

# BiC-Net: Learning Efficient Spatio-Temporal Relation for Text-Video Retrieval

Ning Han, Jingjing Chen, Chuhao Shi, Yawen Zeng, Guangyi Xiao, and Hao Chen

**Abstract**—The task of text-video retrieval aims to understand the correspondence between language and vision, has gained increasing attention in recent years. Previous studies either adopt off-the-shelf 2D/3D-CNN and then use average/max pooling to directly capture spatial features with aggregated temporal information as global video embeddings, or introduce graph-based models and expert knowledge to learn local spatial-temporal relations. However, the existing methods have two limitations: 1) The global video representations learn video temporal information in a simple average/max pooling manner and do not fully explore the temporal information between every two frames. 2) The graph-based local video representations are handcrafted, it depends heavily on expert knowledge and empirical feedback, which may not be able to effectively mine the higher-level fine-grained visual relations. These limitations result in their inability to distinguish videos with the same visual components but with different relations.

To solve this problem, we propose a novel cross-modal retrieval framework, Bi-Branch Complementary Network (BiC-Net), which modifies transformer architecture to effectively bridge text-video modalities in a complementary manner via combining local spatial-temporal relation and global temporal information. Specifically, local video representations are encoded using multiple transformer blocks and additional residual blocks to learn spatio-temporal relation features, calling the module a Spatio-Temporal Residual transformer (SRT). Meanwhile, Global video representations are encoded using a multi-layer transformer block to learn global temporal features. Finally, we align the spatio-temporal relation and global temporal features with the text feature on two embedding spaces for cross-modal text-video retrieval. Extensive experiments are conducted on MSR-VTT, MSVD, and YouCook2 datasets. The results demonstrate the effectiveness of our proposed model. The code is available at: <https://github.com/lionel-hing/BiC-Net>.

**Index Terms**—Text-Video Retrieval, Spatio-Temporal Relation, Bi-Branch Complementary Network.

## I. INTRODUCTION

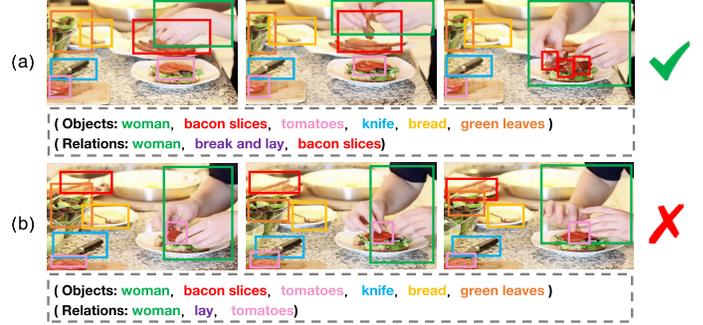
RECENT years have witnessed an exponential growth of multimedia data (e.g., video, image, and text), which increases the demands for effectively retrieving relevant data from another modality, when given a query of one modality. Being one of these challenging tasks, text-video retrieval aims to retrieve the video given a text query, which requires measuring the semantic similarity between a sentence and a video. Video data are distinct from images due to the

Ning Han, Chuhao Shi, Guangyi Xiao, and Hao Chen are with the Department of Information Science and Engineering, Hunan University, Changsha (e-mail: ninghan@hnu.edu.cn; sch8288@hnu.edu.cn; guangyi.xiao@gmail.com; chenhao@hnu.edu.cn).

Jingjing Chen is with the Department of Computer Science, Fudan University, Shanghai (e-mail: chenjingjing@fudan.edu.cn).

Yawen Zeng is with the Bytedance AI Lab, Beijing (e-mail: yawenzeng11@gmail.com).

Textual Query: A woman breaks the two bacon slices into pieces and lays them on the tomatoes.



**Fig. 1: An example of text to video retrieval. Given a textual query, a common pipeline with fine-grained [1] or global semantical visual features [2] will return two videos with the same compositions. The retrieval model with complementary spatio-temporal relation visual features can filter out false-positive without correct interactions.**

temporal dependencies among frames and the additional dynamic relationships among objects, resulting in the inability of existing video retrieval techniques to distinguish videos with the same visual components but with different relations. Figure 1 shows such an example. Given the text query “A woman breaks the two bacon slices into pieces and lay them on the tomatoes”, the existing retrieval systems are likely to consider both (a) and (b) as positive examples, since both of them contain the same motion (“laying”) and objects (“bacon slices”, “tomatoes”) with the text query. However, example (b) is indeed a false positive, as it presents “a woman lays tomatoes on the green leaves” (with bacon slices on the kitchen table). This example suggests that ignoring the visual relations (i.e., object relations) presented in videos could lead to inaccurate retrieval results. Therefore, capturing higher-level spatio-temporal visual relations in videos is crucial to distinguish similar videos.

This paper investigates the problem of cross-modal text-video retrieval. In the literature, many efforts have been devoted to learning better video representations, in order to improve the performance of text-video retrieval. Based on the granularity of feature representations, existing works can be roughly categorized into global and local feature-based methods. Global feature-based methods typically use global representations to represent entire video and sentence, which usually lose part of this temporal information and local details. Such approaches work well in a simple cross-modal retrieval

scenario, where only a single object is presented in the video or text query. For more realistic cases involving complex natural scenes, the performance of these methods is usually unsatisfactory. In contrast, local feature-based methods pay attention to local details and perform matching by detecting objects in videos and texts. With local region modeling, the performance of text-video retrieval has been significantly improved. Nevertheless, the existing efforts can only capture simple visual relations by graph convolutional network (GCN) [1], [3] or utilize an attention mechanism [4], [5] as a cross-modal interaction module to delve into high-level correspondences. As GCN-based video modeling is handcrafted, it depends heavily on expert knowledge and empirical feedback, which may not be able to effectively mine and model the higher-level fine-grained visual relations. Attention-based models, on the other hand, selectively align the key information presented in different modalities. As the fine-grained visual relations are also ignored by attention-based methods, the performance of these methods is still unsatisfactory, and novel modeling solutions are eagerly awaited.

To further improve the performance of text-video retrieval, this paper studies this problem from the perspective of spatio-temporal relation modeling for videos. Generally, there are two major obstacles in modeling the spatio-temporal relation. First, videos contain diverse spatial and temporal information within variations in motion and richer information in local visual details. These objects and interactions increase the difficulties in capturing higher-level fine-grained visual contents. Second, local relation modeling captures considerable fragmented information, which will overlook contextual information. Therefore, the way to comprehensively capture multi-granularity visual information to represent videos from complementary spatial and temporal perspectives is of great importance.

To address the aforementioned problems, we propose a novel Bi-Branch Complementary Network (BiC-Net), which modifies transformer architecture to effectively bridge text-video modalities in a complementary manner via combining local spatial-temporal relation and global temporal information. We present an overview of BiC-Net in Figure 2. Specifically, for videos, our BiC-Net attempts to extract two perspectives of features — global temporal features and local relation features. At the global temporal level, we directly adopt the widely used 2D and 3D-CNN. For local relational features, we use pre-trained Faster-RCNN [6] to extract regional features (i.e., features of bounding boxes). Then, a spatio-temporal residual transformer is employed for learning high-level fine-grained relational features. This module separately captures local spatial relations, and long-term temporal relations among local spatial relations. In addition, a multi-layer transformer block is applied for learning global temporal features. To cover different levels of semantics, we align the global temporal and local relation features with the text feature on two embedding spaces. Lastly, the similarity between videos and texts is measured in both embedding spaces and then summed to obtain the final similarity score. In this way, the global temporal information and local relation information in a video can be utilized for cross-modal text-video retrieval comprehensively.

Our contributions are summarized as below:

- We incorporate feature-split with bi-branch framework called BiC-Net to capture local relations and global temporal features comprehensively, which aligns the global temporal and local relation features with the text feature on two embedding spaces for cross-modal text-video retrieval.
- We first introduce a simple and effective spatio-temporal residual transformer to learn higher-level local relation features, and a multi-layer temporal transformer to further explore global temporal information for global temporal features. In this way, the bi-branch information in video and text can accurately capture cross-modal semantic alignment in a cooperative and complementary manner.
- We conduct extensive experiments on three standard benchmarks and verify the effectiveness of our proposed method by showing that BiC-Net can achieve SOTA performance (86.7% on MSR-VTT 1k-A test set) under similar conditions.

