

Cap4Video: What Can Auxiliary Captions Do for Text-Video Retrieval?

Wenhao Wu^{1,2*} Haipeng Luo^{3*} Bo Fang³ Jingdong Wang² Wanli Ouyang^{4,1}
¹The University of Sydney ²Baidu Inc.
³University of Chinese Academy of Sciences ⁴Shanghai AI Laboratory
whwu.ucas@gmail.com

Abstract

Most existing text-video retrieval methods focus on cross-modal matching between the visual content of videos and textual query sentences. However, in real-world scenarios, online videos are often accompanied by relevant text information such as titles, tags, and even subtitles, which can be utilized to match textual queries. This insight has motivated us to propose a novel approach to text-video retrieval, where we directly generate associated captions from videos using zero-shot video captioning with knowledge from web-scale pre-trained models (e.g., CLIP and GPT-2). Given the generated captions, a natural question arises: what benefits do they bring to text-video retrieval? To answer this, we introduce Cap4Video, a new framework that leverages captions in three ways: i) *Input data*: video-caption pairs can augment the training data. ii) *Intermediate feature interaction*: we perform cross-modal feature interaction between the video and caption to produce enhanced video representations. iii) *Output score*: the Query-Caption matching branch can complement the original Query-Video matching branch for text-video retrieval. We conduct comprehensive ablation studies to demonstrate the effectiveness of our approach. Without any post-processing, Cap4Video achieves state-of-the-art performance on four standard text-video retrieval benchmarks: MSR-VTT (51.4%), VATEX (66.6%), MSVD (51.8%), and DiDeMo (52.0%). The code is available at <https://github.com/whwu95/Cap4Video>.

1. Introduction

Text-video retrieval is a fundamental task in video-language learning. With the rapid advancements in image-language pre-training [15, 30, 46, 47], researchers have focused on expanding pre-trained image-language models, especially CLIP [30], to tackle the text-video retrieval task. The research path has evolved from the most direct global

*Equal contribution.

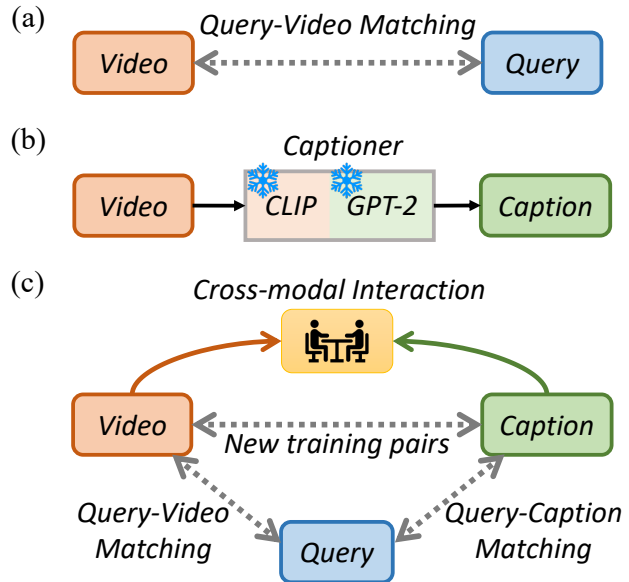


Figure 1. (a) An existing end-to-end learning paradigm for text-video retrieval. (b) Zero-shot video captioning achieved by guiding a large language model (LLM) such as GPT-2 [31] with CLIP [30]. (c) Our Cap4Video framework leverages the generated captions in three aspects: input data augmentation, intermediate feature interaction, and output score fusion.

matching (i.e., video-sentence alignment [11, 24]) to fine-grained matching (e.g., frame-word alignment [36], video-word alignment [13], multi-hierarchical alignment [9, 28], etc.). These studies have demonstrated remarkable performance and significantly outperformed previous models. Two key factors contribute to this improvement. Firstly, CLIP offers powerful visual and textual representations that are pre-aligned in the semantic embedding space, thereby reducing the challenge of cross-modal learning in video-text matching. Secondly, these methods can fine-tune the pre-trained vision and text encoders using sparsely sampled frames in an end-to-end manner. All of these methods aim to learn cross-modal alignment between the visual represen-

tation of videos and the textual representation of the corresponding query, as depicted in Figure 1(a).

However, in real-life scenarios, online videos usually come with related content such as the video’s title or tag on the video website. In addition to the visual signal in the video, the associated textual information can also be used to some extent to describe the video content and match the query (*i.e.*, the common text-to-text retrieval). This raises a pertinent question: *How can we generate associated text descriptions for videos?* One possible solution is to crawl the video title from the video website. However, this method relies on annotations, and there is a risk that the video URL may have become invalid. Another automated solution is to generate captions using zero-shot video caption models. Therefore, we turn our attention to knowledge-rich pre-trained models to handle such challenging open-set scenarios. We find that the recent study ZeroCap [34] provides a good practice to use frozen CLIP [30] and GPT-2 [31] for zero-shot image captioning. Thus, we leverage a video extension [33] of ZeroCap for generating captions in the video domain without any further training.

When provided with auxiliary captions, a natural question naturally arises: *How can we leverage these captions to enhance the text-video retrieval task?* In this paper, we propose the **Cap4Video** learning framework, as illustrated in Figure 1(c), which utilizes captions in three key ways: (i) *Input Data*: One simple approach is to augment the training data with the generated captions. Specifically, the given video and its generated caption can be treated as a matched pair, which serves as an additional positive sample pair for training beyond the query-video pairs. (ii) *Intermediate Feature Interaction*: Cross-modal interaction between the video and captions can be leveraged to improve the video representation. Specifically, we can exploit the complementary information between videos and captions to reduce redundant features from videos and learn more discriminative video representations. (iii) *Output score*: The generated caption can also represent the video’s content, allowing us to employ query-caption matching to complement standard query-video matching for the text-video retrieval task. Moreover, a two-stream architecture can be utilized to reduce model bias and produce more robust results.

We hope that our novel paradigm will encourage further investigation into the video-language learning. In summary, our contributions are as follows:

- We explore a novel problem: leveraging auxiliary captions to further enhance existing text-video retrieval. Besides labor-intensive manual crawling of video website titles, we investigate the potential of rich captions automatically generated by large language models (LLMs) to benefit text-video retrieval.
- We propose the **Cap4Video** learning framework,

which maximizes the utility of the auxiliary captions through three aspects: input data, feature interaction, and output score. Our framework improves the performance of existing query-video matching mechanisms, including global matching and fine-grained matching.

