# TRACKING OBJECTS AND ACTIVITIES WITH ATTENTION FOR TEMPORAL SENTENCE GROUNDING

*Zeyu Xiong*[1*]    *Daizong Liu*[2*]    *Pan Zhou*[1†]    *Jiahao Zhu*[1]

[1]Hubei Key Laboratory of Distributed System Security, Hubei Engineering
Research Center on Big Data Security, School of Cyber Science
and Engineering, Huazhong University of Science and Technology
[2]Wangxuan Institute of Computer Technology, Peking University
zeyuxiong@hust.edu.cn, dzliu@stu.pku.edu.cn, panzhou@hust.edu.cn, jiahaozhu@hust.edu.cn

## ABSTRACT

Temporal sentence grounding (TSG) aims to localize the temporal segment which is semantically aligned with a natural language query in an untrimmed video. Most existing methods extract frame-grained features or object-grained features by 3D ConvNet or detection network under a conventional TSG framework, failing to capture the subtle differences between frames or to model the spatio-temporal behavior of core persons/objects. In this paper, we introduce a new perspective to address the TSG task by tracking pivotal objects and activities to learn more fine-grained spatio-temporal behaviors. Specifically, we propose a novel Temporal Sentence Tracking Network (TSTNet), which contains (A) a Cross-modal Targets Generator to generate multi-modal templates and search space, filtering objects and activities, and (B) a Temporal Sentence Tracker to track multi-modal targets for modeling the targets' behavior and to predict query-related segment. Extensive experiments and comparisons with state-of-the-arts are conducted on challenging benchmarks: Charades-STA and TACoS. And our TSTNet achieves the leading performance with a considerable real-time speed.

*Index Terms*— TSG, tracking, cross-modal, attention

## 1. INTRODUCTION

Temporal sentence grounding (TSG) [1, 2, 3] is an important yet challenging task in multi-modal deep learning due to its complexity of multi-modal interactions and complicated context information. As shown in Fig. 1(a), given an untrimmed video, it aims to determine the segment boundaries including start and end timestamps that contain the interested activity according to a given sentence description.

Most previous works [1, 2, 4, 5, 6, 7, 8] first encode and interact the pair of video-query input, and then employ either a proposal-based or a proposal-free grounding head to predict the target segments. However, these methods extract frame-level video features by a pre-trained 3D ConvNet,
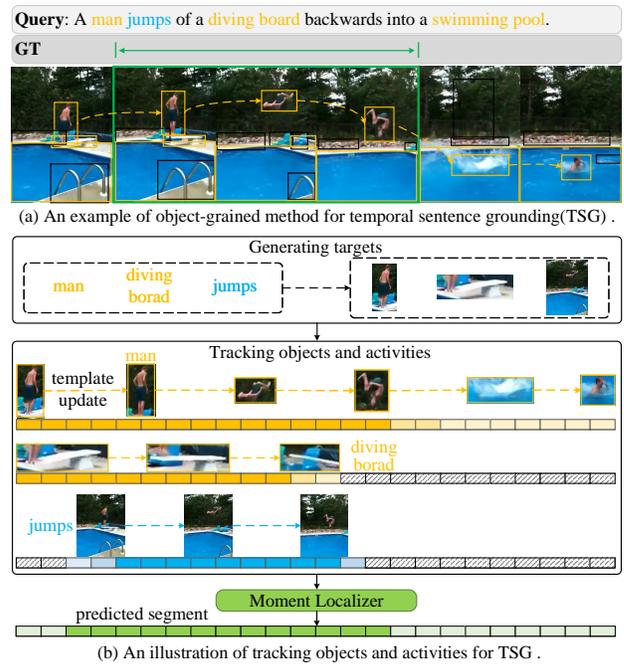


**Fig. 1**. Illustrations of (a) TSG and (b) tracking-view TSG.

which may capture the redundant background appearance in each frame and fails to perceive the subtle differences among video frames with high similarity. Recently, a few detection-based approaches [9, 10, 11] have been proposed to capture fine-grained object appearance features inside each frame for focusing more on the foreground contexts. Although they achieve promising results, these methods directly correlate all spatial-temporal objects in the entire video through a simple graph- or co-attention mechanism, lacking sufficient reasoning on the most query-specific objects.