The rest of this paper is organized as follows. In Section II, we briefly describe a review of related work. In Section III, we describe our proposed BiC-Net model. In Section IV, we provide implementation details and experimental results. In Section V, we finally conclude our paper.

## II. RELATED WORK

### A. Text-video Retrieval

According to the granularity of feature representations, we roughly divide existing works into two groups: global feature-based methods and local relation feature-based methods.

**Global feature-based methods** [2], [7], [8] extract global feature representations of videos and texts and then learn a joint embedding space where visual and textual similarity is measured. For the video representation, they adopt 2D/3D CNN models to extract frame features and aggregate frame features by average-pooling [9]–[11] or max-pooling [2], [12]. For the video representation, they focused only on leveraging the global feature of the video. For instance, Dong et al. [8], [13] employ three levels, i.e., global, temporal, and local to encode videos and texts and learn a hybrid common space for video-text similarity measurement. Miech et al. [2] adopt 2D and 3D CNN to extract frame features and only use max-pooling to obtain global video representation. Yang et al. [7] present a latent semantic tree to encode the text and used a multi-head self-attention mechanism to obtain the temporal-attentive video representation.

**Local feature-based methods** [3]–[5], [14]–[16] use local semantic information from language or video for better text-video alignment from different aspects and then perform text-video retrieval tasks. Wray et al. [14] disentangle action phrases into verbs and nouns for fine-grained video retrieval. The graph-based approaches [1], [3], [16] construct different semantic correlation graphs for videos and learn fine-grained semantic relations for text-video retrieval. Some works [4], [5], [15] also propose fine-grained alignment models that decompose text and video into multiple levels and align text with video at multiple levels for text-video matching.

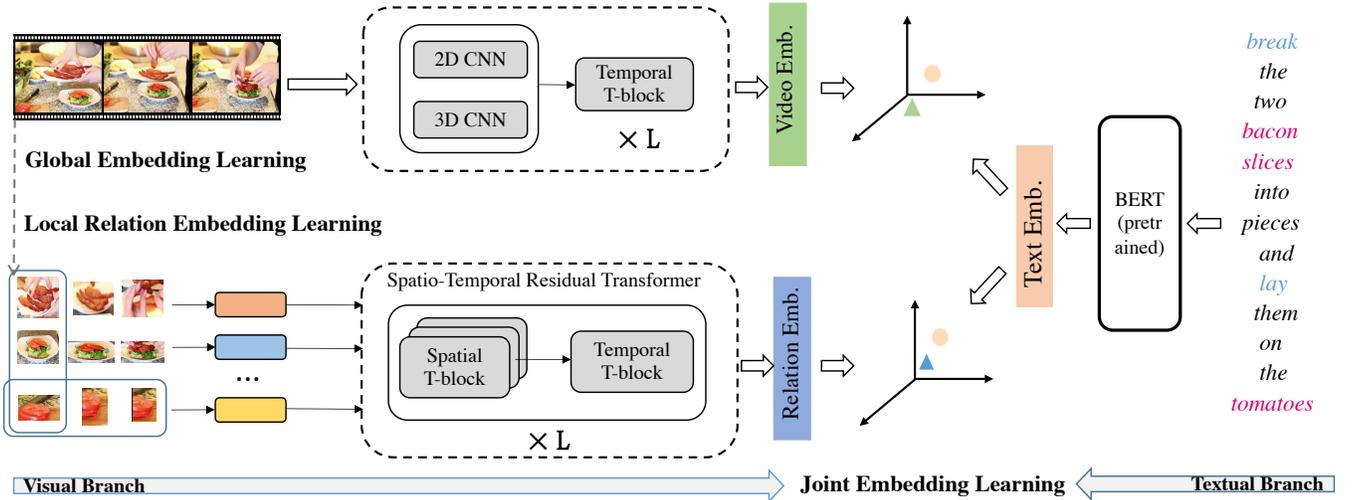


Fig. 2: The overall framework of our proposed BiC-Net. First, we extract local relational and global visual features for videos. The local relations are represented by local regional features using a spatio-temporal residual transformer. The global video features are represented by 2D-CNN and 3D-CNN features via a multi-layer temporal transformer. Then, we extract textual features by BERT. Finally, both video and relational features are leveraged to align with textual features on two embedding spaces for cross-modal text-video retrieval. Among them, T-Block denotes the transformer block.

Recently, some studies have also explored a combination of video experts (e.g., motion, audio, and speech) [10], [17]–[19] or pre-trained video experts [20]–[22] to improve the performance of cross-modal retrieval. Lately, Transformer-based works [23]–[27] have benefited from pre-training models on large-scale language-vision datasets [23], [27]. For example, Bain et al. [27] propose an end-to-end trainable model which adopts a space-time transformer encoder to flexibly train on both video and image datasets. Luo et al. [23] apply the joint language-vision model of CLIP [28], pre-trained on a large-scale text-image dataset as a backbone for text-video retrieval. However, Transformer-based methods have a heavy computational burden due to computational intensive operations and are extremely time-consuming to pre-train on large-scale datasets. Different from these existing works, our study introduces a new spatial-temporal residual transformer to learn higher-level local relation features and a multi-layer transformer to further explore global temporal information for global temporal features. In this way, bi-branch information in video and text can accurately capture cross-modal semantic alignment in a cooperative and complementary manner.

### B. Spatio-Temporal Relation Modeling in Video Understanding

For spatio-temporal relation modeling in video understanding, earlier works adopt 2D/3D CNNs to represent the core operators for spatio-temporal feature learning across downstream video tasks [29]–[32], [32]. However, these video representations focus on learning spatio-temporal features from the entire video and can hardly capture local spatial-temporal relation information. To understand the local relation information in the video, several efforts have demonstrated the effectiveness of incorporating local spatial-temporal relationships into video

understanding in many downstream applications, such as visual relationship detection [33]–[35], action recognition [36]–[38], and video retrieval [1], [3], [16]. For instance, Qian et al. construct a spatio-temporal graph in adjacent video clips to define the relationships between objects. Wang et al. [36] abstract the video as a space-time graphs for action recognition. Song et al. [1] model video as a spatial-temporal graph between object interactions for text-video retrieval. However, modeling object spatio-temporal relations in the video is still not thoroughly investigated. These studies have built visual relation graphs and adopted the GCN [39] to extract visual relation graph features. Massive graph construction and graph feature extraction are hand-crafted, complex, and time-consuming. Recently, the transformer [40] has shown great superiority in understanding 1, 2, and 3-dimensional signals (e.g., natural language processing and computer vision), and has strong interpretability, and strong representation capabilities. Unlike these works, our work designs a spatio-temporal residual transformer to learn the local spatio-temporal relations and further mine the object interactions. Notably, we validate in the experiments that under a strict memory budget, our approach can surpass many related methods.

## III. PROPOSED METHOD

As depicted in Figure 2, the overall pipeline of the proposed method consists of four modules: 1) video embedding learning, which involves extracting video global features; 2) relation embedding learning, which involves extracting local relational features in videos; 3) text embedding learning, which learns the representation for textual sentences by BERT [41]; and 4) joint embedding learning, which optimizes the correspondence between text and video features in a common space with a triple ranking loss.

### A. Video Embedding Learning

Given a long video clip, we sample  $T$  video frames from it with the same temporal duration between every two frames. For frame-level features, we first use 2D-CNN to extract appearance features and 3D-CNN to extract motion features. Then, we concatenate 2D and 3D features and apply a point-wise linear layer to obtain global visual features  $F_g \in \mathbb{R}^{d_g}$ . Finally, we feed the result to standard multi-layer transformer block [40] and an attention-aware feature aggregation layer [42] to obtain its video embedding, which is denoted as  $F_v \in \mathbb{R}^{d^*}$ .

### B. Relation Embedding Learning

In addition to having global visual features, the proposed framework learns local relation features from the video to improve the performance of cross-modal retrieval. The introduction of spatio-temporal relation among objects in the video equips the model with the ability to identify the fine-grained differences of video with similarity. To capture the visual relations from the video, we first adopt the pre-trained Faster RCNN [43] to detect frame-level region proposals and select the top  $N$  region proposals with the highest detection confidence to represent each frame. Prior efforts [36], [38] focus on abstracting frame-level region proposals as fully connected spatial-temporal graphs and using GCN to learn relational features. However, computing all pair-wise relations across all video frames would be inefficient in creating a video as a fully connected graph. In recent years, pure transformer-based models have shown promising performance due to their strong representation capabilities. As a central piece of transformer, self-attention comes with a flexible mechanism to deal with variable-length inputs. It can be understood as a fully connected layer where the weights are dynamically generated from pairwise relations from inputs, which conveys refreshing solutions to process visual relations.