- Extensive experiments conducted on four video benchmarks demonstrate the effectiveness of our method. Our Cap4Video achieves state-of-the-art performance on MSR-VTT [44] (51.4%), VATEX [38] (66.6%), MSVD [43] (51.8%), and DiDeMo [1] (52.0%).

2. Methodology

2.1. Background: Text-Video Matching

Text-video matching aims to evaluate the similarity between a given sentence Q_i and a given video V_j , typically using a similarity function $s(Q_i, V_j)$. In text-to-video retrieval, the goal is to rank all videos based on their similarity scores to a given query sentence. To improve text-video retrieval, recent works [9, 11, 24] have applied CLIP [30] for initialization, leveraging pre-trained knowledge from image-text learning. Our baselines for text-video matching include two typical mechanisms: global matching and fine-grained matching, as illustrated in Figure 2(b).

Global Matching is a commonly used technique in cross-modal contrastive learning [16, 24, 30]. In global matching, each modality is encoded independently to obtain global features, which are then used to calculate similarity. We train the visual encoder to output F frame embeddings for a given video that samples F frames. Similarly, the query encoder returns W word embeddings and the [CLS] embedding as the global representation for a given query sentence that contains W words. The frame embeddings are integrated using average pooling to obtain the global video embedding, which is then compared with the global query embedding to calculate similarity.

Fine-grained Matching focuses on modeling the token-level alignment between two modalities, such as frame-word alignment. In order to achieve token-level patch-word alignment for image-text learning, FILIP [45] and ColBERT [17] employ a *Max-Mean* pipeline. This pipeline finds the token-wise maximum similarity between patch and word tokens, and then averages the maximum similarity of tokens in the image or text to obtain the similarity between an image and a text or vice versa. Moreover, DRL [36] extends the token-wise alignment to text-video retrieval and introduces an attention mechanism to learn weighted pooling instead of mean pooling. We have adopted this mechanism as our enhanced baseline.

2.2. Preprocessing: Caption Generation

To obtain auxiliary captions for a given video, we consider the following two approaches.

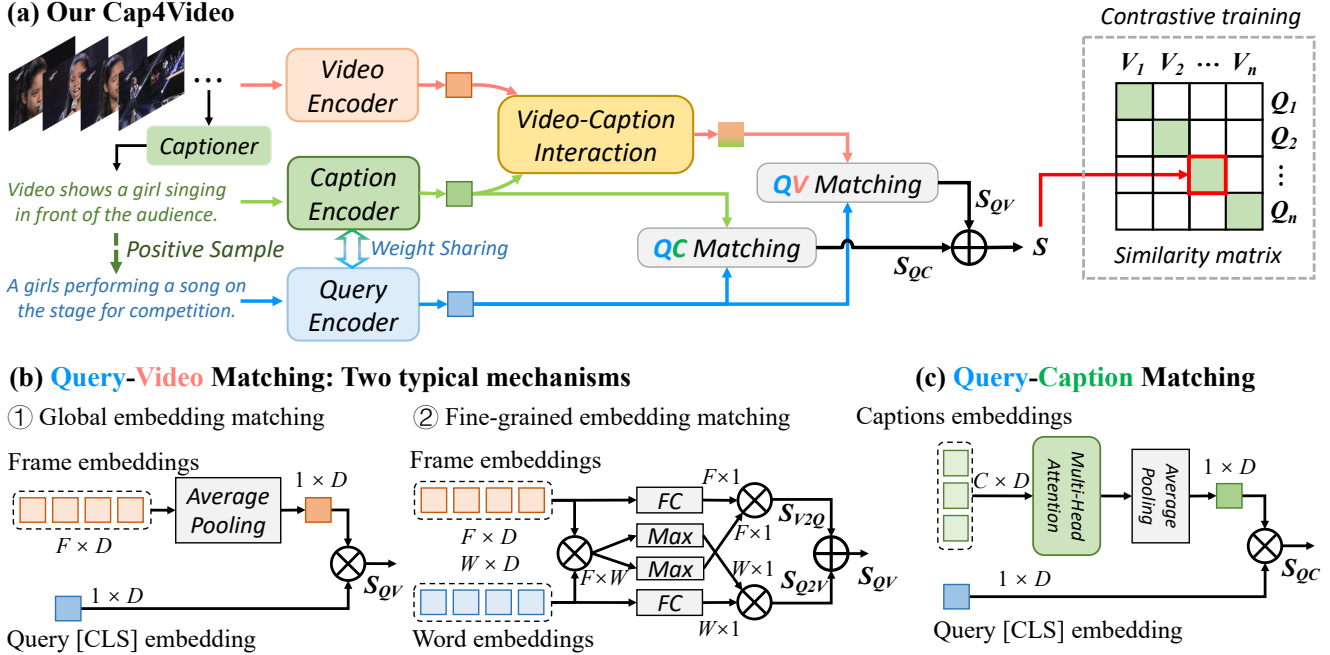


Figure 2. An overview of our **Cap4Video** for text-video retrieval. We first generate captions using a zero-shot video captioner that combines CLIP [30] with GPT-2 [31], leveraging knowledge from both frozen web-scale models. We then utilize the pre-extracted caption information from three different perspectives: i) *Input data*: We use the video and captions to create new positive pairs for data augmentation during training. ii) *Feature interaction*: We perform feature interaction between video and caption to capture intra- and inter-modality context, yielding enhanced video representations. iii) *Output score*: The Query-Caption matching branch can complement the original Query-Video matching branch for text-video retrieval.

Manual Crawling of Video Titles. We extract the video website title by crawling the original links (such as YouTube ID) of each video and utilize it as the caption. However, we skip this step for videos with expired links.

Automatic Video Captioning. In contrast to the manual approach that relies on annotations, we leverage knowledge from the LLM to generate rich and diverse captions. Given the scalability of our framework, we aim to generate captions directly from downstream videos without any additional training, a process known as zero-shot video captioning. To achieve this, we follow [33, 34] and use GPT-2 [31] to predict the next word from an initial prompt, e.g., “Video shows”. A calibrated CLIP [30] loss is then used to drive the model to generate sentences that describe the video, incorporating video-related knowledge into the auto-regressive process. *See Supplementary for more details.*

2.3. Data Augmentation with Auxiliary Captions

Auxiliary captions can be used to augment training data. For example, for a dataset consisting of N videos and their corresponding query sentences, each video and its generated caption can be considered as a positive sample pair for training, in addition to the original query-video pairs. By selecting one caption per video, we can add at least N pairs

as additional data augmentation during training.