To learn more specific spatial-temporal relations among the extracted objects, as shown in Fig. 1(b), we adapt the object tracking perspective into the TSG task to correlate the most query-related objects for activity modeling. Firstly, we generate the multi-modal target templates by selecting core objects and activities (such as "man", "diving broad" and
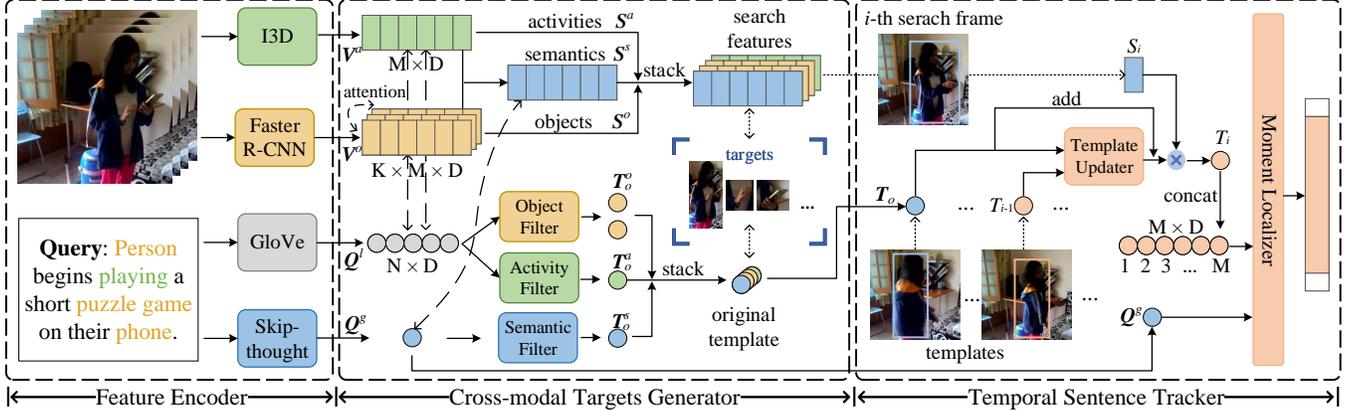
---

**Fig. 2**. The overall architecture of our proposed TSTNet. We first encode the video and query to obtain **V**isual features $\boldsymbol{V}^o, \boldsymbol{V}^a$ and **Q**uery features $\boldsymbol{Q}^l, \boldsymbol{Q}^g$. Then we develop a Cross-modal Targets Generator to generate and filter original **T**emplates $(\boldsymbol{T}_o^o, \boldsymbol{T}_o^a, \boldsymbol{T}_o^s)$ and **S**earch space $(\boldsymbol{S}^o, \boldsymbol{S}^a, \boldsymbol{S}^s)$ for latter **o**bject/**a**ctivity/**s**emantics tracking. The Temporal Sentence Tracker is designed for tracking the query-related target corresponding with sentence semantics and predicting the target segment.

"jumps") and discarding irrelevant ones (black boxes in the video). By aggregating words feature with object-grained visual features, we track the corresponding template in each frame and model the target behavior with continuous target templates. With specific semantics, we deploy a moment localizer to determine the most query-related segment.

However, directly adopting the standard tracking algorithm for TSG will raise two major problems: (1) **Modal gap**: there is a modal gap between vision and language in TSG while utilizing typical tracking algorithm. (2) **Ambiguous target**: there is no specific target provided as supervision in TSG, while standard tracking specifies an object to track. To overcome these two challenges, we develop a novel Temporal Sentence Tracking Network (TSTNet), which contains a Feature Encoder, a Cross-modal Targets Generator and a Temporal Sentence Tracker (Fig. 2). Specifically, we first extract the object-grained features and query features by pre-trained detection model, action recognition 3D ConvNet, Glove [12], Skip-thought [13]. Then, we leverage self and co-attention to establish associations among objects, activities and words for bridging the modal gap and generating search space and templates. And we utilize instance filters to further screening the core targets. A dynamic template updater in temporal sentence tracker is to match and dynamically update templates for each frame in search space, modeling the behavior of targets. Finally, we employ a moment localizer to determine the temporal segment and fine-tune the boundaries of it.

