

# A Survey on Video Moment Localization\*

MENG LIU, Shandong Jianzhu University, China

LIQIANG NIE, Harbin Institute of Technology (Shenzhen), China

YUNXIAO WANG, Shandong University, China

MENG WANG, Hefei University of Technology, China

YONG RUI, Lenovo Company Ltd., China

Video moment localization, also known as video moment retrieval, aiming to search a target segment within a video described by a given natural language query. Beyond the task of temporal action localization whereby the target actions are pre-defined, video moment retrieval can query arbitrary complex activities. In this survey paper, we aim to present a comprehensive review of existing video moment localization techniques, including supervised, weakly supervised, and unsupervised ones. We also review the datasets available for video moment localization and group results of related work. In addition, we discuss promising future directions for this field, in particular large-scale datasets and interpretable video moment localization models.

CCS Concepts: • **Information systems** → **Specialized information retrieval**; **Video search**; **Novelty in information retrieval**.

Additional Key Words and Phrases: Video Moment Localization, Video Moment Retrieval, Vision and Language, Cross-modal Retrieval, Survey

## ACM Reference Format:

Meng Liu, Liqiang Nie, Yunxiao Wang, Meng Wang, and Yong Rui. 2023. A Survey on Video Moment Localization. *ACM Comput. Surv.* 1, 1, Article 188 (June 2023), 36 pages. <https://doi.org/10.1145/3556537>

## 1 INTRODUCTION

With the increasing prevalence of digital cameras and social networks, plenty of videos are recorded, stored, and shared daily [48], especially the surveillance videos. Such large-scale video data prompts video content analysis to become increasingly essential [49]. Compared with static images, one critical characteristic of videos is that they could depict the behavior evolution of a certain object over time. Therefore, understanding and recognizing actions has become the fundamental task to effectively analyze videos. In the past decade, a large body of literature has paid attention to action recognition in videos [85]. And the vast majority of existing studies tackle the action recognition problem as the video classification task [115], where the trimmed videos contain a single action (as shown in Fig. 1(a)). However, in real-world applications, such as surveillance [120], robotics [4], and autonomous driving [14], cameras continuously record video streams. In other words, videos,

\*Several mistakes are corrected from the published version. Liqiang Nie is the corresponding author.

Authors' addresses: Meng Liu, [mengliu.sdu@gmail.com](mailto:mengliu.sdu@gmail.com), Shandong Jianzhu University, China; Liqiang Nie, [nieliqiang@gmail.com](mailto:nieliqiang@gmail.com), Harbin Institute of Technology (Shenzhen), China; Yunxiao Wang, [yunxiao.wang@mail.sdu.edu.cn](mailto:yunxiao.wang@mail.sdu.edu.cn), Shandong University, China; Meng Wang, [eric.mengwang@gmail.com](mailto:eric.mengwang@gmail.com), Hefei University of Technology, China; Yong Rui, [yongrui@lenovo.com](mailto:yongrui@lenovo.com), Lenovo Company Ltd. China.

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

© 2023 Association for Computing Machinery.

0360-0300/2023/6-ART188 \$15.00

<https://doi.org/10.1145/3556537>

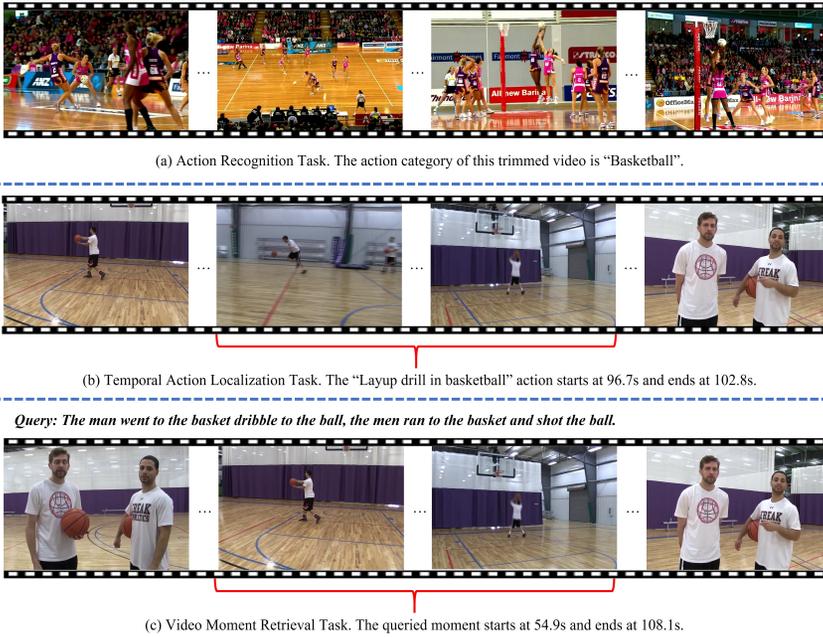


Fig. 1. Illustration of the action recognition, temporal action localization, and video moment localization (retrieval) tasks. (a) The action recognition task aims to classify trimmed videos into different action categories. (b) The goal of temporal action localization is to automatically recognize the interval of the action occurs from the given untrimmed video and judge the category of the detected action. And (c) the purpose of the video moment localization is to localize the temporal moment referenced in the query.

in reality, are mostly untrimmed. Therefore, developing algorithms that simultaneously decide both the temporal intervals and the action categories as they occur is indispensable.

Inspired by this, temporal action localization (as illustrated in Fig. 1(b)), which jointly classifies action instances and localizes them in an untrimmed video, becomes a vital task in video understanding [102]. But existing temporal action localization methods are restricted to a pre-defined list of actions. They are hence inflexible since activities in the wild are even more complex, like "The man went to the basket dribble to the ball, the men ran to the basket and shot the ball". Considering the natural language expression could vary according to the content of videos, it is more natural to utilize the natural language query to localize the desired moment from the given video, as shown in Fig. 1(c). Accordingly, video moment localization, aiming to identify the specific start and end timestamps of the moment in response to the given query, has been a hot research topic recently.

While this task opens up great opportunities to better video understanding, it is substantially more challenging due to the following reasons: 1) The given queries can be arbitrarily complex natural descriptions. Considering the query "The black cat jumps back up after falling" as an example, it describes the temporal relationship between "jump back up" and "fall" actions in a video. Moreover, the language sequence is misaligned with video sequence. Therefore, how to well comprehend the complex query information is one barrier for video moment localization. 2) Different moments in the untrimmed video have varying durations and diverse spatial-temporal characteristics. In addition, a long video often contains multiple moments of interest. For example, given a typical query "A Ferris wheel second comes into view", the model requires to find the second occurrence of "Ferris wheel". Thereby, how to effectively distinguish relevant video moments and precisely localize the moment of interest is very difficult. 3) As a multimodal task, how to model

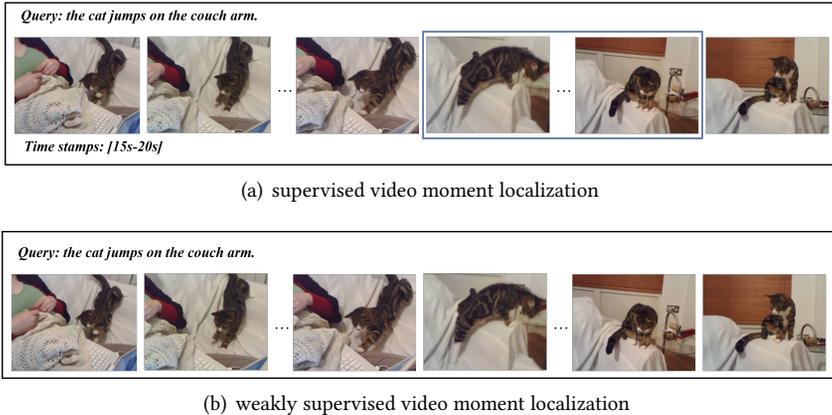


Fig. 2. Visualization of the difference between supervised and weakly supervised video moment localization.

the complex correlations between videos and queries is crucial. And 4) how to boost the efficiency of the localization method is an important issue that matters practical use.

Recently, deep learning techniques have emerged as powerful methods for various tasks, such as video captioning [106] and text-video retrieval [58]. This motivates many researchers to resort to deep learning approaches to tackle the video moment localization task. Particularly, existing work can be divided into three categories: 1) deep supervised learning approaches, 2) deep weakly supervised learning ones, and 3) deep unsupervised learning ones. The main difference between the supervised and weakly supervised learning methods is the availability of the exact temporal annotation of the moment corresponding to the given query, as shown in Fig. 2. To be specific, the former one requires the full annotations of temporal boundaries for each video and the given query, such as video segment of (15s,20s) in Fig. 2(a). In contrast, the latter one merely needs video-sentence pairs for training, and they do not require the temporal boundary annotation information. Different from the aforementioned two categories, unsupervised video moment localization methods require no textual annotations of videos, namely they only require easily available text corpora and a collection of videos to localize.

To give a comprehensive overview of this field, including models, datasets, and future directions, we summarize the work on video moment localization before Dec. 2021 and present this survey. Specifically, to perform a deeper analysis of existing approaches, we establish the fine-grained taxonomy of each category based on their architectures, moment generation strategies, and key characteristics. For instance, for deep supervised learning-based approaches (as summarized in Table 1), we first group them into three categories: two-stage, one-stage, and reinforcement learning methods according to their architectures. Thereinto, two-stage methods commonly adopt separate schemes (e.g., the sliding windows) to generate moment candidates and then match them with the query sentences to obtain the target moment. One-stage ones integrate the moment generation and moment localization into a unified framework, while reinforcement learning paradigms formulate the video moment localization task as a sequential decision making problem. Afterwards, according to the strategies of moment generation, we respectively divide two-stage, one-stage, as well as reinforcement learning methods into different subclasses, and each subclass is further divided into several parts according to the characteristics of corresponding methods. Note that similar fine-grained taxonomies are established for weakly supervised and unsupervised learning methods, as reported in Table 2.

Compared to previous surveys on video moment localization [53][105], the contributions of this survey are as follows:

Table 1. Summarization of supervised video moment localization approaches.

Architecture	Moment Generation	Method	Published	Year	Key Aspect
Two-stage	Hand-crafted Heuristics	MCN [2]	ICCV	2017	Pioneer Work
		MLLC [25]	EMNLP	2018	Visual Modeling
		TCMN [119]	ACM MM	2019	Query Modeling
	Multi-scale Sliding Windows	CTRL [16]	ICCV	2017	Pioneer Work
		ACRN [50]	SIGIR	2018	Visual Modeling
		VAL [76]	PCM	2018	
		SLTA [29]	ICMR	2019	
		MMRG [110]	CVPR	2021	
		ROLE [51]	ACM MM	2018	Query Modeling
		MCF [93]	IJCAI	2018	Inter-modal Interaction Modeling
	ACL [21]	WACV	2019		
	Moment Generation Networks	MIGCN [121]	IEEE TIP	2021	Inter- and Intra-modal Interaction Modeling
QSPN [101]		AAAI	2019	Query-guided Proposal Generation	
SAP [10]		AAAI	2019	Visual Concept-based Proposal Generation	
One-stage	Anchor-based	BPNet [98]	AAAI	2021	VSLNet [114] based Proposal Generation
		TGN [7]	EMNLP	2018	Pioneer Work
		BSSTL [32]	ICIP	2019	Inter-modal Interaction Modeling
		RMN [43]	ACL COLING	2020	
		FIAN [66]	ACM MM	2020	
		I2N [63]	IEEE TIP	2021	
		CMin [122]	SIGIR	2019	Inter- and Intra-modal Interaction Modeling
		CSMGAN [46]	ACM MM	2020	
	PMI-LOC [9]	ECCV	2020		
	CMin-R [41]	IEEE TIP	2020		
	CBP [88]	AAAI	2020	Localization Module	
	Temporal Convolution	SCDM [107]	NIPS	2019	Semantic Modulated Temporal Convolution
		MAN [111]	CVPR	2019	Hierarchical Convolutional Network
		CLEAR [28]	IEEE TIP	2021	Bi-directional Temporal Convolution
		IA-Net [47]	EMNLP	2021	Inter- and Intra-modal Interaction Modeling
	Segment-tree	CMHN [27]	IEEE TIP	2021	Efficiency
		CPN [124]	CVPR	2021	Segment-tree
	Proposal Generation	APGN [44]	EMNLP	2021	Adaptive Proposal Generation Network
		Lp-Net [97]	EMNLP	2021	Learnable Proposal Generation Network
	Enumeration	TMN [42]	ECCV	2018	Pioneer Work
		2D-TAN [118]	AAAI	2020	2D Temporal Map
		MS-2D-TAN [117]	IEEE TPAMI	2021	Multi-scale 2D Temporal Map
		DPIN [86]	ACM MM	2020	Interaction Modeling
		MATN [116]	CVPR	2021	
SMIN [87]		CVPR	2021		
SV-VMR [96]		ICME	2021		
RaNet [17]		EMNLP	2021		
DCM [104]		SIGIR	2021	Location bias	
FVMR [18]		ICCV	2021	Efficiency	
Proposal-free	LNet [8]	AAAI	2019	Pioneer Work	
	ABLR [108]	AAAI	2019	Inter-modal Interaction Modeling	
	ExCL [22]	NAACL-HLT	2019		
	SQAN [60]	CVPR	2020		
	PFGA [64]	WACV	2020		
	VSLNet [114]	ACL	2020		
	VSLNet-L [113]	IEEE TPAMI	2021		
	ACRM [82]	IEEE TMM	2021		
	MIM [37]	ICMR	2021		
	SSMN [54]	ACM TOMM	2021		
	HVTG [11]	ECCV	2020		Inter- and Intra-modal Interaction Modeling
	MQEI [83]	IEE T-ITS	2021		
	CBLN [45]	CVPR	2021	Visual Modeling	
	CP-Net [36]	AAAI	2021		
DORi [70]	WACV	2021			
DEBUG [55]	EMNLP	2019	Imbalance Problem		
DRN [109]	CVPR	2020	Dense Supervision		
IVG [62]	CVPR	2021	Causal Interventions		
DRFT [12]	NIPS	2021	Multi-modal Information		
Reinforcement	Proposal-free	RWM [24]	AAAI	2019	Pioneer Work
		TSP-PRL [95]	AAAI	2020	Policy
		SM-RL [90]	CVPR	2019	Visual Modeling
		STRONG [5]	ACM MM	2020	
		TripNet [23]	BMVC	2020	
		AVMR [6]	ACM MM	2020	Reward
		MABAN [78]	TIP	2021	Inter-modal Interaction Modeling

- Existing surveys pay more attention to analyzing supervised learning based localization methods, disregarding the analysis of weakly supervised and unsupervised ones. Differently, we provide a comprehensive survey of video moment localization approaches, including supervised, weakly supervised, and unsupervised learning based ones.

Table 2. Summarization of weakly-supervised and unsupervised video moment localization approaches.

Paradigm	Architecture	Moment Generation	Method	Published	Year	Key Aspect
Weakly Supervised	Two-stage	Multi-scale Sliding Windows	TGA [59]	CVPR	2019	Pioneer Work
			WSSLN [20]	EMNLP	2019	Proposal Selection
			VLANet [56]	ECCV	2020	
	One-stage	Anchor-based	LoGAN [81]	WACV	2021	Visual Modeling
			SCN [40]	AAAI	2020	Scoring Refinement
		VCA [92]	ACM MM	2021	Visual Modeling	
		Enumeration	RTBPN [123]	ACM MM	2020	Intra-sample Confrontment
			LCNet [103]	TIP	2021	Inter-modal Interaction Modeling
			WSTAN [91]	TMM	2021	Visual Modeling
	MS-2D-RL [91]	ICPR	2021			
Reinforcement	Proposal-free	BAR [94]	ACM MM	2020	Boundary Refinement	
Unsupervised	One-stage	Enumeration	U-VMR [19]	IEEE TCSVT	2021	Knowledge Distillation
		Proposal-free	PSVL [61]	ICCV	2021	Pseudo-supervision Generation

- We establish a holistic and fine-grained taxonomy for existing approaches according to their architectures, moment generation strategies, and key characteristics. This is advantageous to highlight the major strengths and shortcomings of existing methods. Nevertheless, other literature reviews mainly focus on organizing the typical approaches according to the coarse-grained taxonomy, such as localization policies [53] or moment generation schemes [105].
- Our survey covers more papers on the video moment localization topic. To be specific, we summarize the work on video moment localization before Dec. 2021, while existing surveys mainly summarize the work published in 2019 and 2020.
- We introduce the methods of different categories in detail and conduct comparison among them. Nevertheless, these two surveys merely elaborate the typical methods of each category, as well as lack the deep analysis of other approaches.

The remainder of this paper is organized as follows. Section 2 briefly introduces some related research area. Section 3, 4, and 5 review supervised learning, weakly supervised learning, and unsupervised learning approaches on video moment localization, respectively. Section 6 describes the currently available datasets and evaluation metrics used for video moment localization. We separately analyze experimental results of existing approaches and discuss possible future research directions in Section 7 and 8. We conclude the work in Section 9.

## 2 RELATED RESEARCH AREA

There are several research fields closely related to video moment localization, which we now briefly describe.