The automatic video captioner can generate multiple captions (e.g., 20) for each video. However, some of these captions may contain noise and may not be entirely relevant to the video content. To avoid negative effects on training, we use a filtering mechanism that evaluates the semantic similarity between each caption and the ground-truth query of the video using a pre-trained text encoder. The caption with the highest similarity is then chosen for data augmentation. Note that we only use the ground-truth query for caption filtering during the training phase.

2.4. Video-Caption Cross-Modal Interaction

We further consider taking advantage of the complementarity between videos and captions to reduce redundant features and learn more discriminative video representations. To preserve the pre-trained CLIP encoder architecture for efficient transfer learning, we limit the interaction to the final caption and frame embeddings. Specifically, we pass the frame embeddings $e_v = \{v_1, v_2, \dots, v_F\}$ and caption embeddings $e_c = \{c_1, c_2, \dots, c_C\}$ to the interaction module, where F and C represent the number of frames and captions, respectively. Figure 2(b) depicts several ways of interaction between the two modalities.

Sum. To obtain an enhanced frame embedding, an intuitive approach is to compute the sum of the global caption embedding \mathbf{c}_g and each frame embedding:

$$\text{Sum}(\mathbf{v}_i, \mathbf{c}_g) = \mathbf{v}_i + \mathbf{c}_g, \quad i = 1, \dots, F, \quad (1)$$

where $\mathbf{v}_i \in \mathbb{R}^D$ is the i -th frame embedding, and $\mathbf{c}_g \in \mathbb{R}^D$ is computed by averaging the [CLS] embeddings of C generated captions: $\mathbf{c}_g = \frac{1}{C} \sum_{i=1}^C \mathbf{c}_i$.

MLP. To model weighted combinations of each frame embedding and the global caption embedding \mathbf{c}_g , we concatenate them together and pass the result through a learnable Multi-layer Perceptron (MLP):

$$\text{MLP}(\mathbf{v}_i, \mathbf{c}_g) = f_\theta([\mathbf{v}_i, \mathbf{c}_g]), \quad \text{for } i = 1, \dots, F, \quad (2)$$

where $[\cdot, \cdot]$ denotes the concatenation operation, f_θ is the MLP with parameter θ .

Cross Transformer. We also investigate the use of self-attention [35] for interactions. The Cross Transformer operates on a sequence $\{\mathbf{e}_v, \mathbf{e}_c\} = \{\mathbf{v}_1, \dots, \mathbf{v}_F, \mathbf{c}_1, \dots, \mathbf{c}_C\}$ and processes them through L encoder-style transformer blocks to generate final representations:

$$\text{Cross}(\mathbf{e}_v, \mathbf{e}_c) = f_\psi(\{\mathbf{e}_v, \mathbf{e}_c\}), \quad (3)$$

where $\{\cdot\}$ denotes that \mathbf{e}_v and \mathbf{e}_c form a sequence, and f_ψ represents the transformer encoders with parameter ψ .

Co-attention Transformer. Co-attention [23] is another common method for exchanging information between modalities, allowing for mutual attention between video and caption. After this co-attentional transformer layer, we include L transformer layers to model temporal information:

$$\text{CoAttn}(\mathbf{e}_v, \mathbf{e}_c) = f_{\phi_2}(f_{\phi_1}(\{\mathbf{e}_v, \mathbf{e}_c\})), \quad (4)$$

where f_{ϕ_1} is the co-attentional transformer with parameter ϕ_1 and f_{ϕ_2} is the transformer encoders with parameter ϕ_2 .

The video-caption interaction module generates frame embeddings that can be further processed based on the type of matching needed. For global matching, the frame embeddings can be averaged to obtain a single video representation. Alternatively, for fine-grained matching, the individual frame embeddings can be retained.

2.5. Complementary Query-Caption Matching

Besides using the caption for data augmentation and video feature enhancement, it can also directly represent the video content, allowing for text-text retrieval. Specifically, each of the C captions generated by the video is then passed through the caption encoder to obtain its [CLS] text embedding. These caption embeddings are then aggregated to form a global representation, as illustrated in Figure 2(c). The cosine similarity between this global caption embedding and the global query embedding is then calculated to complement the query-video matching.

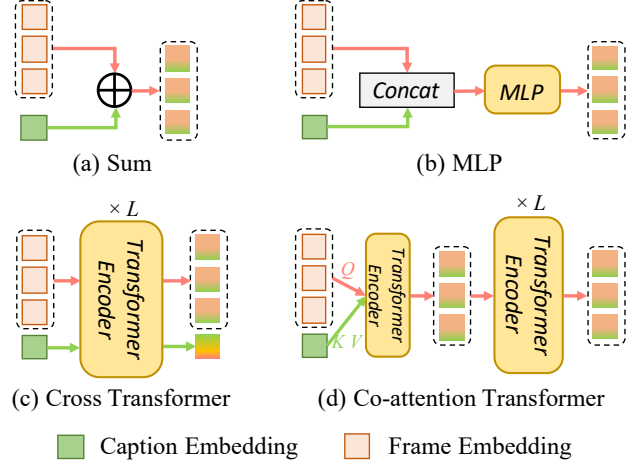


Figure 3. Illustration of four Video-Caption interaction strategies. The enhanced frame embeddings will be followed by a mean pooling for global matching or will remain for fine-grained matching.

Notation. Let $\{\mathbf{e}_{v_i}, \mathbf{e}_{t_i}, \mathbf{e}_{c_i}\}_{i=1}^B$ be a batch of B triples, where \mathbf{e}_{v_i} , \mathbf{e}_{t_i} , and \mathbf{e}_{c_i} denote the i -th video, query, and caption embedding, respectively. Note that the term “embedding” used here is more general for convenience and can vary in meaning depending on the situation. For instance, in query-video global matching, \mathbf{e}_{v_i} and \mathbf{e}_{t_i} represent the averaged video feature and global [CLS] text feature, respectively. In query-video fine-grained matching, \mathbf{e}_{v_i} and \mathbf{e}_{t_i} represent a sequence of frame embeddings and a sequence of word embeddings, respectively. In query-caption matching, \mathbf{e}_{c_i} represents a sequence of caption embeddings, and \mathbf{e}_{t_i} represents a global [CLS] text feature.