Our contributions are summarized as follows:

(1) We provide a new perspective of tracking objects and activities to address the TSG problem, which can focus more on behavior modeling of core targets;

(2) We propose a novel framework TSTNet, which tackles the differences between TSG and standard tracking;

(3) We demonstrate the predominant effectiveness and efficiency of our TSTNet by evaluating on two benchmarks: Charades-STA and TACoS.

## 2. METHODOLOGY

Given an untrimmed video $\mathcal{V}$ and a natural language sentence query $\mathcal{Q}$, the TSG aims at predicting a video segment from time $\tau_s$ to $\tau_e$ corresponding to the same semantic as $\mathcal{Q}$.

### 2.1. Feature Encoder

**Video Encoder.** In order to model both objects and activities, we extract the appearance-aware and motion-aware features of original videos by the pre-trained Faster R-CNN [14] and C3D/I3D [15] networks.

Specifically, for objects, we first uniformly sample fixed $M$ frames from video $\mathcal{V}$, and then extract $K$ objects from each frame using Faster R-CNN with a ResNet-50 FPN backbone. Therefore, we represent the object features as $V^o = \{o_{i,j}, b_{i,j}\}_{i=1,j=1}^{i=M,j=K}$, where $o_{i,j} \in \mathbb{R}^{D_o}$ denotes object features with dimension $D_o$, and $b_{i,j} \in \mathbb{R}^4$ represents the bounding-box coordinate of the $j$-th object in $i$-th frame.

Since the spatial relationship of instances plays a important role in object behavior modeling, we fuse the spatial information $b_{i,j}$ with object features $o_{i,j}$ by *concat* function and Fully Connection (FC), obtaining $\boldsymbol{V}^o = \{\boldsymbol{v}_{i,j}^o\}_{i=1,j=1}^{i=M,j=K}$.

For activity features, we put every 8 frames to a pre-trained 3D ConvNet with stride 4, and sample $M$ output sequence by linear interpolation, which is represented as $\boldsymbol{V}^a = \{\boldsymbol{v}_i^a\}_{i=1}^{i=M}$, where $\boldsymbol{v}_i^a \in \mathbb{R}^D$.

**Query Encoder.** For word-level encoding, following previous works [16, 17], we embed each word in sentence query $\mathcal{Q}$ by Glove [12], obtaining the local semantic of every single word: $\boldsymbol{Q}^l = \{\boldsymbol{q}_i\}_{i=0}^{i=N}$, where $\boldsymbol{q}_i \in \mathbb{R}^D$. To extract the semantic of the whole sentence, the Skip-thought parser [13] is employed to capture the global semantic of the whole query, denoted as $\boldsymbol{Q}^g \in \mathbb{R}^D$.

### 2.2. Cross-modal Targets Generator

To solve the problem of inconsistency of modality and ambiguous targets to retrieved, we developed a Cross-modal Tar-

gets Generator (CTG).

Specifically, we design Search Space Representation (SSR) and Template Generation (TG) across sentences and videos for further cross-modal tracking. As for redundant targets, we develop instance Filters for screening the core semantic-related targets.

**Search Space Representation.** First, we utilize sentence query guidance to represent video search space for retrieval. Because there exist interactions between objects inside each frame, we first learn the self- and inter-correlation between $K$ objects with attention mechanism [18]:

$$\widehat{\boldsymbol{V}}^o = \sigma(\boldsymbol{V}^o\boldsymbol{W}_1(\boldsymbol{V}^o\boldsymbol{W}_2)^\top)\frac{\boldsymbol{V}^o}{\sqrt{D}}, \qquad (1)$$

where $\boldsymbol{W}_1, \boldsymbol{W}_2$ are two learnable matrices, $\sigma$ is an activate function. Next, we associate words with objects in each frame by leveraging query-guide attention to highlight the word-relevant objects while weakening the irrelevant ones:

$$\boldsymbol{w}^q = \sigma(\frac{1}{\sqrt{D}}(\widehat{\boldsymbol{V}}^o\boldsymbol{W}_3)(\boldsymbol{Q}^l\boldsymbol{W}_4)^\top),$$
$$\boldsymbol{V}_q^o = \boldsymbol{w}^q\boldsymbol{Q}^l\boldsymbol{W}_5, \qquad (2)$$

