

EASIER: Relevance-Boosted Captioning and Structural Information Extraction for Zero-Shot Video-Text Retrieval

Anonymous ACL submission

Abstract

While recent progress in video-text retrieval (VTR) has been advanced by the exploration of supervised representation learning, in this paper, we present a novel zero-shot VTR framework, EASIER, to retrieve video/text with off-the-shelf captioning methods, large language models (LLMs), and text retrieval methods. Specifically, we first map videos into captions and then retrieve video captions and text using text retrieval methods, without any model training or fine-tuning. However, due to the limited power of captioning methods, the captions often miss important content in the video, resulting in unsatisfactory retrieval performance. To translate more information into video captions, we designed a novel **relevance-boosted caption generation** method, bringing extra relevant details into video captions by LLMs. Moreover, to emphasize key information and reduce the noise brought by imagination, we extract key visual tokens from captions and design different templates for structuring these tokens with the proposed **structural information extraction**, further boosting the retrieval performance. Benefiting from the enriched captions and structuralized information, extensive experiments on several video-text retrieval benchmarks demonstrate the superiority of EASIER over existing fine-tuned and pretraining methods without any data. A comprehensive study with both human and automatic evaluations shows that the enriched captions capture the key details and barely bring noise to the captions. Codes and data will be released.

1 Introduction

Video-text retrieval (VTR) (Luo et al., 2022; Gao et al., 2021; Ma et al., 2022; Liu et al., 2022a; Zhao et al., 2022; Gorti et al., 2022; Fang et al., 2022) aims to retrieve the corresponding video or text given the query in another modality. Recent years have witnessed the rapid development of VTR with the support from powerful pretraining models (Luo

et al., 2022; Gao et al., 2021; Ma et al., 2022; Liu et al., 2022a), improved retrieval methods (Bertiusius et al., 2021; Dong et al., 2019; Jin et al., 2021), and video-language datasets construction (Xu et al., 2016). However, it remains challenging to precisely match video and language due to the raw data being in heterogeneous spaces and the use of modality-specific encoders.

The most popular paradigm in VTR (Luo et al., 2022; Ma et al., 2022; Liu et al., 2022b) firstly learns a joint feature space across modalities and then compares representations in this space. However, with the discrepancy between different modalities and the design of modality-independent encoders, it is challenging to directly match representations of different modalities generated from different encoders (Liang et al., 2022). On the other side, pioneering works (Wang et al., 2021, 2022e) convert images into captions for better presentation learning on image-language tasks, demonstrating that captioners can mitigate modality discrepancy.

In this work, to take one step forward, we present a zero-shot video-text retrieval framework with our proposed **rElevance-boosted captioning And Structural Information ExtRaction**, EASIER. EASIER first captions videos into video captions. However, we notice that the captions always miss important information in the video, thus leading to bad retrieval performance (see Table 1 without paraphrase and visual tokens). To this end, we propose **relevance-boosted captioning**, which augments video captions by encouraging large language models (LLMs) to add relevant details to captions. Later, to emphasize the key information in the captions, *e.g.*, objects, events, and attributes, we design a **structural information extraction** procedure for extracting and formatting “visual tokens”. Finally, EASIER utilizes off-the-shelf text retrieval methods for zero-shot text retrieval matching video captions with structured visual tokens and text.

084 Finally, to evaluate the effectiveness of our pro- 133
085 posed zero-shot EASIER, we conducted experi- 134
086 ments on three representative video-text bench- 135
087 marks (Chen and Dolan, 2011; Fabian Caba Heil- 136
088 bron and Niebles, 2015; Xu et al., 2016). Results 137
089 show that EASIER outperforms previous methods, 138
090 including fine-tuning methods and few-shot meth- 139
091 ods benefiting from relevance-boosted captioning 140
092 and structural information extraction. 141

093 In summary, our contributions are as follows:

- 094 • We propose a **real zero-shot** video-text re- 142
095 trieval method without requiring any training 143
096 procedure or human-annotated data, only us- 144
097 ing the off-the-shelf captioning method, large 145
098 language models, and text retrieval methods. 146
- 099 • Our proposed EASIER achieves SOTA perfor- 147
100 mance on several metrics across three VTR 148
101 benchmarks. 149
- 102 • Detailed analysis reveals the importance of 150
103 relevance-boosted captioning and structural 151
104 information extraction. We will open-source 152
105 the code and data to facilitate future research. 153

106 2 Related Work

107 **Video-text retrieval**, which involves cross-modal 154
108 alignment and abstract understanding of temporal 155
109 images (videos), has been a popular and fundamen- 156
110 tal task of language-grounding problems (Wang 157
111 et al., 2020a,b, 2021; Yu et al., 2023). Most of 158
112 the existing video-text retrieval frameworks (Yu 159
113 et al., 2017; Dong et al., 2019; Zhu and Yang, 160
114 2020; Miech et al., 2020; Gabeur et al., 2020; Dz- 161
115 abraev et al., 2021; Croitoru et al., 2021) focus 162
116 on learning powerful representations for video and 163
117 text and extracting separated representations. For 164
118 example, in Dong et al. (2019), videos and texts 165
119 are encoded using convolutional neural networks 166
120 and a bi-GRU (Schuster and Paliwal, 1997) while 167
121 mean pooling is employed to obtain multi-level 168
122 representations. MMT (Gabeur et al., 2020) uses 169
123 a cross-modal encoder to aggregate features ex- 170
124 tracted by temporal images, audio, and speech for 171
125 encoding videos. Following that, MDMMT (Dz- 172
126 abraev et al., 2021) further utilizes knowledge 173
127 learned from multi-domain datasets to improve per- 174
128 formance empirically. Further, MIL-NCE (Miech 175
129 et al., 2020) adopts Multiple Instance Learning 176
130 and Noise Contrastive Estimation, addressing the 177
131 problem of visually misaligned narrations from un- 178
132 curated videos. 179

133 Recently, with the success of self-supervised 134
135 pretraining methods (Devlin et al., 2019; Radford 136
137 et al., 2019; Brown et al., 2020), vision-language 137
138 pretraining (Li et al., 2020b; Gan et al., 2020; 138
139 Singh et al., 2022) on large-scale unlabeled cross- 139
140 modal data has shown promising performance in 140
141 various tasks, *e.g.*, image retrieval (Radford et al., 141
142 2021), image captioning (Chan et al., 2023), and 142
143 video retrieval (Luo et al., 2022; Wang and Shi, 143
144 2023a). Recent works (Lei et al., 2021; Cheng 144
145 et al., 2021; Gao et al., 2021; Ma et al., 2022; 145
146 Park et al., 2022a; Wang et al., 2022b,d; Zhao 146
147 et al., 2022; Gorti et al., 2022) have attempted to 147
148 pretrain or fine-tune video-text retrieval models 148
149 in an end-to-end manner. CLIPBERT (Lei et al., 149
150 2021; Bain et al., 2021), as a pioneer, proposes to 150
151 sparsely sample video clips for end-to-end train- 151
152 ing to obtain clip-level predictions and then sum- 152
153 marize them. Frozen in time (Bain et al., 2021) 153
154 uses end-to-end training on both image-text and 154
155 video-text pairs data by uniformly sampling video 155
156 frames. CLIP4Clip (Luo et al., 2022) finetunes 156
157 models and investigates three similarity calculation 157
158 approaches for video-sentence contrastive learn- 158
159 ing on CLIP (Radford et al., 2021). Further, TS2- 159
160 Net (Liu et al., 2022b) proposes a novel token shift 160
161 and selection transformer architecture that adjusts 161
162 the token sequence and selects informative tokens 162
163 in both temporal and spatial dimensions from input 163
164 video samples. While the mainstream of VTR mod- 164
165 els (Xue et al., 2023; Wu et al., 2023) focuses on 165
166 fine-tuning powerful image-text pre-trained mod- 166
167 els, on the other side, as a pioneer, (Tiong et al., 167
168 2022; Wang et al., 2022e) propose to use large lan- 168
169 guage models (LLMs) for zero-shot video question 169
170 answering. 170

171 **Zero-shot cross-modal retrieval.** With the huge 171
172 success of pretrained visual-language model (Rad- 172
173 ford et al., 2021; Luo et al., 2022), zero-shot cross- 173
174 modal retrieval has attracted more and more re- 174
175 search interest recently. Due to the powerful rep- 175
176 resentation learning ability in image and text do- 176
177 mains, CLIP (Radford et al., 2021) achieves sat- 177
178 isfying zero-shot retrieval performance on sev- 178
179 eral representative image-text retrieval bench- 179
180 marks (Huiskes and Lew, 2008; Lin et al., 2014). 180
181 Inspired by this achievement, Liu et al. (2023a,b); 181
182 Chen et al. (2023c); Liu et al. (2024); Guo et al. 182
183 (2024) boost the performance of zero-shot image- 183
184 text retrieval by better representation learning meth- 184