Learning Objectives. For the *Query-Caption* branch, we want the caption embedding \mathbf{e}_c and the query embedding \mathbf{e}_t to be close while they are related and far apart when they are not during training phase. We follow the common practice [24, 36] to consider the bidirectional learning objective. We employ symmetric cross-entropy loss to maximize the similarity between matched *Query-Caption* pairs and minimize the similarity for other pairs:

$$\begin{aligned} \mathcal{L}_{Q2C} &= -\frac{1}{B} \sum_i^B \log \frac{\exp(s_{qc}(\mathbf{e}_{t_i}, \mathbf{e}_{c_i})/\tau)}{\sum_j^B \exp(s_{qc}(\mathbf{e}_{t_i}, \mathbf{e}_{c_j})/\tau)}, \\ \mathcal{L}_{C2Q} &= -\frac{1}{B} \sum_i^B \log \frac{\exp(s_{qc}(\mathbf{e}_{t_i}, \mathbf{e}_{c_i})/\tau)}{\sum_j^B \exp(s_{qc}(\mathbf{e}_{t_j}, \mathbf{e}_{c_i})/\tau)}, \\ \mathcal{L}_{QC} &= \frac{1}{2}(\mathcal{L}_{Q2C} + \mathcal{L}_{C2Q}), \end{aligned} \quad (5)$$

where $s_{qc}(\cdot, \cdot)$ represents the query-caption matching similarity function shown in Figure 2(c), and τ refers to the temperature hyper-parameter for scaling. Similarly, the con-

trastive loss for *Query-Video* branch is formulated as:

$$\begin{aligned}\mathcal{L}_{Q2V} &= -\frac{1}{B} \sum_i^B \log \frac{\exp(s_{qv}(\mathbf{e}_{t_i}, \mathbf{e}_{v_i})/\tau)}{\sum_j^B \exp(s_{qv}(\mathbf{e}_{t_i}, \mathbf{e}_{v_j})/\tau)}, \\ \mathcal{L}_{V2Q} &= -\frac{1}{B} \sum_i^B \log \frac{\exp(s_{qv}(\mathbf{e}_{t_i}, \mathbf{e}_{v_i})/\tau)}{\sum_j^B \exp(s_{qv}(\mathbf{e}_{t_j}, \mathbf{e}_{v_i})/\tau)}, \\ \mathcal{L}_{QV} &= \frac{1}{2}(\mathcal{L}_{Q2V} + \mathcal{L}_{V2Q}),\end{aligned}\quad (6)$$

where $s_{qv}(\cdot, \cdot)$ represents the query-video matching (*e.g.*, global matching, fine-grained matching) similarity function shown in Figure 2(b). The total loss \mathcal{L} is the sum of *Query-Video* loss \mathcal{L}_{QV} and *Query-Caption* loss \mathcal{L}_{QC} :

$$\mathcal{L} = \mathcal{L}_{QV} + \mathcal{L}_{QC}. \quad (7)$$

3. Experiments: Text-Video Retrieval

3.1. Setups

Datasets. We conduct experiment on four popular benchmarks for video-to-text retrieval and text-to-video retrieval tasks. MSR-VTT [44] contains a total of 10K video clips, each having 20 captions. Following the data splits from [10, 24, 27], we train models with associated captions on the Training-9K set and report results on the test 1K-A set. DiDeMo [1] has 10K videos paired with 40K descriptions. Following previous works [2, 19, 24], we concatenate all descriptions of one video to a single query, acting as a *video-paragraph* retrieval task. VATEX [38] collects ~ 35 K videos, each with multiple annotations. There are ~ 26 K videos for training, 1,500 videos for validation and 1,500 videos for testing. MSVD [43] contains 1,970 videos with 80K captions, with ~ 40 captions on average per video. There are 1,200, 100, and 670 videos in the train, validation, and test sets, respectively.

Evaluation Metrics. For brevity, we abbreviate Recall at K to $R@K$ ($K = 1, 5, 10$) upon all datasets, which computes the percentage of correct videos among the top K retrieved videos given textual queries (Text \rightarrow Video, and vice versa). MdR, Median Rank, computes the median of the ground-truth in the retrieval ranking list. MnR, Mean Rank, computes the mean rank of the correct results in the retrieval ranking list. Note that for MdR and MnR, the lower score means the better (indicated as \downarrow).

Implementation Details. All experiments use the visual encoder in CLIP [30] as the video encoder, and the textual encoder in CLIP as both the caption encoder and query encoder. The caption encoder and query encoder share parameters. To reduce conflict between the two branches, the query-video branch is trained first, followed by the query-caption branch. The text length is fixed to 32, and the video length is fixed to 12 for all datasets except DiDeMo (64 max words and 64 frames). The initial learning rate is set to $1e-7$

for the clip parameters and $1e-4$ for the non-clip parameters. The model is trained with a batch size of 128 for 5 epochs, except for DiDeMo (15 epochs), using the Adam [18] optimizer. All learning rates follow the cosine schedule with a linear warmup [14] strategy. For the number of generated captions per video, we set C to 30. The interaction module employs L transformer layers, where L is set to 4 for VATEX and MSR-VTT, and 1 for DiDeMo and MSVD. In the caption branch, the number of transformer layers is set to 2. For caption generation, we directly use the original pre-trained CLIP and GPT-2, without any additional tuning.

3.2. Comparison with State-of-the-Arts

In this section, we compare our Cap4Video with recent state-of-the-art methods on four benchmarks: MSR-VTT [44], MSVD [43], VATEX [38], and DiDeMo [1].

Table 1 shows the comparisons on DiDeMo, where our Cap4Video outperforms CLIP4Clip [24] by a significant margin of **9.2%** in $R@1$ and exceeds DRL [36] by **3.0%**, demonstrating the effectiveness of our method.