### 2.1 Temporal Action Localization

Temporal action localization [77][65], aiming at jointly classifying a pre-defined list of action instances and localizing the timestamps of them in an untrimmed video, is most relevant to the video moment localization task. The main difference between them is that video moment localization aims to localize the specific moment within the untrimmed video via the given natural language query. Apparently, these two tasks encounter some common challenges, e.g., temporal proposal generation and temporal proposal relation modeling. Therefore, early video moment localization approaches directly apply temporal action localization techniques by merely replacing the action classification with cross-modal matching module. Recently, some video moment localization systems explore different grained interactions between the query and the moment as well as the query modeling to further increase both the accuracy and efficiency. Still, temporal action localization techniques will continue to serve as the foundation for the advance of its counterpart in the video moment localization.

## 2.2 Video Object Grounding

Existing video object grounding task can be broken down into the following categories: 1) localizing individual objects mentioned in the language query [125]; 2) localizing a spatio-temporal tube of the target object based on the given language query [13] or person tracklet (i.e., video re-identification [35][34]); and 3) only localizing the referred objects in the language query[73]. The main difference between the first and third category is that the former treats each query word independently and does not distinguish between different instances of the same object, while the latter requires additional disambiguation using object-object relations in both time and space. Compared with the second subtask that requires both spatial and temporal localization, i.e., localizing a sequence of bounding boxes of the queried object, video moment localization merely focuses on determining the temporal boundaries of events corresponding to the given sentence. There are many key technical challenges, such as cross-modality modeling and fine-grained reasoning, that are shared between video object grounding and video moment localization.

## 2.3 Video Corpus Moment Retrieval

Video corpus moment retrieval aims to retrieve a short fraction in a video that semantically corresponds to a text query, from a corpus of untrimmed and unsegmented videos [112]. Differently, video moment localization is to retrieve a short fraction from a single given video that corresponds to a text query. In fact, the task of video corpus moment retrieval was extended from video moment localization by [15], which better matches real-world application scenarios. However, video corpus moment retrieval imposes the additional requirement to distinguish moments from different videos as compared to the task of video moment localization.

## 3 SUPERVISED VIDEO MOMENT LOCALIZATION

**Problem Formulation.** We denote a video as  $V = \{f_t\}_{t=1}^T$ , where  $T$  is the frame number of the video. Each video is associated with temporal annotations  $A = \{(s_j, \tau_j^s, \tau_j^e)\}_{j=1}^M$ , where  $M$  is the annotation number of the video  $V$  and  $s_j$  is a sentence query with respect to a video moment that has  $\tau_j^s$  and  $\tau_j^e$  as start and end time points in the video. In the supervised video moment localization setting, the training data are the annotated query and video pairs (i.e.,  $(V, A)$ ). During the inference, given a video and a natural language query, the trained model would output the location information  $(\tau^s, \tau^e)$  of the target moment regarding a given query.

We summarize existing supervised video moment localization approaches in Table 1. Specifically, they can be roughly grouped into three categories: two-stage, one-stage, and reinforcement learning-based approaches. In what follows, we detail them accordingly.

### 3.1 Two-stage Methods

The two-stage supervised video moment localization approaches commonly utilize a separate scheme (e.g., the sliding windows) to generate moment candidates, and then match them with the query to find the target moment. A general diagram of two-stage video moment localization methods is shown in Fig. 3. In particular, according to the manner of generating moment candidates, existing two-stage methods can be divided into three groups: hand-crafted heuristics, multi-scale sliding windows, and moment generation module based models [101].

**3.1.1 Hand-crafted heuristics.** The work of Hendricks et al. [2] is generally regarded as the pioneer in this branch. It treats the video moment localization task as the moment retrieval problem and presents Moment Context Network (MCN)<sup>1</sup>. To train and evaluate the proposed model, [2]

<sup>1</sup><https://github.com/LisaAnne/LocalizingMoments>.

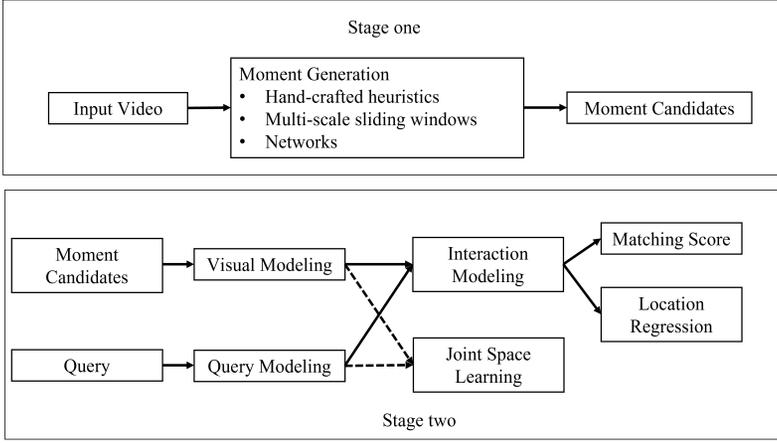


Fig. 3. The schema of the two-stage methods in the whole paper.

collects the Distinct Describable Moments (DiDeMo) dataset, which consists of over 40,000 pairs of natural language queries and localized moments in unedited videos. According to the characteristic of this dataset, [2] generates moment candidates via hand-crafted heuristics. Concretely, a video is first segmented into six five-second clips, and then any contiguous set of clips are utilized to construct a possible moment. Having obtained 21 moment candidates, it learns a shared embedding space for both video temporal context features and natural language queries, adopting the following inter-intra video ranking loss,

$$\begin{cases} L(\theta) = \lambda \sum_i L_i^{intra}(\theta) + (1 - \lambda) \sum_i L_i^{inter}(\theta), \\ L_i^{intra}(\theta) = \sum_{m_j \neq m_i} L^R(D_\theta(s_i, m_i), D_\theta(s_i, m_j)), \\ L_i^{inter}(\theta) = \sum_{j \neq i} L^R(D_\theta(s_i, m_i), D_\theta(s_i, m_j)), \end{cases} \quad (1)$$

where  $L_i^{intra}(\theta)$  is the intra-video ranking loss that encourages queries to be closer to the corresponding video moment than all other possible moments from the same video,  $L_i^{inter}(\theta)$  is the inter-video ranking loss that encourages queries to be closer to corresponding video moments than moments with the same endpoints outside the video,  $L^R(x, y) = \max(0, x - y + b)$  is the ranking loss,  $D$  is the square distance metric,  $\theta$  refers to the learnable parameters of the model, and  $\lambda$  is the weighting parameter. Building on top of MCN, several researchers promoted this branch from different angles, such as visual and query modeling. In the following, we would detail the corresponding studies in sequence.

**Visual Modeling.** [2] considers the entire video as the visual context, which may induce inappropriate information in the context feature, influencing the localization performance. Inspired by this, Hendricks et al. [25] proposed the Moment Localization with Latent Context (MLLC) network<sup>2</sup>, which models visual context as a latent variable. Therefore, MLLC could attend to different video contexts conditioned on the specific query-video pair, overcoming the limitation of query-independent contextual reasoning. Besides, [25] builds the temporal reasoning in video and language dataset (TEMPO) based on the DiDeMo dataset [2], which comprises two parts: TEMPO-TL (Template Language) and TEMPO-HL (Human Language). The former is constructed by the original DiDeMo sentences with language templates, while the latter is built with real videos and newly collected user-provided temporal annotations.

<sup>2</sup><https://people.eecs.berkeley.edu/>.

**Query Modeling.** A common limitation of previous approaches is that they utilize a single textual embedding to represent the query, which may not be precise. To fully exploit temporal dependencies between events within the query, Zhang et al. [119] proposed a Temporal Compositional Modular Network (TCMN)<sup>3</sup>, which leverages a tree attention network to parse the given query into three components concerning the main event, context event, and temporal signal. Besides, a temporal relationship module is designed to measure the similarity between the phrase embedding for the temporal signal and the location feature, and a temporal localization module is designed to measure the similarity between the event-related phrase embedding and the visual feature. The overall matching score between each moment candidate and the given query is calculated by the late fusion strategy.

**3.1.2 Multi-scale sliding windows.** The work of Gao et al. [16] is the pioneer in this branch. They proposed a Cross-modal Temporal Regression Localizer (CTRL) to tackle the video moment localization task<sup>4</sup>. To collect moment candidates for training, CTRL adopts multi-scale temporal sliding windows. In the framework of CTRL, the pre-trained C3D model [84] is first utilized to extract visual features for moment candidates. Meanwhile, a Long Short-term Memory (LSTM) [26] network or Skip-thought [30] is adopted to extract query representations. Afterwards, a cross-modal processing module is designed to jointly model the query and visual features, which calculates element-wise addition, multiplication, and direct concatenation. Finally, the multi-layer perception network (MLP) is designed for visual semantic alignment and moment location regression. To jointly train for visual-semantic alignment and moment location regression, CTRL utilizes a multi-task loss as follows,

$$\begin{cases} L = L_{aln} + \alpha L_{reg}, \\ L_{aln} = \frac{1}{N} \sum_{i=0}^N [\alpha_c \log(1 + \exp(-cs_{i,i})) + \sum_{j=0, j \neq i}^N \alpha_w \log(1 + \exp(cs_{i,j}))], \\ L_{reg} = \frac{1}{N} \sum_{i=0}^N [R(t_{s,i}^* - t_{s,i}) + R(t_{e,i}^* - t_{e,i})], \end{cases} \quad (2)$$

where  $N$  is the batch size,  $cs_{i,j}$  is the matching score between query  $s_j$  and moment candidate  $m_i$ ,  $\alpha_c$  and  $\alpha_w$  are the hyper parameters,  $R(t)$  is the  $L_1$  function,  $t_{s,i}^*$  and  $t_{e,i}^*$  is the ground truth offsets, as well as  $t_{s,i}$  and  $t_{e,i}$  is predicted offsets. More importantly, for evaluation, CTRL adopts the TACoS dataset, and builds a new dataset on top of Charades by adding sentence temporal annotations, called Charades-STA. Hereafter, to further improve the localization performance from different aspects (e.g., visual, query, and interaction modeling), several methods are proposed in the past few years. We will elaborate them in the following paragraphs.

**Visual Modeling.** Although CTRL considers visual context information, it ignores the complex interactions among contexts and fails to identify the importance of each moment, therefore some important cues are missing. Inspired by this, Liu et al. [50] extended the work of [16] and developed an Attentive Cross-Modal Retrieval Network (ACRN)<sup>5</sup>. In particular, they designed a memory attention mechanism to emphasize the visual features mentioned in the query, obtaining the augmented moment representations.

As visual feature extracted from the  $Fc6$  layer of C3D may weaken or ignore the critical visual cues, a Visual-attention Action Localizer (VAL) is introduced in [76], which extracts visual features from the  $CONV5\_3$  layer of C3D. Therefore, it could capture more completed visual information. Similarly, Jiang et al. [29] proposed a Spatial and Language-Temporal Attention (SLTA) method, which takes advantage of object-level local features to enhance the visual representations<sup>6</sup>. To

<sup>3</sup><https://github.com/Sy-Zhang/TCMN-Release>.

<sup>4</sup><https://github.com/jiyanggao/TALL>.

<sup>5</sup><https://sigir2018.wixsite.com/acrn>.

<sup>6</sup><https://github.com/BonnieHuangxin/SLTA>.

be specific, SLTA extracts local features on each frame by Faster R-CNN [69] and introduces the spatial attention to selectively attend to the most relevant local features mentioned in the query. And then it utilizes LSTM to encode the local feature sequence, obtaining the local interaction feature. Finally, SLTA integrates the global motion features extracted by pre-trained C3D with local interaction features as the final visual representations. Recently, to identify the fine-grained differences among similar video moment candidates, Zeng et al. [110] proposed a Multi-Modal Relational Graph (MMRG) framework<sup>7</sup>. It develops a duel-channel relational graph to capture object relations and the phrase relations. Moreover, it considers two self-supervised pre-training tasks: attribute masking and context prediction, to enhance the visual representation and alleviate semantic gap across modalities.

**Query Modeling.** CTRL simply treats queries holistically as one feature vector, which may obfuscate the keywords that have rich temporal and semantic cues. Inspired by this, Liu et al. [51] proposed a cCross-modal mOment Localization nEtwork (ROLE)<sup>8</sup>. It designs a language-temporal attention module, which adaptively reweighs each word's features according to the textual query information and moment context information, therefore deriving useful query representations.

**Interaction Modeling.** Existing scheme CTRL adopts the straightforward multi-modal fusion method, i.e., fusing element-wise addition, multiplication, and concatenation, lacking in-depth analysis. To overcome this drawback, Wu et al. [93] put forward a new multi-modal fusion approach, i.e., Multi-modal Circulant Fusion (MCF)<sup>9</sup>. Particularly, after reshaping feature vectors into circulant matrices, it defines two types of interaction operations between vectors and matrices. The first one is the matrix multiplication between circulant matrix and projection vector, while the second is element-wise product between projection vector and each row vector of circulant matrix. More importantly, the proposed MCF can be integrated into the existing video moment localization models, to further improve the localization accuracy. To adequately exploit rich semantic cues about activities in videos and queries, Ge et al. [21] proposed an Activity Concepts based Localizer (ACL)<sup>10</sup>, which is the first work to mine the activity concepts from both the videos and the sentence queries to facilitate the localization. To explore the interactions between two modalities, ACL separately conducts the multi-modal processing to the pair of activity concepts and the pair of visual features and sentence embeddings. Besides, ACL designs an actionness score generator to calculate the likelihood of the moment candidate containing meaningful activities. Eventually, the alignment score multiplied by the actionness score is set as the final score to predict the alignment confidence between each candidate and query.

Different from aforementioned two methods that focus on inter-modal interaction modeling, Zhang et al. [121] proposed a Multi-modal Interaction Graph Convolutional Network (MIGCN)<sup>11</sup>, which simultaneously explores the complex intra-modal relations and inter-modal interactions residing in the video and sentence query. Particularly, it introduces the multi-modal interaction graph, where two types of nodes are considered (i.e., clip nodes and word nodes) and edges compile both intra-modal relations and inter-modal interactions. Therefore, the graph convolution could refine node representations by jointly considering intra- and inter-modal interaction information.

**3.1.3 Moment generation networks.** Instead of using hand-crafted heuristics or multi-scale sliding windows to generate moment candidates, Xu et al. [101] advanced a Query-guided Segment

<sup>7</sup><https://cvpr-2021.wixsite.com/mmrgr>.

<sup>8</sup><https://acmmm18.wixsite.com/role>.

<sup>9</sup><https://github.com/AmingWu/Multi-modal-Circulant-Fusion>.

<sup>10</sup><https://github.com/runzhouge/MAC>.

<sup>11</sup><https://github.com/zmzhang2000/MIGCN/>.

Proposal Network (QSPN) to generate moment candidates<sup>12</sup>. Concretely, QSPN generalizes the SPN from R-C3D model [100] by introducing query representations as the guidance information to generate moment candidates. And an early fusion retrieval model, instantiated as a two-layer LSTM, is introduced to find the moment that best matches the given query. Besides, a captioning loss is considered to enforce the LSTM to re-generate the query sentence, achieving improved retrieval performance. The combination of retrieval loss and captioning loss is defined as follows,

$$\begin{cases} L = L_{ret} + L_{cap}, \\ L_{ret} = \sum_j \max\{0, \eta + \sigma(s_j, m'_j) - \sigma(s_j, m_j)\}, \\ L_{cap} = -\frac{1}{KT} \sum_{k=1}^K \sum_{t=1}^{T_k} \log P(w_t^k | f(m_k), h_{t-1}^{(2)}, w_1^k, \dots, w_{t-1}^k), \end{cases} \quad (3)$$

where  $\sigma(s_j, m_j)$  is the matching score predicted by the LSTM,  $m'_j$  is the negative moment either comes from the same video or from a different video,  $K$  is the number of queries,  $T_k$  is the number of words in the  $k$ -th query,  $f(m_k)$  represents the visual feature of the moment aligned with the  $k$ -th query. Similarly, Chen et al. [10] proposed a Semantic Activity Proposal (SAP) framework, which also integrates the semantic information in queries into the moment candidate generation process. To be specific, it first trains a visual concept detection CNN with paired query-clip training data. Subsequently, the visual concepts extracted from the query and video frames are utilized to calculate the visual-semantic correlation score for every frame. By grouping frames with high visual-semantic correlation scores, moment candidates could be generated. Different from the above two methods, Xiao et al. [98] proposed a Boundary Proposal Network (BPNet), which utilizes VSLNet [114] to generate several high-quality moment proposals, avoiding redundant candidates.

**Summarization.** In what follows, we summarize the two-stage supervised video moment retrieval methods.

- *Hand-crafted heuristics:* In this branch, the pioneer work MCN focuses on learning the shared embedding space for the moment representation and the query representation. To further improve the performance of MCN, MLLC enforces the model to attend to different video contexts, while TCMN exploits temporal dependencies between events in the query.
- *Multi-scale sliding windows:* To enhance the performance of the pioneer work CTRL in this branch, ACRN, VAL, SLTA, and MMRG focus on the visual modeling. Specifically, ACRN emphasizes the visual features mentioned in the query, VAL captures more completed visual information from the feature map of C3D, SLTA takes advantage of object-level local features, and MMRG considers self-supervised pre-training tasks. Differently, ROLE focuses on query modeling and designs a language-temporal attention module to derive useful query representations. Moreover, MCF, ACL, and MIGCN pay attention to the interaction modeling. Particularly, MCF and ACL are devoted to modeling inter-modal interactions, while MIGCN jointly considers the inter-modal and intra-modal interaction modeling.
- *Moment generation networks:* To generate moment candidates, QSPN generalizes the SPN by introducing query representations as the guidance information, while SAP generates moment candidates by grouping frames with high visual-semantic correlation scores. Differently, BPNet directly utilizes the VSLNet [114] to generate moment proposals.