where $\boldsymbol{W}_3, \boldsymbol{W}_4, \boldsymbol{W}_5$ are the transform matrices, $\boldsymbol{w}^q \in \mathbb{R}^{K \times N}$ represents the correlation between each word-object pair. Since too many pre-extracted objects will interfere with tracking and modeling the key targets, we are supposed to filter out redundant backgrounds or instances, and accurately select core objects/activities for tracking. Therefore, we deploy $k$ adaptive Object Filters to select $k$ core object search space from $K$ object features. In detail, we implement the filter with a linear layer, followed by a Leaky_ReLU function and a 1d-maxpool layer to activate and filter targets. At last, we obtain $k$ object search space $\boldsymbol{S}^o = \{_i S^o\}_{i=1}^{i=k}$, where $_i S^o \in \mathbb{R}^{M \times D}$ represents the $i$-th search space for tracking.

For the activity search space representation, similarly, we gain the query-guide activity features $\boldsymbol{V}_q^a$ by replacing $\widehat{\boldsymbol{V}}^o$ with $\boldsymbol{V}^a$ in Eq. (2). Then a linear layer followed by a Leaky_ReLU function is employed to generate the search space $\boldsymbol{S}^a \in \mathbb{R}^{M \times D}$.

Considering object or activity alone is not enough to model the semantics and relationships between them. Therefore, we learn the semantic features $\boldsymbol{V}^s$ with activity features $\boldsymbol{V}^a$ and object features $\widehat{\boldsymbol{V}}^o$ by Eq. (2), and then calculate element-wise multiplication with $\boldsymbol{Q}^g$ to get semantic search space $\boldsymbol{S}^s \in \mathbb{R}^{M \times D}$

**Template Generation.** After getting the search space of the video, we need to determine an initial template as the target for matching and tracking. In this case, we consider deeming the sentence query $\boldsymbol{Q}^l, \boldsymbol{Q}^g$ as initial matching template by combining it with instances in videos $\boldsymbol{V}^o, \boldsymbol{V}^a$.

In details, we first calculate the cosine similarity between each word and object at each frame, and obtain the object-aware query feature by equation:

$$\boldsymbol{w}_{i,j,t}^s = \frac{\boldsymbol{q}_{i,t}(\boldsymbol{v}_{j,t}^o)^\top}{\|\boldsymbol{q}_{i,t}\|\|\boldsymbol{v}_{j,t}^o\|},$$
$$\widehat{\boldsymbol{q}}_{i,t}^o = \sum_{j=1}^K \boldsymbol{w}_{i,j,t}^s\boldsymbol{v}_{j,t}^o, \widehat{\boldsymbol{Q}}^o = \{\widehat{\boldsymbol{q}}_{i,t}^o\}_{i=1,t=1}^{i=N,t=M}, \qquad (3)$$

where $\boldsymbol{w}_{i,j,t}^s$ is the similarity between $i$-th word and $j$-th object at frame $t$.

Then, we pick $k$ original object templates $\boldsymbol{T}_o^o = \{_i \boldsymbol{T}_o^o\}_{i=1}^{i=k} \in \mathbb{R}^{k \times D}$ from $\widehat{\boldsymbol{Q}}^o$ via $k$ object filters, which corresponds to $k$ search spaces $\boldsymbol{S}^o$.

For activity template generation, similarly, we obtain an activity-aware enhanced query feature $\widehat{\boldsymbol{Q}}^a \in \mathbb{R}^{N \times M \times D}$ by Eq. (3) and filter an activity original template $\boldsymbol{T}_o^a \in \mathbb{R}^{1 \times D}$.

The semantic original template $\boldsymbol{T}_o^s \in \mathbb{R}^{1 \times D}$ is generated from $\boldsymbol{Q}^g$ and $\boldsymbol{V}^s$ by Eq. (3) followed a semantic filter.

### 2.3. Temporal Sentence Tracker

In order to model the targets' behavior for text-visual alignment, we develop a Dynamic Template Updater (DTU) to track targets in the search space with original templates and then deploy a Moment Localizer to localize the most query-related moment in videos.