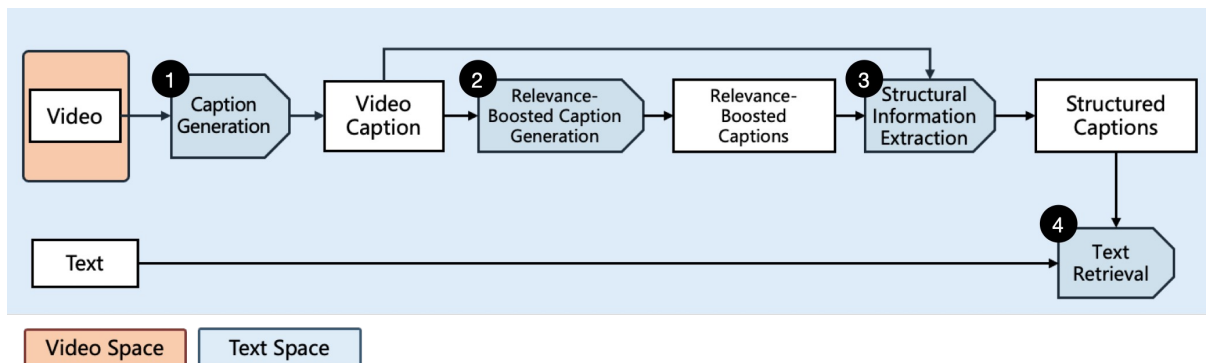


Figure 1: The illustration of our proposed EASIER. EASIER includes four steps. First, we generate video captions for video using off-the-shelf video captioning methods. Second, to enrich the captions, we propose the relevance-boosted caption-generation method using LLMs. Third, to emphasize the important information in the captions, we propose a novel structural information extraction. Finally, after obtaining structured video captions, we employ off-the-shelf text retrieval methods to perform zero-shot video-text retrieval.

Dolan, 2011; Fabian Caba Heilbron and Niebles, 2015), video-language pre-trained models (Wang et al., 2022c; Chen et al., 2023a; Xu et al., 2023; Chen et al., 2023c; Li et al., 2023b; Liu et al., 2023c; Zhu et al., 2024) also achieve satisfying zero-shot video-text retrieval results.

In this paper, inspired by these pioneering works, to explore zero-shot video-text retrieval, we step forward and propose a simple but effective zero-shot video-text retrieval method, EASIER, by utilizing off-the-shelf captioning, large language models, and text retrieval methods.

3 EASIER - Zero-Shot Video Text Retrieval

In this section, we present the details of our proposed method, EASIER. Specifically, we first generate captions for videos using video caption generation methods. Then, to cover most of the details in videos, with our proposed **relevance-boosted caption generation**, we obtain a detailed caption containing almost all the details. Finally, we propose the **structural information extraction** to emphasize important information in the captions for better video-text retrieval performance. The whole procedure is summarized in Figure 1.

3.1 Step 1 - Video Caption Generation

Video captioning with off-the-shelf captioners. Specifically, we employ [Tewel et al. \(2021, 2022\)](#) to generate video captions and then use GPT-2 ([Radford et al., 2019](#)) to enrich sentences using the prompts, *i.e.*, “Video presents”.

3.2 Step 2 - Relevance-Boosted Caption Generation

As shown in Figure 3, we notice that the generated captions always miss some important information, leading to unsatisfying retrieval performance. A simple solution to this problem is to fine-tune the captioning models, which will improve their caption-generation abilities. However, this approach needs a huge amount of annotated video-caption data and expensive computation resources, and the fine-tuned models are always not able to be transferred to other benchmarks. To this end, we propose the **relevance-boosted caption generation**, which is training-free and generates detailed captions that contain almost every detail of the video.

Specifically, we use large language models (LLMs) ([Brown et al., 2020](#); [Touvron et al., 2023](#)) to conduct the hallucination-based generation using the following prompt template.

The following is a caption from a video: [" + <Video Caption> + "]. Based on this caption, generate two paraphrased captions capturing the key information and main themes, each of which should be in one sentence with up to twenty words. Meanwhile, please be creative, you can have some imagination and add the necessary details. Generated sentences should be in the number list. Also please generate text without any comment.

Our proposed method generates multiple captions (*e.g.*, 1, 2, and 3). However, some of these

captions might introduce noise or lack strong relevance to the video’s content. To mitigate potential negative impacts, we apply a filtering method to assess the semantic similarity between relevance-boosted captions and the original video caption by leveraging a pre-trained text encoder (Reimers and Gurevych, 2019). Then we concatenate the filtered captions along with the original video caption to obtain the final captions.

3.3 Step 3 - Structural Information Extraction

To understand which kind of information is essential to VTR, we analyze the contextual text of video captions by breaking down the video captions into four different visual tokens using NLTK (Bird et al., 2009), *i.e.*, phrase, object, event, and attribute. Then we structure the information into the following structure,

<p><Caption> <Phrases> <Attributes> <Events> <Objects></p>
--

3.4 Step 4 - Video (Video Caption)-Text Retrieval

Finally, after obtaining structured video caption data, we are ready to perform the retrieval step. Specifically, we compute the similarity score at the video level between text and video caption using off-the-shelf retrieval methods, *i.e.*, BM25 (Robertson and Walker, 1994) and Sentence transformers (Reimers and Gurevych, 2019).

4 Experiments

4.1 Benchmarks, Baselines, and Evaluation Metrics

Benchmarks. Following previous work (Luo et al., 2022; Ma et al., 2022), we use three representative benchmarks for evaluating EASIER, *i.e.*, MSR-VTT (Xu et al., 2016), MSVD (Chen and Dolan, 2011), and ActivityNet (Fabian Caba Heilbron and Nieves, 2015). Details of the dataset split are presented in appendix A.1.

Baselines. To show the empirical efficiency of our EASIER, we compare it with fine-tuned models, pre-trained methods, and few-shot methods. Details are presented in Appendix A.2.

Evaluation metric. To evaluate the retrieval performance of our proposed model, we use recall at Rank K (R@K, higher is better), median rank (M_dR, lower is better), and mean rank (M_nR, lower is better) as retrieval metrics, which are widely used

in previous retrieval works (Radford et al., 2021; Luo et al., 2022; Ma et al., 2022).

Implementation details and related model details are defferd to Appendix A.3.

4.2 Quantitative Results

In this part, we present the qualitative results of EASIER on three VTR benchmarks.

MSR-VTT. We found that the contextual video text obtained directly through video captioning methods generally have mediocre performance (R@1: 20.3) compared to other baseline Text-Video Retrieval method. However, after using LLM to do relevance boosting from the video caption, the R@1 of our method nearly doubled (R@1 = 40.9) shown in Table 5 . Therefore, we further boosted each sentence and expanded it into two sentences. From the results presented in Table 1, it can be seen that this approach outperforms the second-best method by 9.9. This indicates the significant impact of relevance boosting and expanding captions on enhancing the performance of Text-Video Retrieval systems. Compared to DiscreteCodebook (Liu et al., 2022a), which aligns modalities in an unsupervised manner, EASIER outperforms DiscreteCodebook on every metric. Meanwhile, EASIER also outperforms VidIL (Wang et al., 2022e), which uses few-shot prompting, demonstrating the usability of integrating zero-shot LLM on text-to-video retrieval. This suggests that leveraging zero-shot on LLMs is a promising approach to enhance text-to-video retrieval performance. Also, we notice that EASIER has bad results on mean rank. To understand why this happens, we visualize the distribution of rank in Figure 2. It is obvious that though most of the videos have very good rank, *e.g.*, lower than 10, there are still some captions ranked in the last. This might be due to the failure of caption generation for some videos, where the generated captions do not contain any information from the video.

MSVD and ActivityNet. The results on MSVD and ActicityNet are shown in Table 3. EASIER achieves the best R@1 on text-to-video retrieval on two datasets compared to the previous methods.

4.3 Ablation Studies

In this part, we present a series of ablation experiments on MSR-VTT to better understand the effectiveness of different components of EASIER, using LLaMA2-7b-chat-hf and BM25. Due to space limitations, we present the ablation study on retrieval