### 3.2 One-stage Methods

Although two-stage approaches have achieved promising performance, they are suffering from inferior efficiency. For instance, to handle the diverse temporal scales and locations of the moment candidates, exhaustive matching between a large number of overlapping moments and the query is

<sup>12</sup>[https://github.com/VisionLearningGroup/Text-to-Clip\\_Retrieval](https://github.com/VisionLearningGroup/Text-to-Clip_Retrieval).

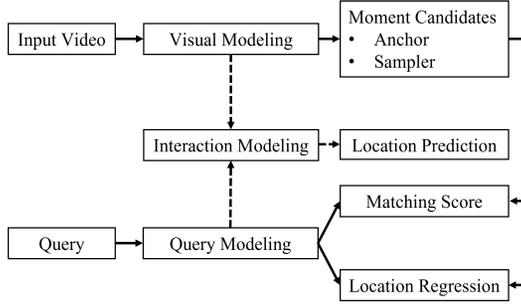


Fig. 4. A block diagram of one-stage video moment localization.

required, which is very computationally expensive. Therefore, developing one-stage methods that localize the golden moment (i.e., the moment that best matches the given query) in one single pass is critical.

As summarized in Table 1, according to the processing of moment candidates, existing one-stage or single-pass video moment localization methods can be grouped into three categories: anchor-based, sampler-based, and proposal-free ones. A block diagram of typical one-stage methods is shown in Fig. 4. In what follows, we review them in sequence.

**3.2.1 Anchor-based models.** Chen et al. [7] proposed the first dynamic single-stream end-to-end Temporal GroundNet (TGN) model for video moment localization<sup>13</sup>. It sets each time step as an anchor and generates  $K$  moment candidates that end at the current time step. Moreover, it utilizes the hidden states generated by the interaction module to yield confidence scores for moments ending at the corresponding time step. The objective of TGN is defined as,

$$L = \sum_{(V,S)} \sum_{t=1}^T \left\{ - \sum_{k=1}^K w_0^k y_t^k \log c_t^k + w_1^k (1 - y_t^k) \log(1 - c_t^k) \right\}, \quad (4)$$

where  $(V, S)$  denotes a training pair from the training set,  $y_t^k$  is interpreted as whether the  $k$ -th moment candidate at time step  $t$  corresponds to the given query,  $c_t^k$  denotes the prediction result,  $w_0^k$  and  $w_1^k$  are calculated according to the frequencies of positive and negative samples in the training set with length  $k\delta$ .

According to the primary contributions of existing anchor-based one-stage approaches, they can be further classified into two categories, i.e., focusing on localization module building and interaction modeling.

**Localization Module.** To yield more precise localization, Wang et al. [88] proposed an end-to-end Contextual Boundary-aware Prediction (CBP) model<sup>14</sup>, which jointly predicts temporal anchors and boundaries at each time step. Its localization module consists of an anchor sub-module as [7] did and a boundary sub-module that is utilized to decide whether the current step is a semantic boundary corresponding to the start/end time of the moment.

**Interaction Modeling.** As simply combining the video and query information is not expressive enough to explore the interaction information, Li et al. [32] proposed a Bidirectional Single-Stream Temporal Localization (BSSTL) model. It introduces an attentive cross-modal fusion model, consisting of dynamic attention weighted query, joint gating, and fusion, to capture the interaction between the queries and videos. Liu et al. [43] designed a deep Rectification-Modulation Network

<sup>13</sup><https://github.com/JaywongWang/TGN>.

<sup>14</sup><https://github.com/JaywongWang/CBP>.

(RMN), which adopts a multi-step reasoning framework to gradually capture the higher-order multi-modal interaction. To be specific, it utilizes multi-step rectification-modulation layers to progressively refine video and query interactions, improving localization accuracy. Different from previous approaches that merely focus on exploiting unidirectional interaction information from the video to query, Qu et al. [66] proposed a Fine-grained Iterative Attention Network (FIAN) to capture bilateral query-video interaction information. Specifically, it leverages a symmetrical iterative attention module to generate video-aware query and query-aware video representations, where two attention branches share the same architecture but opposite input. To model multi-level vision-language cross-modal relation and long-range context, Ning et al. [63] proposed an Interaction-Integrated Network (I2N). To be specific, it introduces an interaction-integrated cell, which integrates both cross-modal and contextual interactions between vision and language. By stacking a few interaction-integrated cells in a chain, the network could explore more detailed vision-language relations on multiple granularities, therefore achieving more accurate semantics understanding and more accurate localization results.

The aforementioned single-pass methods mainly focus on exploring the cross-modal relations between the video and query, ignoring the importance of jointly investigating the cross- and self-modal relations. In light of this, Zhang et al. [122] introduced a Cross-Modal Interaction Network (CMIN)<sup>15</sup>, which jointly considers multiple crucial factors. Specifically, the syntactic graph convolution network is adopted to enhance the query understanding. And the multi-head self-attention module is built to capture long-range semantic dependencies from the video context. Moreover, a multi-stage cross-modal interaction module is presented to exploit the potential relations of video and query contents. Recently, inspired by the fact that captioning can improve the performance of image-based grounding of textual phrases [71], Zhang et al. further improved the model CMIN by integrating a query reconstruction module [41], dubbed as CMIN-R. Likewise, Liu et al. [46] presented a Cross- and Self-Modal Graph Attention Network (CSMGAN)<sup>16</sup>. It introduces a cross-modal relation graph to highlight relevant instances across the video and query. Meanwhile, a self-modal relation graph is utilized to model the pairwise relation inside each modality. In the same year, Chen et al. [9] proposed a Pairwise Modality Interaction module (dubbed PMI-LOC), which feeds multimodal features to a channel-gated modality interaction module, to exploit both intra- and inter-modality information.

**3.2.2 Sampler-based models.** Different from anchor-based approaches that generate moment candidates via multi-scale windows at each time step, sampler-based ones output moment candidates via different sampler networks, such as temporal convolution, segment-tree, and enumeration.

**Temporal Convolution.** Yuan et al. [107] proposed a Semantic Conditioned Dynamic Modulation (SCDM) mechanism based approach<sup>17</sup>, which leverages query semantic information to modulate the temporal convolution processes in a hierarchical temporal convolutional network. Particularly, similar to [39][52] for object/action detection, the position prediction module is introduced to output location offsets and overlap scores of moment candidates based on the modulated features. Similarly, Zhang et al. [111] proposed Moment Alignment Network (MAN)<sup>18</sup> to deal with semantic misalignment. In particular, it treats the encoded word features as efficient dynamic filters to convolve with input visual representations, and then a hierarchical convolutional network is applied to directly produce multi-scale moment candidates. Different from MAN, Hu et al. [28]

<sup>15</sup><https://github.com/ikuinen/CMIN>.

<sup>16</sup><https://github.com/liudaizong/CSMGAN>.

<sup>17</sup><https://github.com/ytytzy/SCDM>.

<sup>18</sup><https://github.com/dazhang-cv/MAN>.

presented an end-to-end Coarse-to-fine cross-modal sEmantic Alignment netwoRk (CLEAR)<sup>19</sup>, which utilizes bi-directional temporal convolution network followed by multi-scale temporal pooling to generate moment candidates. And to explore cross-modal semantic correlation and improve the localization efficiency, it respectively advances a multi-granularity interaction module and a semantic pruning strategy.

Different from that existing methods mainly leverage vanilla soft attention to perform the alignment in a single-step process, Liu et al. [47] presented an Iterative Alignment Network (IA-Net), which captures complicated relations between inter- and intra-modality through multi-step reasoning. To be specific, it proposes an improved co-attention mechanism that utilizes learnable paddings to address nowhere-to-attend problem with deep latent clues, and a calibration module to refine the alignment knowledge of inter- and intra-modal relations.

To boost the efficiency of moment localization, Hu et al. [27] proposed an end-to-end Cross-Modal Hashing Network (CMHN)<sup>20</sup>. It designs a cross-model hashing module to project cross-modal heterogeneous representations into a shared isomorphic Hamming space for compact hash code learning. Therefore, with the well-trained model at hand, the hash codes of any upcoming videos could be obtained offline and independently, improving the localization efficiency and scalability.

**Segment-tree.** Unlike the aforementioned sampler-based approaches, Zhao et al. [124] formulated the video moment localization task as a multi-step decision problem and proposed a Cascaded Prediction Network (CPN). To be specific, it adopts a segment tree-based structure to generate moment candidates in different temporal scales and refine the representation of them in a message-passing way via graph neural network.

**Proposal Generation Network.** To well exploit moment-level interaction and speed up the localization efficiency, Liu et al. [44] proposed an Adaptive Proposal Generation Network (APGN). It first leverages a binary classification module to predict foreground frames, and then it utilizes the boundary regression module to generate proposals on each foreground frame. In this way, the redundant proposals are decreased. Meanwhile, Xiao et al. [97] proposed a Learnable Proposal Network (LP-Net) for video moment localization<sup>21</sup>. In the LP-Net, proposal boxes are represented by 2-d parameters ranging from 0 to 1, denoting normalized center coordinates and lengths, which are randomly initialized. And these learnable proposal boxes are updated by dynamic adjustor during training.

**Enumeration.** Given that the input video has temporal length  $n$ , the Temporal Modular Network (TMN) proposed in [42] will regress  $n$  correspondence scores for each temporal segment, and then combine the scores of consecutive segments to produce  $\frac{n(n+1)}{2}$  scores for all possible moment candidates. Finally, the sub-video with the maximum score is predicted as the golden moment.

To model temporal relations between video moments, Zhang et al. [118] presented a 2D temporal map, where one dimension indicates the start time of a moment and the other indicates the end time. Based on the 2D map, [118] introduces a Temporal Adjacent Network (2D-TAN)<sup>22</sup> to generate the 2D score map, i.e., the matching scores of moment candidates on the 2D temporal map with the given query. Recently, Zhang et al. [117] extended the 2D-TAN to a multi-scale version and proposed a Multi-Scale Temporal Adjacency Network (MS-2D-TAN), which models the temporal context between video moments by a set of predefined two-dimensional maps under different temporal scales<sup>23</sup>.

<sup>19</sup><https://github.com/Huyp777/CSUN>.

<sup>20</sup><https://github.com/Huyp777/CMHN>.

<sup>21</sup><https://github.com/xiaoneil/LPNet/>.

<sup>22</sup><https://github.com/microsoft/2D-TAN>.

<sup>23</sup><https://github.com/microsoft/2D-TAN>.

Different from aforementioned methods that use alignment information to find out the best-matching candidate, Wang et al. [86] proposed a unified Dual Path Interaction Network (DPIN). It jointly considers the alignment and discrimination information to make the prediction. Particularly, the frame-level representation path extracts the discriminative boundary information from the fused features, while the candidate-level path arranges the moment candidate features in a 2D temporal map as [118] did and extract the alignment information. Finally, DPIN fuses the two kinds of representations to make prediction. Likewise, Wang et al. [87] developed a Structured Multi-level Interaction Network (SMIN) by jointly considering multiple levels of visual-textual interaction and moment structured interaction. [116] presents the Multi-stage Aggregated Transformer Network (MATN) to enhance the cross-modal alignment and localization accuracy. Particularly, it designs a new visual-language transformer backbone by using different parameters to process different modality contents, and introduces a multi-stage aggregation module to calculate the moment representation via considering three stage-specific representations. Wu et al. [96] proposed a video moment retrieval model, named SV-VMR, which jointly models the fine-grained and comprehensive relations by using both semantic and visual structures. Unlike previous work, Gao et al. [17] formulated the video moment localization task as the video reading comprehension and presented a Relation-aware Network (RaNet)<sup>24</sup>. To distinguish similar moment candidates in the visual modality, it introduces a coarse-and-fine cross-modal interaction module to simultaneously capture the sentence-moment and token-moment level interaction, obtaining sentence-aware and token-aware moment representations. Meanwhile, it leverages graph convolutional network to capture moment-moment relations, which further strengthens the discriminative of moment representations.

To against the temporal location biases of video moment localization, Yang et al. [104] advanced the Deconfounded Cross-modal Matching (DCM) method<sup>25</sup>, which considers the moment temporal location as a hidden confounding variable. To remove the confounding effects of moment location, it disentangles the moment representation to infer the core feature of visual content, and then applies causal intervention on the multi-modal input based on backdoor adjustment. To well balance the localization accuracy and speed, Gao et al. [18] proposed the Fast Video Moment Retrieval (FVMR) model, which replaces the complex cross-modal interaction module in existing methods with a cross-modal common space.

**3.2.3 Proposal-free models.** Unlike the anchor- and sampler-based methods that depend on moment candidates, the proposal-free ones directly predict the start and end time of the target moment. This removes the need to retrieve and re-rank multiple moment candidates.

As the first work in this subbranch, Chen et al. [8] proposed a localizing network (LNet), which works in an end-to-end fashion, to tackle the video moment localization task. It first matches the query and video sequence by cross-gated attended recurrent networks, to exploit their fine-grained interactions and generate a query-aware video representation. Afterwards, a self interactor is designed to perform cross-frame matching, which dynamically encodes and aggregates the matching evidences. Finally, a boundary model is introduced to locate the positions of video moments corresponding to the query by directly predicting the start and end points. To further improve the localization accuracy of LNet, the follow-up research of this branch is dedicated to exploring how to well exploit interactive information between the video and the given query.

To adequately explore the cross-modal interactions between the query and video, Yuan et al. [108] proposed an end-to-end Attention Based Location Regression (ABLR) model<sup>26</sup>. It leverages a multi-modal co-attention mechanism to learn both video and query attentions. Moreover, a multi-layer

<sup>24</sup><https://github.com/Hunterssxx/RaNet>.

<sup>25</sup>[https://github.com/Xun-Yang/Causal\\_Video\\_Moment\\_Retrieval](https://github.com/Xun-Yang/Causal_Video_Moment_Retrieval).

<sup>26</sup>[https://github.com/ytytzy/ABLR\\_code](https://github.com/ytytzy/ABLR_code).

attention-based location prediction network is proposed to regress the temporal coordinates for the target video moment. Similarly, Ghosh et al. [22] proposed an Extractive Clip Localization (ExCL) approach to predict the start and end frames of the target moment. Particularly, they compared three variants of span predictor, i.e., MLP, Tied-LSTM, and Conditioned-LSTM, to predict the start and end probabilities for each frame. Liu et al. [54] designed a Single-shot Semantic Matching Network (SSMN), which properly explores the cross-modal relationship and the temporal information via an enhanced cross-modal attention module. For in-depth relationship modeling between semantic phrases and video segments, Mun [60] proposed a Sequential Query Attention Network (SQAN) based method by performing local-global video-query interactions<sup>27</sup>. Tang et al. [82] proposed an Attentive Cross-modal Relevance Matching (ACRM) model<sup>28</sup>, which uses element-wise multiplication and subtraction functions to model the interaction between the frame and frame-specific query features. Liang et al. [37] proposed a Multi-branches Interaction Model (MIM), which re-weights the video features according to relevance of query and video in multiple sub-spaces. By treating the video as a text passage and the target moment as the answer span, Zhang et al. [114] proposed a video span localizing network (VSLNet) on top of the standard span-based QA framework<sup>29</sup>. In particular, a context-query attention (CQA) [99] is designed to capture the cross-modal interactions between visual and textural features. Recently, to address the performance degradation on long videos, Zhang et al. [113] further extended VSLNet to VSLNet-L by introducing a multi-scale split-and-concatenation strategy. It first partition long video into clips of different lengths, and then locate the target moment in the clips that are more likely to contain it. Unlike the above methods, Opazo et al. [64] introduced a Proposal-free Temporal Moment Localization model using Guided Attention (PFGA)<sup>30</sup>, which utilizes an attention-based dynamic filter to transfer query information to the video. Moreover, a new loss is designed to enforce the model to focus on the most relevant part of the video. Since humans have difficulty agreeing on the start and end time of an action inside a video [1][74], PFGA uses soft-labels [80] to model the uncertainty associated to the labels. The final loss for training PFGA method is defined as,

$$\begin{cases} L = L_{KL} + L_{att}, \\ L_{KL} = D_{KL}(\hat{\tau}^s || \tau^s) + D_{KL}(\hat{\tau}^e || \tau^e), \\ L_{att} = - \sum_{i=1}^n (1 - \delta_{\tau^s \leq i \leq \tau^e} \log(1 - a_i)), \end{cases} \quad (5)$$

where  $L_{KL}$  loss aims to minimize the Kullback-Leibler divergence between the predicted and ground truth probability distributions,  $L_{att}$  loss is designed to encourage the model to attend relevant features,  $\delta$  is the Kronecker delta,  $n$  denotes the length of visual sequence, and  $a_i$  is the attention of the corresponding location.