Inspired by these pioneering efforts, to capture higher-level visual relations from the video, we design a new architecture to learn the relation embeddings, named Spatio-Temporal Residual Transformer (SRT), that exploits all the variants of transformer blocks and residual connections but composes each in different placement of SRT. In the following, the basic components used in the transformer block and the transformer block used in the SRT module are presented in detail.

**Transformer Block.** The Transformer consists of multi-head self-attention (MSA), multi-layer perceptron (MLP), and layer-norm (LN). In the self-attention module, the inputs  $X \in \mathbb{R}^{n \times d}$  are linearly transformed to three parts, i.e., queries  $Q \in \mathbb{R}^{n \times d_k}$ , keys  $K \in \mathbb{R}^{n \times d_k}$  and values  $V \in \mathbb{R}^{n \times d_v}$ , where  $n$  is the sequence length,  $d$ ,  $d_k$ ,  $d_v$  are the dimensions of inputs, queries (keys) and values, respectively. The scaled dot-product attention is applied on  $Q$ ,  $K$ ,  $V$ :

$$SA(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V. \quad (1)$$

With  $SA(Q, K, V)$ ,  $MSA$  is defined as:

$$MSA(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_M) W^O, \quad (2)$$

where  $\text{head}_i = SA(QW_i^Q, KW_i^K, VW_i^V)$ .

Where  $QW_i^Q, KW_i^K, VW_i^V$  are projections of different heads,  $W^O$  is another mapping function. The MLP is applied between self-attention layers for feature transformation and non-linearity:

$$MLP(X) = GELU(XW_1 + b_1)W_2 + b_2, \quad (3)$$

where  $W_1 \in \mathbb{R}^{d \times d_m}$  and  $W_2 \in \mathbb{R}^{d_m \times d}$  are weights of the two fully-connected layers respectively,  $b_1 \in \mathbb{R}^{d_m}$  and  $b_2 \in \mathbb{R}^d$  are the bias terms, and GELU [44] is the activation function. Layer normalization [45] is a key part in transformer for stable training and faster convergence, and LN is applied over each sample  $x \in \mathbb{R}^d$  as follows:

$$LN(x) = \frac{x - \mu}{\eta} \odot \gamma + \beta, \quad (4)$$

where  $\mu \in \mathbb{R}$ ,  $\eta \in \mathbb{R}$  are the mean and standard deviation of the feature respectively,  $\odot$  is the element-wise dot, and  $\gamma \in \mathbb{R}^d$ ,  $\beta \in \mathbb{R}^d$  are learnable affine transform parameters.

**SRT for relation embedding learning.** We propose a spatio-temporal residual transformer architecture to learn local relation information in a video. In this spatio-temporal residual transformer, we have two data flows in which one flow operates across the frame and the other processes the object proposals inside each frame. Suppose that a set of object proposals  $Y^t = \{y^t\}_{n=1}^N$  are in frame  $t$ , where  $y^t \in \mathbb{R}^{d_r}$  is the feature vector of the  $n$ -th proposal and  $N$  is the top  $N$  region proposals. We view each frame tensor  $Y_0^t$  as a sequence of object proposal embeddings:

$$Y_0^t = [y_0^{t,1}, y_0^{t,2}, \dots, y_0^{t,N}]. \quad (5)$$

For the object proposal embeddings, to capture spatial relations among visual objects, we utilize a transformer block to explore the interaction pattern in spatial (frame) between object proposals. Then, a residual connection is used to aggregate spatial information and original local information:

$$Y_l^{t'} = Y_{l-1}^t + MSA(LN(Y_{l-1}^t)), \quad (6)$$

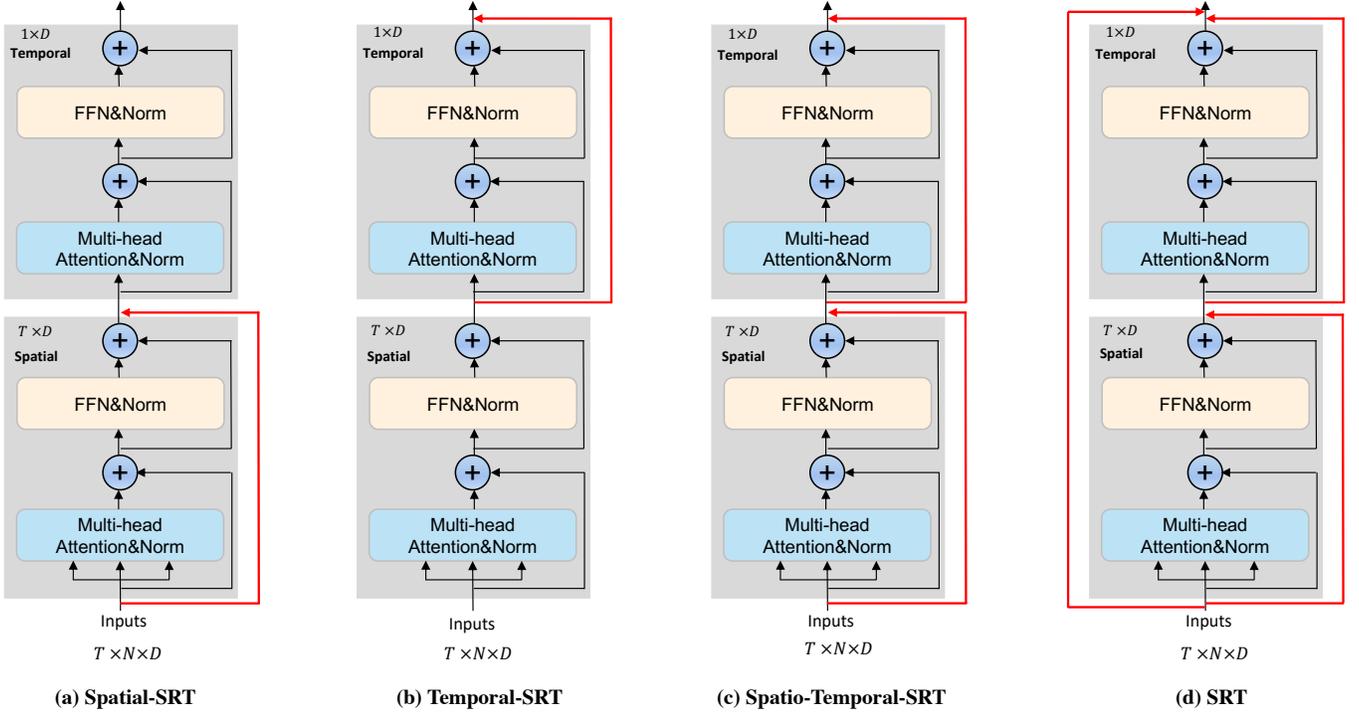
$$Y_l^{t''} = Y_l^{t'} + MLP(LN(Y_l^{t'})), \quad (7)$$

$$Y_l^{t'''} = Y_l^{t''} + Y_{l-1}^t. \quad (8)$$

where  $l = 1, 2, \dots, L$  is index of the  $l$ -th layer, and  $L$  is the total number of layers. The updated features after multi-layer transformer block are forwarded to an average pooling layer, which calculates the mean of all the proposal features and leads to a  $1 \times d^r$  dimensions representation. All frame tensors after transformation are

$$Z_0 = [Y_L^{t'''}, Y_L^{t'''}, \dots, Y_L^{t'''}]. \quad (9)$$

This process builds the relationship among proposals by computing interactions between any two proposals. For the frame level, we create the object proposal embedding memories to store the sequence of frame-level representations  $Z_0$ . Similar to the object proposal level processing, we use a transformer block for transforming the frame embeddings. Then, a residual connection is used to aggregate temporal information with spatial information and spatio-temporal information with original



**Fig. 3: Designs of SRT and its variants: (a) spatial residual (Spatial-SRT); (b) temporal residual (Temporal-SRT); (c) spatio-temporal residual (Spatio-Temporal-SRT); (d) our SRT.**

local information, respectively. Our final relation embedding is defined as:

$$Z'_l = Z_{l-1} + MSA(LN(Z_{l-1})), \quad (10)$$

$$Z''_l = Z'_l + MLP(LN(Z'_l)), \quad (11)$$

$$Z'''_l = Z''_l + Z_{l-1}, \quad (12)$$

$$F_r = Z'''_l + Y_{l-1}^t. \quad (13)$$

The temporal transformer block is used for modeling temporal relation among frame embeddings. Finally, we apply an attention-aware feature aggregation layer [42] to obtain the final relation embedding, denoted as  $F_r \in \mathbb{R}^{d_*}$ .