Table 2 provides a comparison of our approach with recent state-of-the-art models on MSR-VTT. Our method achieves new state-of-the-art performance on text-to-video retrieval for both ViT-B/32 and ViT-B/16 backbones, significantly surpassing previous works. For instance, we achieve a **+4.8%** higher $R@1$ than CLIP4Clip with the same ViT-B/32 on text-to-video retrieval. Additionally, our Cap4Video outperforms the recent TS2-Net [22] by **2.3%** and **2.0%** with ViT-B/32 and ViT-B/16, respectively.

Table 3 and Table 4 show the results for the MSVD and VATEX datasets, respectively, where we use ViT-B/16 as our backbone. For MSVD, our Cap4Video achieves a remarkable performance of 51.8% $R@1$ and outperforms CLIP-based models CLIP4Clip [24] and X-Pool [25] by **6.6%** and **4.6%** on text-to-video retrieval, respectively. For VATEX, our approach also outperforms the recent state-of-the-art methods and achieves a **+7.5%** $R@1$ improvement over TS2-Net [22] for text-to-video retrieval.

In recent studies, several methods have been proposed

Method	R@1	R@5	R@10	MdR	MnR
CE [21]	15.6	40.9	-	8.2	-
ClipBERT [19]	21.1	47.3	61.1	6.3	-
Frozen [2]	31.0	59.8	72.4	3.0	-
TMVM [20]	36.5	64.9	75.4	3.0	-
CLIP4Clip [24]	42.8	68.5	79.2	2.0	18.9
TS2-Net [22]	41.8	71.6	82.0	2.0	14.8
HunYuan [28]	45.0	75.6	83.4	2.0	12.0
DRL [36]	49.0	76.5	84.5	2.0	-
Cap4Video	52.0	79.4	87.5	1	10.5

Table 1. Results of text-to-video retrieval on the DiDeMo [1].

Method	Venue	Text → Video					Video → Text				
		R@1	R@5	R@10	MdR↓	MnR↓	R@1	R@5	R@10	MdR↓	MnR↓
ClipBERT [19]	CVPR'20	22.0	46.8	59.9	6.0	-	-	-	-	-	-
MMT [10]	ECCV'20	26.6	57.1	69.6	4.0	-	27.0	57.5	69.7	3.7	21.3
T2VLAD [39]	CVPR'21	29.5	59.0	70.1	4.0	-	31.8	60.0	71.1	3.0	-
SupportSet [29]	ICLR'21	30.1	58.5	69.3	3.0	-	28.5	58.6	71.6	3.0	-
Frozen [2]	ICCV'21	32.5	61.5	71.2	3.0	-	-	-	-	-	-
BridgeFormer [12]	CVPR'22	37.6	64.8	75.1	-	-	-	-	-	-	-
TMVM [20]	NeurIPS'22	36.2	64.2	75.7	3.0	-	34.8	63.8	73.7	3.0	-
<i>CLIP-ViT-B/32</i>											
CLIP4Clip [24]	arXiv'21	44.5	71.4	81.6	2.0	15.3	42.7	70.9	80.6	2.0	11.6
CenterCLIP [48]	SIGIR'22	44.2	71.6	82.1	2.0	15.1	42.8	71.7	82.2	2.0	10.9
CAMoE [7]	arXiv'21	44.6	72.6	81.8	2.0	13.3	45.1	72.4	83.1	2.0	10.0
CLIP2Video [9]	arXiv'21	45.6	72.6	81.7	2.0	14.6	43.5	72.3	82.1	2.0	10.2
X-Pool [13]	CVPR'22	46.9	72.8	82.2	2.0	14.3	-	-	-	-	-
QB-Norm [4]	CVPR'22	47.2	73.0	83.0	2.0	-	-	-	-	-	-
TS2-Net [22]	ECCV'22	47.0	74.5	83.8	2.0	13.0	45.3	74.1	83.7	2.0	9.2
DRL [36]	arXiv'22	47.4	74.6	83.8	2.0	-	45.3	73.9	83.3	2.0	-
Cap4Video		49.3	74.3	83.8	2.0	12.0	47.1	73.7	84.3	2.0	8.7
<i>CLIP-ViT-B/16</i>											
CLIP2TV [11]	arXiv'21	48.3	74.6	82.8	2.0	14.9	46.5	75.4	84.9	2.0	10.2
CenterCLIP [48]	SIGIR'22	48.4	73.8	82.0	2.0	13.8	47.7	75.0	83.3	2.0	10.2
TS2-Net [22]	ECCV'22	49.4	75.6	85.3	2.0	13.5	46.6	75.9	84.9	2.0	8.9
DRL [36]	arXiv'22	50.2	76.5	84.7	1.0	-	48.9	76.3	85.4	2.0	-
Cap4Video		51.4	75.7	83.9	1.0	12.4	49.0	75.2	85.0	2.0	8

Table 2. Retrieval results on the validation set of MSR-VTT 1K [44]. Here we report results **without** any post-processing operations (e.g., DSL [7] or QB-Norm [4]) during inference.

Method	R@1	R@5	R@10	MdR	MnR
CE [21]	19.8	49.0	63.8	6.0	-
SUPPORT [29]	28.4	60.0	72.9	4.0	-
CLIP [30]	37.0	64.1	73.8	3.0	-
Frozen [2]	33.7	64.7	76.3	3.0	-
TMVM [20]	36.7	67.4	81.3	2.5	-
CLIP4Clip [24]	45.2	75.5	84.3	2.0	10.3
X-Pool [13]	47.2	77.4	86.0	2.0	9.3
Cap4Video	51.8	80.8	88.3	1	8.3

Table 3. Results of text-to-video retrieval on the MSVD [43].

Method	R@1	R@5	R@10	MdR	MnR
HGR [6]	35.1	73.5	83.5	2.0	-
CLIP [30]	39.7	72.3	82.2	2.0	12.8
SUPPORT [29]	44.9	82.1	89.7	1.0	-
CLIP4Clip [24]	55.9	89.2	95.0	1.0	3.9
Clip2Video [9]	57.3	90.0	95.5	1.0	3.6
QB-Norm [4]	58.8	88.3	93.8	1.0	-
TS2-Net [22]	59.1	90.0	95.2	1.0	3.5
Cap4Video	66.6	93.1	97.0	1	2.7

Table 4. Results of text-to-video retrieval on the VATEX [38].

to improve text-video retrieval performance by adjusting similarity during inference using other query information. Notably, our results adheres to the standard retrieval logic, where the most relevant video is retrieved for each query from a set of videos, without any knowledge of the relationship between other queries and videos. Therefore, all results reported in our tables do not involve any post-processing procedures such as DSL [7] and QB-Norm [4]. Overall, the consistent state-of-the-art performance across four benchmarks demonstrates the effectiveness of our Cap4Video.

