

Text query: the baby takes something out of his mouth.



Figure 1. An illustration of localizing clips of interest by natural language description. In this example, we are looking for the clip that the baby takes something out of his mouth. The video clip above the yellow bar is the clip we are interested in.

pressive than a pre-defined label described by a single word or a short phrase. For example, in the description of “the baby takes something out of his mouth”, there are different object categories (“baby”, “mouth”), attributes (“his”) and human-object interactions (“take something out of”). To localize the described clip in a video sequence, the visual model has to recognize objects, actions, spatial relations, and also the context of the event. Standard action detection algorithms only localize actions at a coarse level. In this new task, the model has to understand both language descriptions and video content in finer scales. These datasets contain a great number of descriptions, covering a large vocabulary and a variety of complex sentence structures. With such rich information, though challenging, this task enables us to learn a model that generalizes to new natural language descriptions.

The existing methods for video clip localization by natural language description follow traditional cross-modal retrieval approaches [8], embed visual information and language information individually [2, 9]. Hendricks *et al.* [2] used a ranking loss to train the embedding functions, so that the language feature is closer to the ground-truth clip than all negative clips. Gao *et al.* [9] merged clip visual embedding and language embedding by a multi-modal processing module, followed by alignment score estimation and location regression. In this paper, we propose a novel attentive sequence to sequence translator (ASST) for this challenging task. Specifically, we make the following contributions.

First, we introduce a finely calibrated video-language attentive mechanism for understanding rich semantic descriptions. A bi-directional Recurrent Neural Network (biRNN) is utilized to parse text descriptions in a forward and a backward way. The method in [2] generates a vector for each fixed-length clip by average pooling, during which it loses the local information of each individual frame. The method in [9] keeps temporal structures of sampled video clips by temporal convolutions, but ignores detailed relation with natural language descriptions. Likewise, the methods in [2] and [9] only generate one vector for each sentence, without word-level information. Our method works in a much finer manner and attends every meaningful word or phrase to each frame. For example, given a sentence of 6 words and a video of 100 frames, we generate an attention

weight matrix of size 6×100 . Richer information as represented by the attention weight matrix enables our method to perform a finer video-language matching. Our attention mechanism aligns language descriptions and video content, which benefits the understanding of both semantic description and video content. For example, when localizing “we can see someone in a mask”, the target video clip will attend to the word “mask” as opposed to the entire sentence. In this way, a more accurate video-language representation can be generated.

Second, we design a hierarchical network architecture for jointly learning video-language representations with multiple granularities. The hierarchical visual network models videos at multiple temporal scales to extract subtle details of video content as well as global context. Lower layers in the hierarchy exploit local video content, which are then integrated at a higher layer to model temporal dependencies in a longer duration. In other words, lower layers provide finer exploration while higher layers provide context information. For example, given a video of birthday party, we first model each frame with details, then we model three sub-events including “cut a birthday cake”, “sing a song”, and “eat the cake” at a higher layer, and lastly we model the entire event of birthday party in the highest layer. A robust joint video-language representation is then generated, which preserves fine-scale details, frame and clip dependencies as well as temporal context of video content.

We evaluate our ASST on two large-scale datasets for localizing clips by natural language descriptions. On the DiDeMo dataset, it achieves 32.38% in Rank@1 and outperforms the state-of-the-art by 4.28%. On the Charades-STA dataset, we obtain 37.04% in Rank@1, $IoU = 0.5$. Notably, it outperforms the state-of-the-art by a large margin of 15.62%.

2. Related work

Action detection in videos. In action detection, the goal is to localize interested video clips of given action labels. This task is usually divided into two branches: temporal detection and spatio-temporal detection. For temporal detection [17], the model has to output the start and end time of each action clip. For spatio-temporal detection [29], be-

sides temporally localizing the actions, the model should also generate the spatial locations of the actions for each frame. We do not consider spatial localization in this paper.