Next, we discuss several variants for SRT, as illustrated in Figure 3. **Spatial-SRT** only utilizes Eq.(8) to aggregate spatial information and original local information by a residual connection (i.e., Figure (3a)). **Temporal-SRT** only adopts Eq.(12) to aggregate temporal information and spatial information by a residual connection (i.e., Figure (3b)). **Spatio-Temporal-SRT** only uses Eq.(8) and Eq.(12) to aggregate temporal information with spatial information and spatial information with original local information by a residual connection, respectively (i.e., Figure (3c)). Besides, we use **Non-SRT** as a base variant, and the module indicates that no residuals are added between the transformer blocks. We compare the above five variants of SRT on a standard benchmark in Section IV-B and observe the SRT achieves the best performance. Moreover, we find that SRT introduces minor modifications of the residual connection but grants maximum benefits.

### C. Text Embedding Learning

For learning the contextual relations between the words in the video description sentence  $s_i$ , we adopt a BERT language representation model to encode the word sequence, and it applies the bidirectional training of transformer [40] to language modeling. It includes 12 layers of transformer blocks. Each block has 12 attention heads, and the hidden size is 768. Here, we take the hidden state of the per-token outputs of the last 2 layers to represent the information of the entire input sentence  $F_s \in \mathbb{R}^{d_t}$ . Finally, we transform each sentence representation  $F_s \in \mathbb{R}^{d_t}$  into a text embedding feature  $F_t \in \mathbb{R}^{d_*}$  by using a pointwise linear layer and an attention-aware feature aggregation layer [42].

### D. Joint Embedding Learning

The purpose of joint embedding learning between video and textual features is to perform similarity comparisons. For a given video  $V_i$ , the proposed framework extracts two types of embedding features — video embeddings  $F_v$  and relation embeddings  $F_r$ . We calculate the similarity between videos and sentences in both embedding spaces. Specifically, for a given sentence  $T_i$ , the similarity score with  $V_i$  is obtained by summing the cosine similarities between its text embedding features  $F_t$  and such two types of video embedding features,

$$S(V_i, T_i) = \lambda \cdot \text{cosine}(F_r, F_t) + (1 - \lambda) \cdot \text{cosine}(F_v, F_t). \quad (14)$$

where  $0 \leq \lambda \leq 1$  is a hyper-parameter to balance the importance of two similarity scores. Based on the defined similarity score, we use a hinge-based triplet ranking loss to

**TABLE I: Performance of introducing visual relations (VR) for cross-modal retrieval. The evaluations are done on 1k-A test set (Training-9k) [18] for MSR-VTT.**

Method	Text-to-Video				Video-to-Text			
	R@1	R@5	R@10	MedR	R@1	R@5	R@10	MedR
dataset split from [18]								
VG	32.9	65.8	79.7	3	32.1	65.2	77.6	3
VR <sub>st</sub>	29.8	64.4	77.2	3	29.2	63.6	76.8	3
VG + VR <sub>m</sub>	32.7	66.0	79.4	3	32.9	67.0	79.6	3
VG + VR <sub>s</sub>	33.8	69.3	82.9	2	36.2	72.4	84.3	2
VG + VR <sub>t</sub>	33.7	67.5	81.7	3	34.0	70.2	82.7	2
BiC-Net	<b>39.4</b>	<b>75.5</b>	<b>86.7</b>	<b>2</b>	<b>39.4</b>	<b>76.5</b>	<b>85.9</b>	<b>2</b>

**TABLE II: Performance of the variants for SRT for cross-modal retrieval. The evaluations are done on 1k-A test set (Training-9k) [18] for MSR-VTT.**

Method	Text-to-Video				Video-to-Text			
	R@1	R@5	R@10	MedR	R@1	R@5	R@10	MedR
dataset split from [18]								
Non-SRT	36.2	73.9	84.4	2	38.2	74.3	87.5	2
Spatial-SRT	37.8	71.7	85.0	2	39.2	74.5	85.2	2
Temporal-SRT	37.8	71.9	85.2	2	40.1	73.9	85.8	2
Spatio-Temporal-SRT	38.2	73.2	85.8	2	39.3	74.2	85.5	2
SRT (BiC-Net)	<b>39.4</b>	<b>75.5</b>	<b>86.7</b>	<b>2</b>	<b>39.4</b>	<b>76.5</b>	<b>85.9</b>	<b>2</b>

**TABLE III: Cross-modal retrieval comparison with state-of-the-art methods on MSR-VTT.**

Method	Text-to-Video				Video-to-Text			
	R@1	R@5	R@10	MedR	R@1	R@5	R@10	MedR
Full test set [46]								
STG [1]	8.3	23.7	33.9	28	–	–	–	–
HGR [4]	9.2	26.2	36.5	24	15.0	36.7	48.8	11
DualEncoding [8]	11.6	30.3	41.3	17	<b>22.5</b>	47.1	58.9	7
T2VLAD [19]	12.7	34.8	47.1	12	20.7	48.9	62.1	6
BiC-Net	<b>19.2</b>	<b>47.0</b>	<b>62.5</b>	<b>6</b>	20.6	<b>49.3</b>	<b>63.7</b>	<b>6</b>
1k-B test set [12]								
CE [10]	18.2	46.0	60.7	7	18.0	46.0	60.3	6.5
DualEncoding [8]	23.0	50.6	62.5	5	25.1	52.1	64.6	5
MMT [18]	20.3	49.1	63.9	6	21.1	49.4	63.2	6
T2VLAD [19]	26.1	54.7	68.1	4	26.7	56.1	70.4	4
BiC-Net	<b>34.0</b>	<b>71.1</b>	<b>84.1</b>	<b>3</b>	<b>37.9</b>	<b>73.4</b>	<b>85.3</b>	<b>2</b>
1k-A test set (training-7k) [2]								
Miech et al. [2]	12.1	35.0	48.0	12	–	–	–	–
STG [1]	15.5	39.2	50.4	10	–	–	–	–
TCE [7]	17.1	39.9	53.7	9	–	–	–	–
DualEncoding [8]	21.6	49.5	62.3	6	27.8	48.7	58.7	6
BiC-Net	<b>32.8</b>	<b>68.2</b>	<b>82.4</b>	<b>3</b>	<b>36.8</b>	<b>71.5</b>	<b>83.5</b>	<b>2</b>
1k-A test set (Training-9k) [18]								
MMT [18]	24.6	54.0	67.1	4	24.4	56.0	67.8	4
SUPPORT-SET [22]	27.4	56.3	67.7	3	26.6	55.1	67.5	3
Frozen [27]	31.0	59.5	70.5	3	–	–	–	–
CLIP4Clip [23]	<b>44.5</b>	71.4	81.6	2	–	–	–	–
BiC-Net	39.4	<b>75.5</b>	<b>86.7</b>	<b>2</b>	<b>39.4</b>	<b>76.5</b>	<b>85.9</b>	<b>2</b>

encourage the similarity score of matched video and sentence to be larger than those of mismatched ones:

$$\mathcal{L}_r = [\delta - S(V_i, T_i) + S(V_i, T_j)]_+ + [\delta - S(V_i, T_i) + S(V_j, T_i)]_+, \quad (15)$$

where  $0 < \delta \leq 1$  is the margin, the operator  $[x]_+ = \max(x, 0)$ , and  $S(\cdot, \cdot)$  is the similarity function.  $(V_i, T_i)$  represents the positive pair, while  $(V_i, T_j)$  and  $(V_j, T_i)$  represent the negative pairs available in the mini-batch.

## IV. EXPERIMENTS

### A. Experimental Setup

1) Dataset: We evaluated the proposed BiC-Net model on three benchmarks: MSR-VTT, MSVD, and YouCook2. The MSR-VTT dataset [46] is the most widely-used dataset for text-video retrieval. It contains 10,000 Youtube video clips with 20 different text captions. Following the settings in [12], [46], [47], we adopt three kinds of evaluation settings. For the 1k-A test set [47], we using 7k train+val videos [2] and 9k train+val videos for training [18] and report results. The MSVD dataset [48] contains 1,970 video clips from YouTube. Each video clip has around 40 descriptions in multiple lan-

TABLE IV: Cross-modal retrieval comparison with state-of-the-art methods on MSVD.