To further explore fine-grained interaction information within the video and query sentence, Chen et al. [11] advanced a Hierarchical Visual-Textual Graph (HVTG) model for video moment localization<sup>31</sup>. To be specific, it first builds an object-sentence subgraph to obtain sentence-aware object features for each frame. And then it feeds sentence-aware object features into an object-object subgraph, to capture object-object interactions inside each frame. Afterwards, a sentence-guided attention followed by a Bi-LSTM is introduced to establish temporal relations among frames, outputting final visual representations. To explore multi-level query semantics (both word- and phrase-level) as well as model multi-level interactions between two modalities, Tang et al. [83] proposed a Multi-level Query Exploration and Interaction (MQEI) model. To be specific, it leverages

<sup>27</sup><https://github.com/JonghwanMun/LGI4temporalgrounding>.

<sup>28</sup><https://github.com/tanghaoyu258/ACRM-for-moment-retrieval>.

<sup>29</sup><https://github.com/IsaacChanghau/VSLNet>.

<sup>30</sup><https://github.com/crodriguez/TMLGA>.

<sup>31</sup><https://github.com/forwchen/HVTG>.

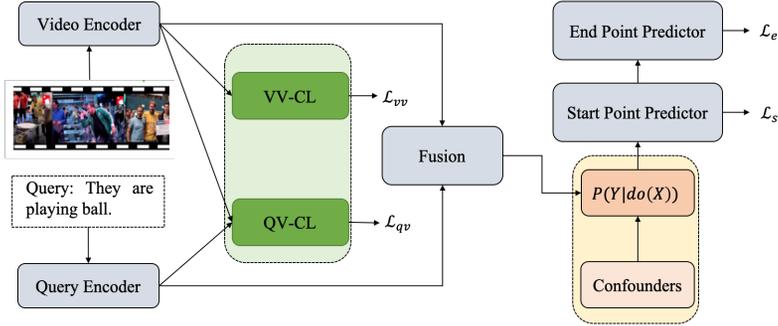


Fig. 5. The architecture of IVG. VV-CL and QV-CL respectively refer to two contrastive modules with losses expressed as  $L_{vv}$  and  $L_{qv}$ .  $L_s$  and  $L_e$  denote the cross-entropy losses for predicting the boundary of the target span, figure from [62].

a stack of syntactic GCNs to model the syntactic dependencies of the query, obtaining the word-level query features. And a sequential attention module is adopted to learn phrase-level query representations. With such multi-level query representations, the context-query attention is adopted for the word-level and phrase-level cross-modal feature fusion.

Different from most proposal-free methods that are devoted to multi-modal modeling (i.e., fusion or interaction), Li et al. [36] focused on exploiting the fine-grained temporal clues in videos and presented a Contextual Pyramid Network (CP-Net) to mine rich temporal contexts through fine-grained hierarchical correlation at different 2D temporal scales. Similarly, a language-conditioned message-passing algorithm is proposed in [70] to well learn the video feature embedding, which could capture the relationships between humans, objects and activities in the video. Inspired by the successful application of biaffine mechanism in dependency parsing, Liu et al. [45] applied the biaffine mechanism to video moment localization and presented the Context-aware Biaffine Localizing Network (CBLN)<sup>32</sup>. It leverages a multi-context biaffine localization module to aggregate both multi-scale local and global contexts for each frame representation, and then scores all possible pairs of start and end frames for moment localization.

To alleviate the imbalance problem between positive and negative samples to some extent, Lu et al. [55] proposed a DEense Bottom-Up Grounding (DEBUG) framework by regarding all frames in the ground truth moment as positive samples. For each positive frame, DEBUG has a classification subnet to predict its relatedness with the query, and a boundary regression subnet to regress the unique distances from its location to bi-directional ground truth boundaries. This helps to avoid falling into the local optimum caused by independent predictions since each pair of boundary predictions is based on the same frame feature. Note that DEBUG can be seamlessly incorporated into any backbone to boost performance.

Through utilizing the distances between the frame within the ground truth and the start (end) frame as dense supervisions, Zeng et al. [109] designed a dense regression network (DRN)<sup>33</sup>. To be specific, DRN first forwards the video frames and the query to the video-query interaction module for extracting the multi-scale feature maps. Afterwards, each feature map is processed by the grounding module to predict a temporal bounding box, a semantic matching score, and an IoU score at each temporal location for ranking. Finally, combining the matching score and the IoU score, DRN could find the best grounding result.

Despite the enormous success of the aforementioned models, they may suffer from the spurious correlations between textual and visual features due to the selection bias of the dataset. To address

<sup>32</sup><https://github.com/liudaizong/CBLN>.

<sup>33</sup><https://github.com/Alvin-Zeng/DRN>.

this problem, [62] first presents the Interventional Video Grounding (IVG) by introducing the causal interventions, as shown in Fig. 5. Moreover, it introduces a dual contrastive learning to learn more informative visual and textual representations. To be specific, the loss  $L_{vv}$  is expressed as follows,

$$L_{vv} = -I_{\theta}^{vv(s)} - I_{\theta}^{vv(e)}, \quad (6)$$

where  $I_{\theta}^{vv(s)}$  and  $I_{\theta}^{vv(e)}$  respectively denote the mutual information between the start and end boundaries of the video and the other clips. And the contrastive loss  $L_{vq}$  is defined as follow,

$$L_{vq} = -I_{\theta}^{vq} = -E_{V'_a}[sp(C(q, V'))] + E_{V'_b}[sp(C(q, V'))], \quad (7)$$

where  $V'_a$  denote the features that reside within a target moment,  $V'_b$  denote the features are located outside of the target moment,  $C$  refers to the mutual information discriminator, and  $sp(z) = \log(1 + e^z)$ .

Different from most existing methods that only consider RGB images as visual features, Chen et al. [12] proposed a multi-modal learning framework for video temporal grounding using (D)epth, (R)GB, and optical (F)low with the (T)ext as the query (DRFT). To model the interactions between modalities, this paper develops a dynamic fusion mechanism across modalities via co-attentional transformer. Moreover, to enhance intra-modal feature representations, it leverages self-supervised contrastive learning across videos for each modality.

**Summarization.** Below, we summarize the one-stage video moment localization methods.

- *Anchor-based models:* In this branch, on top of TGN, CBP designs a boundary prediction module to predict temporal boundaries at each time step. To capture the interaction information, BSSTL, RMN, FIAN, and I2N focus on exploring the inter-modal interactions, while CMIN, CSMGAN, PMI-LOC and CMIN-R focus on exploring both intra- and inter-modal interactions.
- *Sampler-based models:* In this branch, SCDM, MAN, CLEAR, IA-Net, and CMHN adopt the temporal convolution network to generate moment candidates. Specifically, SCDM designs semantic modulated temporal convolution, MAN utilizes hierarchical convolution, and CLEAR leverages bi-directional temporal convolution. Although IA-Net and CMHN also utilizes temporal convolution to generate moment candidates, they respectively are devoted to tackling the interaction modeling and location efficiency issue. CPN innovatively integrates the segment-tree into the video moment localization for moment generation. Differently, APGN and Lp-Net introduce different proposal generation networks for moment candidates generation. TMN, 2D-TAN, MS-2D-TAN, DPIN, MATN, SMIN, SV-VMR, RaNet, DCM, and FVMR enumerate all possible moment candidates. More specifically, 2D-TAN and MS-2D-TAN respectively utilize the temporal adjacent network and multi-scale temporal adjacent network to gradually perceive more context of adjacent moment candidates. DPIN, MATN, SMIN, SV-VMR, and RaNet focus on developing different mechanisms to capturing interaction information. Differently, DCM is devoted to addressing the location bias issue, while FVMR aims to well balance the location accuracy and efficiency.
- *Proposal-free models:* In this branch, ABLR, ExCL, SQAN, PFGA, VSLNet, VSLNet-L, ACRM, MIM, and SSMN focus on exploring cross-modal interaction information between videos and queries. Thereinto, VSLNet and VSLNet-L tackle video moment localization by using the standard span-based QA framework. HVTG and MQEI are devoted to modeling both intra- and inter-modal interactions, while CBLN, CP-Net, and DORi merely pay attention to enhance the visual representation. DEBUG aims to tackle the imbalance problem between positive and negative samples, while DRN aims to promote the accuracy of localization by considering frame-level dense supervisions. IVG targets to address the spurious correlations

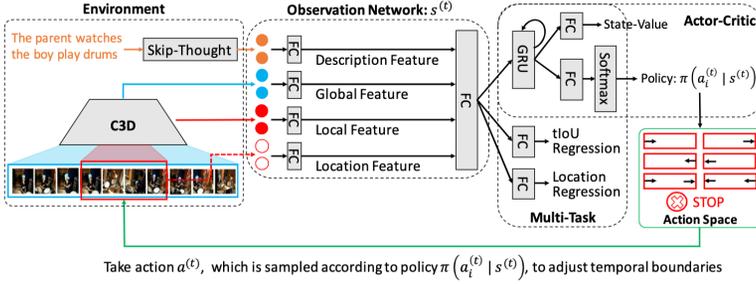


Fig. 6. The overall pipeline of RWM model, figure from [24].

between textual and visual features, while DRFT tackles the video moment localization by adopting multi-modal visual features.

### 3.3 Reinforcement Learning

As aforementioned, the one- and two-stage video moment localization studies inevitably suffer from inefficiency and unintelligent issues, supervised reinforcement learning based localization approaches are then proposed recently, as summarized in Table 1.

He et al. [24] presented the first end-to-end reinforcement learning based framework, named RWM, which formulates the video moment localization as a problem of sequential decision making (as shown in Fig. 6). Concretely, its action space consists of 7 predefined actions<sup>34</sup>. At each time step, the observation network concatenates and fuses sentence embedding, global video feature, local video feature, and normalized temporal boundaries, generating the state vector. The obtained state vector is then fed into the actor-critic [79] to learn the policy and state-value. To learn more representative state vectors, RWM combines reinforcement learning and supervised learning into a multi-task learning framework.

Inspired by the coarse-to-fine decision-making paradigm of human, Wu et al. [95] formulated a Tree-Structured Policy based Progressive Reinforcement Learning (TSP-PRL) framework<sup>35</sup>. Different from the previous work RWM [24], it builds a hierarchical action space, containing all primitive actions in this task. Moreover, TSP-PRL designs a tree-structured policy to decompose complex action policies and obtains a reasonable primitive action via two stages selection. To optimize the tree-structured policy, the progressive reinforcement learning (PRL) is designed on the basis of [79].

To bridge the huge visual-semantic gap between videos and queries, Wang et al. [90] proposed a semantic matching reinforcement learning (SM-RL) model, which improves video representation by introducing visual semantic concepts. Specifically, in the observation module, unlike previous methods merely leverage the global image information, the sentence query, as well as the location information to output the action and state values, SM-RL also integrates semantic concept features of videos. To jointly explore crucial clues hidden in the temporal and spatial information, Cao et al. [5] contributed a Spatio-Temporal ReinfOrcement learniNG (STRONG) framework for video moment localization<sup>36</sup>. It first exploits a temporal-level reinforcement learning to dynamically adjust the boundary of localized video moment. Thereafter, a spatial-level reinforcement learning is proposed to track the scene on consecutive image frames, therefore filtering out less relevant information. Concretely, the state is defines as the combination of the query sentence feature, the

<sup>34</sup>Moving start/end point ahead by  $\delta$  (a predefined video length), moving start/end point backward by  $\delta$ , shifting both start and end point backward/forward by  $\delta$ , and a STOP action.

<sup>35</sup><https://github.com/WuJie1010/TSP-PRL>.

<sup>36</sup><https://github.com/yawenzeng/STRONG>.

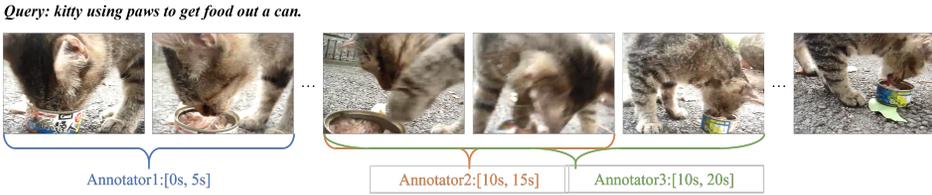


Fig. 7. An annotation example from the DiDeMo dataset. Each color refers to one annotation result of the given query.

local video feature, and the spatial video feature. To simultaneously tackle the challenge of efficient search and query-video alignment, Hahn et al. [23] presented a Tripping through time Network (TripNet). Compared with [24] and [95], the main contribution of TripNet is that it integrates the reinforcement learning (RL) and fine-grained video analysis. Concretely, in the state processing module, it designs a novel gated attention mechanism to model fine-grained textual and visual representations for video-text alignment. Besides, its reward function is defined as the difference of potentials between the previous state and current state. Therefore, the agent could find the target without looking at all frames of the video.

Considering both limited moment selection and insufficient structural comprehension, Sun et al. [78] proposed a Multi-Agent Boundary-Aware Network (MABAN)<sup>37</sup>, which utilizes multi-agent reinforcement learning to obtain the two temporal boundary points for the target moment. At each step of temporal boundary adjustment, the start point agent and end point agent receive the state vector from the observation network and adjust the temporal boundaries in variable directions and scales, making moment selection more flexible and goal-oriented. To overcome the latter issue, a cross-modal interaction that considers semantic fusion in global and local phases is introduced to explore rich contextual information.

Different from pioneer methods that consider the video moment localization as either a ranking issue or a localization problem, Cao et al. [6] proposed an Adversarial Video Moment Retrieval model (AVMR), which combines ranking and localization into a unified framework<sup>38</sup>. To accelerate the convergence and increase the diversity of generated video moments, AVMR employs deep deterministic policy gradient (DDPG) [38] to learn the policy. In addition, as utilizing indistinguishable values of IoU as the reward is unstable and difficult to convergence, AVMR introduces a discriminator to provide flexible reward, i.e., the Bayesian personalized ranking model.

**Summarization.** In this category, TSP-PRL build a hierarchical action space and designs the progressive reinforcement learning to optimize its tree-structured policy. Different from TSP-PRL, TripNet focuses on enhancing visual representations, SM-RL introduces visual semantic concepts to calculate state values, while STRONG jointly considers the temporal-level and spatial-level reinforcement learning to improve the localization performance. Differently, AVMR utilizes the Bayesian personalized ranking model to provide reward and MABAN is devoted to capturing cross-modal interactions.

#### 4 WEAKLY SUPERVISED VIDEO MOMENT LOCALIZATION

**Problem Statement.** It is difficult for human to mark the start and end locations of a certain moment. As shown in Fig. 7, three annotators give completely different annotation results for the same query and video. These inconsistent temporal annotations may introduce ambiguity in the training data. Moreover, acquiring dense annotations of query-temporal boundaries is often tedious. Intuitively, video-sentence pairs may be obtained with minimum human intervention as compared

<sup>37</sup><https://mic.tongji.edu.cn/e5/23/c9778a189731/page.htm>.

<sup>38</sup><https://github.com/yawenzeng/AVMR>.

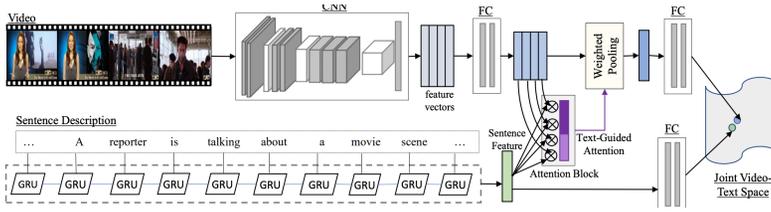


Fig. 8. The overall framework of TGA, figure from [59]

to temporal sentence annotations. Therefore, weakly supervised video moment localization task is introduced. Unlike supervised video moment localization that has access to the start and end time points of the queries, weakly supervised video moment localization only utilizes video-sentence pairs for training, namely it does not rely on temporal annotation information.

As summarized in Table 2, weakly supervised video moment localization studies can also be categorized into tree groups: two-stage, one-stage, and reinforcement learning based methods.

#### 4.1 Two-stage Methods

To the best of our knowledge, the first weakly supervised video moment localization model is proposed in [59]<sup>39</sup>. As shown in Fig. 8, it aims to train a joint embedding network to project video and query features into the same joint space. In particular, following [16], it first obtains the moment candidates via multi-scale sliding windows. After the feature extraction process, it introduces the Text-Guided Attention (TGA) module to obtain the query-wise video feature. Thereafter, two FC layers are utilized in [59] to project the query-specific video feature and paired query feature into the same joint space.

Inspired by the success of two-stream structure in weakly supervised object and action detection tasks [3][89], [20] presents a Weakly Supervised Natural Language Localization Network (WSSLN), consisting of two branches: alignment branch and selection branch. The former branch is utilized to estimate the consistency score, while the latter one is adopted to calculate the selection score. Scores from both branches are then merged to produce the final score of each moment-query pair. Similarly, Ma et al. [56] proposed the Video-Language Alignment Network (VLANet) to prune out spurious moment candidates. To be specific, a surrogate proposal selection module is designed to select the best-matched moment from each moment group based on the cosine similarity to the query embedding. Moreover, VLANet is trained using contrastive loss that enforces semantically similar videos and queries to cluster in the joint embedding space.