Two-Stream ConvNets [39, 44] and 3D-ConvNets [42] have been proposed to model video sequences for action recognition. These methods have been widely used in temporal action detection [22, 35, 47, 49] and other video analysis tasks, *e.g.*, video captioning [43]. Shou *et al.* [35] applied frame-level classification via a modified 3D-ConvNet, and then generated temporal window prediction by merging frame-level prediction. Lin *et al.* [22] adopted Single Shot Detector [26] to detect interested actions in one shot. Xu *et al.* [47] tried to tackle the problem of temporal proposal generation. They incorporated an RoI-pooling layer, followed by a classification branch and a temporal coordinate regression branch. Zhao *et al.* [49] achieved a good performance by hard coding context regions, which enables the detector to model temporal context. They also proposed a sophisticated region merging strategy. We aggregate contextual temporal structures via dilated convolutional layers, followed by squeezing-expansion layers. Our model jointly encodes neighboring and all observable visual information, which is more flexible for temporal structure modeling.

Natural language object localization. Natural language object localization [14, 27] is to spatially localize objects of interest in images by natural language descriptions. Mao *et al.* [27] encoded each possible region by a Long Short-Term Memory (LSTM), and trained the model in a discriminative way by using a softmax loss over all encoded regions. Hu *et al.* [14] ranked each object proposal by considering spatial configuration, local visual information and global visual information. Rohrbach *et al.* [33] used attention mechanism to choose the region that could be best used for description reconstruction. Liu *et al.* [25] extracted image attributes from image regions proposals, and measured the similarity between image region attribute embeddings and textual embeddings. In this paper, we focus on natural language video clip localization. We propose a novel ASST to jointly learn a video-language representation.

Video clip localization by natural language descriptions. Video clip localization by natural language descriptions is a new task introduced recently [2, 9]. Hendricks *et al.* [2] proposed a Moment Context Network (MCN) for this task. MCN is a ranking based method. It encodes natural language description and video clip information individually, and ranks video clips by measuring distance between video clip embeddings and sentence embeddings. MCN considers representations from two modals separately, ignores detailed relevance between two modals. Gao *et al.* [9] fused video feature and language feature by a multi-modal processing module which consists of addition, multiplication and fully connected operations. This Cross-modal Temporal Regression Localizer (CTRL) builds an alignment score

estimator and a location regressor on the fused representation to produce clip prediction. Our ASST explores cross-modal relevance at finer granularities, which is able to better exploit detailed information.

3. Our approach

In our work, we leverage natural language information as a flexible network module to translate visual information. We exploit detailed connections between natural language information and visual information on multiple hierarchies as well. We directly sample possible clips from the final video-language joint representation sequence generated by our ASST for clip localization.

As illustrated in Figure 2, the proposed attentive sequence to sequence translator (ASST) consists of three components: a video subnet for modeling video content, a language subnet for modeling natural language description and a multi-stage cross-modal attention module.

3.1. Video subnet

Our video subnet is shown on the right part in Figure 2. The video subnet first translates the input visual feature sequence into a representation sequence. The length of the representation sequence and the length of the input feature sequence are identical. Each element in this final representation sequence contains temporal visual contextual information, ranging from subtle details to global context, as well as clues from natural language descriptions. Then, video clip representations are obtained by directly sampling clips on the final representation sequence.

3.1.1 Visual feature learning

The video subnet is a 1D convolutional network spanned on temporal domain. We take a sequence of vectors $\{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{m-1}\}$ ($\{\mathbf{x}\}$ for simplicity) with dimension d_v as the input of the video subnet, where m is the number of frame-level features in the input sequence. We also denote the temporal length in seconds of the input video sequence as τ . Visual feature modeling consists of two stages. One is local temporal feature modeling via stacked dilated convolution. The other one is global temporal context modeling via a squeezing-expansion process. The video subnet transforms an input sequence $\{\mathbf{x}\}$ to a representation sequence with same length $\{\mathbf{x}'_0, \mathbf{x}'_1, \dots, \mathbf{x}'_{m-1}\}$ ($\{\mathbf{x}'\}$ for simplicity) via these two stages.