Method	Text-to-Video				Video-to-Text			
	R@1	R@5	R@10	MedR	R@1	R@5	R@10	MedR
Mithun et al. [17]	16.1	41.1	53.5	9	23.4	45.4	53.0	8
CE [10]	19.8	49.0	63.8	6	-	-	-	-
ViSERN [3]	18.1	48.4	61.3	6	24.3	46.2	59.5	7
SUPPORT-SET [22]	23.0	52.8	65.8	5	<b>27.3</b>	50.7	60.8	5
BiC-Net	<b>24.6</b>	<b>57.0</b>	<b>70.3</b>	<b>4</b>	24.2	<b>58.7</b>	<b>70.1</b>	<b>4</b>

TABLE V: Cross-modal retrieval comparison with state-of-the-art methods on YouCook2. TS: trained from scratch on YouCook2.

Method	Text-to-Video				Video-to-Text			
	R@1	R@5	R@10	MedR	R@1	R@5	R@10	MedR
HGLMM FV CCA [2]	4.6	14.3	21.6	75	-	-	-	-
Miech et al. [2]	4.2	13.7	21.5	65	-	-	-	-
COOT [42]	5.9	16.7	24.8	49.7	-	-	-	-
AME-Net [16]	7.6	21.5	32.8	<b>28</b>	7.9	22.5	32.2	<b>28</b>
BiC-Net (TS)	<b>8.7</b>	<b>23.9</b>	<b>33.5</b>	31	<b>8.3</b>	<b>23.6</b>	<b>32.6</b>	31

guages. We only adopt English annotations in experiments. Following prior work [49], we separate the dataset into 1,200 clips for training, 100 clips for validation, and 670 clips for testing. The YouCook2 dataset [50] contains 2,000 cooking videos with 14,000 video clips. It covers 89 types of recipes. Each video clip is described by a textual sentence. Referring to [2], we evaluate the text-video clip retrieval task on the validation clips.

2) Evaluation Metrics: We employ the widely used median retrieval rank (MedR) and recall rate at top  $\mathcal{K}$  ( $R@K$ ) for assessing retrieval accuracy. MedR measures the median rank position among where true positives are returned.  $R@K$  measures the fraction of true positives being ranked at top  $\mathcal{K}$  returned results. Therefore, lower MedR scores indicate higher performance; in contrast, higher  $R@K$  scores indicate better performance.

3) Implementation Details: We sample 26 video frames from it with the same temporal duration between every two frames. In our experiments, the ILSVRC-2012-CLS [51] pre-trained InceptionResNetV2 [52] is adopted to extract 1536-D 2D features and the Kinetics [53] pre-trained I3D [54] to extract 1024-D 3D features. The number  $N$  of regions within a frame is 36, identical to [43]. The dimension  $d$  of region features extracted from ResNet-101 is 2048-D. The dimensionality of video-embedding vectors  $F_v$  and relation-embedding vectors  $F_r$  are set as 1024-D. For each sentence, we use pre-trained BERT to extract 1536-D word embedding and apply a point-wise linear layer and an attention-aware feature aggregation layer [42] to obtain 1024-D text-embedding vectors.

We implement our proposed model using PyTorch<sup>1</sup> and train it on 4 Tesla V100 GPUs. We train for 60 epochs using Adam optimizer [55] with a mini-batch size of 64. On the MSR-VTT, MSVD, and YouCook2, the learning rates are set to 0.0002, 0.0004, and 0.0004, respectively. As for the layer number  $L$  of transformer block, we set it to 4, 2 and 4 on the MSR-VTT, MSVD, and YouCook2 datasets, respectively. In addition, the trade-off parameter  $\lambda$  in Eq. (14), the margin  $\delta$  in Eq. (15) are set to 0.5 and 0.2, respectively.

## B. Ablation Studies

1) **Experiments with spatio-temporal Relation.** We experimented with variants of our model to verify the effectiveness of introducing spatio-temporal relation for text-video retrieval:

- **VG.** We only utilize the pre-trained 2D and 3D CNNs to extract the global features of the whole video as video embedding learning.
- **VG + VR<sub>m</sub>.** We apply the average-pooling features of all regions without using the features extracted by spatio-temporal residual transformer as relation embedding learning.
- **VG + VR<sub>s</sub>.** We only utilize regional spatial relation features as relation embedding learning and global features as video embedding learning.
- **VG + VR<sub>t</sub>.** We only utilize regional temporal relation features as relation embedding learning and global features as video embedding learning.
- **VR<sub>st</sub>.** We only utilize regional spatio-temporal relation features as relation embedding learning.

We explore these model variants on the MSR-VTT, as shown in Table II. We omit the results on MSVD and YouCook2 because of space limitations, but they show similar trends to MSR-VTT. From the results, we have the following observations. First, as expected, on both text-to-video and video-to-text, our BiC-Net, **VG + VR<sub>s</sub>**, and **VG + VR<sub>t</sub>** significantly outperforms **VG** alone. The result verifies the significance of introducing spatio-temporal relation representation. Second, compared with **VG + VR<sub>m</sub>**, the performance of our BiC-Net verifies that the spatio-temporal residual transformer can capture the fine-grained local relational features. Third, compared with two variants of models (i.e., **VG** and **VR<sub>st</sub>**) that only use either global visual features or local relational features, our model considers both global and local relational features to achieve the best performance. This verifies the effectiveness of aligning the global visual and local relational features with text features on two embedding spaces. Notably, global visual features and local relational features are highly complementary, and their combination leads to an improvement far beyond the performance of the global visual features alone. Moreover, compared with **VG + VR<sub>s</sub>**

<sup>1</sup><http://www.pytorch.org>

and  $\mathbf{VG} + \mathbf{VR}_t$ , our BiC-Net achieves substantially better performance, which reveals the complementary of the spatial and temporal relation features.

2) **Evaluation of of SRT.** We test the effectiveness of our proposed SRT and its variants on relation embedding learning. As shown in Table IV, our SRT and its variants achieve better performance than Non-SRT, which indicates the effectiveness of spatio-temporal relation modeling by adding residual blocks. The difference between Spatial-SRT and Temporal-SRT is that a residual block is added at different positions. We can see that Temporal-SRT significantly surpasses Spatial-SRT, which indicates the importance of temporal relation modeling. Spatio-Temporal-SRT adds a residual block based on Spatial-SRT/Temporal-SRT, which achieves better performance than Spatial-SRT/Temporal-SRT by aggregating temporal information with spatial information. In the end, compared to the other variants, we observed that our proposed SRT achieves the best performance when three residual blocks are added, indicating that simultaneously adding residual blocks in our model performs better than adding only one of them. To sum up, the contribution of each component enables our SRT to learn higher-level spatio-temporal relation information.

### C. Comparison with State-of-the-art Methods

To demonstrate the effectiveness of the BiC-Net solution, we compared it to several state-of-the-art baselines: (1) RNN-based methods: DualEncoding [8], TCE [7], (2) Multimodal Fusion methods: Mithun et al. [17], CE [10], MMT [18], (3) GCN-based methods: ViSERN, [3], STG [1] and AME-Net [16], (4) Transformer-based methods: COOT [42], CLIP4Clip [23] and Frozen [27], (5) other methods: HGLMM FV CCA [2], Miech et al. [2], SUPPORT-SET [22], T2VLAD [19].

1) Experiments on MSR-VTT: The experimental results are presented in Table III. We can observe that for all data partitions, our proposed method consistently outperforms all, compared to traditional RNN-based methods and multimodal Fusion methods in all evaluation metrics by a large margin, including CE [10], MMT [18], and T2VLAD [19], which use expert features (e.g., object, motion, face, scene, sound, and speech). Moreover, our BiC-Net significantly outperforms recent spatio-temporal relation-based method (STG [1]) in all evaluation metrics, especially, it boosts the text-video retrieval quality by a margin of 28.6% in R@10 on full test set. This condition reveals the effectiveness of our BiC for modeling video global and relational information. The obvious performances are shown on full test set and 1k-A test set (training-7k).

