.2846.70Table 19: Ablation study on for the negative samples on the NExT-QA VideoQA dataset with 2 as our base model.ModelTemporal↑Causal↑DescriptiModelTemporal↑Causal↑Descriptiw/o negative samples61.363.471.5w/ negative samples69.270.178.4		45.02	73.65	27.40	33.92	w/o LL aMA
Table 19: Ablation study on for the negative samples on the NExT-QA VideoQA dataset with 2 as our base model.ModelTemporal \uparrow Causal \uparrow Descripting 1.5w/o negative samples61.363.471.5w/ negative samples69.270.178.4	-	45.59	73.33	27.40	34.00	
Table 19: Ablation study on for the negative samples on the NExT-QA VideoQA dataset with 2 as our base model. Model Temporal↑ Causal↑ Descripting w/o negative samples 61.3 63.4 71.5 w/ negative samples 69.2 70.1 78.4	_	40.70	/4.20	27.90	34.05	<u> </u>
Table 19: Ablation study on for the negative samples on the NExT-QA VideoQA dataset with 2 as our base model. Model Temporal↑ Causal↑ Descripting w/o negative samples 61.3 63.4 71.5 w/ negative samples 69.2 70.1 78.4						
Table 19: Ablation study on for the negative samples on the NExT-QA VideoQA dataset with 2 as our base model.ModelTemporal↑Causal↑Descriptingw/o negative samples61.363.471.5w/ negative samples69.270.178.4						
Table 19: Ablation study on for the negative samples on the NEx I-QA VideoQA dataset with 2 as our base model.2 as our base model					·	10. Al-1-4
ModelTemporal↑Causal↑Descriptiw/o negative samples61.363.471.5w/ negative samples69.270.178.4	aset with BLIP-	ideoQA dataset wi	NEXT-QA V	mples on the	ne negative sa	19: Ablation study on for th
w/o negative samples 61.3 63.4 71.5 w/ negative samples 69.2 70.1 78.4	escription	Descripti	Causal↑	oral↑	Tempo	Model
w/ negative samples 69.2 70.1 78.4	71.5	71.5	63.4	3	61.	w/o negative samples
	78.4	78.4	70.1	2	69.	w/ negative samples
Table 20: Ablation study on for the negative samples on the NExT-QA VideoQA dataset with 2 as our base model.	aset with BLIP-	ideoQA dataset wi	NExT-QA V	mples on the	he negative sa	20: Ablation study on for th ur base model.
Model Temporal↑ Causal↑ Description↑	otion↑	Description [↑]	ausal↑	Ca	Temporal↑	Model
w/o attribute 60.9 62.5 67.0	0	67.0	62.5	(60.9	w/o attribute
w/ attribute 69.2 70.1 78.4	4	78.4	70.1		69.2	w/ attribute
	aset with BLIP-	ideoQA dataset wi	NExT-QA V	mples on the	he negative sa	21: Ablation study on for th
Table 21: Ablation study on for the negative samples on the NExT-QA VideoQA dataset with 2 as our base model.			emporal↑	Te		Model
Table 21: Ablation study on for the negative samples on the NExT-QA VideoQA dataset with 2 as our base model. Model Temporal↑	sal↑ Descriptio	Causal↑	· ·		ge alignment	w/o corse-grained language
Table 21: Ablation study on for the negative samples on the NExT-QA VideoQA dataset with 2 as our base model. 2 as our base model.	sal↑ Descriptio .5 72.3	Causal↑ 67.5	64.3		6 6	
Table 21: Ablation study on for the negative samples on the NExT-QA VideoQA dataset with 2 as our base model. Model Temporal↑	sal↑ Descriptio	Causal↑	÷ · ·		ge alignment	w/o corse-grained language



Figure 7: The t-SNE visualizations of "before language alignment" and "after language alignment". Green circles denote the original sentences, while purple triangles denote denote the augmented sentences.