Dilation Layers. We first model the frame-level input features with their neighboring frames. This process builds temporal hierarchical structure of the input sequence via stacked dilated convolution layers [18, 21, 28, 48]. The dilation rates are doubled from each layer forward. Stacked dilated convolution increases receptive field in an exponential speed as the network goes deeper, while the sizes of

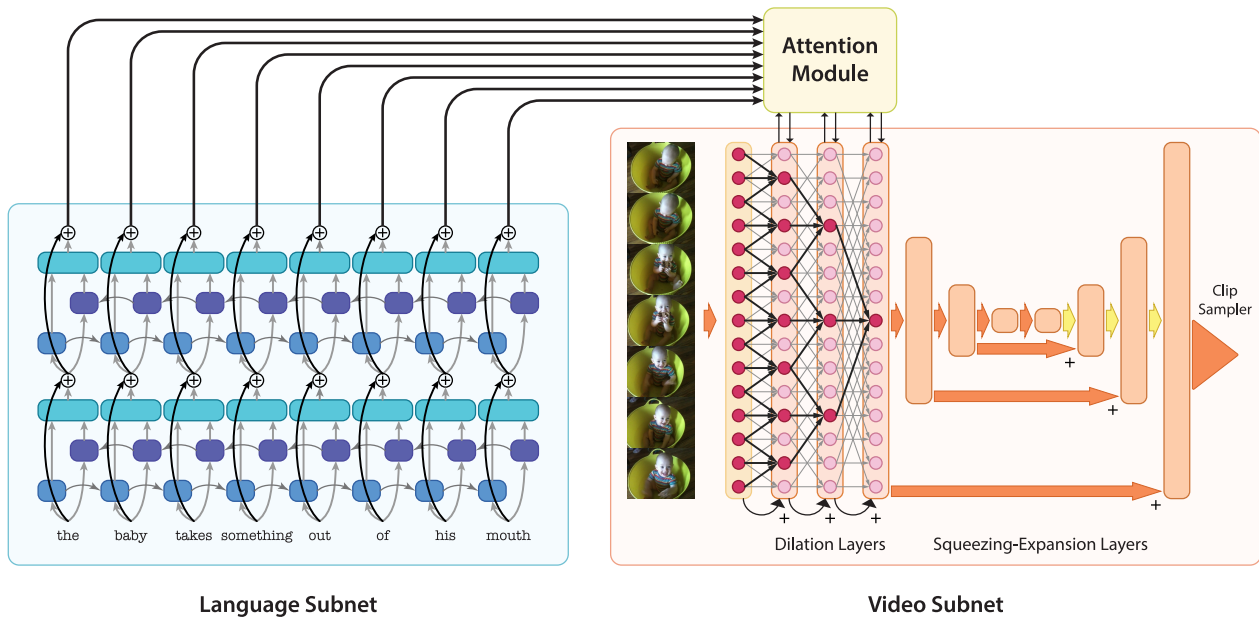


Figure 2. An illustration of our model. Our model consists of an RNN-based language subnet (left) and a ConvNet-based video subnet (right). The two subnets are integrated via an attention module (top) on every dilation layer and squeezing layer. In this example, the language subnet has one stacked LSTM layer. For the video subnet, there are three dilation layers, followed by three squeezing layers and three expansion layers. The output of the last dilation layer is used as the input of first squeezing layer. Best viewed in color.

feature maps remain unchanged. These layers model temporal contextual structures at multiple time scales for every input frame. In Figure 2, we show a model with three dilation layers. We additionally apply shortcut connections [13] between every two adjacent layers. Shortcut connections shorten the information path from the bottom to the top of the network, which makes the training process easier. After the above operations, every element in the last dilation layer perceives a large temporal window.