To improve the latent alignment between videos and natural language, Tan et al. [81] designed a latent co-attention network (named LoGAN), which learns contextualized visual semantic representations from fine-grained frame-by-word interactions. Particularly, it builds a word-conditioned visual graph module to learn contextualized visual semantic representations by integrating temporal contextual information into the visual features.

#### 4.2 One-stage Methods

Recently, Lin et al. [40] presented a Semantic Completion Network (SCN), which scores all the moments sampled at different scales in a single pass. Particularly, it introduces the semantic completion module to measure the semantic similarity between moments and the query, as well as compute rewards. To train the semantic completion module, SCN designs a reconstruction loss to force it to extract key information from the visual context. Different from previous works using paired sentences as supervision information, Wang et al. [92] proposed the visual co-occurrence

<sup>39</sup>[https://github.com/niluthpol/weak\\_supervised\\_video\\_moment](https://github.com/niluthpol/weak_supervised_video_moment).

Table 3. Statistics of datasets for video moment localization task.

Dataset	#Videos	#Queries	Duration	Domain	Source
DiDeMo [2]	10,464	40,543	30s	Open	Flickr
TEMPO [25]	10,464	-	30s	Open	Flickr
Charades-STA [16]	9,848	16,128	31s	Daily activities	Homes
TACoS [67]	127	18,818	296s	Cooking	Lab Kitchen
ActivityNet-Captions [31]	19,209	71,957	180s	Open	YouTube

Table 4. Performance comparison between two supervised localization models on TEMPO-HL.

Methods	Year	DiDeMo		Before		After		Then		While		Average		
		R@1	mIoU	R@1	mIoU	R@1	mIoU	R@1	mIoU	R@1	mIoU	R@1	R@5	mIoU
MLLC [25]	2018	27.38	<b>42.45</b>	32.33	56.91	14.43	37.33	19.58	50.39	10.39	35.95	20.82	71.68	44.57
TCMN [119]	2019	<b>28.77</b>	42.37	<b>35.47</b>	<b>59.28</b>	17.91	<b>40.79</b>	20.47	<b>50.78</b>	<b>18.81</b>	<b>42.95</b>	<b>24.29</b>	<b>76.98</b>	<b>47.24</b>

Table 5. Performance comparison between two supervised localization models on TEMPO-TL.

Methods	Year	DiDeMo		Before		After		Then		Average		
		R@1	mIoU	R@1	mIoU	R@1	mIoU	R@1	mIoU	R@1	R@5	mIoU
MLLC [25]	2018	27.46	<b>41.20</b>	35.31	41.81	29.38	38.90	26.83	54.97	29.74	76.76	44.22
TCMN [119]	2019	<b>28.90</b>	41.03	<b>37.68</b>	<b>44.78</b>	<b>32.61</b>	<b>42.77</b>	<b>31.16</b>	<b>55.46</b>	<b>32.85</b>	<b>78.73</b>	<b>46.01</b>

alignment (VCA) method, which targets to learn more discriminative and robust visual features by mining visual supervision information. Concretely, it utilizes the similarity among sentences to mine positive pairs of relevant video moments from different videos as well as negatives pairs for contrastive learning.

Existing weak-supervised methods ignore the influence of intra-sample confrontment between semantically similar moments, they hence fail to distinguish the target moment from plausible negative moments. Inspired by this, Zhang et al. [123] advanced a Regularized Two-Branch Proposal Network (RTBPN) that simultaneously considers the inter-sample and intra-sample confrontments. To well model the fine-grained video-text local correspondences, Yang et al. [103] presented a Local Correspondence Network (LCNet) for weakly supervised temporal sentence grounding. Specifically, it leverages the hierarchical feature representation module to model the one-one, one-many, many-one, and many-many correspondences between the video and text. Differently, Wang et al. [91] treated the weakly supervised task as a multiple instance learning (MIL) problem, and introduced a Weakly Supervised Temporal Adjacent Network (WSTAN) by integrating the temporal adjacent network, a complementary branch, and the self-discriminating loss into a unified framework. Therein, the temporal adjacent network is leveraged to model the relationships among candidate proposals; the complementary branch is utilized to refine the predictions and rediscover more semantically meaningful clips; and the self-discriminating loss is designed to force the model to be more temporally discriminative. Similarly, Li et al. [33] proposed a Multi-Scale 2D Representation Learning (MS-2D-RL) method, which conducts convolution over multi-scale 2D temporal maps to capture the moment relations. Besides, it designs a moment evaluation module to generate pseudo label for training.

### 4.3 Reinforcement Learning

To the best of our knowledge, Boundary Adaptive Refinement (BAR) [94] framework is the first work that extends reinforcement learning to weakly supervised video moment localization. Particularly, it consists of a context-aware feature extractor, an adaptive action planner, and a cross-modal alignment evaluator. Therein, the context-aware feature extractor is utilized to encode the current

Table 6. Performance comparison among supervised localization models on DiDeMo.

Type	Method	Year	R@1	R@5	mIoU
Two-stage	MCN [2]	2017	28.10	78.21	41.08
	MLLC [25]	2018	28.37	78.64	<b>43.65</b>
	TCMN [119]	2019	<b>28.90</b>	<b>79.00</b>	41.03
One-stage	TGN [7]	2018	28.23	79.26	42.97
	TMN [42]	2018	22.92	76.08	35.17
	LNet [8]	2019	-	-	41.43
	MAN [111]	2019	27.02	<b>81.70</b>	41.16
	I2N [63]	2021	<b>29.00</b>	73.09	<b>44.32</b>
RL	SM-RL [90]	2019	31.06	80.45	43.94

and contextualized environment states. The adaptive action planner is designed to infer action sequences to refine the temporal boundary, and the advantage actor-critic (A2C) algorithm is chosen to train the adaptive action planner. Moreover, the cross-modal alignment evaluator is utilized to estimate the alignment score between each segment-query pair and assign a target-oriented reward to each action.

**Summarization.** Among two-stage methods, TGA aims to learn a joint embedding space for video and query representations. Moreover, WSLLN and VLANet focus on proposal selection. To be specific, WSLLN leverages a two-stream structure to measure moment-query consistency and conduct moment selection simultaneously, while VLANet proposes a surrogate proposal selection module to prune out spurious moment candidates. Differently, LoGAN is devoted to learning contextualized visual representation. For one-stage methods, according to the processing of moment candidates, they are roughly divided into anchor-based and enumeration based approaches. Particularly, among anchor-based methods, SCN leverages the semantic completion module to refine the matching score, while VCA targets to learn more discriminative and robust visual features. As to enumeration based methods, RTBPN aims to tackle the inter-sample and intra-sample confrontation issue. LCNet focuses on exploring one-one, one-many, many-one, and many-many interaction information between the video and text. WSTAN designs the temporal adjacent network to model the relationships among moment candidates, while MS-2D-RL conducts convolution over multi-scale 2D temporal maps to capture moment relations.

## 5 UNSUPERVISED VIDEO MOMENT LOCALIZATION

Although weakly-supervised video moment localization approaches have achieved favorable performance, a certain amount of paired video-sentence data is still indispensable for model learning. Moreover, it is difficult to collect a large-scale video-sentence dataset for an arbitrary domain in the wild. Inspired by this, Gao et al. [19] proposed the unpaired video moment retrieval approach, namely U-VMR, which requires no textual annotations of videos and instead leverages the existing visual concept detectors and a pre-trained image-sentence embedding space. To be specific, it first designs a video-conditioned sentence generator to generate suitable sentence representations by leveraging the mined visual concepts. And then it develops a relation-aware moment localizer, which leverages a sentence-guided graph neural network to effectively select a portion of video clips as the moment representation. To obtain different candidate moments, U-VMR follows 2D-TAN [118] to utilize a sparse sampling strategy for generating candidate moments. Finally, the pre-trained image-sentence embedding model is adopted to align the generated sentence and moment representations for model learning. Similarly, Nam et al. [61] proposed the zero-shot natural language video localization task, which only requires easily available text corpora, off-the-shelf

Table 7. Performance comparison among supervised localization models on TACoS.

Type	Model	Year	R(1,0.5)	R(1,0.3)	R(1,0.1)	R(5,0.5)	R(5,0.3)	R(5,0.1)
Two-stage	CTRL [16]	2017	13.30	18.32	24.32	25.42	36.69	48.73
	ACRN [50]	2018	14.62	19.52	24.22	24.88	34.97	47.42
	VAL [76]	2018	14.74	19.76	25.74	26.52	38.55	51.87
	MCF [93]	2018	12.53	18.64	25.84	24.73	37.13	52.96
	SLTA [29]	2019	11.92	17.07	23.13	20.86	32.90	46.52
	ACL [21]	2019	20.01	24.17	31.64	30.66	42.15	57.85
	SAP [10]	2019	18.24	-	31.15	28.11	-	53.51
	BPNet [98]	2021	20.96	25.96	-	-	-	-
	MMRG [110]	2021	<b>39.28</b>	<b>57.83</b>	<b>85.34</b>	<b>56.34</b>	<b>78.38</b>	<b>84.37</b>
	One-stage	TGN [7]	2018	18.90	21.77	41.87	31.02	39.06
BSSTL [32]		2019	18.73	22.31	43.11	29.89	40.87	54.67
CMIN [122]		2019	18.05	24.64	32.48	27.02	38.46	62.13
ABLR [108]		2019	9.40	19.50	34.70	-	-	-
ExCL [22]		2019	28.00	45.50	-	-	-	-
SCDM [107]		2019	21.17	26.11	-	32.18	40.16	-
DEBUG [55]		2019	-	23.45	41.15	-	-	-
VSLNet [114]		2020	24.27	29.61	-	-	-	-
DRN [109]		2020	23.17	-	-	33.36	-	-
CBP [88]		2020	24.79	27.31	-	37.40	43.64	-
CMIN-R [41]		2020	19.57	27.33	36.88	28.53	43.35	64.93
CSMGAN [46]		2020	27.09	33.90	42.74	41.22	53.98	68.97
FIAN [66]		2020	28.58	33.87	39.55	39.16	47.76	56.14
RMN [43]		2020	25.61	32.21	42.17	40.58	54.20	68.75
2D-TAN [118]		2020	25.32	37.29	47.59	45.04	57.81	70.31
DPIN [86]		2020	32.92	46.74	59.04	50.26	62.16	75.78
CLEAR [28]		2021	30.27	42.18	-	51.76	63.61	-
IA-Net [17]		2021	26.27	37.91	47.18	46.39	57.62	71.75
APGN [44]		2021	27.86	40.47	-	47.12	59.98	-
CMHN [27]		2021	25.58	30.04	-	35.23	44.05	-
SMIN [87]		2021	35.24	48.01	-	53.36	65.18	-
RaNet [17]		2021	33.54	43.34	-	55.09	67.33	-
IVG [62]		2021	29.07	38.84	49.36	-	-	-
CPN [124]		2021	36.58	48.29	<b>61.24</b>	-	-	-
CBLN [45]		2021	27.65	38.98	49.16	46.24	59.96	73.12
MS-2D-TAN [117]		2021	34.29	41.74	49.24	56.76	67.01	<b>78.33</b>
I2N [63]		2021	29.25	31.47	-	46.08	52.65	-
FVMR [18]		2021	29.12	41.48	53.12	50.00	64.53	78.12
ACRM [82]		2021	33.84	47.11	-	-	-	-
MIM [37]		2021	26.54	35.44	44.69	-	-	-
CP-Net [36]	2021	29.29	42.61	-	-	-	-	
DORi [70]	2021	28.69	31.80	-	-	-	-	
MATN [116]	2021	<b>37.57</b>	<b>48.79</b>	-	<b>57.91</b>	<b>67.63</b>	-	
RL	SM-RL [90]	2019	15.95	20.25	26.51	27.84	38.47	50.01
	TripNet [23]	2020	<b>19.17</b>	<b>23.95</b>	-	-	-	-

object detector, and a collection of videos to localize. To address the task, a Pseudo-Supervised Video Localization (PSVL) method is advanced. To be specific, it first generates temporal event proposal by compositing discovered atomic event, and then generates simplified sentence for each temporal event proposal via considering nouns detected by the object detector and verbs predicted from text corpora. Finally, it utilizes a attentive cross-modal neural network to predict the temporal boundary regions.

Table 8. Performance comparison among supervised localization models on Charades-STA.

Type	Model	Year	R(1,0.7)	R(1,0.5)	R(1,0.3)	R(5,0.7)	R(5,0.5)	R(5,0.3)
Two-stage	CTRL [16]	2017	7.15	21.42	-	26.91	59.11	-
	ROLE [51]	2018	7.82	21.74	37.68	30.06	70.37	92.79
	VAL [76]	2018	9.16	23.12	-	27.98	61.26	-
	SLTA [29]	2019	8.25	22.81	38.96	31.46	72.39	94.01
	ACL [21]	2019	12.20	30.48	-	35.13	64.84	-
	QSPN [101]	2019	15.80	35.60	54.70	45.40	<b>79.40</b>	<b>95.60</b>
	SAP [10]	2019	13.36	27.42	-	38.15	66.37	-
	MIGCN [121]	2021	22.04	42.26	-	-	-	-
	BPNNet [98]	2021	20.51	38.25	55.46	-	-	-
	MMRG [110]	2021	-	<b>44.25</b>	<b>71.60</b>	-	<b>60.22</b>	78.67
One-stage	SCDM [107]	2019	33.43	54.44	-	58.08	74.43	-
	MAN [111]	2019	22.72	46.53	-	53.72	86.23	-
	DEBUG [55]	2019	17.69	37.39	54.95	-	-	-
	ExCL [22]	2019	22.40	44.10	61.50	-	-	-
	CBP [88]	2020	18.87	36.80	-	50.19	70.94	-
	PMI-LOC [9]	2020	19.27	39.37	55.48	-	-	-
	FIAN [66]	2020	37.72	58.55	-	63.52	87.80	-
	RMN [43]	2020	36.98	59.13	-	61.02	87.51	-
	2D-TAN [118]	2020	23.25	39.81	-	52.15	79.33	-
	DPIN [86]	2020	26.96	47.98	-	55.00	85.53	-
	PFGA [64]	2020	33.74	52.02	67.53	-	-	-
	VSLNet [114]	2020	30.19	47.31	64.30	-	-	-
	DRN [109]	2020	23.68	42.90	-	54.87	87.80	-
	HVTG [11]	2020	23.30	47.27	61.37	-	-	-
	APGN [44]	2021	38.86	62.58	-	62.11	91.24	-
	SMIN [87]	2021	<b>40.75</b>	<b>64.06</b>	-	<b>68.09</b>	89.49	-
	RaNet [17]	2021	39.65	60.40	-	64.54	89.57	-
	IA-Net [17]	2021	37.91	61.29	-	62.04	<b>89.78</b>	-
	SSMN [54]	2021	28.49	51.51	66.13	-	-	-
	IVG [62]	2021	32.88	50.24	67.63	-	-	-
	CPN [124]	2021	36.67	59.77	75.53	-	-	-
	CBLN [45]	2021	28.22	47.94	-	57.47	88.20	-
	DRFT [12]	2021	40.15	63.03	<b>76.68</b>	-	-	-
	LP-Net [97]	2021	34.03	54.33	66.59	-	-	-
	MQEI [83]	2021	38.04	59.35	73.67	-	-	-
	MS-2D-TAN [117]	2021	23.25	41.10	-	48.55	81.53	-
	I2N [63]	2021	22.88	41.69	-	46.85	75.73	-
	FVMR [18]	2021	18.22	38.16	-	44.96	82.18	-
	ACRM [82]	2021	22.60	44.10	65.10	-	-	-
	SV-VMR [96]	2021	19.98	38.09	-	38.15	66.37	-
MIM [37]	2021	35.51	53.23	70.08	-	-	-	
CP-Net [36]	2021	22.47	40.32	-	-	-	-	
DORi [70]	2021	40.56	59.65	72.72	-	-	-	
RL	RWM [24]	2019	-	36.70	-	-	-	-
	SM-RL [90]	2019	-	11.17	24.36	-	32.08	61.25
	TSP-PRL [95]	2020	24.75	45.45	-	-	-	-
	TripNet [23]	2020	16.07	38.29	<b>54.64</b>	-	-	-
MABAN [78]	2021	<b>32.26</b>	<b>56.29</b>	-	-	-	-	

## 6 DATASETS AND EVALUATION

### 6.1 Datasets

The statistics of the datasets for video moment localization task are summarized in Table 3.

Table 9. Performance comparison among supervised localization models on ActivityNet-Captions.