**Dynamic Template Updater.** For any search space and template tuple $(\boldsymbol{S}, \boldsymbol{T})$, we fuse the original template $\boldsymbol{T}_o$ and the $(i-1)$-th template $T_{i-1}$ as a new template for aligning the $i$-th search frame $S_i$, as is shown in Fig. 2 (right). We utilize a template updater $\phi$ to update the templates and then concatenate all templates by sequence, formulated as:

$$T_i = (\phi(\boldsymbol{T}_o, T_{i-1}) + \boldsymbol{T}_o) \cdot S_i,$$
$$\boldsymbol{F}_T = [T_1, T_2, ..., T_M], \qquad (4)$$

where $\phi(\cdot)$ is the Feedforward Neural Network (FNN) followed by a GRU [19] unit. $\boldsymbol{F}_T$ contains the behavior information composed of continuous templates. In practice, we feed different tuples of search space and template into different template updaters, and obtain $\boldsymbol{F}_T^o \in \mathbb{R}^{k \times M \times D}, \boldsymbol{F}_T^a \in \mathbb{R}^{M \times D}, \boldsymbol{F}_T^s \in \mathbb{R}^{M \times D}$. Then we concatenate them followed by a FC layer and get $\widehat{\boldsymbol{F}}_T$:

Noting that reversed trace of the target also provides rich behavior information, we track the target from the last search frame as Eq. (4) and obtain the reversed features $\widehat{\boldsymbol{F}}_T^r$. At last, we concatenate the forward $\widehat{\boldsymbol{F}}_T$ and the reversed $\widehat{\boldsymbol{F}}_T^r$ as $\widetilde{\boldsymbol{F}}_T$.

**Moment Localizer.** As many temporal localizers are plug-and-play, we follow the previous work [16] to predict the target moment for fair comparison.

## 3. EXPERIMENTS

**Datasets. Charades-STA** dataset was built on Charades by [1], including 9,848 videos of indoor scenarios. By convention, we use 12,408 and 3,720 video-sentence pairs for training and testing. **TACoS** is collected from the MPII Cook-

| Charades-STA | | | | | |
|---|---|---|---|---|---|
| Methods | Feature | R@1, IoU = | | | mIoU |
| | | 0.3 | 0.5 | 0.7 | |
| CBP | C3D | - | 36.80 | 18.87 | 35.74 |
| 2DTAN | VGG | - | 39.81 | 23.31 | - |
| VSLNet | I3D | 70.46 | 54.19 | 35.22 | 50.02 |
| LGI | I3D | 72.96 | 59.46 | 35.48 | 51.38 |
| CPN | I3D | 75.53 | 59.77 | 36.67 | 53.14 |
| IA-Net | I3D | - | 61.29 | 37.91 | - |
| DRFT | I3D+F+D | *76.68* | 63.03 | 40.15 | *54.89* |
| MARN | I3D+Obj | - | *66.43* | *44.80* | - |
| **TSTNet** | C3D+Obj | 76.26 | 65.34 | 43.61 | 56.76 |
| | I3D+Obj | **77.62** | **67.49** | **45.21** | **57.82** |

**Table 1**. Comparison with SOTAs on Charades-STA.

| TACoS | | | | | |
|---|---|---|---|---|---|
| Methods | Feature | R@1, IoU = | | | mIoU |
| | | 0.3 | 0.5 | 0.7 | |
| CMIN | C3D | 24.64 | 18.05 | - | - |
| CBP | C3D | 27.31 | 24.79 | 19.10 | 21.59 |
| 2DTAN | C3D | 37.29 | 25.32 | - | - |
| VSLNet | I3D | 29.61 | 24.27 | 20.03 | 24.11 |
| IA-Net | I3D | 37.91 | 26.27 | - | - |
| CPN | I3D | 48.29 | 36.58 | *21.25* | *34.63* |
| MARN | I3D+Obj | *48.47* | *37.25* | - | - |
| **TSTNet** | C3D+Obj | 50.21 | 38.47 | 23.12 | 35.26 |
| | I3D+Obj | **53.39** | **41.23** | **26.62** | **37.83** |