Squeezing-Expansion Layers. The dilation layers demonstrated above contain only local temporal contextual structures. The elements in the final dilation layer do not observe the entire input visual sequence. However, the understanding of global context from the entire video is crucial for localizing clips in videos. For example, someone may be interested in localizing the clip when a certain target event happens the second time. To distinguish the positive clip from visually similar negative clips, the model has to understand the global context to make accurate predictions. To model global context information, we generate our final representation sequence through a squeezing phase followed by an expansion phase (Figure 2). The output of the last dilation layer is used as the input of the first squeezing layer. The architecture of “Squeezing-Expansion Layers” is inspired by top-down module for object detection [23, 24, 37]. The model further encodes global context information into the final representation sequence through this squeezing-expansion process.

The squeezing phase consists of several convolution layers with stride 2 and kernel size 3, and generates a global representation vector. This vector summarizes visual information and temporal structures across all the input visual features. It will be used as global context background for further representation modeling.

The expansion phase and the squeezing phase are connected by a convolution operation. During the expansion phase, we expand the global representation vector in a reverse order as we did during the squeezing phase. Expansion stops when it reaches the size of the last dilation layer. We apply a convolution with kernel size 1 on each feature map from the squeezing phase, while the output is added to the expanded feature map with the same size. The result feature map includes both global contextual information (from expanded feature maps) and local contextual information (from convolved feature maps in the squeezing phase). We use linear interpolation as our expansion operation. The output of the last layer is our final representation sequence.

With dilation layers and squeezing-expansion layers, our translator is able to transform input visual sequence into a representation sequence, exploiting both global and local temporal contextual structures.

3.1.2 Clip sampling

We sample our target clips on the final representation sequence $\{\mathbf{x}'\}$ by RoI-pooling [12] in temporal domain with linear interpolation. For each clip sampler, we sample 7 elements from the final representation sequence $\{\mathbf{x}'\}$ as the input.

The clip sampler is a submodule consists of two stacked convolution layers with kernel size 3. All clip samplers share weights. Under different circumstances, clip samplers can be built differently. We show the following two cases.

If there are pre-defined temporal segments, we address this problem as a classification problem by enumerating pre-defined segments. We build one clip sampler for each possible segment. Each sampler generates one scalar, which represents the confidence score for the corresponding video clip. Finally, a softmax loss function is applied over all clip samplers to train our ASST in a discriminative manner following [27]. We denote this model as *classification model*.

If there are no pre-defined temporal segments, our model will have to generate clip proposals. We build clip samplers similar to [31]. In our model, we have six groups of samplers. The i -th clip sampler group consists of clips with length $l_i = \frac{\tau}{2^i}$, where $0 \leq i \leq 5, i \in \mathbb{N}$, and τ is the temporal length of the final representation sequence $\{\mathbf{x}'\}$ as we mentioned before. We sample clips densely by placing every two adjacent samplers with distance $\frac{1}{3}l_i$. Therefore, there are $2^{i+2} - 3$ samplers in the i -th group. For each clip sampler, we predict a tuple $(\hat{y}, \hat{d}_c, \hat{d}_l)$, where \hat{y} is confidence score of the presence of target video clip, \hat{d}_c represents clip center deviation and \hat{d}_l is the log difference from the sampled clip length. During inference, for each sampled clip with temporal center coordinate c_s and clip length l_s , the predicted temporal window is calculated by,

$$(\hat{c}, \hat{l}) = (c_s + \hat{d}_c l_s, e^{\hat{d}_l} l_s), \quad (1)$$

where \hat{c} is the center coordinate of the predicted window and \hat{l} is the length of the window. We use 2-class softmax loss to train the confidence score of presence \hat{y} . Smooth-L1 loss is used to train temporal coordinate regression factors \hat{d}_c and \hat{d}_l following most object detection literatures [5, 15, 24, 26, 31]. We denote this model as *detection model*.