We also compare with some typical transformer-based methods, such as CLIP4Clip [23] and Frozen [27]. CLIP4Clip adopts the language-vision transformer model of CLIP [28] pre-trained on a large-scale text-image dataset as a backbone. Frozen uses a transformer-based video model [56] as a backbone. In contrast, we design a new transformer-based backbone to model spatio-temporal relations and global temporal information. Our BiC outperforms most of the compared methods on 1k-A test set (Training-9k), e.g., BiC 86.7% vs CLIP4Clip 81.6% w.r.t. text-to-video R@10. This indicates

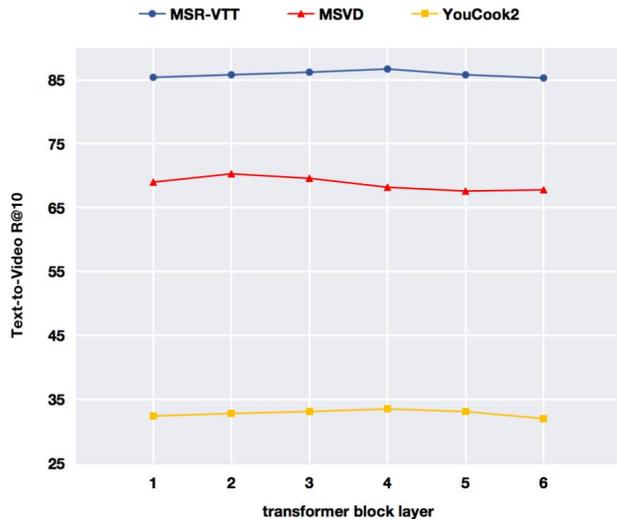


Fig. 4: Performance of text-to-video retrieval with different layers number of transformer blocks.

that learning cross-modal complementarity in a cooperative and complementary manner takes effect.

2) Experiments on MSVD: Table IV summarizes the performance comparison results. We also observe that our proposed BiC outperforms recent state-of-the-art methods in terms of most indicators. Note that among all these methods, ViSERN [3] uses only local video features to compute the similarity between the video and text. Analogously, we also observe that BiC-Net outperforms the local feature-based method ViSERN [3] by a great margin. This reveals that jointly modeling the global and local video representation plays a significant role in text-video retrieval, contributing to more powerful representation. To ensure a fair comparison, we compare the previous SOTA method, SUPPORT-SET [22] without pre-training on HowTo100M [2]. Under the full fair comparison, our BiC outperforms the previous best method SUPPORT-SET by 9.3% on video-text retrieval R@10. Notably, on MSVD, the performance of our model is not as outstanding. The reason for small gains is that the transformer has the property of lacking structural bias making it prone to overfitting for small-scale data.

3) Experiments on YouCook2: As shown in Table V, our method achieves the best performance, which is 8.7% absolute gains in the evaluation metric of text-video retrieval R@10 better than COOT [42]. In Miech et al. [2] and COOT [42], the global video features are used for video representation, whose performances are worse than most methods. AME-Net [16] adopts global features and handcrafted graph-based relation features. AME-Net achieves better performances than Miech et al. [2], while their performances are worse than ours, which indicates that the global and local information learned by our method can be mutually promoted in a complementary manner. This observation indicates that in addition to the global video features, local relation features are also important for video representation.

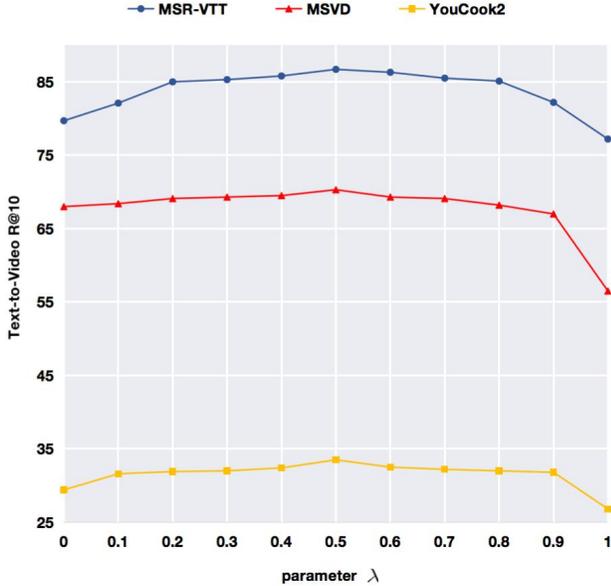


Fig. 5: Evaluation of different weight combination of the global and relation similarities.

#### D. Parametric Sensitivity Analysis

We carry out experiments to explore how the layer number of transformer block  $L$  and the trade-off parameter  $\lambda$  affect the retrieval performance. Notably, we omit the video-text retrieval results on three datasets due to space limitations, which show similar trends to text-video retrieval. First, we analyze the influence of the hyperparameter, that is, the layer number of transformer block on the MSR-VTT 1k-A test set [18], MSVD, and YouCook2 datasets. Figure 4 presents the results across the layer number of transformer block on the two datasets by R@10; note that R@1 and R@5 present the same trend, in which performances increase until certain numbers (4, 2, and 4 for MSR-VTT, MSVD and YouCook2 datasets, respectively) and then become stable. This result is due to the model’s capability of capturing the spatio-temporal relations of the deepest layer numbers.

Moreover, the influence of the hyperparameter  $\lambda$  in Eq. (11) is revealed in Figure 5. We assign different trade-off parameter  $\lambda$  to the two scores (i.e.,  $\mathbf{VR}_{st}$  and  $\mathbf{VG}$ ) to observe their influence on the matching performance on the three datasets. By analyzing the results shown in Figure 5, we have the following observations: 1) The leftmost part of Figure 5 shows the results when  $\mathbf{VR}_{st}$  accounts for 0, that is, when the proportion of  $\mathbf{VG}$  is 1, which means that we remove the visual relations module from our model (i.e.,  $\mathbf{VG}$ ). We can observe that when the spatio-temporal relations module is removed, the retrieval performance is reduced by a large margin over the three datasets. This condition shows the positive effect of comprehensively introducing spatio-temporal relations for text-video retrieval. 2) Increasing the proportion of  $\mathbf{VR}_{st}$  substantially boosts the performance of model. Our model performs best performance on the three datasets when  $\lambda = 0.5$ . Therefore, we argue that  $\mathbf{VR}_{st}$  and  $\mathbf{VG}$  occupy the same contribution to the overall similarity. We conclude that the two similarities work together to obtain the best retrieval

TABLE VI: Comparison with different models in terms of model size and computation overhead at the inference stage.

Model	Parameters (M)	FLOPs (G)
MMT [18]	133.4	12.64
DualEncoding [8]	95.9	<b>3.64</b>
BiC-Net	<b>31.48</b>	10.33

performance in a cooperative manner.

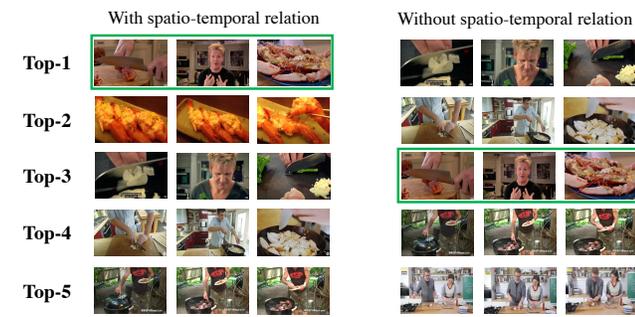
#### E. Model Complexity

We compare our method with open-source methods in terms of model size and computation overhead at the inference stage. As shown in Figure 4, since the performance of using one layer of transformer block outperforms MMT by a large margin, we only calculate the model size and computational overhead of using one layer of transformer block. Analogously, we also observe that our BiC-Net with one layer of transformer block outperforms DualEncoding by a great margin on the MSR-VTT 1k-A test set [2]. Notably, we omit the text-video retrieval results of VSR with a layer of transformer block on the MSR-VTT 1k-A test set [2] due to space limitations, which show similar trends to the MSR-VTT 1k-A test set [18]. In addition, we conclude that for each additional layer of transformer block, the computational cost will increase by 8.29 GFLOPs and the parameters will increase by 25.19M. Following [8], we measure the number of FLOPs required for a text-video pair. As shown in Table VI and Figure 4, we have two main observations: 1) our BiC-Net with one layer of transformer block achieves 85.6% text-to-video R@10 accuracy on the MSR-VTT 1k-A test set [18], which is 18.5% higher than MMT, with fewer parameters and lower computational cost. 2) our BiC-Net with one layer of transformer block is smaller and slightly slower than DualEncoding.