**Table 2**. Comparison with SOTAs on TACoS.

ing [20], which contains 127 long videos of cooking scenarios. Following [1], we obtain 10,146, 4,589 and 4,083 clip-sentence pairs as training ,validation and testing dataset.
**Evaluation Metrics.** We adopt "R@$n$, IoU=$\mu$" and "mIoU" metrics for evaluation. The "R@$n$, IoU=$\mu$" denotes the percentage of at least one of top-$n$ predictions having IoU larger than $\mu$. "mIoU" represents the mean average IoU.
**Implementation Details.** For object feature extraction, we utilize Faster R-CNN [14] with a ResNet-50 FPN [21] backbone to obtain object features. The number $K$ of extracted objects is set to 15 and the number $k$ of object filter is set to 5. The length of frame sequences $M$ in our model is 64, 200 for Charades-STA and TACoS. For query encoding, we utilize GloVe 840B 300d [12] to embed each word as word features. For model setting, the activate function $\sigma$ is Sigmoid. The hidden dimension $D$ is 512. We sample 800 segment proposals for TACoS and 384 for Charades-STA similar to [16]. We train our model by an Adam optimizer with the learning rate of 0.0008 for 60 epoches. Batch size is 64.

### 3.1. Experimental Results and Analysis

We compare the proposed TSTNet with the following state-of-the-arts: (1) *Proposal-based* methods: CBP [22], 2DTAN [23], CMIN [16], IA-Net [17]; (2) *Proposal-free* methods: VSLNet [24], LGI [25] , CPN [26]; (3) *Multi-stream* methods: DRFT [27], MARN [9]. The best results are in **bold** and

| Methods | TGN | 2DTAN | CMIN | TSTNet |
|---|---|---|---|---|
| V-QPS | 2.23 | 3.89 | 86.29 | **103.27** |
| Parameters | 166 | 363 | 78 | **67** |
| Accuracy | 18.89 | 25.32 | 18.05 | **41.23** |

**Table 3**. Effyciency comparision in terms of video-query pairs per second (V-QPS), Parameters (Mb) and Accuracy (R@1, IoU=0.5 metric) on TACoS dataset.

| Components | Changes | R@1, IoU = | | |
|---|---|---|---|---|
| | | 0.3 | 0.5 | 0.7 |
| Cross-modal | w/o SSR | 72.59 | 62.21 | 39.79 |
| Target | w/o TG | 74.25 | 63.82 | 41.59 |
| Generator | w/o Filter | 73.12 | 62.74 | 40.27 |
| Temporal | w/o DTU | 70.12 | 60.49 | 38.84 |
| Sentence | w/o GRU | 72.31 | 61.92 | 39.65 |
| Tracker | w/o $\widehat{\boldsymbol{F}}_T^r$ | 73.61 | 62.64 | 40.35 |
| **Full** | | **77.62** | **67.49** | **45.21** |

**Table 4**. Ablation study on Charades-STA dataset.

the second bests are in *italic*.
**Quantitative Comparion.** As summarized in Table 1 and 2, our proposed TSTNet surpasses all existing methods on two datasets. Observe that the performance improvements of TSTNet are more significant under more strict metrics (R@1, IoU=0.7), indicating that TSTNet can predict more precise moment boundaries of untrimmed videos. As multi-stream methods, DRFT integrates the three modalities of visual information RGB (I3D), optical flow (F) and depth maps (D), and MARN fuses activity (I3D) with object (Obj). They perform well and imply that combining multi-dimension sources helps the model learn more accurate semantics. Differing from DRFT and MARN, our TSTNet traces objects from a finer granularity and mines more explicit behavior of core targets, thus outperforming better in results.
**Efficiency Comparison.** We compare the inference speed and effectiveness of our TSTNet with previous methods on a single Nvidia Quadro RTX5000 GPU on TACoS dataset. Table 3 shows that TSTNet achieves a significantly faster inference speed and a lightweight model size.
**Ablation Study.** As shown in Table 4, we verify the contribution of several modules in our TSTNet: we remove SSR, TG, Filter in Cross-modal Target Generator, and DTU, GRU, reversed feature $\widehat{\boldsymbol{F}}_T^r$ in Temporal Sentence Tracker (mentioned in Sec. 2). The result manifests each above component provides a positive contribution.

### 4. CONCLUSION

In this work, we solve the TSG task with a multi-modal instance tracking framework and propose the TSTNet. With this effective and efficient framework, TSTNet outperforms state-of-the-arts on two challenging benchmarks.