3.2. Language subnet

The architecture of our language subnet is shown on the left part in Figure 2. We use pre-trained word embedding sequence $\{\mathbf{w}_0^0, \mathbf{w}_1^0, \dots, \mathbf{w}_{n-1}^0\}$ as linguistic inputs. Stacked bi-directional LSTM layers are built on top of word embedding sequence for sentence context modeling [46]. We also add shortcut connections from the input to the output of each LSTM layer for more efficient training.

$$\begin{aligned} \{\mathbf{w}_0^i, \mathbf{w}_1^i, \dots, \mathbf{w}_{n-1}^i\} = & biRNN_i(\{\mathbf{w}_0^{i-1}, \mathbf{w}_1^{i-1}, \dots, \mathbf{w}_{n-1}^{i-1}\}) \\ & + \{\mathbf{w}_0^{i-1}, \mathbf{w}_1^{i-1}, \dots, \mathbf{w}_{n-1}^{i-1}\}, \end{aligned} \quad (2)$$

where $\{\mathbf{w}^i\}$ is the result word-level feature sequence of the i -th layer. After p layers of bi-directional LSTMs, the output sequence $\{\mathbf{w}_0^p, \mathbf{w}_1^p, \dots, \mathbf{w}_{n-1}^p\}$ ($\{\mathbf{w}^p\}$ for simplicity) is the output of the language subnet which will be used as the input for the cross-modal attention mechanism.

We choose Gated Linear Unit (GLU) [7] as the non-linearity function of convolutions for language modeling. GLU processes a input vector $[\mathbf{q}_0, \mathbf{q}_1]$ by,

$$GLU([\mathbf{q}_0, \mathbf{q}_1]) = \mathbf{q}_0 \odot \sigma(\mathbf{q}_1), \quad (3)$$

where \mathbf{q}_0 and \mathbf{q}_1 are the inputs to the GLU with dimension d , σ is the sigmoid function, and \odot is element-wise multiplication. GLU has been used in context modeling in natural language processing [7, 11]. The gate $\sigma(\mathbf{q}_1)$ dynamically controls data flow via the gating mechanism. GLU shows superior performance than ReLU in our preliminary experiments.

3.3. Cross-modal attention module

To combine our video subnet and language subnet, we introduce a cross-modal attention module that combines visual and linguistic information. The attention module attends word-level representations for each frame, discovers connection between word-level linguistic information and frame-level visual information. It yields more detailed understanding of video content and description semantics. Comparing to using sentence-level linguistic feature, using word-level linguistic feature also shortens the path of gradient flow back to each word. In addition, this attention module is applied on every dilation layer and squeezing layer. Since different layers aggregate temporal contextual information at different temporal scales, feeding language feature on every layer enables our model to exploit video-language relevance at multiple temporal scales. Through this attentive cross-modal fusion process, our ASST translates the input visual sequence into the final representation sequence that consists of joint language-video representation, which can be easily recognized by clip samplers. In Figure 3, we show our attention module for one dilation layer.

For each visual layer after temporal convolution $\{\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_{m-1}\}$ ($\{\mathbf{v}\}$ for simplicity), the attention module takes language feature $\{\mathbf{w}^p\}$ and $\{\mathbf{v}\}$ as inputs. We split language feature $\{\mathbf{w}^p\}$ into two equal-length sequences by convolutions with kernel size 1: $\{\mathbf{w}_0^a, \mathbf{w}_1^a, \dots, \mathbf{w}_{n-1}^a\}$ for attention weight matrix computation and $\{\mathbf{w}_0^v, \mathbf{w}_1^v, \dots, \mathbf{w}_{n-1}^v\}$ for further feature computation. We also derive $\{\mathbf{v}_0^a, \mathbf{v}_1^a, \dots, \mathbf{v}_{m-1}^a\}$ from $\{\mathbf{v}\}$ for at-