Type	Model	Year	R(1,0.7)	R(1,0.5)	R(1,0.3)	R(5,0.7)	R(5,0.5)	R(5,0.3)
Two-stage	QSPN [101]	2019	13.60	27.70	45.30	38.30	59.20	75.70
	MIGCN [121]	2021	<b>44.94</b>	<b>60.03</b>	-	-	-	-
	BPNet [98]	2021	24.69	42.07	<b>58.98</b>	-	-	-
One-stage	TGN [7]	2018	-	28.47	45.51	-	43.33	57.32
	BSSL [32]	2019	-	47.68	55.32	-	57.53	70.53
	CMIN [122]	2019	23.88	43.40	63.61	50.73	67.95	80.54
	SCDM [107]	2019	19.86	36.75	54.80	41.53	64.99	77.29
	ABLR [108]	2019	-	36.79	55.67	-	-	-
	ExCL [22]	2019	23.60	43.60	63.00	-	-	-
	DEBUG [55]	2019	-	39.72	55.91	-	-	-
	CBP [88]	2020	17.80	35.76	54.30	46.20	65.89	77.63
	CMIN-R [41]	2020	24.48	44.62	64.41	52.96	69.66	82.39
	PMI-LOC [9]	2020	17.83	38.28	59.69	-	-	-
	CSMGAN [46]	2020	29.15	<b>49.11</b>	<b>68.52</b>	59.63	77.43	87.68
	FIAN [66]	2020	29.81	47.90	64.10	59.66	77.64	87.59
	RMN [43]	2020	27.21	47.41	67.01	56.76	75.64	87.03
	2D-TAN [118]	2020	27.38	44.05	58.75	62.26	76.65	85.65
	DPIN [86]	2020	28.31	47.27	62.40	60.03	77.45	87.52
	PFGA [64]	2020	19.26	33.04	51.28	-	-	-
	VSLNet [114]	2020	26.16	43.22	63.16	-	-	-
	SQAN [60]	2020	23.07	41.51	58.52	-	-	-
	DRN [109]	2020	24.36	45.45	-	50.30	77.97	-
	HVTG [11]	2020	18.27	40.15	57.60	-	-	-
	APGN [44]	2021	28.64	48.92	-	63.19	78.87	-
	CLEAR [28]	2021	28.05	45.33	59.96	62.13	77.26	85.83
	IA-Net [17]	2021	27.95	48.57	67.14	63.12	78.99	87.21
	SMIN [87]	2021	30.34	48.46	-	62.11	<b>81.16</b>	-
	RaNet [17]	2021	28.67	45.59	62.97	<b>75.93</b>	-	-
	MATN [116]	2021	<b>31.78</b>	48.02	-	63.18	78.02	-
	SSMN [54]	2021	20.03	35.38	52.76	-	-	-
	IVG [62]	2021	27.10	43.84	63.22	-	-	-
	CPN [124]	2021	28.10	45.10	62.81	-	-	-
	CBLN [45]	2021	27.60	48.12	66.34	63.41	79.32	<b>88.91</b>
	DRFT [12]	2021	27.79	45.72	62.91	-	-	-
	LP-Net [97]	2021	25.39	45.92	64.29	-	-	-
	MQEI [83]	2021	24.58	45.86	64.39	-	-	-
	MS-2D-TAN [117]	2021	29.21	46.16	61.04	60.85	78.80	87.30
	FVMR [18]	2021	26.85	45.00	60.63	61.04	77.42	86.11
	CMHN [27]	2021	24.02	43.47	62.49	53.16	73.42	85.37
SV-VMR [96]	2021	27.32	45.21	61.39	63.44	77.10	85.98	
CP-Net [36]	2021	21.63	40.56	-	-	-	-	
DORi [70]	2021	26.41	41.35	57.89	-	-	-	
RL	RWM [24]	2019	-	36.90	-	-	-	-
	TSP-PRL [95]	2020	-	38.82	<b>56.02</b>	-	-	-
	TripNet [23]	2020	13.93	32.19	48.42	-	-	-
	MABAN [78]	2021	<b>23.05</b>	<b>40.01</b>	-	-	-	-

**DiDeMo [2].** The Distinct Describable Moments (DiDeMo) dataset is recently proposed in [2]. It contains 10,464 videos with 40,543 annotated queries. To annotate moment-query pairs, videos are trimmed to a maximum of 30 seconds and then divided into 6 segments with 5 seconds long each. Besides, each moment in this dataset is constructed by one or more consecutive segments.

Therefore, there are 21 moment candidates in each video, and the task is shifted to selecting the moment that best matches the query.

**TEMPO [25].** The TEMPO dataset is collected based on the DiDeMo dataset [2]. Specifically, it further extends the language descriptions that involve multiple events, while keeping its videos the same. The extended language expressions are collected based on four commonly used temporal words, before, after, while, and then. Simple sentences that come from DiDeMo are also included in this dataset. There are two parts in TEMPO, i.e., TEMPO-TL (Template Language) which is constructed by the original DiDeMo sentences with language templates and TEMPO-HL (Human Language) which is built by human annotations. Note that since the download link of this dataset is invalid currently, we hence cannot obtain the statistic result of the queries.

**Charades-STA [16].** The Charades dataset [75] is first collected from daily indoor activities for activity understanding. Each video contains temporal activity annotation (from 157 activity categories) and multiple video-level descriptions. To make it suitable for language-based temporal location task, Gao et al. [16] decomposed the original video-level descriptions into shorter sub-sentences, and performed keyword matching to assign them to temporal segments in videos. The alignment annotations are further verified manually. In total, there are 9, 848 videos with 16,128 annotated queries in this dataset, and these videos are 31 seconds long on average.

**TACoS [67].** This dataset is constructed by [68]. It was built on the top of MPII Compositive dataset [72], which contains different activities in the cooking domain. In TACoS, each video is associated with two type of annotations. The first one is fine-grained activity labels with temporal location (start and end time). The second is natural language descriptions with temporal locations. Note that the natural language descriptions are obtained by crowd-sourcing annotators, who are asked to describe the content of the video clips by sentences. In total, there are 127 videos picturing people who perform cooking tasks with 18,818 queries.

**ActivityNet-Captions [31].** This dataset is proposed by Krishna et al. [31] for the dense video captioning task, which contains annotations from the open domain. In this dataset, each video contains at least two ground truth segments, and each segment is paired with one ground truth caption. In total, this dataset contains 19, 209 videos and 71,957 queries.

## 6.2 Evaluation Metrics

In existing studies, “ $R@n, IoU=m$ ” proposed by [26] is commonly adopted as the evaluation metric to measure their performance. It is defined as the percentage of at least one of the top- $n$  predicted moments which have IoU with ground-truth moment larger than  $m$  [18]. In the following, we use  $R(n, m)$  to denote “ $R@n, IoU=m$ ”. Note that for DiDeMo and TEMPO datasets, existing methods measure their performance with Rank@1 ( $R@1$ ), Rank@5 ( $R@5$ ), and mean intersection over union (mIoU).

## 7 EXPERIMENTAL RESULTS

### 7.1 Supervised Localization Models

Supervised video moment localization models are evaluated on the DiDeMo, TEMPO, TACoS, Charades-STA, and ActivityNet-Captions datasets. We directly summarized their experimental results from the corresponding papers in Tables 4-9.

*7.1.1 Two-stage ones.* Among two-stage supervised video moment localization models, the hand-crafted heuristics based ones (i.e., MCN, MLLC, and TCMN) are mainly evaluated on DiDeMo and TEMPO. Their experimental results are reported in Table 4, 5, and 6. As one of the earliest introduced methods, MCN is considered the de facto baseline result. MLLC outperforms the baseline model MCN, suggesting that learning to reason about which context moment is correct (as opposed

to considering global video context) is beneficial. TCMN exhibits the promising performance across all the metrics of the complex sentence and comparative results in simple sentences on Tempo-HL. The results show that the compositional modeling of complex queries can improve the localization performance.

As to video moment localization approaches that utilize multi-scale sliding windows to generate moment candidates (i.e., CTRL, ACRN, VAL, SLTA, MMRG, ROLE, MCF, ACL, and MIGCN), experiments are primarily executed on TACoS and Charades-STA. The experimental results are reported in Table 7 and 8. We can see that: **1)** ROLE achieves better performance as compared to CTRL on Charades-STA, verifying the importance of visual-aware query modeling. **2)** Although ACRN, VAL, SLTA, and MMRG all focus on improving visual modeling, MMRG achieves superior performance on both TACoS and Charades-STA. This mainly because it jointly considers the object relations and the phrase relations modeling for video moment localization, meanwhile, it introduces pre-training tasks to enhance the visual representation. **3)** Compared with MCF that also focus on intra-modal interaction modeling, ACL achieves remarkable performance improvements on both TACoS and Charades-STA. This reflects that the pair of activity concepts extracted from both videos and queries play a vital role in improving the cross-modal alignment. **4)** MIGCN outperforms ACL on Charades-STA, justifying the necessity of capturing both intra- and inter-modal interaction information. And **5)** MMRG achieves superior performance, and the results are competitive to other multi-scale sliding windows based methods. One possible reason is that it adopts two self-supervised pre-training tasks: attribute masking and context prediction to alleviate semantic gap across modalities.

In contrast, the methods that leverages moment generation networks are mainly evaluated on Charades-STA and ActivityNet-Captions. As reported in Table 8 and 9, BPNet outperforms both QSPN and SAP. Because it utilizes VSLNet to generate high-quality moment candidates, therefore boosting the localization accuracy. However, the localization accuracy of BPNet is significantly worse than the MMRG. The reasons may be that **1)** Compared with VSLNet, the multi-scale sliding windows generate higher quality moment candidates. And **2)** MMRG could identify the fine-grained differences among similar video moment candidates. Despite achieving promising performance, these two-stage approaches should pre-process the untrimmed videos to obtain moment candidates, therefore they are inferior in efficiency.

**7.1.2 One-stage ones.** We summarized experimental results of all discussed one-stage methods in Table 6-9. From these results, we could find that: **1)** Among anchor-based one-stage methods, FIAN achieves the best performance on all datasets. Particular, the performance of FIAN significantly surpasses TGN, demonstrating the importance of inter-modal interaction modeling. Moreover, compared with other inter-modal interaction modeling methods, FIAN achieves better performance. This verifies the effectiveness of the iterative attention mechanism. More importantly, FIAN even achieves superior performance as compared to methods that jointly consider intra- and inter-modal interaction modeling. This may be because it iteratively captures bilateral query-video interaction information. Furthermore, CMIN-R performs better than CMIN, verifying that reconstructing the natural language queries could indeed enhance the cross-modal representations. **2)** For sampler-based localization methods, SMIN achieves the best performance on Charades-STA. However, MATN achieves the best performance on TACoS and ActivityNet-Captions and significantly surpasses the work SMIN. On one hand, this reflects that the enumeration strategy could generate higher quality moment candidates as compared to other strategies. On the other hand, it demonstrates the powerful ability of the visual-language transformer for learning modality representations. **3)** For proposal-free methods, CP-Net achieves the best performance on TACoS, but its performance is slightly worse on the other two datasets. This is mainly because the video length of TACoS is

Table 10. Performance Comparison among weakly supervised localization models on DiDeMo.

Type	Model	Year	R@1	R@5	mIoU
Two-stage	TGA [59]	2019	12.19	39.74	24.92
	WSLLN [20]	2019	18.40	54.40	27.40
	VLANet [56]	2020	19.32	<b>65.68</b>	25.33
	LoGAN [81]	2021	<b>39.20</b>	64.04	<b>38.28</b>
One-stage	RTBPN [123]	2020	<b>20.79</b>	<b>60.26</b>	29.81
	WSTAN [91]	2021	19.40	54.64	<b>31.94</b>

Table 11. Performance comparison among weakly supervised and unsupervised localization models on Charades-STA.

Type	Model	Year	R(1,0.7)	R(1,0.5)	R(1,0.3)	R(5,0.7)	R(5,0.5)	R(5,0.3)
Two-stage	TGA [59]	2019	8.84	19.94	32.14	33.51	65.52	86.58
	VLANet [56]	2020	14.17	31.83	45.24	33.09	<b>82.85</b>	<b>95.70</b>
	LoGAN [81]	2021	<b>14.54</b>	<b>34.68</b>	<b>51.67</b>	<b>39.11</b>	74.30	92.74
One-stage	SCN [40]	2020	9.97	23.58	42.96	38.87	71.80	95.56
	RTBPN [123]	2020	13.24	32.36	<b>60.04</b>	41.18	71.85	97.48
	LCNet [103]	2021	18.87	<b>39.19</b>	59.60	<b>45.24</b>	<b>80.56</b>	94.78
	VCA [92]	2021	<b>19.57</b>	38.13	58.58	37.75	78.75	<b>98.08</b>
	WSTAN [91]	2021	12.28	29.35	43.39	41.53	76.13	93.04
	MS-2D-RL [33]	2021	17.31	30.38	-	34.92	69.60	-
RL	BAR [94]	2020	12.23	27.04	44.97	-	-	-
Unsupervised	U-VMR [19]	2021	8.27	20.14	<b>46.69</b>	32.45	72.07	91.18
	PSVL [61]	2021	<b>14.17</b>	<b>31.29</b>	46.47	-	-	-

longer than the other two datasets, namely there are lots of visual similar moment candidates. CP-Net focuses on exploiting the fine-grained temporal clues to enhance the discriminative of different moments, therefore achieving the better performance on TACoS. As DORi could capture complex relationships between humans, objects and activities in the video, it outperforms CP-Net on Charades-STA which is collected for activity understanding. Compared to TACoS and Charades-STA datasets, ActivityNet-Captions is collected from Youtube, of which the videos contain multi-modal information. Thereby, DRFT achieves the best performance on ActivityNet-Captions since it utilizes multimodal information and adequately exploit interactions between modalities. And 4) among all one-stage methods, MATN achieves superior performance as compared to FIAN, CP-Net, DRFT, and DORi. This reflects that the performance gap between moment generation based methods and proposal-free methods is still large. Thereby, for proposal-free methods, it is worth to explore more effective interaction strategies to further improve localization accuracy.

**7.1.3 Reinforcement learning.** Reinforcement learning based localization models alleviate the efficiency issue to a certain extent, yet their performance is inferior as reported in Table 6-9. The main reason may be that they mostly focus on the design of policy and rewards, ignoring the importance of multiple crucial factors, such as query representation learning, video context modeling, and multimodal fusion.

## 7.2 Weakly Supervised and Unsupervised Localization Models

Weakly supervised video moment localization approaches are mainly evaluated on the DiDeMo, Charadea-STA, and ActivityNet-Captions datasets. We directly summarized their experimental results from the corresponding papers in Table 10, 11, and 12. TGA is the first weakly supervised video moment localization model, considered as the de facto baseline result. Among the two-stage

Table 12. Performance comparison among weakly supervised and unsupervised localization models on ActivityNet-Captions.

Type	Model	Year	R(1,0.5)	R(1,0.3)	R(1,0.1)	R(5,0.5)	R(5,0.3)	R(5,0.1)
Two-stage	WSLLN [20]	2019	22.70	42.80	75.40	-	-	-
One-stage	SCN [40]	2020	29.22	47.23	71.48	55.69	71.45	90.88
	RTBPN [123]	2020	-	29.63	49.77	-	60.59	79.89
	LCNet [103]	2021	26.33	48.49	78.58	62.66	<b>82.51</b>	93.95
	VCA [92]	2021	<b>31.00</b>	50.45	67.96	53.83	71.79	92.14
	WSTAN [91]	2021	30.01	<b>52.45</b>	<b>79.78</b>	<b>63.42</b>	79.38	<b>93.15</b>
	MS-2D-RL [33]	2021	-	29.68	49.79	-	58.66	72.57
RL	BAR [94]	2020	30.73	49.03	-	-	-	-
Unsupervised	U-VMR [19]	2021	11.64	26.38	<b>46.15</b>	30.83	54.27	73.13
	PSVL [61]	2021	<b>14.74</b>	<b>30.08</b>	44.74	-	-	-

Table 13. Efficiency comparison among some video localization models on three datasets. TE: time cost of query Embedding. CML: time cost of the cross-modal learning. ALL: the total time cost of TE and CML.

Methods	TACoS			ActivityNet Captions			Charades-STA		
	TE	CML	ALL	TE	CML	ALL	TE	CML	ALL
PFGA [64]	1.14	11.37	12.51	1.24	8.97	10.21	1.15	4.37	5.52
VSLNet [114]	3.58	5.02	8.59	3.87	4.86	8.74	3.90	4.27	8.18
SQAN [60]	-	-	-	1.53	7.03	8.56	1.23	4.76	5.99
DRN [109]	4.67	22.13	26.81	4.86	18.46	23.32	4.52	12.39	16.91
CTRL [16]	4.32	534.23	538.55	4.75	398.25	403.00	4.53	12.20	16.73
SCDM [107]	3.65	780.00	783.65	3.27	359.76	363.03	2.97	23.77	26.07
CBP [88]	3.17	2659.01	2662.18	2.44	522.65	525.09	2.87	266.08	268.95
2D-TAN [118]	1.72	135.84	137.56	1.69	80.35	403.10	1.59	16.78	18.37
FVMR [18]	3.51	0.14	3.65	3.14	0.09	3.23	2.86	0.01	2.87

approaches (i.e., WSLLN, LoGAN, and VLANet), LoGAN achieves the best performance on DiDeMo and Charadea-STA, as compared to WSLLN and VLANet. This reflects the importance of learning contextualized visual semantic representations for weakly-supervised moment localization. As to the one-stage models, VCA achieves the best performance in terms of “R(1, 0.7)” on both Charadea-STA and ActivityNet-Captions. This is mainly because enumeration based methods would generate overmuch moment candidates, while they cannot well learn discriminative and robust visual features under video-level supervision. Differently, VCA utilizes the similarity among sentences to mine positive pairs as well as negatives pairs for contrastive learning, therefore learning more discriminative visual features. Besides, BAR as the first work that extends reinforcement learning to weakly supervised video moment localization also achieves promising performance. Particular, it outperforms TGA and SCAN on Charadea-STA and ActivityNet-Captions.