Methods	Venue	Text-to-Video Retrieval					Video-to-Text Retrieval				
		R@1 \uparrow	R@5 \uparrow	R@10 \uparrow	MdR \downarrow	MnR \downarrow	R@1 \uparrow	R@5 \uparrow	R@10 \uparrow	MdR \downarrow	MnR \downarrow
<i>Training-based</i>											
LiteVL-S	EMNLP'22	46.7	71.8	81.7	2.0	-	-	-	-	-	-
X-Pool	CVPR'22	46.9	72.8	82.2	2.0	14.3	-	-	-	-	-
CenterCLIP	SIGIR'22	44.2	71.6	82.1	2.0	15.1	42.8	71.7	82.2	2.0	10.9
TS2-Net	ECCV'22	47.0	74.5	83.8	2.0	13.0	45.3	74.1	83.7	2.0	9.2
X-CLIP	ACM MM'22	46.1	74.3	83.1	2.0	13.2	46.8	73.3	84.0	2.0	9.1
NCL	EMNLP'22	43.9	71.2	81.5	2.0	15.5	44.9	71.8	80.7	2.0	12.8
TABLE	AAAI'23	47.1	74.3	82.9	2.0	13.4	47.2	74.2	84.2	2.0	11.0
VOP	CVPR'23	44.6	69.9	80.3	2.0	16.3	44.5	70.7	80.6	2.0	11.5
DiscreteCodebook	ACL'22	43.4	72.3	81.2	-	14.8	42.5	71.2	81.1	-	12.0
VCM	AAAI'22	43.8	71.0	-	2.0	14.3	45.1	72.3	82.3	2.0	10.7
CenterCLIP	SIGIR'22	48.4	73.8	82.0	2.0	13.8	47.7	75.0	83.3	2.0	10.2
HiSE	ACM MM'22	45.0	72.7	81.3	2.0	-	46.6	73.3	82.3	2.0	-
TS2-Net	ECCV'22	49.4	75.6	85.3	2.0	13.5	46.6	75.9	84.9	2.0	8.9
S3MA	EMNLP'23	53.1	78.2	86.2	1.0	10.5	52.7	79.2	86.3	1.0	8.2
UCOFIA	ICCV'23	49.4	72.1	-	-	12.9	47.1	74.3	-	-	-
ProST	ICCV'23	49.5	75.0	84.0	2.0	11.7	48.0	75.9	85.2	2.0	8.3
UATVR	ICCV'23	49.8	76.1	85.5	2.0	12.9	51.1	74.8	85.1	1.0	8.3
MV-Adapter	CVPR'24	46.2	73.2	82.7	-	-	47.2	74.8	83.9	-	-
<i>Zero-Shot (Pretrained Models)</i>											
VLM	ACL'21	28.1	55.5	67.4	4.0	-	-	-	-	-	-
HERO	EMNLP'21	16.8	43.3	57.7	-	-	-	-	-	-	-
VideoCLIP	EMNLP'21	30.9	55.4	66.8	-	-	-	-	-	-	-
EvO	CVPR'22	23.7	52.1	63.7	4.0	-	-	-	-	-	-
OA-Trans	CVPR'22	35.8	63.4	76.5	3.0	-	-	-	-	-	-
RaP	EMNLP'22	40.9	67.2	76.9	2.0	-	-	-	-	-	-
OmniVL	NeurIPS'22	34.6	58.4	66.6	-	-	-	-	-	-	-
mPLUG-2	ICML'23	48.3	<u>75.0</u>	<u>83.2</u>	-	-	-	-	-	-	-
InternVL	Arxiv'23	42.4	65.9	75.4	-	-	46.3	70.5	79.6	-	-
LanguageBind	ICLR'24	42.6	65.4	75.5	-	-	-	-	-	-	-
<i>Few-Shot</i>											
VidIL	NeurIPS'22	40.8	65.2	-	-	-	<u>39.6</u>	<u>64.5</u>	-	-	-
<i>Zero-Shot</i>											
EASIER w/o paraphrase and visual tokens		20.3	40.9	51.7	9.0	60.3	17.5	36.7	47.3	12.0	<u>82.3</u>
EASIER w/o visual tokens		<u>54.0</u>	73.9	80.2	1.0	<u>24.5</u>	27.9	41.3	47.3	15.0	136.2
EASIER		58.2	75.8	83.5	1.0	18.9	36.4	56.5	<u>63.8</u>	3.0	75.7

Table 1: Video-Text retrieval results on MSR-VTT. The best results are marked in **bold**. The second best results are underlined. “NC” refers to Neurocomputing.

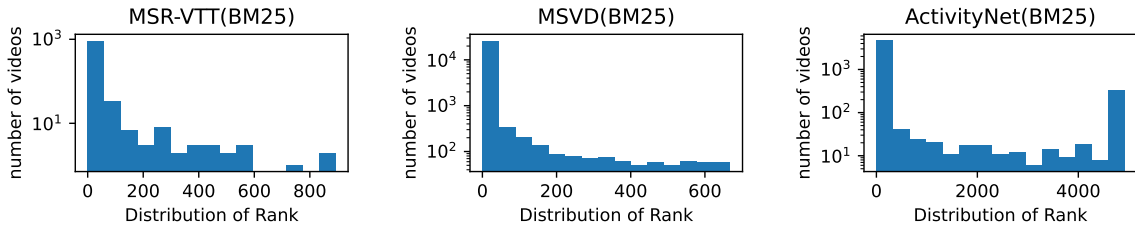


Figure 2: These figures illustrate the distribution of the rank of each (test) gallery video (captions) retrieved by (test) text queries.

352 methods and the investigation on the need for struc- 363
353 tural information extraction and relevance-boosted 364
354 captions in appendix A.4 and appendix A.6. 365

355 **Impact of combination of structural information** 366
356 **(visual tokens).** To choose the best combination 367
357 method for the extracted visual tokens (phrases, 368
358 attributes, objects, and events), we conduct experi- 369
359 ments using different arrangements of these visual 370
360 tokens, as shown in Table 2. By reducing the inclu- 371
361 sion of visual tokens, the retrieval performance of 372
362 EASIER decreases, thereby proving the superiority 373

of integrating these four visual tokens together.

The order of different structural information. 364
365 Another important factor to consider is the order 366
367 of these visual tokens. To this end, we systemat- 368
369 ically evaluate which specific order of <phrase>, 370
371 <object>, <attribute>, and <event> maximizes the 372
373 efficiency and accuracy of the retrieval process. The results are shown in Table 4. We discover that among various arrangements, the model performs best when either phrases or objects are placed at the end of the sequence. This superior performance

Caption	Phrase	Object	Event	Attribute	Text-to-Video Retrieval					Video-to-Text Retrieval				
					R@1↑	R@5↑	R@10↑	MdR↓	MnR↓	R@1↑	R@5↑	R@10↑	MdR↓	MnR↓
✓					54.0	73.9	80.2	1.0	24.5	27.9	41.3	47.3	15.0	136.2
✓	✓				57.4	76.2	83.0	1.0	19.3	29.9	45.6	52.3	8.0	115.4
✓		✓			56.9	77.5	83.8	1.0	18.6	35.8	56.9	64.8	3.0	73.9
✓			✓		54.2	73.2	79.6	1.0	24.9	28.4	42.7	49.1	12.0	130.3
✓				✓	55.0	74.2	80.2	1.0	24.1	28.6	43.3	48.9	11.0	132.2
✓	✓	✓			57.4	76.2	83.5	1.0	18.7	34.5	54.0	62.5	4.0	79.9
✓	✓		✓		57.3	76.3	82.6	1.0	19.8	31.5	47.3	54.2	7.0	109.0
✓	✓			✓	57.6	76.3	83.5	1.0	19.1	31.0	47.4	54.7	7.0	110.5
✓		✓	✓		56.9	76.6	83.2	1.0	19.3	35.9	57.9	65.6	3.0	71.2
✓		✓		✓	57.6	77.4	83.8	1.0	18.2	37.4	58.5	66.3	3.0	71.8
✓			✓	✓	54.0	73.3	79.6	1.0	24.9	30.0	44.3	51.1	10.0	126.1
✓	✓	✓	✓		58.0	75.9	83.7	1.0	19.3	35.1	55.1	63.0	4.0	77.3
✓	✓	✓		✓	57.8	76.3	84.1	1.0	18.3	35.7	55.5	63.1	3.0	78.3
✓	✓		✓	✓	57.8	76.0	82.5	1.0	19.5	31.8	48.5	55.2	6.0	106.6
✓		✓	✓	✓	57.3	76.7	83.2	1.0	18.9	37.5	59.4	67.0	3.0	69.2
✓	✓	✓	✓	✓	58.2	75.8	83.5	1.0	18.9	36.4	56.5	63.8	3.0	75.7

Table 2: Retrieval performance with different combinations of four visual tokens (Phrase, Object, Event, Attribute) on MSR-VTT using EASIER. Best in **Bold**.

Methods	Venue	Text-to-Video Retrieval			
		R@1↑	R@5↑	R@10↑	MnR↓
<i>MSVD</i>					
RaP	EMNLP'22	35.9	64.3	73.7	-
LanguageBind	ICLR'24	52.2	79.4	87.3	-
EASIER		57.2	80.0	88.2	15.6
<i>ActivityNet</i>					
LanguageBind	ICLR'24	35.1	63.4	76.6	-
EASIER		59.0	71.4	77.0	387.4

Table 3: Text-Video retrieval results on MSVD and ActivityNet. The best results are marked in **bold**.

might be due to the detailed and specific information that phrases and objects offer, enhancing the model’s ability to accurately match and retrieve relevant video content.

Number of relevance-boosted captions. In this part, we aim to explore how many relevance-boosted captions work the best. More captions have the potential to offer more detailed descriptions, which may enhance the viewer’s comprehension of the visual content. Previous studies (Biten et al., 2019; Tang et al., 2023) have demonstrated that longer captions tend to be more descriptive and semantically rich, achieving improved comprehension and retrieval performance. However, more relevance-boosted captions also mean more noises are injected. So balancing the number of relevance-boosted captions would be highly important. From the results shown in Table 5, we notice that paraphrasing into two or three sentences significantly improved R@1, R@5, and R@10. Considering computational constraints and the similar effectiveness of paraphrasing into two and three sentences, we decide to boost into two sentences.

Complexity of prompt templates for structural

information extraction. The complexity of the prompt plays a pivotal role in shaping the output generated by the model, influencing the depth of analysis and the richness of information conveyed. An intricate prompt may provide the model with additional context and guidance, enabling it to produce more detailed responses. Specifically, we compare four templates (Basic, Structured, Detailed Description, and Narrative Format) offering different levels of complexity for organizing video content as shown in Appendix A.5. The results are shown in Table 6. We notice that intricate prompts provide the model with additional context and guidance, enabling it to produce more detailed responses. However, it may also lead to a loss of information, which is important for the retrieval performance. The results show that with the narrative format template, R@1, R@5, and R@10 on text-to-video retrieval are improved. This might be because the simplest format template provides insights into storytelling elements such as objects, events, and phrases, which leads to more precise keyword assignments.