#### F. Qualitative Results

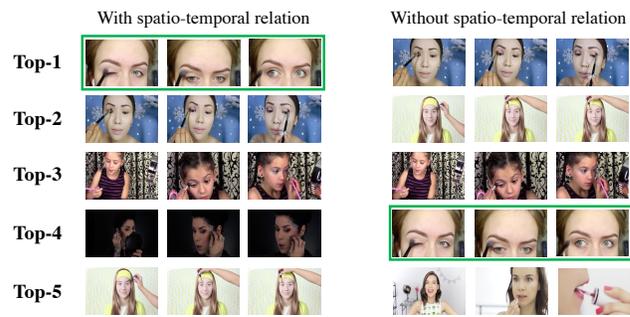
Figure 6 shows two examples of text-to-video retrieval results between the model with and without visual spatio-temporal relations. We specifically choose two text-video retrieval examples that include complex spatio-temporal relationships. Figure 6 (a) shows multiple sub-actions retrieval example; the sentence describes two objects (“man” and “crab”), and two actions (“cutting open a crab” and “taking the meat out”) in a short-term segment, which requires accurate spatio-temporal grounding. Comparing BiC-Net to its variant  $\mathbf{VG}$ , our model successfully retrieves the correct video, which contains all spatio-temporal relationships and entities described in the sentence. The second video only contains “taking the meat out from carb” actions. The third and fourth videos only involve a “cutting” action and similar objects (e.g., “man” and “knife”). The fifth video also only contains an action (“taking the meat”). In the left example, the  $\mathbf{VG}$  model also retrieves similar scenes (e.g., similar man and cutting action) in the video. However, we observe that videos involving related elements are only ranked as the true positive in the top-3 positions. The performance of  $\mathbf{VG}$  indicates that removing the fine-grained spatio-temporal relationships hurts the expressiveness

**Query:** a man cutting open a crab and taking the meat out to prepare a food dish.



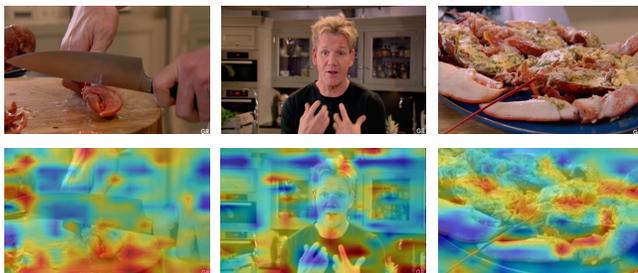
(a)

**Query:** a woman applies eye shadow to her right eye with a makeup brush.



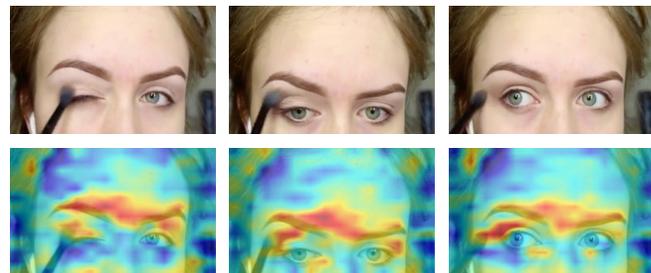
(b)

**Fig. 6: Qualitative examples of the text-video tasks:** In (a), (b), we show retrieval ranks of BiC-Net and the variant VG on MSR-VTT dataset test set [2]. Given a textual description as a query, we retrieve the most relevant video ranked from top to bottom. True positives are bounded in green boxes.



**Caption:** a man cutting open a crab and taking the meat out to prepare a food dish.

(a)



**Caption:** a woman applies eye shadow to her right eye with a makeup brush.

(b)

**Fig. 7: Visualization of attention map on sample clips from the MSR-VTT. The top row presents original frames, and the bottom presents corresponding attention maps.**

of the video representation and further degrades the retrieval performance. Another example is showed in Figure 6 (b), which requires fine-grained spatio-temporal relation grounding. The positive example contains a scenario involving two objects (“woman”, and “makeup brush”) and a fine-grained action (“applies eye shadow”). Comparing these results, we observe that the variant VG retrieves a list of similar action videos, which cannot capture the fine-grained action (“applies eye shadow to her right eye”) in the video. Our model not only identifies the relevant objects “woman” and “makeup brush” but also captures the fine-grained relations between them. Again, it verifies the effectiveness of introducing spatio-temporal relation features to distinguish videos with the same visual components but with different relations.

### G. Visualization Results

To intuitively observe the effectiveness of introducing spatio-temporal relations, we visualize the attention map to infer the value of the spatio-temporal relation features. We select 2 videos, including two positive example in text-to-video retrieval from MSR-VTT. In Figure 7, we show the original frames and attention maps. As can be seen, our BiC-Net learns

to value core parts with intense semantic relations such as “man + crab” in “cutting open a crab and taking the meat out”, “woman + makeup brush” in “applies eye shadow to her right eye”. Furthermore, we find that the salient regions (e.g., man, crab, woman’s right eye) are highlighted separately in Figure 7. This also verifies that our model can learn fine-grained relational information with the corresponding text sentences.

## V. CONCLUSIONS

This work contributes to a novel modeling method for cross-modal text-video retrieval. We claim that video representation should learn not only from global features but also from local spatio-temporal relationships. To fulfill this target, we design the Bi-Branch Complementary Network (BiC-Net) to capture local relational and global visual information for modeling comprehensively. Extensive experimental results on three benchmarks have demonstrated the effectiveness and superiority of our proposed method. Besides, we still face an inherent computational burden of attention in processing long-length video with more complex local relations. Therefore, we leave computational optimization of the multi-layer spatio-temporal transformer as future works.