It is noteworthy that unsupervised video moment localization approaches achieve favorable performance on both Charades-STA and ActivityNet-Captions. Specifically, PSVL obtains comparable results against some supervised and weakly-supervised approaches on Charades-STA, such as ACL, SAP, VSA-STV, RTBPN, and WSTAN. This verifies that even without any moment level or video-level annotations, we can also obtain acceptable localization results by carefully designing the effective video-conditioned sentence generator.

### 7.3 Efficiency Comparison

[18] conducts efficiency comparison among some video moment localization methods during inference, the results are summarized in Table 13. We can see that: 1) TE is not the test-time computational bottleneck since all methods cost similar time (~3ms). 2) CBP spends the most

time in CML ( $\sim 2600$ ms), demonstrating that self-attention technique as the contextual integration module for CML is time consuming. And 3) FVMR is  $35\times$  to  $20,000\times$  faster than other methods, especially one-stage methods VSLNet and SQAN. This reflects that learning cross-modal common space is much more efficient than cross-modal interaction. Moreover, FVMR outperforms others in both speed and accuracy metrics, verifying that introducing common space learning strategy could well balance the efficiency and accuracy.

## 8 DISCUSSION AND FUTURE DIRECTIONS

The introduction of video moment localization task has aroused great interest, yet it is a challenging task since it requires joint reasoning over the visual and textual information. To better localize the target moment within the video, the video moment localization model should understand both the video and query information, such as objects, attributes, and actions, and then identify a particular moment via reasoning. Existing methods can be divided into supervised, weakly-supervised, and unsupervised learning paradigm based methods. Particular, the former two can be further classified into two-stage, one-stage, as well as reinforcement learning methods. We first discuss the main differences among them and point the possible technological trend of each type.

Two-stage video moment localization approaches commonly utilize a separate scheme to generate moment candidates, and then match them with the query. Therefore, their localization efficiency is relatively low, and the localization accuracy is restricted by the moments generated by the first stage. *In the future, it is a possible trend to design a more powerful moment generation model for the first-stage or a more efficient location regression module for the second-stage, to improve the localization efficiency and accuracy.* Compared with two-stage methods, one-stage ones do not need a separate stage for moment candidate generation, therefore they locate target moments faster. Some one-stage methods also utilize moment generation strategies to generate moment candidates, such as temporal convolution, but these processes are optimized together with the localization accuracy. This ensures the quality of generated moment candidates. However, existing models typically resort to various attention mechanisms to estimate the correlation between “video-query” pair, resulting in inefficiency and low scalability. *In the future, building a effective and efficient cross-modal learning module, replacing the cross-modal interaction module, is deserved to exploring.* More importantly, both simple and complex input data are processed by the same interaction module, which is inappropriate. In light of this, *building a novel interaction module, which could dynamically adopt different mechanisms to explore correlation information for different “video-query” pairs, is a possible technological trend.* Different from the two- and one-stage methods, the reinforcement learning methods formulate this task as a problem of sequential decision making by learning an agent which regulates the temporal grounding boundaries progressively based on its policy. In other words, they can localize more accurate temporal locations in a few glimpses of video. Although they are more efficient, their localization accuracy is still far from satisfactory due to the insufficient structural comprehension. *One possible solution may be construct more efficient cross-modal interaction module to explore rich contextual information.*

Afterwards, we further point out several possible future directions to advance this research area. Current state-of-the-art models suffer serious dataset bias problem [57]. To be specific, on one hand, many queries describe the moments that appear at the start or end of the video. This imbalance makes the model can only predict the location of moments according to the statistical correlations without understanding the visual contents. On the other hand, as stated by Mithun et al. [59], it is difficult to mark the start and end locations of a certain moment. This may introduce ambiguity in the training data and influence the alignment between visual and textual information. More importantly, sentence queries are typically simple and short, they mainly focus on one object and one action. In addition, many videos in existing datasets contain limited information, which require

no complex reasoning. Therefore, *building a large-scale dataset with long videos and semantic-rich queries would be an interesting research direction for this task*. Moreover, existing models tend to give a direct location prediction without an intermediate reasoning process. Therefore, it is difficult for people to evaluate the reasoning capability of models and analyze the localization results of the model. *A promising research direction is to construct interpretable video moment localization models*. Although recent approaches have achieved significant progress, their performance is still far below that of humans. This is because humans have extensive domain knowledge. In other words, they can employ the corresponding background knowledge to successfully localize the target moment from the video, given the specific query. Thereby, *it would be an interesting and promising direction to augment the video moment localization models with external knowledge*.

## 9 CONCLUSION

In this paper, we comprehensively review the state-of-the-art methods on video moment localization. To be more specific, we first review supervised learning based approaches, including two-stage, one-stage, and reinforcement learning models. We then describe the recently emerging weakly supervised and unsupervised video moment localization methods. After reviewing different video moment localization datasets, we group results according to the datasets. Finally, we figure out a number of promising directions for future research.

## ACKNOWLEDGMENTS

This work is supported by the National Natural Science Foundation of China, No.: U1936203 and No.: 62006142; the Shandong Provincial Natural Science Foundation for Distinguished Young Scholars, No.: ZR2021JQ26; the Major Basic Research Project of Natural Science Foundation of Shandong Province, No.: ZR2021ZD15; Science and Technology Innovation Program for Distinguished Young Scholars of Shandong Province Higher Education Institutions, No.: 2021KJ036; as well as the special fund for distinguished Professors of Shandong Jianzhu University.

## REFERENCES

- [1] Humam Alwassel, Fabian Caba Heilbron, Victor Escorcia, and Bernard Ghanem. 2018. Diagnosing Error in Temporal Action Detectors. In *Proceedings of the European Conference on Computer Vision*. 256–272.
- [2] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. 2017. Localizing Moments in Video with Natural Language. In *Proceedings of the IEEE International Conference on Computer Vision*. 5803–5812.
- [3] Hakan Bilen and Andrea Vedaldi. 2016. Weakly Supervised Deep Detection Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2846–2854.
- [4] J. Burgner-Kahrs, D. C. Rucker, and H. Choset. 2015. Continuum Robots for Medical Applications: A Survey. *IEEE Transactions on Robotics* 31, 6 (2015), 1261–1280.
- [5] Da Cao, Yawen Zeng, Meng Liu, Xiangnan He, Meng Wang, and Zheng Qin. 2020. STRONG: Spatio-Temporal Reinforcement Learning for Cross-Modal Video Moment Localization. In *Proceedings of the ACM International Conference on Multimedia*. 4162–4170.
- [6] Da Cao, Yawen Zeng, Xiaochi Wei, Liqiang Nie, Richang Hong, and Zheng Qin. 2020. Adversarial Video Moment Retrieval by Jointly Modeling Ranking and Localization. In *Proceedings of the ACM International Conference on Multimedia*. 898–906.
- [7] Jingyuan Chen, Xinpeng Chen, Lin Ma, Zequn Jie, and Tat-Seng Chua. 2018. Temporally Grounding Natural Sentence in Video. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. 162–171.
- [8] Jingyuan Chen, Lin Ma, Xinpeng Chen, Zequn Jie, and Jiebo Luo. 2019. Localizing Natural Language in Videos. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 8175–8182.
- [9] Shaoxiang Chen, Wenhao Jiang, Wei Liu, and Yu-Gang Jiang. 2020. Learning Modality Interaction for Temporal Sentence Localization and Event Captioning in Videos. In *Proceedings of the European Conference on Computer Vision*. 333–351.
- [10] Shaoxiang Chen and Yugang Jiang. 2019. Semantic Proposal for Activity Localization in Videos via Sentence Query. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 8199–8206.

- [11] Shaoxiang Chen and Yu-Gang Jiang. 2020. Hierarchical Visual-Textual Graph for Temporal Activity Localization via Language. In *European Conference on Computer Vision*. 601–618.
- [12] Yi-Wen Chen, Yi-Hsuan Tsai, and Ming-Hsuan Yang. 2021. End-to-end multi-modal video temporal grounding. *Advances in Neural Information Processing Systems* 34 (2021), 28442–28453.
- [13] Zhenfang Chen, Lin Ma, Wenhan Luo, and Kwan-Yee Kenneth Wong. 2019. Weakly-Supervised Spatio-Temporally Grounding Natural Sentence in Video. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. 1884–1894.
- [14] L. Claussmann, M. Revilloud, D. Gruyer, and S. Glaser. 2020. A Review of Motion Planning for Highway Autonomous Driving. *IEEE Transactions on Intelligent Transportation Systems* 21, 5 (2020), 1826–1848.
- [15] Victor Escorcia, Mattia Soldan, Josef Sivic, Bernard Ghanem, and Bryan C. Russell. 2019. Temporal Localization of Moments in Video Collections with Natural Language. *CoRR* (2019), 1–17.
- [16] Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. 2017. TALL: Temporal Activity Localization via Language Query. In *Proceedings of the IEEE International Conference on Computer Vision*. 5267–5275.
- [17] Jialin Gao, Xin Sun, Mengmeng Xu, Xi Zhou, and Bernard Ghanem. 2021. Relation-aware Video Reading Comprehension for Temporal Language Grounding. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. 3978–3988.
- [18] Junyu Gao and Changsheng Xu. 2021. Fast Video Moment Retrieval. In *Proceedings of the IEEE International Conference on Computer Vision*. 1523–1532.
- [19] Junyu Gao and Changsheng Xu. 2021. Learning Video Moment Retrieval Without a Single Annotated Video. *IEEE Transactions on Circuits and Systems for Video Technology* 32, 3 (2021), 1646–1657.
- [20] Mingfei Gao, Larry Davis, Richard Socher, and Caiming Xiong. 2019. WSLN: Weakly Supervised Natural Language Localization Networks. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. 1481–1487.
- [21] Runzhou Ge, Jiyang Gao, Kan Chen, and Ram Nevatia. 2019. MAC: Mining Activity Concepts for Language-Based Temporal Localization. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*. 245–253.
- [22] Soham Ghosh, Anuva Agarwal, Zarana Parekh, and Alexander G. Hauptmann. 2019. ExCL: Extractive Clip Localization Using Natural Language Descriptions. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 1984–1990.
- [23] Meera Hahn, Asim Kadav, James M. Rehg, and Hans Peter Graf. 2020. Tripping through time: Efficient Localization of Activities in Videos. In *Proceedings of the British Machine Vision Conference*. 1–16.
- [24] Dongliang He, Xiang Zhao, Jizhou Huang, Fu Li, Xiao Liu, and Shilei Wen. 2019. Read, Watch, and Move: Reinforcement Learning for Temporally Grounding Natural Language Descriptions in Videos. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 8393–8400.
- [25] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan C. Russell. 2018. Localizing Moments in Video with Temporal Language. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. 1380–1390.
- [26] Ronghang Hu, Huazhe Xu, Marcus Rohrbach, Jiashi Feng, Kate Saenko, and Trevor Darrell. 2016. Natural Language Object Retrieval. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4555–4564.
- [27] Yupeng Hu, Meng Liu, Xiaobin Su, Zan Gao, and Liqiang Nie. 2021. Video Moment Localization via Deep Cross-Modal Hashing. *IEEE Transactions on Image Processing* 30 (2021), 4667–4677.
- [28] Yupeng Hu, Liqiang Nie, Meng Liu, Kun Wang, Yinglong Wang, and Xian-Sheng Hua. 2021. Coarse-to-fine Semantic Alignment for Cross-modal Moment Localization. *IEEE Transactions on Image Processing* 30 (2021), 5933–5943.
- [29] Bin Jiang, Xin Huang, Chao Yang, and Junsong Yuan. 2019. Cross-Modal Video Moment Retrieval with Spatial and Language-Temporal Attention. In *Proceedings of the International Conference on Multimedia Retrieval*. 217–225.
- [30] Ryan Kiros, Yukun Zhu, Russ R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Skip-Thought Vectors. In *Proceedings of the Advances in Neural Information Processing Systems*. 3294–3302.
- [31] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. 2017. Dense-Captioning Events in Videos. In *Proceedings of the IEEE International Conference on Computer Vision*. 706–715.
- [32] Cheng Li, Yuming Zhao, Shihao Peng, and Jinting Chen. 2019. Bidirectional Single-Stream Temporal Sentence Query Localization in Untrimmed Videos. In *Proceedings of the IEEE International Conference on Image Processing*. 270–274.
- [33] Ding Li, Rui Wu, Yongqiang Tang, Zhizhong Zhang, and Wensheng Zhang. 2021. Multi-scale 2D Representation Learning for Weakly-supervised Moment Retrieval. In *2020 25th International Conference on Pattern Recognition*. 8616–8623.
- [34] Jianing Li, Jingdong Wang, Qi Tian, Wen Gao, and Shiliang Zhang. 2019. Global-local Temporal Representations for Video Person Re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*. 3958–3967.
- [35] Jianing Li, Shiliang Zhang, and Tiejun Huang. 2020. Multi-scale Temporal Cues Learning for Video Person Re-identification. *IEEE Transactions on Image Processing* 29 (2020), 4461–4473.

- [36] Kun Li, Dan Guo, and Meng Wang. 2021. Proposal-Free Video Grounding with Contextual Pyramid Network. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 1902–1910.
- [37] Guoqiang Liang, Shiyu Ji, and Yanning Zhang. 2021. Local-enhanced Interaction for Temporal Moment Localization. In *Proceedings of the International Conference on Multimedia Retrieval*. 201–209.
- [38] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. 2016. Continuous Control with Deep Reinforcement Learning. In *International Conference on Learning Representations*. 1–14.
- [39] Tianwei Lin, Xu Zhao, and Zheng Shou. 2017. Single Shot Temporal Action Detection. In *Proceedings of the ACM International Conference on Multimedia*. 988–996.
- [40] Zhijie Lin, Zhou Zhao, Zhu Zhang, Qi Wang, and Huasheng Liu. 2020. Weakly-Supervised Video Moment Retrieval via Semantic Completion Network. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 11539–11546.
- [41] Z. Lin, Z. Zhao, Z. Zhang, Z. Zhang, and D. Cai. 2020. Moment Retrieval via Cross-Modal Interaction Networks With Query Reconstruction. *IEEE Transactions on Image Processing* 29 (2020), 3750–3762.
- [42] Bingbin Liu, Serena Yeung, Edward Chou, Dean Huang, Li Feifei, and Juan Carlos Nieves. 2018. Temporal Modular Networks for Retrieving Complex Compositional Activities in Videos. In *Proceedings of the European Conference on Computer Vision*. 569–586.
- [43] Daizong Liu, Xiaoye Qu, Jianfeng Dong, and Pan Zhou. 2020. Reasoning Step-by-Step: Temporal Sentence Localization in Videos via Deep Rectification-Modulation Network. In *Proceedings of the International Conference on Computational Linguistics*. 1841–1851.
- [44] Daizong Liu, Xiaoye Qu, Jianfeng Dong, and Pan Zhou. 2021. Adaptive Proposal Generation Network for Temporal Sentence Localization in Videos. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. 9292–9301.
- [45] Daizong Liu, Xiaoye Qu, Jianfeng Dong, Pan Zhou, Yu Cheng, Wei Wei, Zichuan Xu, and Yulai Xie. 2021. Context-aware Biaffine Localizing Network for Temporal Sentence Grounding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 11235–11244.
- [46] Daizong Liu, Xiaoye Qu, Xiao-Yang Liu, Jianfeng Dong, Pan Zhou, and Zichuan Xu. 2020. Jointly Cross- and Self-Modal Graph Attention Network for Query-Based Moment Localization. In *Proceedings of the ACM International Conference on Multimedia*. 4070–4078.
- [47] Daizong Liu, Xiaoye Qu, and Pan Zhou. 2021. Progressively Guide to Attend: An Iterative Alignment Framework for Temporal Sentence Grounding. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. 9302–9311.
- [48] M. Liu, L. Nie, X. Wang, Q. Tian, and B. Chen. 2019. Online Data Organizer: Micro-Video Categorization by Structure-Guided Multimodal Dictionary Learning. *IEEE Transactions on Image Processing* 28, 3 (2019), 1235–1247.
- [49] M. Liu, L. Qu, L. Nie, M. Liu, L. Duan, and B. Chen. 2020. Iterative Local-Global Collaboration Learning Towards One-Shot Video Person Re-Identification. *IEEE Transactions on Image Processing* 29 (2020), 9360–9372.
- [50] Meng Liu, Xiang Wang, Liqiang Nie, Xiangnan He, Baoquan Chen, and Tat-Seng Chua. 2018. Attentive Moment Retrieval in Videos. In *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval*. 15–24.
- [51] Meng Liu, Xiang Wang, Liqiang Nie, Qi Tian, Baoquan Chen, and Tat-Seng Chua. 2018. Cross-Modal Moment Localization in Videos. In *Proceedings of the ACM International Conference on Multimedia*. 843–851.
- [52] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C. Berg. 2016. SSD: Single Shot MultiBox Detector. In *European Conference on Computer Vision*. 21–37.
- [53] Xinfang Liu, Xiushan Nie, Zhifang Tan, Jie Guo, and Yilong Yin. 2021. A Survey on Natural Language Video Localization. *arXiv preprint arXiv:2104.00234* (2021), 1–13.
- [54] Xinfang Liu, Xiushan Nie, Junya Teng, Li Lian, and Yilong Yin. 2021. Single-shot Semantic Matching Network for Moment Localization in Videos. *ACM Transactions on Multimedia Computing, Communications, and Applications* 17, 3 (2021), 1–14.
- [55] Chujie Lu, Long Chen, Chile Tan, Xiaolin Li, and Jun Xiao. 2019. Debug: Dense Bottom-Up Grounding Approach for Natural Language Video Localization. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing*. 5144–5153.
- [56] Minuk Ma, Sunjae Yoon, Junyeong Kim, Youngjoon Lee, Sunghun Kang, and Chang D. Yoo. 2020. VLANet: Video-Language Alignment Network for Weakly-Supervised Video Moment Retrieval. In *Proceedings of the European Conference on Computer Vision*. 156–171.
- [57] Esa Rahtu Mayu Otani, Yuta Nakahima and Janne Heikkilä. 2020. Uncovering Hidden Challenges in Query-Based Video Moment Retrieval. In *The British Machine Vision Conference*. 1–12.
- [58] A. Miech, D. Zhukov, J. Alayrac, M. Tapaswi, I. Laptev, and J. Sivic. 2019. HowTo100M: Learning a Text-Video Embedding by Watching Hundred Million Narrated Video Clips. In *Proceedings of the IEEE International Conference*

on *Computer Vision*. 2630–2640.