5 Analysis on Quality of Relevance-Boosted Captions

As the details brought by relevance-boosted generation might bring irrelevant information, we analyze the quality of relevance-boosted captions.

5.1 Automatic Evaluation

Inspired by Li et al. (2023a), we generate video captions with varying levels of relevant details by using different prompts to control the level of relevance generation. Specifically, we generate cap-

Order List	Text-to-Video Retrieval					Video-to-Text Retrieval				
	R@1↑	R@5↑	R@10↑	MdR↓	MnR↓	R@1↑	R@5↑	R@10↑	MdR↓	MnR↓
Order List 1	58.2	75.8	83.5	1.0	18.9	36.4	56.5	63.8	3.0	75.7
Order List 2	57.9	75.9	83.4	1.0	18.7	36.7	56.4	64.4	3.0	75.3
Order List 3	58.0	75.7	83.2	1.0	19.1	36.3	56.6	64.2	3.0	75.0

Table 4: Retrieval performance with different order of four visual tokens (Phrase, Object, Event, Attribute) on MSR-VTT using EASIER. Best in **Bold**. We discovered three unique sequencing methods for visual tokens, each producing distinct outcomes based on their specific arrangements. Order List 1 places objects or phrases to the end, *i.e.*, {Caption}, ..., {Phrase/Object}, Order List 2 represents {Caption}, ..., {Event}, and Order List 3 represents {Caption}, ..., {Attribute}.

# of Hallucinated Captions	Text-to-Video Retrieval					Video-to-Text Retrieval				
	R@1↑	R@5↑	R@10↑	MdR↓	MnR↓	R@1↑	R@5↑	R@10↑	MdR↓	MnR↓
1	40.9	55.5	60.9	3.0	227.3	34.3	54.2	62.6	4.0	114.0
2	58.2	75.8	83.5	1.0	18.9	36.4	56.5	63.8	3.0	75.7
3	55.7	73.9	82.2	1.0	21.2	35.1	52.8	62.4	4.0	87.1

Table 5: Retrieval performance with different numbers of relevance-boosted captions on MSR-VTT using EASIER. Best in **Bold**.

Template	Text-to-Video Retrieval					Video-to-Text Retrieval				
	R@1↑	R@5↑	R@10↑	MdR↓	MnR↓	R@1↑	R@5↑	R@10↑	MdR↓	MnR↓
Basic Template	58.2	75.8	83.5	1.0	18.9	36.4	56.5	63.8	3.0	75.7
Structured Template	55.7	74.6	81.2	1.0	21.1	31.0	45.4	51.5	9.0	77.2
Template with Detailed Description	55.9	74.6	81.7	1.0	21.2	34.3	53.9	61.7	4.0	84.7
Narrative Format Template	56.5	74.7	81.7	1.0	20.9	26.9	43.4	49.0	11.0	129.7

Table 6: Retrieval performance with different template formats on MSR-VTT using EASIER. Best in **Bold**.

Relevance	Automatic Evaluation Metric		Human Evaluation				Text-to-Video Retrieval				Video-to-Text Retrieval				
	HHEM	Factual Accuracy	Relevance	Coherence	Specificity	R@1↑	R@5↑	R@10↑	MdR↓	MnR↓	R@1↑	R@5↑	R@10↑	MdR↓	MnR↓
High-level	16.1%	0.33	0.42	0.24	-0.95	56.9	75.1	82.6	1.0	21.4	23.1	37.0	43.2	22.0	147.6
Medium-level	14.7%	0.52	0.78	1.21	0.07	57.3	75.2	82.4	1.0	18.1	25.0	37.7	43.5	19.0	150.1
Low-level	9.6%	0.85	0.81	1.38	0.68	57.6	74.9	83.3	1.0	19.1	34.7	52.6	64.2	3.0	88.6
EASIER	10.9%	0.87	0.86	1.28	0.52	58.2	75.8	83.5	1.0	18.9	36.4	56.5	63.8	3.0	75.7

Table 7: Retrieval performance with different level of Relevance Boosting on MSR-VTT. Best in **Bold**.

tions at three levels: high, medium, and low (see Appendix B). We used the HHEM model (Honovich et al., 2022) to compute the hallucination rate by comparing the relevance-boosted captions and original video captions. As shown in Table 7, lower levels of generation do not significantly change retrieval results. We also observe that captions with a lower boosting rate perform worse than captions with higher levels.

5.2 Human Evaluation

We also conduct a human evaluation to further evaluate the relevance-boosted captions.

Participants: Our human evaluation task involves reading relevance-boosted captions from different levels, video captions without relevance-boosting, and rating those relevance-boosted captions from them. We recruited 10 participants (7M, 3F). We conducted a rigorous qualification process, evaluating their English proficiency, to ensure high-quality annotations. We hired them by sending invited emails to graduate students. We allocated up to 30 minutes for each participant to complete the study, and for their valuable time and input, each

participant received a compensation of \$15.

Task: We randomly selected 50 pairs of relevance-boosted captions and original video captions from EASIER. Note that each pair has only one relevance-boosted caption and one original video caption. Each participant is assigned 50 pairs. Each pair is evaluated by 10 individuals. In each trial, a participant reads 4 relevance-boosted captions for the original video caption: one by high-level boosting, one by medium-level boosting, one by low-level boosting, and one from EASIER. The order of these four is also randomized, so participants do not know which generated caption is from which method. The participant is asked to rate the 4 captions along four dimensions using a five-point Likert scale from -2 to 2.

- *Factual Accuracy:* The relevance-boosted caption is factually correct to convey the content from the video caption.
- *Relevance:* The relevance-boosted caption is relevant to the video caption.
- *Coherence:* The relevance-boosted caption is coherent to the video caption.
- *Specificity:* The relevance-boosted caption is spe-

LLM	Text-to-Video Retrieval					Video-to-Text Retrieval				
	R@1↑	R@5↑	R@10↑	MdR↓	MnR↓	R@1↑	R@5↑	R@10↑	MdR↓	MnR↓
LLaMA	58.2	75.8	83.5	1.0	18.9	36.4	56.5	63.8	3.0	75.7
GPT 3.5	61.2	80.4	86.8	1.0	15.0	35.5	58.3	65.5	3.0	77.6

Table 8: Retrieval performance with different LLM models on MSR-VTT using EASIER. Best in **Bold**.

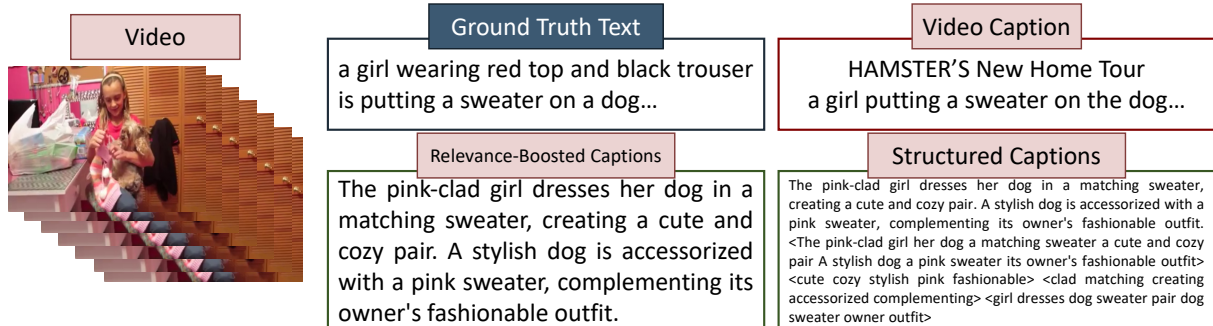


Figure 3: An example. Relevance-boosted captions contain more information compared to vanilla video captions in the video though some noises are also injected.

cific and detailed to the video caption.

Evaluation Results: We conducted *Wilcoxon tests* (Woolson, 2007) with a significance level of 0.05 to compare the performance of high-level, medium-level, low-level boosting, and EASIER in the Factual Accuracy, Relevance, Coherence, and Specificity dimensions. The Wilcoxon test is a non-parametric statistical test used to compare two paired groups of data. The obtained p-values indicate the probability of observing the reported differences if there were no true differences between the models.

The results indicate significant differences in the Factual Accuracy dimension, where EASIER outperforms High-level boosting ($V = 4836$, $p = 1.45e-30$), Medium-level boosting ($V = 4819$, $p = 7.22e-31$). For the Coherence dimension, we notice that they are almost at the same level, likely because captions refined by the LLM are already sufficiently coherent for users. In the Relevance dimension, EASIER surpasses high-level boosting ($V = 3247$, $p = 1.44e-21$), medium-level boosting ($V = 3693$, $p = 1.69e-20$), low-level boosting ($V = 3188$, $p = 1.53e-20$). For the Specificity dimension which considers whether the relevance-boosted caption is detailed and specified, Low-level boosting outperforms all methods: High-level boosting ($V = 4463$, $p = 1.25e-7$), Medium-level boosting ($V = 3830$, $p = 3.48e-14$), EASIER ($V = 2260$, $p = 2.63e-7$). It is worth noting that while low-level boosting is more detailed than EASIER, it performs slightly worse in VTR, possibly due to the higher importance of factual accuracy in evaluating the effectiveness of relevance-boosted captions. Future

work can focus on designing an innovative framework for the relevance-boosted captioning method to integrate useful dimensions.