## REFERENCES

- [1] X. Song, J. Chen, Z. Wu, and Y.-G. Jiang, "Spatial-temporal graphs for cross-modal text2video retrieval," *IEEE Transactions on Multimedia*, 2021.
- [2] A. Miech, D. Zhukov, J. Alayrac, M. Tapaswi, I. Laptev, and J. Sivic, "Howto100m: Learning a text-video embedding by watching hundred million narrated video clips," in *Proceedings of the IEEE/CVF International Conference on Computer Vision, (ICCV)*, 2019, pp. 2630–2640.
- [3] Z. Feng, Z. Zeng, C. Guo, and Z. Li, "Exploiting visual semantic reasoning for video-text retrieval," in *International Joint Conference on Artificial Intelligence, (IJCAI)*, 2020, pp. 1005–1011.
- [4] S. Chen, Y. Zhao, Q. Jin, and Q. Wu, "Fine-grained video-text retrieval with hierarchical graph reasoning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2020, pp. 10 635–10 644.
- [5] P. Wu, X. He, M. Tang, Y. Lv, and J. Liu, "Hanet: Hierarchical alignment networks for video-text retrieval," in *Proceedings of the ACM international conference on Multimedia*, 2021.
- [6] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *Advances in neural information processing systems*, vol. 28, 2015.
- [7] X. Yang, J. Dong, Y. Cao, X. Wang, M. Wang, and T.-S. Chua, "Tree-augmented cross-modal encoding for complex-query video retrieval," in *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*, 2020, pp. 1339–1348.
- [8] J. Dong, X. Li, C. Xu, X. Yang, G. Yang, X. Wang, and M. Wang, "Dual encoding for video retrieval by text," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [9] J. Dong, X. Li, and C. G. Snoek, "Predicting visual features from text for image and video caption retrieval," *IEEE Transactions on Multimedia*, vol. 20, no. 12, pp. 3377–3388, 2018.
- [10] Y. Liu, S. Albanie, A. Nagrani, and A. Zisserman, "Use what you have: Video retrieval using representations from collaborative experts," in *British Machine Vision Conference*, 2019, p. 279.
- [11] X. Li, F. Zhou, C. Xu, J. Ji, and G. Yang, "Sea: Sentence encoder assembly for video retrieval by textual queries," *IEEE Transactions on Multimedia*, vol. 23, pp. 4351–4362, 2020.
- [12] A. Miech, I. Laptev, and J. Sivic, "Learning a text-video embedding from incomplete and heterogeneous data," *arXiv preprint arXiv:1804.02516*, 2018.
- [13] J. Dong, X. Li, C. Xu, S. Ji, Y. He, G. Yang, and X. Wang, "Dual encoding for zero-example video retrieval," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2019, pp. 9346–9355.
- [14] M. Wray, D. Larlus, G. Csurlus, and D. Damen, "Fine-grained action retrieval through multiple parts-of-speech embeddings," in *Proceedings of the IEEE/CVF International Conference on Computer Vision, (ICCV)*, 2019, pp. 450–459.
- [15] N. Han, J. Chen, G. Xiao, H. Zhang, Y. Zeng, and H. Chen, "Fine-grained cross-modal alignment network for text-video retrieval," in *Proceedings of the ACM international conference on Multimedia*, 2021, pp. 3826–3834.
- [16] N. Han, J. Chen, H. Zhang, H. Wang, and H. Chen, "Adversarial multi-grained embedding network for cross-modal text-video retrieval," *ACM Transactions on Multimedia Computing, Communications, and Applications*, vol. 18, no. 2, pp. 1–23, 2022.
- [17] N. C. Mithun, J. Li, F. Metzke, and A. K. Roy-Chowdhury, "Learning joint embedding with multimodal cues for cross-modal video-text retrieval," in *Proceedings of the 2018 ACM on International Conference on Multimedia Retrieval*, 2018, pp. 19–27.
- [18] V. Gabeur, C. Sun, K. Alahari, and C. Schmid, "Multi-modal transformer for video retrieval," in *Proceedings of the European Conference on Computer Vision, (ECCV)*, vol. 5, 2020.
- [19] X. Wang, L. Zhu, and Y. Yang, "T2vлад: global-local sequence alignment for text-video retrieval," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2021, pp. 5079–5088.
- [20] A. Miech, J.-B. Alayrac, L. Smaira, I. Laptev, J. Sivic, and A. Zisserman, "End-to-end learning of visual representations from uncurated instructional videos," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2020, pp. 9876–9886.
- [21] A. Rouditchenko, A. Boggust, D. Harwath, D. Joshi, S. Thomas, K. Audhkhasi, R. Feris, B. Kingsbury, M. Picheny, A. Torralba *et al.*, "Avlnet: Learning audio-visual language representations from instructional videos," *arXiv preprint arXiv:2006.09199*, 2020.
- [22] M. Patrick, P.-Y. Huang, Y. Asano, F. Metzke, A. G. Hauptmann, J. F. Henriques, and A. Vedaldi, "Support-set bottlenecks for video-text representation learning," in *International Conference on Learning Representations, (ICLR)*, 2021.
- [23] H. Luo, L. Ji, M. Zhong, Y. Chen, W. Lei, N. Duan, and T. Li, "Clip4clip: An empirical study of clip for end to end video clip retrieval," *arXiv preprint arXiv:2104.08860*, 2021.
- [24] J. Lei, L. Li, L. Zhou, Z. Gan, T. L. Berg, M. Bansal, and J. Liu, "Less is more: Clipbert for video-and-language learning via sparse sampling," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2021, pp. 7331–7341.
- [25] H. Fang, P. Xiong, L. Xu, and Y. Chen, "Clip2video: Mastering video-text retrieval via image clip," *arXiv preprint arXiv:2106.11097*, 2021.
- [26] S. Liu, H. Fan, S. Qian, Y. Chen, W. Ding, and Z. Wang, "Hit: Hierarchical transformer with momentum contrast for video-text retrieval," in *Proceedings of the IEEE/CVF International Conference on Computer Vision, (CVPR)*, 2021, pp. 11 915–11 925.
- [27] M. Bain, A. Nagrani, G. Varol, and A. Zisserman, "Frozen in time: A joint video and image encoder for end-to-end retrieval," in *Proceedings of the IEEE/CVF International Conference on Computer Vision, (ICCV)*, 2021, pp. 1728–1738.
- [28] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International Conference on Machine Learning, (ICML)*, 2021, pp. 8748–8763.
- [29] W. Lu, D. Li, L. Nie, P. Jing, and Y. Su, "Learning dual low-rank representation for multi-label micro-video classification," *IEEE Transactions on Multimedia*, 2021.
- [30] Y. Zhang, W. Min, L. Nie, and S. Jiang, "Hybrid-attention enhanced two-stream fusion network for video venue prediction," *IEEE Transactions on Multimedia*, vol. 23, pp. 2917–2929, 2020.
- [31] M. Qi, J. Qin, Y. Yang, Y. Wang, and J. Luo, "Semantics-aware spatial-temporal binaries for cross-modal video retrieval," *IEEE Transactions on Image Processing*, vol. 30, pp. 2989–3004, 2021.
- [32] W. Wang, J. Gao, X. Yang, and C. Xu, "Many hands make light work: Transferring knowledge from auxiliary tasks for video-text retrieval," *IEEE Transactions on Multimedia*, 2022.
- [33] X. Qian, Y. Zhuang, Y. Li, S. Xiao, S. Pu, and J. Xiao, "Video relation detection with spatio-temporal graph," in *Proceedings of the ACM international conference on Multimedia*, 2019, pp. 84–93.
- [34] J. Xiao, X. Shang, X. Yang, S. Tang, and T.-S. Chua, "Visual relation grounding in videos," in *European conference on computer vision, (ECCV)*, 2020, pp. 447–464.
- [35] Y. Li, X. Yang, X. Shang, and T.-S. Chua, "Interventional video relation detection," in *Proceedings of the 29th ACM International Conference on Multimedia*, 2021, pp. 4091–4099.
- [36] X. Wang and A. Gupta, "Videos as space-time region graphs," in *Proceedings of the European Conference on Computer Vision, (ECCV)*, 2018, pp. 399–417.
- [37] S. Yan, Y. Xiong, and D. Lin, "Spatial temporal graph convolutional networks for skeleton-based action recognition," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.
- [38] J. Wu, L. Wang, L. Wang, J. Guo, and G. Wu, "Learning actor relation graphs for group activity recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2019, pp. 9964–9974.
- [39] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *International Conference on Learning Representations, (ICLR)*, 2017.
- [40] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances on Neural Information Processing Systems, (NeurIPS)*, 2017, pp. 6000–6010.
- [41] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "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 Computational Linguistics: Human Language Technologies*, 2019, p. 4171–4186.
- [42] S. Geng, M. Zolfaghari, H. Pirsiavash, and T. Brox, "Coot: Cooperative hierarchical transformer for video-text representation learning," in *Advances on Neural Information Processing Systems, (NeurIPS)*, 2020.
- [43] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang, "Bottom-up and top-down attention for image captioning and visual question answering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2018, pp. 6077–6086.

- [44] D. Hendrycks and K. Gimpel, “Gaussian error linear units (gelus),” *arXiv preprint arXiv:1606.08415*, 2016.
- [45] J. L. Ba, J. R. Kiros, and G. E. Hinton, “Layer normalization,” *arXiv preprint arXiv:1607.06450*, 2016.
- [46] J. Xu, T. Mei, T. Yao, and Y. Rui, “Msr-vtt: A large video description dataset for bridging video and language,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2016, pp. 5288–5296.
- [47] Y. Yu, J. Kim, and G. Kim, “A joint sequence fusion model for video question answering and retrieval,” in *Proceedings of the European Conference on Computer Vision, (ECCV)*, 2018, pp. 487–503.
- [48] D. Chen and W. B. Dolan, “Collecting highly parallel data for paraphrase evaluation,” in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, 2011, pp. 190–200.
- [49] S. Venugopalan, M. Rohrbach, J. Donahue, R. Mooney, T. Darrell, and K. Saenko, “Sequence to sequence-video to text,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision, (ICCV)*, 2015, pp. 4534–4542.
- [50] L. Zhou, C. Xu, and J. J. Corso, “Towards automatic learning of procedures from web instructional videos,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.
- [51] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, “Imagenet large scale visual recognition challenge,” *International journal of computer vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [52] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. Alemi, “Inception-v4, inception-resnet and the impact of residual connections on learning,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 31, no. 1, 2017.
- [53] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev *et al.*, “The kinetics human action video dataset,” *arXiv preprint arXiv:1705.06950*, 2017.
- [54] J. Carreira and A. Zisserman, “Quo vadis, action recognition? a new model and the kinetics dataset,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2017, pp. 6299–6308.
- [55] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *International Conference on Learning Representations, (ICLR)*, 2015.
- [56] G. Bertasius, H. Wang, and L. Torresani, “Is space-time attention all you need for video understanding,” *arXiv preprint arXiv:2102.05095*, vol. 2, no. 3, p. 4, 2021.