- [59] Niluthpol Chowdhury Mithun, Sujoy Paul, and Amit K. Roy-Chowdhury. 2019. Weakly Supervised Video Moment Retrieval From Text Queries. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 11592–11601.
- [60] Jonghwan Mun, Minsu Cho, and Bohyung Han. 2020. Local-Global Video-Text Interactions for Temporal Grounding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 10807–10816.
- [61] Jinwoo Nam, Daechul Ahn, Dongyeop Kang, Seong Jong Ha, and Jonghyun Choi. 2021. Zero-shot Natural Language Video Localization. In *Proceedings of the IEEE International Conference on Computer Vision*. 1470–1479.
- [62] Guoshun Nan, Rui Qiao, Yao Xiao, Jun Liu, Sicong Leng, Hao Zhang, and Wei Lu. 2021. Interventional Video Grounding with Dual Contrastive Learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2765–2775.
- [63] Ke Ning, Lingxi Xie, Jianzhuang Liu, Fei Wu, and Qi Tian. 2021. Interaction-Integrated Network for Natural Language Moment Localization. *IEEE Transactions on Image Processing* 30 (2021), 2538–2548.
- [64] Cristian Rodriguez Opazo, Edison Marrese-Taylor, Fatemeh Sadat Saleh, Hongdong Li, and Stephen Gould. 2020. Proposal-free Temporal Moment Localization of a Natural-Language Query in Video using Guided Attention. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*. 2453–2462.
- [65] Zhiwu Qing, Haisheng Su, Weihao Gan, Dongliang Wang, Wei Wu, Xiang Wang, Yu Qiao, Junjie Yan, Changxin Gao, and Nong Sang. 2021. Temporal Context Aggregation Network for Temporal Action Proposal Refinement. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 485–494.
- [66] Xiaoye Qu, Pengwei Tang, Zhikang Zou, Yu Cheng, Jianfeng Dong, Pan Zhou, and Zichuan Xu. 2020. Fine-Grained Iterative Attention Network for Temporal Language Localization in Videos. In *Proceedings of the ACM International Conference on Multimedia*. 4280–4288.
- [67] Michaela Regneri, Marcus Rohrbach, Dominikus Wetzels, Stefan Thater, Bernt Schiele, and Manfred Pinkal. 2013. Grounding Action Descriptions in Videos. *Transactions of the Association for Computational Linguistics* 1 (2013), 25–36.
- [68] Michaela Regneri, Marcus Rohrbach, Dominikus Wetzels, Stefan Thater, Bernt Schiele, and Manfred Pinkal. 2013. Grounding Action Descriptions in Videos. *Transactions of the Association for Computational Linguistics* 1 (2013), 25–36.
- [69] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. In *Proceedings of the Advances in Neural Information Processing Systems*. 91–99.
- [70] Cristian Rodriguez-Opazo, Edison Marrese-Taylor, Basura Fernando, Hongdong Li, and Stephen Gould. 2021. DORI: Discovering Object Relationships for Moment Localization of a Natural Language Query in a Video. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*. 1079–1088.
- [71] Anna Rohrbach, Marcus Rohrbach, Ronghang Hu, Trevor Darrell, and Bernt Schiele. 2016. Grounding of Textual Phrases in Images by Reconstruction. In *European Conference on Computer Vision*. 817–834.
- [72] Marcus Rohrbach, Michaela Regneri, Mykhaylo Andriiuka, Sikandar Amin, Manfred Pinkal, and Bernt Schiele. 2012. Script Data for Attribute-Based Recognition of Composite Activities. In *Proceedings of the European Conference on Computer Vision*. 144–157.
- [73] Arka Sadhu, Kan Chen, and Ram Nevatia. 2020. Video Object Grounding Using Semantic Roles in Language Description. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 10417–10427.
- [74] Gunnar A. Sigurdsson, Olga Russakovsky, and Abhinav Gupta. 2017. What Actions Are Needed for Understanding Human Actions in Videos?. In *Proceedings of the IEEE International Conference on Computer Vision*. 2137–2146.
- [75] Gunnar A Sigurdsson, Gul Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav Gupta. 2016. Hollywood in homes: Crowdsourcing Data Collection for Activity Understanding. In *Proceedings of the European Conference on Computer Vision*. 510–526.
- [76] Xiaomeng Song and Yahong Han. 2018. VAL: Visual-Attention Action Localizer. In *Proceedings of the Pacific-Rim Conference on Multimedia*, Vol. 11165. 340–350.
- [77] Haisheng Su, Weihao Gan, Wei Wu, Yu Qiao, and Junjie Yan. 2021. BSN++: Complementary Boundary Regressor with Scale-Balanced Relation Modeling for Temporal Action Proposal Generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 2602–2610.
- [78] Xiaoyang Sun, Hanli Wang, and Bin He. 2021. MABAN: Multi-Agent Boundary-Aware Network for Natural Language Moment Retrieval. *IEEE Transactions on Image Processing* 30 (2021), 5589–5599.
- [79] Richard S Sutton and Andrew G Barto. 2018. *Reinforcement Learning: An Introduction*. MIT press.
- [80] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2016. Rethinking the Inception Architecture for Computer Vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2818–2826.

- [81] Reuben Tan, Huijuan Xu, Kate Saenko, and Bryan A Plummer. 2021. Logan: Latent Graph Co-attention Network for Weakly-Supervised Video Moment Retrieval. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*. 2083–2092.
- [82] Haoyu Tang, Jihua Zhu, Meng Liu, Zan Gao, and Zhiyong Cheng. 2021. Frame-wise Cross-modal Matching for Video Moment Retrieval. *IEEE Transactions on Multimedia* 24 (2021), 1338–1349.
- [83] Haoyu Tang, Jihua Zhu, Lin Wang, Qinghai Zheng, and Tianwei Zhang. 2021. Multi-Level Query Interaction for Temporal Language Grounding. *IEEE Transactions on Intelligent Transportation Systems* (2021), 1–10.
- [84] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. [n.d.]. Learning Spatiotemporal Features With 3D Convolutional Networks. In *Proceedings of the IEEE International Conference on Computer Vision*. 4489–4497.
- [85] Z. Tu, H. Li, D. Zhang, J. Dauwels, B. Li, and J. Yuan. 2019. Action-Stage Emphasized Spatiotemporal VLAD for Video Action Recognition. *IEEE Transactions on Image Processing* 28, 6 (2019), 2799–2812.
- [86] Hao Wang, Zheng-Jun Zha, Xuejin Chen, Zhiwei Xiong, and Jiebo Luo. 2020. Dual Path Interaction Network for Video Moment Localization. In *Proceedings of the ACM International Conference on Multimedia*. 4116–4124.
- [87] Hao Wang, Zheng-Jun Zha, Liang Li, Dong Liu, and Jiebo Luo. 2021. Structured Multi-Level Interaction Network for Video Moment Localization via Language Query. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 7026–7035.
- [88] Jingwen Wang, Lin Ma, and Wenhao Jiang. 2020. Temporally Grounding Language Queries in Videos by Contextual Boundary-Aware Prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 12168–12175.
- [89] Limin Wang, Yuanjun Xiong, Dahua Lin, and Luc Van Gool. 2017. Untrimmednets for Weakly Supervised Action Recognition and Detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4325–4334.
- [90] Weining Wang, Yan Huang, and Liang Wang. 2019. Language-Driven Temporal Activity Localization: A Semantic Matching Reinforcement Learning Model. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 334–343.
- [91] Yuechen Wang, Jiajun Deng, Wengang Zhou, and Houqiang Li. 2021. Weakly Supervised Temporal Adjacent Network for Language Grounding. *IEEE Transactions on Multimedia* (2021), 1–13.
- [92] Zheng Wang, Jingjing Chen, and Yu-Gang Jiang. 2021. Visual Co-Occurrence Alignment Learning for Weakly-Supervised Video Moment Retrieval. In *Proceedings of the ACM International Conference on Multimedia*. 1459–1468.
- [93] Aming Wu and Yahong Han. 2018. Multi-modal Circulant Fusion for Video-to-Language and Backward. In *Proceedings of the International Joint Conference on Artificial Intelligence*. 1029–1035.
- [94] Jie Wu, Guanbin Li, Xiaoguang Han, and Liang Lin. 2020. Reinforcement Learning for Weakly Supervised Temporal Grounding of Natural Language in Untrimmed Videos. In *Proceedings of the ACM International Conference on Multimedia*. 1283–1291.
- [95] Jie Wu, Guanbin Li, Si Liu, and Liang Lin. 2020. Tree-Structured Policy Based Progressive Reinforcement Learning for Temporally Language Grounding in Video. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 12386–12393.
- [96] Ziyue Wu, Junyu Gao, Shucheng Huang, and Changsheng Xu. 2021. Diving Into The Relations: Leveraging Semantic and Visual Structures For Video Moment Retrieval. In *IEEE International Conference on Multimedia and Expo*. 1–6.
- [97] Shaoning Xiao, Long Chen, Jian Shao, Yueting Zhuang, and Jun Xiao. 2021. Natural Language Video Localization with Learnable Moment Proposals. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. 4008–4017.
- [98] Shaoning Xiao, Long Chen, Songyang Zhang, Wei Ji, Jian Shao, Lu Ye, and Jun Xiao. 2021. Boundary Proposal Network for Two-Stage Natural Language Video Localization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 2986–2994.
- [99] Caiming Xiong, Victor Zhong, and Richard Socher. 2017. Dynamic Co-attention Networks for Question Answering. In *Proceedings of the International Conference on Learning Representations*. 1–14.
- [100] Huijuan Xu, Abir Das, and Kate Saenko. 2017. R-C3D: Region Convolutional 3D Network for Temporal Activity Detection. In *Proceedings of the IEEE International Conference on Computer Vision*. 5783–5792.
- [101] Huijuan Xu, Kun He, Bryan A Plummer, Leonid Sigal, Stan Sclaroff, and Kate Saenko. 2019. Multilevel Language and Vision Integration for Text-to-Clip Retrieval. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 9062–9069.
- [102] L. Yang, H. Peng, D. Zhang, J. Fu, and J. Han. 2020. Revisiting Anchor Mechanisms for Temporal Action Localization. *IEEE Transactions on Image Processing* 29 (2020), 8535–8548.
- [103] Wenfei Yang, Tianzhu Zhang, Yongdong Zhang, and Feng Wu. 2021. Local Correspondence Network for Weakly Supervised Temporal Sentence Grounding. *IEEE Transactions on Image Processing* 30 (2021), 3252–3262.

- [104] Xun Yang, Fuli Feng, Wei Ji, Meng Wang, and Tat-Seng Chua. 2021. Deconfounded Video Moment Retrieval with Causal Intervention. In *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1–10.
- [105] Yulan Yang, Zhaohui Li, and Gangyan Zeng. 2020. A Survey of Temporal Activity Localization via Language in Untrimmed Videos. In *2020 International Conference on Culture-oriented Science & Technology*. 596–601.
- [106] Y. Yang, J. Zhou, J. Ai, Y. Bin, A. Hanjalic, H. T. Shen, and Y. Ji. 2018. Video Captioning by Adversarial LSTM. *IEEE Transactions on Image Processing* 27, 11 (2018), 5600–5611.
- [107] Yitian Yuan, Lin Ma, Jingwen Wang, Wei Liu, and Wenwu Zhu. 2019. Semantic Conditioned Dynamic Modulation for Temporal Sentence Grounding in Videos. In *Proceedings of the Conference on Neural Information Processing Systems*. 534–544.
- [108] Yitian Yuan, Tao Mei, and Wenwu Zhu. 2019. To Find Where You Talk: Temporal Sentence Localization in Video with Attention based Location Regression. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 9159–9166.
- [109] Runhao Zeng, Haoming Xu, Wenbing Huang, Peihao Chen, Mingkui Tan, and Chuang Gan. 2020. Dense Regression Network for Video Grounding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 10284–10293.
- [110] Yawen Zeng, Da Cao, Xiaochi Wei, Meng Liu, Zhou Zhao, and Zheng Qin. 2021. Multi-Modal Relational Graph for Cross-Modal Video Moment Retrieval. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2215–2224.
- [111] Da Zhang, Xiyang Dai, Xin Wang, Yuanfang Wang, and Larry S. Davis. 2019. MAN: Moment Alignment Network for Natural Language Moment Retrieval via Iterative Graph Adjustment. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1247–1257.
- [112] Hao Zhang, Aixin Sun, Wei Jing, Guoshun Nan, Liangli Zhen, Joey Tianyi Zhou, and Rick Siow Mong Goh. 2021. Video Corpus Moment Retrieval with Contrastive Learning. In *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval*. 685–695.
- [113] Hao Zhang, Aixin Sun, Wei Jing, Liangli Zhen, Joey Tianyi Zhou, and Rick Siow Mong Goh. 2022. Natural Language Video Localization: A Revisit in Span-based Question Answering Framework. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44, 8 (2022), 4252–4266.
- [114] Hao Zhang, Aixin Sun, Wei Jing, and Joey Tianyi Zhou. 2020. Span-based Localizing Network for Natural Language Video Localization. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. 6543–6554.
- [115] J. Zhang, K. Mei, Y. Zheng, and J. Fan. 2019. Exploiting Mid-Level Semantics for Large-Scale Complex Video Classification. *IEEE Transactions on Multimedia* 21, 10 (2019), 2518–2530.
- [116] Mingxing Zhang, Yang Yang, Xinghan Chen, Yanli Ji, Xing Xu, Jingjing Li, and Heng Tao Shen. 2021. Multi-Stage Aggregated Transformer Network for Temporal Language Localization in Videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 12669–12678.
- [117] Songyang Zhang, Houwen Peng, Jianlong Fu, Yijuan Lu, and Jiebo Luo. 2021. Multi-Scale 2D Temporal Adjacency Networks for Moment Localization with Natural Language. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021), 1–14.
- [118] Songyang Zhang, Houwen Peng, Jianlong Fu, and Jiebo Luo. 2020. Learning 2D Temporal Adjacent Networks for Moment Localization with Natural Language. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 12870–12877.
- [119] Songyang Zhang, Jinsong Su, and Jiebo Luo. 2019. Exploiting Temporal Relationships in Video Moment Localization with Natural Language. In *Proceedings of the ACM International Conference on Multimedia*. 1230–1238.
- [120] S. Zhang, Y. Zhu, and A. K. Roy-Chowdhury. 2016. Context-Aware Surveillance Video Summarization. *IEEE Transactions on Image Processing* 25, 11 (2016), 5469–5478.
- [121] Zongmeng Zhang, Xianjing Han, Xuemeng Song, Yan Yan, and Liqiang Nie. 2021. Multi-Modal Interaction Graph Convolutional Network for Temporal Language Localization in Videos. *IEEE Transactions on Image Processing* 30 (2021), 8265–8277.
- [122] Zhu Zhang, Zhijie Lin, Zhou Zhao, and Zhenxin Xiao. 2019. Cross-Modal Interaction Networks for Query-Based Moment Retrieval in Videos. In *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval*. 655–664.
- [123] Zhu Zhang, Zhijie Lin, Zhou Zhao, Jieming Zhu, and Xiuqiang He. 2020. Regularized Two-Branch Proposal Networks for Weakly-Supervised Moment Retrieval in Videos. In *Proceedings of the ACM International Conference on Multimedia*. 4098–4106.
- [124] Yang Zhao, Zhou Zhao, Zhu Zhang, and Zhijie Lin. 2021. Cascaded Prediction Network via Segment Tree for Temporal Video Grounding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4197–4206.
- [125] Luowei Zhou, Nathan Louis, and Jason J Corso. 2018. Weakly-Supervised Video Object Grounding from Text by Loss Weighting and Object Interaction. In *British Machine Vision Conference*. 1–12.