5.3 Qualitative Results

To qualitatively validate the effectiveness of EASIER, we present an example in fig. 3. The retrieval results show that relevance-boosted captions have more information in the video than vanilla video captions. Besides, our proposed structural information extraction methods clearly emphasize the important visual tokens, *i.e.*, phrase, object, event, and attribute, further boosting the performance.

6 Conclusion

In this paper, we present an innovative zero-shot framework, EASIER, which revolutionizes video-text retrieval by capitalizing on existing captioning methods, large language models (LLMs), and text retrieval techniques. By sidestepping the need for model training or fine-tuning, our framework offers a streamlined approach to retrieval. To overcome the shortcomings of traditional captioning methods, we propose a groundbreaking relevance-boosted caption generation technique that incorporates LLMs' generated information into video captions. Moreover, our introduction of structural information extraction further enhances retrieval performance by highlighting key visual tokens. Through extensive experimentation across diverse benchmarks, we demonstrate the superior efficacy of EASIER compared to conventional fine-tuned and pretraining methods, even in the absence of training data.

544
545
546
547
548
549
550
551
552
553

554

555
556
557
558
559
560

561
562
563
564
565
566
567

568
569
570
571

572
573
574
575
576
577

578
579
580
581
582
583
584
585
586
587
588
589
590
591
592

593
594
595
596
597

Limitations

In the future, it would be interesting to explore more detailed methods for zero-shot video-text retrieval, such as incorporating the audio modality and corresponding off-the-shelf foundation models. Moreover, as a pioneering work, our work mainly focuses on establishing the paradigm. It would be great if we could explore more text retrieval methods, video captioning methods, and LLMs for relevance-boosted caption generation.

References

Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. 2021. [Frozen in time: A joint video and image encoder for end-to-end retrieval](#). In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pages 1708–1718. IEEE.

Gedas Bertasius, Heng Wang, and Lorenzo Torresani. 2021. [Is space-time attention all you need for video understanding?](#) In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 813–824. PMLR.

Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural language processing with Python: analyzing text with the natural language toolkit*. " O'Reilly Media, Inc."

Ali Furkan Biten, Lluís Gomez, Marçal Rusinol, and Dimosthenis Karatzas. 2019. Good news, everyone! context driven entity-aware captioning for news images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

Shuqiang Cao, Bairui Wang, Wei Zhang, and Lin Ma. 2022. [Visual consensus modeling for video-text retrieval](#). In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence,*

IAAI 2022, The Twelveth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022, pages 167–175. AAAI Press.

David M. Chan, Austin Myers, Sudheendra Vijayanarasimhan, David A. Ross, and John Canny. 2023. [\\$IC^3\\$: Image Captioning by Committee Consensus](#). ArXiv:2302.01328 [cs].

David Chen and William Dolan. 2011. [Collecting highly parallel data for paraphrase evaluation](#). In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 190–200, Portland, Oregon, USA. Association for Computational Linguistics.

Dongsheng Chen, Chaofan Tao, Lu Hou, Lifeng Shang, Xin Jiang, and Qun Liu. 2022. [LiteVL: Efficient video-language learning with enhanced spatial-temporal modeling](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7985–7997, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Sihan Chen, Handong Li, Qunbo Wang, Zijia Zhao, Mingzhen Sun, Xinxin Zhu, and Jing Liu. 2023a. [VAST: A vision-audio-subtitle-text omni-modality foundation model and dataset](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.

Yizhen Chen, Jie Wang, Lijian Lin, Zhongang Qi, Jin Ma, and Ying Shan. 2023b. [Tagging before Alignment: Integrating Multi-Modal Tags for Video-Text Retrieval](#). In *AAAI Conference on Artificial Intelligence*. arXiv. ArXiv:2301.12644 [cs].

Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. 2023c. [Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks](#). *arXiv preprint arXiv:2312.14238*.

Xing Cheng, Hezheng Lin, Xiangyu Wu, Fan Yang, and Dong Shen. 2021. [Improving video-text retrieval by multi-stream corpus alignment and dual softmax loss](#). *CoRR*, abs/2109.04290.

Ioana Croitoru, Simion-Vlad Bogolin, Marius Litordeanu, Hailin Jin, Andrew Zisserman, Samuel Albanie, and Yang Liu. 2021. [Teachtext: Crossmodal generalized distillation for text-video retrieval](#). In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pages 11563–11573. IEEE.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for*

654		Qingpei Guo, Furong Xu, Hanxiao Zhang, Wang Ren,	711
655	<i>Computational Linguistics: Human Language Tech-</i>	Ziping Ma, Lin Ju, Jian Wang, Jingdong Chen, and	712
656	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages	Ming Yang. 2024. M2-encoder: Advancing bilingual	713
657	4171–4186, Minneapolis, Minnesota. Association for	image-text understanding by large-scale efficient pre-	714
	Computational Linguistics.	training.	715
658	Jianfeng Dong, Xirong Li, Chaoxi Xu, Shouling Ji,	Or Honovich, Roei Aharoni, Jonathan Herzig, Ha-	716
659	Yuan He, Gang Yang, and Xun Wang. 2019. Dual	gai Taitelbaum, Doron Kukliansy, Vered Cohen,	717
660	encoding for zero-example video retrieval. In <i>IEEE</i>	Thomas Scialom, Idan Szpektor, Avinatan Has-	718
661	<i>Conference on Computer Vision and Pattern Recogni-</i>	sidim, and Yossi Matias. 2022. True: Re-evaluating	719
662	<i>tion, CVPR 2019, Long Beach, CA, USA, June 16-20,</i>	factual consistency evaluation. <i>arXiv preprint</i>	720
663	<i>2019</i> , pages 9346–9355. Computer Vision Founda-	<i>arXiv:2204.04991.</i>	721
664	tion / IEEE.		
665	Maksim Dzabraev, Maksim Kalashnikov, Stepan	Siteng Huang, Biao Gong, Yulin Pan, Jianwen Jiang,	722
666	Komkov, and Aleksandr Petiushko. 2021. MDMMT:	Yiliang Lv, Yuyuan Li, and Donglin Wang. 2023.	723
667	multidomain multimodal transformer for video re-	VoP: Text-Video Co-Operative Prompt Tuning for	724
668	trieval. In <i>IEEE Conference on Computer Vision and</i>	Cross-Modal Retrieval. In <i>Proceedings of the</i>	725
669	<i>Pattern Recognition Workshops, CVPR Workshops</i>	<i>IEEE/CVF Conference on Computer Vision and Pat-</i>	726
670	<i>2021, virtual, June 19-25, 2021</i> , pages 3354–3363.	<i>tern Recognition</i> , pages 6565–6574.	727
671	Computer Vision Foundation / IEEE.		
672	Bernard Ghanem Fabian Caba Heilbron, Victor Escor-	Mark J. Huiskes and Michael S. Lew. 2008. The MIR	728
673	cia and Juan Carlos Niebles. 2015. Activitynet: A	Flickr Retrieval Evaluation. In <i>Proceedings of the 1st</i>	729
674	large-scale video benchmark for human activity un-	<i>ACM International Conference on Multimedia Infor-</i>	730
675	derstanding. In <i>Proceedings of the IEEE Conference</i>	<i>mation Retrieval, MIR '08</i> , pages 39–43, New York,	731
676	<i>on Computer Vision and Pattern Recognition</i> , pages	NY, USA. Association for Computing Machinery.	732
677	961–970.	Event-place: Vancouver, British Columbia, Canada.	733
678	Bo Fang, Wenhao Wu, Chang Liu, Yu Zhou, Yuxin	Weike Jin, Zhou Zhao, Pengcheng Zhang, Jieming Zhu,	734
679	Song, Weiping Wang, Xiangbo Shu, Xiangyang Ji,	Xiuqiang He, and Yueting Zhuang. 2021. Hierar-	735
680	and Jingdong Wang. 2023. Uatvr: Uncertainty-	chical cross-modal graph consistency learning for	736
681	adaptive text-video retrieval.	video-text retrieval. In <i>SIGIR '21: The 44th Inter-</i>	737
682	Xiang Fang, Daizong Liu, Pan Zhou, and Yuchong Hu.	<i>national ACM SIGIR Conference on Research and</i>	738
683	2022. Multi-modal cross-domain alignment network	<i>Development in Information Retrieval, Virtual Event,</i>	739
684	for video moment retrieval. <i>IEEE Transactions on</i>	<i>Canada, July 11-15, 2021</i> , pages 1114–1124. ACM.	740
685	<i>Multimedia.</i>		
686	Valentin Gabeur, Chen Sun, Karteek Alahari, and	Xiaojie Jin, Bowen Zhang, Weibo Gong, Kai Xu, Xue-	741
687	Cordelia Schmid. 2020. Multi-modal transformer	Qing Deng, Peng Wang, Zhao Zhang, Xiaohui Shen,	742
688	for video retrieval. In <i>Computer Vision - ECCV 2020</i>	and Jiashi Feng. 2024. Mv-adapter: Multimodal	743
689	<i>- 16th European Conference, Glasgow, UK, August</i>	video transfer learning for video text retrieval.	744
690	<i>23-28, 2020, Proceedings, Part IV</i> , volume 12349 of		
691	<i>Lecture Notes in Computer Science</i> , pages 214–229.	Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L.	745
692	Springer.	Berg, Mohit Bansal, and Jingjing Liu. 2021. Less is	746
693	Zhe Gan, Yen-Chun Chen, Linjie Li, Chen Zhu,	more: Clipbert for video-and-language learning via	747
694	Yu Cheng, and Jingjing Liu. 2020. Large-scale adversar-	sparse sampling. In <i>IEEE Conference on Computer</i>	748
695	ial training for vision-and-language representation	<i>Vision and Pattern Recognition, CVPR 2021, virtual,</i>	749
696	learning. In <i>Advances in Neural Information Process-</i>	<i>June 19-25, 2021</i> , pages 7331–7341. Computer Vi-	750
697	<i>ing Systems 33: Annual Conference on Neural</i>	sion Foundation / IEEE.	751
698	<i>Information Processing Systems 2020, NeurIPS 2020,</i>		
699	<i>December 6-12, 2020, virtual.</i>	Junyi Li, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun	752
700	Zijian Gao, Jingyu Liu, Sheng Chen, Dedan Chang,	Nie, and Ji-Rong Wen. 2023a. Halueval: A large-	753
701	Hao Zhang, and Jinwei Yuan. 2021. CLIP2TV: an	scale hallucination evaluation benchmark for large	754
702	empirical study on transformer-based methods for	language models. In <i>Proceedings of the 2023 Con-</i>	755
703	video-text retrieval. <i>CoRR</i> , abs/2111.05610.	<i>ference on Empirical Methods in Natural Language</i>	756
704	Satya Krishna Gorti, Noël Vouitsis, Junwei Ma, Key-	<i>Processing</i> , pages 6449–6464.	757
705	van Golestan, Maksims Volkovs, Animesh Garg, and	Kunchang Li, Yali Wang, Yizhuo Li, Yi Wang, Yinan	758
706	Guangwei Yu. 2022. X-pool: Cross-modal language-	He, Limin Wang, and Yu Qiao. 2023b. Unmasked	759
707	video attention for text-video retrieval. In <i>IEEE/CVF</i>	teacher: Towards training-efficient video foundation	760
708	<i>Conference on Computer Vision and Pattern Recogni-</i>	models. In <i>Proceedings of the IEEE/CVF Interna-</i>	761
709	<i>tion, CVPR 2022, New Orleans, LA, USA, June</i>	<i>tional Conference on Computer Vision (ICCV)</i> , pages	762
710	<i>18-24, 2022</i> , pages 4996–5005. IEEE.	19948–19960.	763
		Linjie Li, Yen-Chun Chen, Yu Cheng, Zhe Gan, Licheng	764
		Yu, and Jingjing Liu. 2020a. HERO: Hierarchical	765
		encoder for Video+Language omni-representation	766
		pre-training. In <i>Proceedings of the 2020 Conference</i>	767