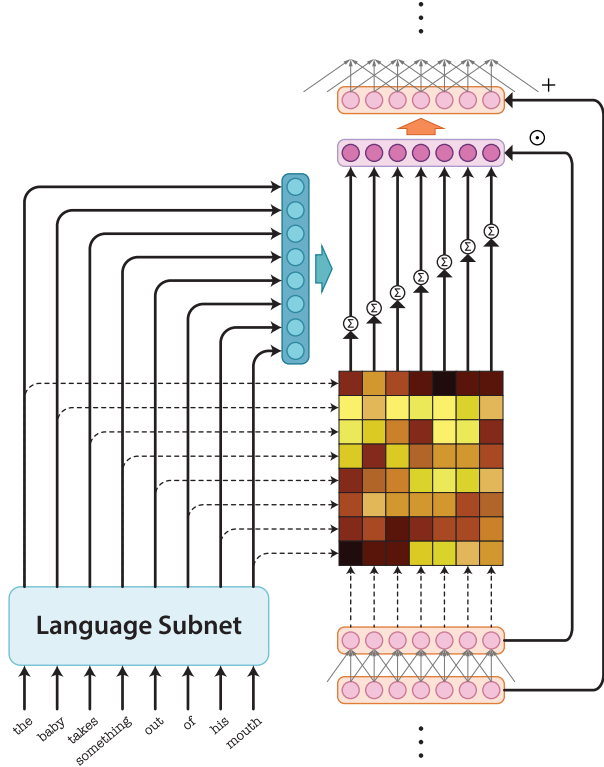


Figure 3. An illustration of the attention module for combining the video subnet and the language subnet on one dilation layer. This module takes language feature (left) and video feature (bottom right) as inputs. It first generates an attention weight matrix from both input features, and then attends language feature for every video frame. This attended language feature multiplies with input video feature. After a convolution operation, shortcut connection is used to connect visual feature maps across adjacent layers. Best viewed in color.

tention weight matrix computation. The attention weight matrix is computed as,

$$A'_{ij} = \frac{\mathbf{w}_i^a \cdot \mathbf{v}_j^a}{d_a},$$

$$A_{ij} = \frac{e^{A'_{ij}}}{\sum_i e^{A'_{ij}}},$$
(4)

where d_a is the dimension of \mathbf{w}_i^a and \mathbf{v}_j^a , and A is the attention weight matrix with size $n \times m$. We use this matrix to compute vision attended language feature,

$$\mathbf{u}'_j = \sum_i A_{ij} \mathbf{w}_i^v.$$
(5)

We then apply element-wise multiplication on $\{\mathbf{u}'\}$ and $\{\mathbf{v}\}$,

$$\mathbf{u}_j = \mathbf{u}'_j \odot \mathbf{v}_j.$$
(6)

$\{\mathbf{u}\}$ is our language augmented context for this convolution layer. We apply a BatchNorm [16] followed by a

ReLU to normalize $\{\mathbf{u}\}$. Finally, we add normalized $\{\mathbf{u}\}$ to the feature map before the last temporal convolution, as the input for the next temporal convolution.

4. Experiments

4.1. Datasets

We use the DiDeMo dataset [2] and the Charades-STA dataset [9, 38] to evaluate our model for the natural language video clip localization task.

DiDeMo. DiDeMo is a dataset for natural language description localization in open-world videos. There are 33,005, 4,180 and 4,021 video-description pairs in the training, validation and testing subsets, respectively. To annotate video clips, videos in DiDeMo are segmented every 5 seconds. The max available video length in this dataset is 30 seconds. The ground truth of each video-description pair is annotated by multiple people on these segments. Therefore, for each video, there are 21 possible clips. As we described previously, we address this task as a 21-way classification task. Following [2], we use Rank@1, Rank@5 and mean Intersection over Union ($mIoU$) as evaluation metrics.