3.3. Ablation Study

In this section, we provide detailed ablation studies to clarify the effects of each part of our design.

Auxiliary Caption as Data Augmentation. We begin by investigating the impact of captions on data augmentation for training. In a real-world scenario, the original video title would naturally serve as an additional auxiliary caption. Therefore, we manually extracted the title from the video’s original webpage and compared it to the caption generated by the GPT-2 model. Table 5 presents the re-

Method	Global Matching					Fine-grained Matching				
	R@1	R@5	R@10	MdR↓	MnR↓	R@1	R@5	R@10	MdR↓	MnR↓
Baseline	42.8	70.4	79.0	2	16.6	45.7	73.7	82.6	2	13.1
<i>+Different Sources of Caption as Data Augmentation</i>										
Video Title from Source URL	43.8	71.1	80.9	2	15.1	44.3	72.7	83.5	2	13.1
Zero-shot Video Captioning	44.2	70.7	81.5	2	16.2	46.3	72.5	81.7	2	12.9
<i>+Different Number of Captions for Data Augmentation</i>										
Top-1	44.2	70.7	81.5	2	16.2	46.3	72.5	81.7	2	12.9
Top-3	43.3	71.7	81.6	2	15.0	45.5	73.8	82.4	2	12.7
Top-5	43.4	70.6	80.4	2	16.2	45.6	72.7	82.7	2	12.9
<i>+Video-Caption Feature Interaction</i>										
Video Only	44.2	70.7	81.5	2	16.2	46.3	72.5	81.7	2	12.9
Sum	43.8	71.5	80.3	2	16.1	47.2	73.3	82.8	2	13.1
Concat-MLP	37.5	66.1	78.4	3	15.7	40.0	68.7	79.9	2	12.7
Cross Transformer	44.6	71.6	80.3	2	14.6	47.9	75.4	83.0	2	11.5
Co-attention Transformer	45.3	71.2	80.9	2	15.0	48.5	74.0	82.5	2	12.7
<i>+Query-Caption Matching Score</i>										
Query-Video Only	45.3	71.2	80.9	2	15.0	48.5	74.0	82.5	2	12.7
Query-Caption Only	30.3	55.2	67.5	4	26.4	30.3	55.2	67.5	4	26.4
Query-Video + Query-Caption	45.6	71.7	81.2	2	14.8	49.3	74.2	83.4	2	12.1

Table 5. Component-wise evaluation of our framework on the MSR-VTT 1K validation set. With the ViT-B/32 backbone, we report the text-to-video retrieval results for two representative Query-Video matching mechanisms: global matching and fine-grained matching. The consistent improvement on two typical matching mechanisms demonstrates the generalization ability and effectiveness of our method.

sults of using different sources of the caption, from which we observe that using captions generated by the web-scale model as data augmentation for training can lead to direct improvements in R@1 (+1.4%, +0.6%) under both matching mechanisms. Using video titles can also bring a 1% improvement for global matching.

We also explore the impact of the number of generated captions used for augmentation. We used the caption filtering mechanism mentioned in Sec. 2.3 to rank the relevance of captions and the ground-truth query, and selected different numbers of captions for training. The results demonstrated that using only one caption is sufficient.

Benefit on Both Online and Offline Videos Scenarios. Our method is applicable for both online and offline videos. Offline videos are local videos with no title, while online videos have a title on the video website. Therefore, for online videos, the step of generating a caption with CLIP+GPT-2 is skipped, and the website title is used directly as a caption. In Table 6, the results show that our method has significantly improved over the global matching baseline for **both** offline and online videos.




Video-Caption Feature Interaction. As mentioned in Sec. 2.4, we have designed four approaches for Video-Caption feature interaction. Based on the results presented in Table 5, we can conclude the following: 1) The basic approach of *Sum* has been shown to enhance fine-grained matching by 0.9% R@1, but there is no noticeable improvement in global matching. 2) The *MLP* approach is difficult

	Caption Source	R@1
Baseline	N/A	42.8
Cap4Video	Original Website Title	45.8
	Captioner (CLIP+GPT-2)	45.6

Table 6. Exploring the effectiveness of captions from different sources on MSR-VTT 1k-A. Setting: ViT-B/32, global matching.

to optimize and performs poorly in both matching scenarios. We speculate that the *MLP*'s operation in a black-box environment, despite creating a nonlinear metric space, may lead to degradation. 3) The *Cross Transformer* approach has demonstrated improvements of +0.4% and +1.6% in two matching settings, respectively. These enhancements may be attributed to the self-attention mechanism's ability to capture the inter-modal relationship between the video and caption. 4) Moreover, the *Co-attention Transformer* approach has significantly boosted performance, with gains of +1.1% and +2.2% for these two matching mechanisms. In summary, the results indicate that proper interaction between the video and generated caption can lead to better video representation and improved Query-Video matching.

Query-Caption matching. We also investigate the Query-Caption matching branch for text-video retrieval. We aggregate caption embeddings using mean pooling to yield a global embedding. As shown in Table 5, the single Query-Caption matching branch achieves 30.3% R@1 on text-to-

Query7765 : a person is discussing a car.			
Video	Rank	+ Caption	Rank
	6	video of a car camera recording the driver's voice.	1
	4	video showing the car in a parking spot.	2
	5	video of SUV in the video below shows a salesman talking to an audience.	3


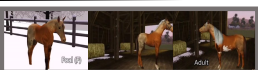

Query9616 : person is recording the brown horse which is having fun.			
Video	Rank	+ Caption	Rank
	2	video of the horse jumping over a fence at Ranch in Nevada was captured on camera.	1
	1	video showing animation of a horse's simulation, which simulates the game.	2
	4	video showing a horse simulation video game in which you could see your avatar being animated by the camera.	3

Figure 4. The text-video results on the MSR-VTT 1K-A test set. **Left:** The ranking results of the query-video matching model. **Right:** The ranking results of Cap4Video, which incorporates generated captions to enhance retrieval. Please zoom in for best view.

video retrieval, outperforming previous query-video matching methods such as ClipBERT [19] (22.0%) and MMT [10] (26.6%). This suggests that the Query-Caption matching branch can complement the regular Query-Video matching branch for improved text-video retrieval. Combining the score of Query-Caption matching branch with Query-Video matching branch further improves performance (+0.8%).