768				
769				
770				
771	Pandeng Li, Chen-Wei Xie, Liming Zhao, Hongtao			
772	Xie, Jiannan Ge, Yun Zheng, Deli Zhao, and Yong-			
773	dong Zhang. 2023c. Progressive spatio-temporal pro-			
774	totype matching for text-video retrieval . In <i>2023</i>			
775	<i>IEEE/CVF International Conference on Computer</i>			
776	<i>Vision (ICCV)</i> , pages 4077–4087.			
777	Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang,			
778	Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong			
779	Hu, Li Dong, Furu Wei, Yejin Choi, and Jianfeng			
780	Gao. 2020b. Oscar: Object-semantics aligned pre-			
781	training for vision-language tasks . In <i>Computer Vi-</i>			
782	<i>sion - ECCV 2020 - 16th European Conference, Glas-</i>			
783	<i>gow, UK, August 23-28, 2020, Proceedings, Part</i>			
784	<i>XXX</i> , volume 12375 of <i>Lecture Notes in Computer</i>			
785	<i>Science</i> , pages 121–137. Springer.			
786	Weixin Liang, Yuhui Zhang, Yongchan Kwon, Serena			
787	Yeung, and James Zou. 2022. Mind the gap: Under-			
788	standing the modality gap in multi-modal contrastive			
789	representation learning . In <i>Advances in neural infor-</i>			
790	mation processing systems .			
791	Tsung-Yi Lin, Michael Maire, Serge Belongie, James			
792	Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,			
793	and C. Lawrence Zitnick. 2014. Microsoft COCO:			
794	Common Objects in Context . In <i>Computer Vision –</i>			
795	<i>ECCV 2014</i> , pages 740–755, Cham. Springer Inter-			
796	national Publishing.			
797	Alexander Liu, SouYoung Jin, Cheng-I Lai, Andrew			
798	Rouditchenko, Aude Oliva, and James Glass. 2022a.			
799	Cross-modal discrete representation learning . In <i>Pro-</i>			
800	<i>ceedings of the 60th Annual Meeting of the Associa-</i>			
801	<i>tion for Computational Linguistics (Volume 1: Long</i>			
802	<i>Papers)</i> , pages 3013–3035, Dublin, Ireland. Associa-			
803	tion for Computational Linguistics.			
804	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae			
805	Lee. 2023a. Improved baselines with visual instruc-			
806	tion tuning.			
807	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan			
808	Zhang, Sheng Shen, and Yong Jae Lee. 2024. Llava-			
809	next: Improved reasoning, ocr, and world knowledge .			
810	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae			
811	Lee. 2023b. Visual instruction tuning.			
812	Ruyang Liu, Chen Li, Yixiao Ge, Ying Shan, Thomas H			
813	Li, and Ge Li. 2023c. One for all: Video conversation			
814	is feasible without video instruction tuning. <i>arXiv</i>			
815	<i>preprint arXiv:2309.15785</i> .			
816	Yuqi Liu, Pengfei Xiong, Luhui Xu, Shengming Cao,			
817	and Qin Jin. 2022b. Ts2-net: Token shift and selec-			
818	tion transformer for text-video retrieval . In <i>Computer</i>			
819	<i>Vision - ECCV 2022 - 17th European Conference, Tel</i>			
820	<i>Aviv, Israel, October 23-27, 2022, Proceedings, Part</i>			
821	<i>XIV</i> , volume 13674 of <i>Lecture Notes in Computer</i>			
822	<i>Science</i> , pages 319–335. Springer.			
	Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen			823
	Lei, Nan Duan, and Tianrui Li. 2022. Clip4clip: An			824
	empirical study of CLIP for end to end video clip			825
	retrieval and captioning . <i>Neurocomputing</i> , 508:293–			826
	304.			827
	Yiwei Ma, Guohai Xu, Xiaoshuai Sun, Ming Yan,			828
	Ji Zhang, and Rongrong Ji. 2022. X-CLIP: End-to-			829
	end multi-grained contrastive learning for video-text			830
	retrieval . In <i>ACM international conference on mul-</i>			831
	<i>timedia, MM '22</i> , pages 638–647, New York, NY,			832
	USA. Association for Computing Machinery. Num-			833
	ber of pages: 10 Place: Lisboa, Portugal.			834
	Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira,			835
	Ivan Laptev, Josef Sivic, and Andrew Zisserman.			836
	2020. End-to-end learning of visual representa-			837
	tions from uncurated instructional videos . In <i>2020</i>			838
	<i>IEEE/CVF Conference on Computer Vision and Pat-</i>			839
	<i>tern Recognition, CVPR 2020, Seattle, WA, USA,</i>			840
	<i>June 13-19, 2020</i> , pages 9876–9886. Computer Vi-			841
	sion Foundation / IEEE.			842
	Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac,			843
	Makarand Tapaswi, Ivan Laptev, and Josef Sivic.			844
	2019. Howto100m: Learning a text-video embed-			845
	ding by watching hundred million narrated video			846
	clips . In <i>2019 IEEE/CVF International Confer-</i>			847
	<i>ence on Computer Vision, ICCV 2019, Seoul, Ko-</i>			848
	<i>rea (South), October 27 - November 2, 2019</i> , pages			849
	2630–2640. IEEE.			850
	Jae Sung Park, Sheng Shen, Ali Farhadi, Trevor Darrell,			851
	Yejin Choi, and Anna Rohrbach. 2022a. Exposing			852
	the limits of video-text models through contrast sets .			853
	In <i>Proceedings of the 2022 Conference of the North</i>			854
	<i>American Chapter of the Association for Computa-</i>			855
	<i>tional Linguistics: Human Language Technologies,</i>			856
	pages 3574–3586, Seattle, United States. Associa-			857
	tion for Computational Linguistics.			858
	Yookoon Park, Mahmoud Azab, Seungwhan Moon,			859
	Bo Xiong, Florian Metze, Gourab Kundu, and Kir-			860
	mani Ahmed. 2022b. Normalized contrastive learn-			861
	ing for text-video retrieval . In <i>Proceedings of the</i>			862
	<i>2022 Conference on Empirical Methods in Natural</i>			863
	<i>Language Processing</i> , pages 248–260, Abu Dhabi,			864
	United Arab Emirates. Association for Computa-			865
	tional Linguistics.			866
	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya			867
	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-			868
	try, Amanda Askell, Pamela Mishkin, Jack Clark,			869
	Gretchen Krueger, and Ilya Sutskever. 2021. Learn-			870
	ing transferable visual models from natural language			871
	supervision . In <i>Proceedings of the 38th International</i>			872
	<i>Conference on Machine Learning, ICML 2021, 18-24</i>			873
	<i>July 2021, Virtual Event</i> , volume 139 of <i>Proceedings</i>			874
	<i>of Machine Learning Research</i> , pages 8748–8763.			875
	PMLR.			876
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,			877
	Dario Amodei, Ilya Sutskever, et al. 2019. Language			878
	models are unsupervised multitask learners. <i>OpenAI</i>			879
	<i>blog</i> , 1(8):9.			880