# 5. REFERENCES

[1] Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia, "Tall: Temporal activity localization via language query," in *ICCV*, 2017. 1, 3, 4

[2] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell, "Localizing moments in video with natural language," in *ICCV*, 2017. 1

[3] Daizong Liu, Xiaoye Qu, Xiao-Yang Liu, Jianfeng Dong, Pan Zhou, and Zichuan Xu, "Jointly cross-and self-modal graph attention network for query-based moment localization," in *ACM MM*, 2020. 1

[4] Daizong Liu, Xiaoye Qu, Jianfeng Dong, Pan Zhou, Yu Cheng, Wei Wei, Zichuan Xu, and Yulai Xie, "Context-aware biaffine localizing network for temporal sentence grounding," in *CVPR*, 2021. 1

[5] Daizong Liu, Xiaoye Qu, and Wei Hu, "Reducing the vision and language bias for temporal sentence grounding," in *ACM MM*, 2022. 1

[6] Daizong Liu and Wei Hu, "Skimming, locating, then perusing: A human-like framework for natural language video localization," in *ACM MM*, 2022. 1

[7] Jiahao Zhu, Daizong Liu, Pan Zhou, Xing Di, Yu Cheng, et al., "Rethinking the video sampling and reasoning strategies for temporal sentence grounding," in *EMNLP Findings*, 2022. 1

[8] Daizong Liu, Xiang Fang, Pan Zhou, Xing Di, Weining Lu, and Yu Cheng, "Hypotheses tree building for one-shot temporal sentence localization," in *AAAI*, 2023. 1

[9] Daizong Liu, Xiaoye Qu, Pan Zhou, and Yang Liu, "Exploring motion and appearance information for temporal sentence grounding," in *AAAI*, 2022. 1, 4

[10] Zeyu Xiong, Daizong Liu, and Pan Zhou, "Gaussian kernel-based cross modal network for spatio-temporal video grounding," in *ICIP*, 2022. 1

[11] Daizong Liu, Xiang Fang, Wei Hu, and Pan Zhou, "Exploring optical-flow-guided motion and detection-based appearance for temporal sentence grounding," *TMM*, 2023. 1

[12] Jeffrey Pennington, Richard Socher, and Christopher D Manning, "Glove: Global vectors for word representation," in *EMNLP*, 2014. 2, 4

[13] Ryan Kiros, Yukun Zhu, Russ R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler, "Skip-thought vectors," in *NeurIPS*, 2015. 2

[14] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *NeurIPS*, 2015. 2, 4

[15] Joao Carreira and Andrew Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *CVPR*, 2017. 2

[16] Zhu Zhang, Zhijie Lin, Zhou Zhao, and Zhenxin Xiao, "Cross-modal interaction networks for query-based moment retrieval in videos," in *SIGIR*, 2019. 2, 3, 4

[17] Daizong Liu, Xiaoye Qu, and Pan Zhou, "Progressively guide to attend: An iterative alignment framework for temporal sentence grounding," in *EMNLP*, 2021. 2, 4

[18] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin, "Attention is all you need," in *NeurIPS*, 2017. 3

[19] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," in *NeurIPS*, 2014. 3

[20] Marcus Rohrbach, Michaela Regneri, Mykhaylo Andriluka, Sikandar Amin, Manfred Pinkal, and Bernt Schiele, "Script data for attribute-based recognition of composite activities," in *ECCV*, 2012. 4

[21] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie, "Feature pyramid networks for object detection," in *CVPR*, 2017. 4

[22] Jingwen Wang, Lin Ma, and Wenhao Jiang, "Temporally grounding language queries in videos by contextual boundary-aware prediction," in *AAAI*, 2020. 4

[23] Songyang Zhang, Houwen Peng, Jianlong Fu, and Jiebo Luo, "Learning 2d temporal adjacent networks for moment localization with natural language," in *AAAI*, 2020. 4

[24] Hao Zhang, Aixin Sun, Wei Jing, and Joey Tianyi Zhou, "Span-based localizing network for natural language video localization," in *ACL*, 2020. 4

[25] Jonghwan Mun, Minsu Cho, and Bohyung Han, "Local-global video-text interactions for temporal grounding," in *CVPR*, 2020. 4

[26] Yang Zhao, Zhou Zhao, Zhu Zhang, and Zhijie Lin, "Cascaded prediction network via segment tree for temporal video grounding," in *CVPR*, 2021. 4

[27] Yi-Wen Chen, Yi-Hsuan Tsai, and Ming-Hsuan Yang, "End-to-end multi-modal video temporal grounding," in *NeurIPS*, 2021. 4