Overall, Cap4Video utilizes generated captions in three ways: input data augmentation, intermediate feature interaction, and output score fusion, leading to consistent improvements (+2.8% / +3.6%) in both matching mechanisms.

3.4. Visualization

We provide two examples of videos retrieved by our Cap4Video method and a model without auxiliary captions. As illustrated in Figure 4, our approach successfully retrieves the ground-truth video with the assistance of the caption, while the video-only model returns multiple videos that are somewhat relevant to the query but not precise. See more qualitative results in Supplementary.

4. Related Works

Zero-shot Image Captioning. In the field of natural language processing, transformer-based GPT models [5, 31] have been successful in generating text from prompts by training on large-scale text corpora. Similarly, CLIP [30], a vision-language alignment model trained on 400 million

image-text pairs, has demonstrated impressive zero-shot performance on vision tasks. However, research on transferring web-scale models to zero-shot video captioning remains limited. Recently, ZeroCap [34] proposed a method of using CLIP and the GPT-2 language model to generate textual descriptions of input images, leveraging knowledge from both models in a truly zero-shot manner without re-training or fine-tuning model parameters. MAGIC [32] has also used CLIP scores to align GPT-2 logits with corresponding images but requires fine-tuning on the MS-COCO caption text corpus. More recently, a study [33] extended the zero-shot capability of ZeroCap to the video domain. In this paper, we employ this video extension to generate auxiliary captions without any additional training.

Text-Video Retrieval aims to retrieve relevant video content based on natural language descriptions. Early studies [6, 10, 21, 37, 39] focused on knowledge transfer from “expert” models and captured intra-modal and cross-modal interactions based on pre-extracted features. However, the performance of these methods is limited since they cannot perform end-to-end optimization. Recently, more methods have involved end-to-end training for text-video retrieval. One typical approach [2, 26, 27] is to first perform large-scale text-video pre-training, then transfer the model to downstream text-video retrieval tasks. With the emergence of pre-trained Vision-Language Models (VLMs), there have been increased efforts to leverage them for improving video understanding [41, 42]. Thus, another training-efficient line is to directly expand the pre-trained VLM to the text-video retrieval task. CLIPBERT [19] enables affordable pioneering end-to-end training with a sparse sampling strategy. After that, recent works [3, 4, 8, 9, 11, 13, 22, 24, 40, 48] focus on transferring knowledge from CLIP models that have been pre-trained on 400M image-text pairs. The research path has evolved from the most direct global matching (*i.e.*, video-sentence alignment [11, 24]) to fine-grained matching (*e.g.*, frame-word alignment [36], video-word alignment [13], multi-hierarchical alignment [9, 28]). Unlike above CLIP-Based efforts on query-video matching, we propose to generate auxiliary captions from videos to improve text-video retrieval. Thus our method is compatible with both global and fine-grained matching.

5. Conclusion

We introduce Cap4Video, a novel framework that leverages captions generated by web-scale language models to enhance text-video matching in three key ways: 1) Input data augmentation for training, 2) Intermediate video-caption feature interaction for compact video representations, and 3) Output score fusion for improved text-video retrieval. Our approach demonstrates consistent performance gains on four standard text-video retrieval benchmarks, outperforming state-of-the-art methods by a clear margin.

References

- [1] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. Localizing moments in video with natural language. In *ICCV*, pages 5803–5812, 2017. [2](#), [5](#)
- [2] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *ICCV*, pages 1728–1738, 2021. [5](#), [6](#), [8](#)
- [3] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. A clip-hitchhiker’s guide to long video retrieval. *arXiv preprint arXiv:2205.08508*, 2022. [8](#)
- [4] Simion-Vlad Bogolin, Ioana Croitoru, Hailin Jin, Yang Liu, and Samuel Albanie. Cross modal retrieval with querybank normalisation. In *CVPR*, pages 5194–5205, 2022. [6](#), [8](#)
- [5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *NeurIPS*, 33:1877–1901, 2020. [8](#)
- [6] Shizhe Chen, Yida Zhao, Qin Jin, and Qi Wu. Fine-grained video-text retrieval with hierarchical graph reasoning. In *CVPR*, pages 10638–10647, 2020. [6](#), [8](#)
- [7] Xing Cheng, Hezheng Lin, Xiangyu Wu, Fan Yang, and Dong Shen. Improving video-text retrieval by multi-stream corpus alignment and dual softmax loss. *arXiv preprint arXiv:2109.04290*, 2021. [6](#)
- [8] 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. *arXiv preprint arXiv:2301.06309*, 2023. [8](#)
- [9] Han Fang, Pengfei Xiong, Luhui Xu, and Yu Chen. Clip2video: Mastering video-text retrieval via image clip. *arXiv preprint arXiv:2106.11097*, 2021. [1](#), [2](#), [6](#), [8](#)
- [10] Valentin Gabeur, Chen Sun, Karteek Alahari, and Cordelia Schmid. Multi-modal transformer for video retrieval. In *ECCV*, pages 214–229. Springer, 2020. [5](#), [6](#), [8](#)
- [11] Zijian Gao, Jingyu Liu, Sheng Chen, Dedan Chang, Hao Zhang, and Jinwei Yuan. Clip2tv: An empirical study on transformer-based methods for video-text retrieval. *arXiv preprint arXiv:2111.05610*, 2021. [1](#), [2](#), [6](#), [8](#)
- [12] Yuying Ge, Yixiao Ge, Xihui Liu, Dian Li, Ying Shan, Xiaohu Qie, and Ping Luo. Bridging video-text retrieval with multiple choice questions. In *CVPR*, pages 16167–16176, 2022. [6](#)
- [13] 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*, pages 5006–5015, 2022. [1](#), [6](#), [8](#)
- [14] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large mini-batch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017. [5](#)
- [15] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, pages 4904–4916. PMLR, 2021. [1](#)
- [16] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, pages 4904–4916. PMLR, 2021. [2](#)
- [17] Omar Khattab and Matei Zaharia. Colbert: Efficient and effective passage search via contextualized late interaction over bert. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*, pages 39–48, 2020. [2](#)
- [18] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. [5](#)
- [19] Jie Lei, Linjie Li, Luwei Zhou, Zhe Gan, Tamara L Berg, Mohit Bansal, and Jingjing Liu. Less is more: Clipbert for video-and-language learning via sparse sampling. In *CVPR*, pages 7331–7341, 2021. [5](#), [6](#), [8](#)
- [20] Chengzhi Lin, Ancong Wu, Junwei Liang, Jun Zhang, Wenhong Ge, Wei-Shi Zheng, and Chunhua Shen. Text-adaptive multiple visual prototype matching for video-text retrieval. In *NeurIPS*, 2022. [5](#), [6](#)
- [21] Yang Liu, Samuel Albanie, Arsha Nagrani, and Andrew Zisserman. Use what you have: Video retrieval using representations from collaborative experts. *arXiv preprint arXiv:1907.13487*, 2019. [5](#), [6](#), [8](#)
- [22] Yuqi Liu, Pengfei Xiong, Luhui Xu, Shengming Cao, and Qin Jin. Ts2-net: Token shift and selection transformer for text-video retrieval. In *ECCV*, pages 319–335. Springer, 2022. [5](#), [6](#), [8](#)
- [23] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *NeurIPS*, 32, 2019. [4](#)
- [24] 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. [1](#), [2](#), [4](#), [5](#), [6](#), [8](#)
- [25] Yiwei Ma, Guohai Xu, Xiaoshuai Sun, Ming Yan, Ji Zhang, and Rongrong Ji. X-clip: End-to-end multi-grained contrastive learning for video-text retrieval. In *ACM MM*, pages 638–647, 2022. [5](#)
- [26] Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman. End-to-end learning of visual representations from uncurated instructional videos. In *CVPR*, pages 9879–9889, 2020. [8](#)
- [27] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. Howto100m: Learning a text-video embedding by watching hundred million narrated video clips. In *ICCV*, pages 2630–2640, 2019. [5](#), [8](#)
- [28] Shaobo Min, Weijie Kong, Rong-Cheng Tu, Dihong Gong, Chengfei Cai, Wenzhe Zhao, Chenyang Liu, Sixiao Zheng, Hongfa Wang, Zhifeng Li, et al. Hunyuan_tvr for text-video retrieval. *arXiv preprint arXiv:2204.03382*, 2022. [1](#), [5](#), [8](#)
- [29] Mandela Patrick, Po-Yao Huang, Yuki Asano, Florian Metzger, Alexander G Hauptmann, Joao F Henriques, and An-