881	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. <i>arXiv preprint arXiv:1908.10084</i> .	937
882		938
883		939
884	Stephen E Robertson and Steve Walker. 1994. Some simple effective approximations to the 2-poisson model for probabilistic weighted retrieval. In <i>SIGIR'94: Proceedings of the Seventeenth Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval, organised by Dublin City University</i> , pages 232–241. Springer.	940
885		941
886		942
887		943
888		944
889		945
890		946
891	Mike Schuster and Kuldeep K. Paliwal. 1997. Bidirectional recurrent neural networks. <i>IEEE Trans. Signal Process.</i> , 45(11):2673–2681.	947
892		948
893		949
894	Nina Shvetsova, Brian Chen, Andrew Rouditchenko, Samuel Thomas, Brian Kingsbury, Rogerio Feris, David Harwath, James Glass, and Hilde Kuehne. 2022. Everything at Once – Multi-modal Fusion Transformer for Video Retrieval . In <i>2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 19988–19997, New Orleans, LA, USA. IEEE.	950
895		951
896		952
897		953
898		954
899		955
900		956
901		
902	Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. 2022. FLAVA: A foundational language and vision alignment model . In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022</i> , pages 15617–15629. IEEE.	957
903		958
904		959
905		960
906		961
907		962
908		
909	Benny J Tang, Angie Boggust, and Arvind Satyanarayan. 2023. Vistext: A benchmark for semantically rich chart captioning. <i>arXiv preprint arXiv:2307.05356</i> .	963
910		964
911		965
912	Yoad Tewel, Yoav Shalev, Roy Nadler, Idan Schwartz, and Lior Wolf. 2022. Zero-shot video captioning with evolving pseudo-tokens. <i>arXiv preprint arXiv:2207.11100</i> .	966
913		967
914		968
915		969
916	Yoad Tewel, Yoav Shalev, Idan Schwartz, and Lior Wolf. 2021. Zero-shot image-to-text generation for visual-semantic arithmetic. <i>arXiv preprint arXiv:2111.14447</i> , 1(3):6.	970
917		971
918		972
919		973
920	Anthony Meng Huat Tiong, Junnan Li, Boyang Li, Silvio Savarese, and Steven C.H. Hoi. 2022. Plug-and-play VQA: Zero-shot VQA by conjoining large pre-trained models with zero training . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 951–967, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	974
921		975
922		976
923		977
924		978
925		979
926		
927	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura,	980
928		981
929		982
930		983
931		984
932		985
933		986
934		987
935		988
936		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

994	Yimu Wang, Bo Xue, Quan Cheng, Yuhui Chen, and Lijun Zhang. 2021. Deep unified cross-modality hashing by pairwise data alignment. In <i>Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21</i> , pages 1129–1135.	1050
995		1051
996		1052
997		1053
998		1054
999	Zhenhailong Wang, Manling Li, Ruochen Xu, Luowei Zhou, Jie Lei, Xudong Lin, Shuohang Wang, Ziyi Yang, Chenguang Zhu, Derek Hoiem, Shih-Fu Chang, Mohit Bansal, and Heng Ji. 2022e. Language models with image descriptors are strong few-shot video-language learners. In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .	1055
1000		1056
1001		1057
1002		1058
1003		1059
1004		
1005		
1006		
1007		
1008		
1009	Ziyang Wang, Yi-Lin Sung, Feng Cheng, Gedas Bertasius, and Mohit Bansal. 2023. Unified coarse-to-fine alignment for video-text retrieval.	1060
1010		1061
1011		1062
1012	Robert F Woolson. 2007. Wilcoxon signed-rank test. <i>Wiley encyclopedia of clinical trials</i> , pages 1–3.	1063
1013		1064
1014	Wenhao Wu, Haipeng Luo, Bo Fang, Jingdong Wang, and Wanli Ouyang. 2023. Cap4Video: What Can Auxiliary Captions Do for Text-Video Retrieval? In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 10704–10713.	1065
1015		1066
1016		1067
1017		1068
1018		1069
1019		1070
1020	Xing Wu, Chaochen Gao, Zijia Lin, Zhongyuan Wang, Jizhong Han, and Songlin Hu. 2022. RaP: Redundancy-aware video-language pre-training for text-video retrieval. In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 3036–3047, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	1071
1021		1072
1022		1073
1023		1074
1024		1075
1025		1076
1026		1077
1027	Haiyang Xu, Qinghao Ye, Ming Yan, Yaya Shi, Jiabo Ye, Yuanhong Xu, Chenliang Li, Bin Bi, Qi Qian, Wei Wang, Guohai Xu, Ji Zhang, Songfang Huang, Fei Huang, and Jingren Zhou. 2023. mplug-2: A modularized multi-modal foundation model across text, image and video. <i>ArXiv</i> , abs/2302.00402.	1078
1028		1079
1029		1080
1030		1081
1031		1082
1032		1083
1033	Hu Xu, Gargi Ghosh, Po-Yao Huang, Prahal Arora, Masoumeh Aminzadeh, Christoph Feichtenhofer, Florian Metze, and Luke Zettlemoyer. 2021a. VLM: Task-agnostic video-language model pre-training for video understanding. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 4227–4239, Online. Association for Computational Linguistics.	1084
1034		1085
1035		1086
1036		1087
1037		1088
1038		1089
1039		1090
1040		1091
1041	Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metze, Luke Zettlemoyer, and Christoph Feichtenhofer. 2021b. Video-CLIP: Contrastive pre-training for zero-shot video-text understanding. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 6787–6800, Punta Cana, Dominican Republic. Association for Computational Linguistics.	1092
1042		1093
1043		1094
1044		1095
1045		1096
1046		1097
1047		1098
1048		1099
1049		
	Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. MSR-VTT: A Large Video Description Dataset for Bridging Video and Language. In <i>2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 5288–5296, Las Vegas, NV, USA. IEEE.	1050
		1051
		1052
		1053
		1054
	Hongwei Xue, Yuchong Sun, Bei Liu, Jianlong Fu, Ruihua Song, Houqiang Li, and Jiebo Luo. 2023. CLIP-ViP: Adapting Pre-trained Image-Text Model to Video-Language Alignment. In <i>The Eleventh International Conference on Learning Representations</i> .	1055
		1056
		1057
		1058
		1059
	Qiyang Yu, Yang Liu, Yimu Wang, Ke Xu, and Jingjing Liu. 2023. Multimodal federated learning via contrastive representation ensemble. In <i>The Eleventh International Conference on Learning Representations</i> .	1060
		1061
		1062
		1063
		1064
	Youngjae Yu, Jongseok Kim, and Gunhee Kim. 2018. A joint sequence fusion model for video question answering and retrieval. In <i>Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part VII</i> , volume 11211 of <i>Lecture Notes in Computer Science</i> , pages 487–503. Springer.	1065
		1066
		1067
		1068
		1069
		1070
		1071
	Youngjae Yu, Hyungjin Ko, Jongwook Choi, and Gunhee Kim. 2017. End-to-end concept word detection for video captioning, retrieval, and question answering. In <i>2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017</i> , pages 3261–3269. IEEE Computer Society.	1072
		1073
		1074
		1075
		1076
		1077
		1078
	Shuai Zhao, Linchao Zhu, Xiaohan Wang, and Yi Yang. 2022. Centerclip: Token clustering for efficient text-video retrieval. In <i>Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '22</i> , page 970–981, New York, NY, USA. Association for Computing Machinery.	1079
		1080
		1081
		1082
		1083
		1084
		1085
	Bin Zhu, Bin Lin, Munan Ning, Yang Yan, Jiaxi Cui, WANG HongFa, Yatian Pang, Wenhao Jiang, Junwu Zhang, Zongwei Li, Cai Wan Zhang, Zhifeng Li, Wei Liu, and Li Yuan. 2024. Languagebind: Extending video-language pretraining to n-modality by language-based semantic alignment. In <i>The Twelfth International Conference on Learning Representations</i> .	1086
		1087
		1088
		1089
		1090
		1091
		1092
		1093
	Linchao Zhu and Yi Yang. 2020. Actbert: Learning global-local video-text representations. In <i>2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020</i> , pages 8743–8752. Computer Vision Foundation / IEEE.	1094
		1095
		1096
		1097
		1098
		1099

A Experiments

A.1 Details of Benchmarks

- **MSR-VTT** (Xu et al., 2016) contains 10,000 videos with length varying from 10 to 32 seconds, each paired with about 20 human-labeled captions. Following the evaluation protocol from previous works (Yu et al., 2018; Miech et al., 2019), we use the training-9k / test 1k-A splits for training and testing respectively.
- **MSVD** (Chen and Dolan, 2011) contains 1,970 videos with a split of 1200, 100, and 670 as the train, validation, and test set, respectively. The duration of videos varies from 1 to 62 seconds. Each video is paired with 40 English captions.
- **ActivityNet** (Fabian Caba Heilbron and Niebles, 2015) is consisted of 20,000 Youtube videos with 100,000 densely annotated descriptions. For a fair comparison, following the previous setting (Luo et al., 2022; Gabeur et al., 2020), we concatenate all captions together as a paragraph to perform a video-paragraph retrieval task by concatenating all the descriptions of a video. Performances are reported on the “val1” split of the ActivityNet.