- drea Vedaldi. Support-set bottlenecks for video-text representation learning. In *ICLR*, 2021. 6
- [30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763. PMLR, 2021. 1, 2, 3, 5, 6, 8
- [31] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019. 1, 2, 3, 8
- [32] Yixuan Su, Tian Lan, Yahui Liu, Fangyu Liu, Dani Yogatama, Yan Wang, Lingpeng Kong, and Nigel Collier. Language models can see: Plugging visual controls in text generation. *arXiv preprint arXiv:2205.02655*, 2022. 8
- [33] Yoad Towel, Yoav Shalev, Roy Nadler, Idan Schwartz, and Lior Wolf. Zero-shot video captioning with evolving pseudo-tokens. *arXiv preprint arXiv:2207.11100*, 2022. 2, 3, 8
- [34] Yoad Towel, Yoav Shalev, Idan Schwartz, and Lior Wolf. Zerocap: Zero-shot image-to-text generation for visual-semantic arithmetic. In *CVPR*, pages 17918–17928, 2022. 2, 3, 8
- [35] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017. 4
- [36] Qiang Wang, Yanhao Zhang, Yun Zheng, Pan Pan, and Xian-Sheng Hua. Disentangled representation learning for text-video retrieval. *arXiv preprint arXiv:2203.07111*, 2022. 1, 2, 4, 5, 6, 8
- [37] Wenzhe Wang, Mengdan Zhang, Runnan Chen, Guanyu Cai, Penghao Zhou, Pai Peng, Xiaowei Guo, Jian Wu, and Xing Sun. Dig into multi-modal cues for video retrieval with hierarchical alignment. In *IJCAI*, pages 1113–1121, 2021. 8
- [38] Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. Vatex: A large-scale, high-quality multilingual dataset for video-and-language research. In *ICCV*, pages 4581–4591, 2019. 2, 5, 6
- [39] Xiaohan Wang, Linchao Zhu, and Yi Yang. T2vlad: global-local sequence alignment for text-video retrieval. In *CVPR*, pages 5079–5088, 2021. 6, 8
- [40] Xiaohan Wang, Linchao Zhu, Zhedong Zheng, Mingliang Xu, and Yi Yang. Align and tell: Boosting text-video retrieval with local alignment and fine-grained supervision. *IEEE TMM*, 2022. 8
- [41] Wenhao Wu, Zhun Sun, and Wanli Ouyang. Revisiting classifier: Transferring vision-language models for video recognition. In *AAAI*, 2023. 8
- [42] Wenhao Wu, Xiaohan Wang, Haipeng Luo, Jingdong Wang, Yi Yang, and Wanli Ouyang. Bidirectional cross-modal knowledge exploration for video recognition with pre-trained vision-language models. In *CVPR*, 2023. 8
- [43] Zuxuan Wu, Ting Yao, Yanwei Fu, and Yu-Gang Jiang. Deep learning for video classification and captioning. In *Frontiers of multimedia research*, pages 3–29. 2017. 2, 5, 6
- [44] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *CVPR*, pages 5288–5296, 2016. 2, 5, 6
- [45] Lewei Yao, Runhui Huang, Lu Hou, Guansong Lu, Minzhe Niu, Hang Xu, Xiaodan Liang, Zhenguo Li, Xin Jiang, and Chunjing Xu. Filip: Fine-grained interactive language-image pre-training. In *ICLR*, 2021. 2
- [46] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*, 2022. 1
- [47] Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, et al. Florence: A new foundation model for computer vision. *arXiv preprint arXiv:2111.11432*, 2021. 1
- [48] Shuai Zhao, Linchao Zhu, Xiaohan Wang, and Yi Yang. Centerclip: Token clustering for efficient text-video retrieval. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 970–981, 2022. 6, 8