A.2 Baselines

To show the empirical efficiency of our EASIER, we compare it with fine-tuned models (LiteVL (Chen et al., 2022), NCL (Park et al., 2022b), TABLE (Chen et al., 2023b), VOP (Huang et al., 2023), X-CLIP (Ma et al., 2022), Discrete-Codebook (Liu et al., 2022a), TS2-Net (Liu et al., 2022b), VCM (Cao et al., 2022), HiSE (Wang et al., 2022b), CenterCLIP (Zhao et al., 2022), X-Pool (Gorti et al., 2022), S3MA (Wang and Shi, 2023b)), and MV-Apapter (Jin et al., 2024), pre-trained methods (VLM (Xu et al., 2021a), HERO (Li et al., 2020a), VideoCLIP (Xu et al., 2021b), EvO (Shvetsova et al., 2022), OA-Trans (Wang et al., 2022a), RaP (Wu et al., 2022), OmniVL (Wang et al., 2022c), mPLUG-2 (Xu et al., 2023), InternVL (Chen et al., 2023c), Languange-Bind (Zhu et al., 2024), UCOFIA (Wang et al., 2023), ProST (Li et al., 2023c), and UATVR (Fang et al., 2023),), and a few-shot method, *i.e.*, VidIL (Wang et al., 2022e).

A.3 Implementation Details

For video caption generation, we use [Tewel et al. \(2021, 2022\)](#) to generate video captions and [GPT-2 \(Radford et al., 2019\)](#) to enrich sentences. For relevance-boosted caption generation, we employ [LLaMA2-7b-chat-hf \(Touvron et al., 2023\)](#) and get two boosted captions. For structural information extraction, we use [NLTK \(Bird et al., 2009\)](#). For text retrieval, we use [BM25 \(Robertson and Walker, 1994\)](#).

We use [GPT2 \(Radford et al., 2019\)](#) for sentence enrichment during video caption generation. [GPT-2 \(Radford et al., 2019\)](#), developed by OpenAI, is a large-scale transformer-based language model renowned for its ability to generate coherent and contextually relevant text. With 1.5 billion parameters, GPT-2 can be fine-tuned for a variety of natural language processing tasks, such as text generation, summarization, and captioning. In our task, we enrich image captions with GPT-2 with one NVIDIA A100 GPU using around 20 hours.

We use [Llama \(Touvron et al., 2023\)](#)(version: Llama-2-7b-chat-hf) to conduct the relevance-boosted caption generation task. [Llama \(Touvron et al., 2023\)](#) is an advanced language model with approximately 65 billion parameters. Its default backend is designed for efficiency and scalability. The computational budget for LLaMA in our task is approximately 23 hours with one NVIDIA A100 GPU. Its ability to understand context, generate coherent and contextually relevant responses, and perform a wide range of language-related tasks is significantly enhanced. LLaMA is a powerful and accessible tool, widely used in various applications. Therefore, it is included as an advanced baseline.

A.4 Choice of Retrieval Methods

In this part, we investigate the impact of different retrieval methods, *i.e.*, [BM25 \(Robertson and Walker, 1994\)](#) and sentence transformers ([Reimers and Gurevych, 2019](#)). The results are shown in section 6. It shows that [BM25](#) outperforms the sentence transformer.

A.5 Prompts for Structural Information Extraction

1. Basic Template: the simplest, providing a straightforward list of video elements, the one shown in Section 3.3.
2. Structured Template: It adds categorized elements, making the information easier to navi-

Retrieval Methods	Text-to-Video Retrieval					Video-to-Text Retrieval				
	R@1↑	R@5↑	R@10↑	MdR↓	MnR↓	R@1↑	R@5↑	R@10↑	MdR↓	MnR↓
BM25	58.2	75.8	83.5	1.0	18.9	36.4	56.5	63.8	3.0	75.7
Sentence Transformer	41.2	62.1	70.5	2.0	33.5	34.7	57.8	67.5	3.0	36.0

Table 9: Retrieval performance with different retrieval models on MSR-VTT using EASIER. Best in **Bold**.

Retrieval Methods	Text-to-Video Retrieval					Video-to-Text Retrieval				
	R@1↑	R@5↑	R@10↑	MdR↓	MnR↓	R@1↑	R@5↑	R@10↑	MdR↓	MnR↓
EASIER	58.2	75.8	83.5	1.0	18.9	36.4	56.5	63.8	3.0	75.7
EASIER (video caption only repeats to the same length as structured caption)	54.0	73.9	80.2	1.0	24.6	27.9	41.3	47.3	491.5	136.2
EASIER (structured information extraction only repeats to the same length as video caption)	18.6	25.1	27.1	15.0	444.6	25.1	40.2	44.9	21.0	287.6

Table 10: Comparative Analysis of Caption Repetition and Extracted Structural Information Repetition on Retrieval Performance

1196	gate for the retrieval method.	offering insights into the	1234
1197		unfolding story.	1235
1198			
1199	Video Caption : <Caption>. Key		
1200	Phrases: <{Phrases}>. Main	A.6 Ablation Study: Are Structural	1237
1201	Objects: <Objects>. Notable	Information Extraction and	1238
1202	Features: <{Attributes}>. Key	Relevance-Boosted Caption Generation	1239
1203	Events: <Events>	necessary?	1240
1204	3. Template with Detailed Description: This fur-	We also conduct another ablation study to investi-	1241
1205	ther elaborates on each element, offering in-	gate the effect of the video caption repeating itself	1242
1206	depth insights.	several times to form text that is the same length	1243
1207		as the structured caption stage. According to the	1244
1208	Detailed Video Description:	Table section 6, we find that our EASIER method	1245
1209	Caption: <{Caption}> Objects and	outperforms the others, indicating that a blend of	1246
1210	Attributes Overview: Each	relevance boosting (imagined or generated content)	1247
1211	object, <{Objects}>, is detailed	and structured information significantly improves	1248
1212	with attributes such as <{	retrieval results. Specifically, in text-to-video re-	1249
1213	Attributes}> to provide a	trieval, EASIER achieves much higher recall rates	1250
1214	clearer image. Event Analysis:	and lower median and mean ranks than the other	1251
1215	The video's narrative is driven	methods, which rely solely on caption repetition	1252
1216	by events like <{Events}>, which	or structured information extraction. Also, we find	1253
1217	are elaborated for better	that caption repetition outperforms structured in-	1254
1218	understanding. Phrases Insight:	formation extraction repetition. This suggests that	1255
1219	Phrases like <{Phrases}> are	incorporating relevance boosting is crucial for en-	1256
1220	explained for their significance	hancing retrieval effectiveness.	1257
1221	to the content.		
1222		B Prompt to Generate Captions in	1258
1223	4. Narrative Format Template: it weaves the ele-	Different Levels of Relevance Boosting	1259
1224	ments into a cohesive story, enhancing engage-		
1225	ment and providing a thematic understand-	B.1 Low-level Relevance	1260
1226	ing.		
1227		The following is a caption from a	1261
1228	Caption: <Caption> In this video	video: [" + text + "]. Based on this	1262
1229	, we observe <{Objects}> with <{	caption, generate two paraphrased	1263
1230	Attributes}>, a vivid	captions capturing the key	1264
1231	representation of <{Events}>.	information and main themes, each of	1265
1232	Phrases such as <{Phrases}>	which should be in one sentence	1266
1233	punctuate the narrative,	with up to twenty words (Do not	1267
		include any details not mentioned in	1268
			1269

1270 the text. Focus on the main points
1271 and key details.). Also Please
1272 generate text without any comment.
1273

1274 **B.2 Medium-level Relevance**

1275 The following is a caption from a
1276 video: [" + text + "]. Based on this
1277 caption, generate two paraphrased
1278 captions capturing the key
1279 information and main themes, each of
1280 which should be in one sentence
1281 with up to twenty words (Feel free
1282 to elaborate on points that seem
1283 important, even if not explicitly
1284 mentioned.). Also Please generate
1285 text without any comment.
1286
1287

1288 **B.3 High-level Relevance**

1289 The following is a caption from a
1290 video: [" + text + "]. Based on this
1291 caption, generate two paraphrased
1292 captions capturing the key
1293 information and main themes, each of
1294 which should be in one sentence
1295 with up to twenty words (Feel free
1296 to add any details or
1297 interpretations that you think
1298 enhance the summary, even if they
1299 are not directly mentioned in the
1300 text.). Also Please generate text
1301 without any comment.
1302
1303