Charades-STA. Charades-STA is a dataset for natural language description localization in indoor videos. There are 12,408 and 3,720 video-description pairs in the training and testing subsets, respectively¹. There are no pre-defined segments provided as in the DiDeMo dataset. We consider this task as a 2-way temporal detection task. Following [9], we use Rank@1, $IoU = 0.5$, Rank@1, $IoU = 0.7$, Rank@5, $IoU = 0.5$ and Rank@5, $IoU = 0.7$ as evaluation metrics.

4.2. Implementation Details

We implement our ASST by TensorFlow [1]. The code is publicly available at <https://github.com/NeonKrypton/ASST>.

In our experiments, we build our ASST with one bi-directional LSTM layer for language subnet, four dilation layers, six squeezing layers and six expansion layers for video subnet. Our language subnet takes GloVe word embeddings [30] as input. Following previous literatures [9, 2], we use pre-trained Two-Stream ConvNets [39] feature and 3D-ConvNets [42] (C3D) feature as the visual input. Two-stream ConvNets consist of RGB network and optical flow network, which take original RGB images and optical flows as inputs, recognize static objects and motions in videos respectively. C3D takes a sequence of consecutive RGB frames as input. It models static objects and motions simultaneously by applying hierarchical 3D convolution operations. The channel size of LSTM states, dilation layers and

¹The authors did some cleaning to the dataset. Updated dataset and results can be found at <https://github.com/jiyanggao/TALL>

Method	Rank@1	Rank@5	<i>mIoU</i>
Frequency Prior [2]	19.40	66.38	26.65
CCA [2]	18.11	52.11	37.82
MCN [2]	28.10	78.21	41.08
Ours	32.38	78.44	47.49

Table 1. Natural language localization in videos results on test subset of the DiDeMo dataset.

Method	R@1 <i>IoU</i> =0.5	R@1 <i>IoU</i> =0.7	R@5 <i>IoU</i> =0.5	R@5 <i>IoU</i> =0.7
CTRL [9]	21.42	7.15	59.11	26.91
Ours	37.04	18.04	68.12	38.28
Ours w/ Two-Stream	42.72	24.06	71.32	43.98

Table 2. Natural language localization in videos results on test subset of the Charades-STA dataset.

squeezing-expansion layers are 512, 1, 024 and 512. We apply dropout with rate 0.5 on input visual feature, and 0.8 for the rest of the model. We use Adam [20] as our optimizer. The learning rate starts from 5×10^{-4} , and multiplies 0.9 every 2,500 steps. The batch size we used for training is 128. We use different training strategies for two datasets. The details are as follows.

DiDeMo. Following [2], we use VGG16 [40] trained on ImageNet [34] training set as our RGB network. Inception-BN [16] trained on UCF101 [41] split 1 training set [45] is used as our optical flow network. The frame sample rate is $\frac{128}{30}$, so that our ASST can observe entire 30-second video in one observation. There are more than one annotation for each video-description pair. We first filter out annotations which has no overlap with other annotations from the same pair. Then, we randomly choose one annotation from remaining annotations for each iteration. To stabilize predictions, we average predictions from multiple model checkpoints.

Charades-STA. Following [9], we use C3D [42] trained on Sports-1M [19] as our visual input. The frame sample rate is 4. During training, we randomly stretch or compress our model by a scale within [0.8, 1.25], then sample a clip of random time from training videos. Any clip sampler has an overlap $IoU \geq 0.5$ with ground truth clip is considered as a positive training sample. The positive-negative sampling ratio is set to 1 : 1. During inference, we slide our model along testing videos. Non-Maximum Suppression (NMS) with threshold 0.8 is applied during post processing.

4.3. Comparison with other methods

DiDeMo. We perform experiments on the DiDeMo dataset for natural language description localization in videos. Table 1 shows our results and the baselines provided in [2].

Our model outperforms MCN in all three metrics. In Rank@5, our ASST is better than MCN by 0.23%. In

Rank@1 and *mIoU*, our ASST significantly outperforms MCN by 4.28% and 6.41%, respectively. In Figure 4, we show some of our localization results. In the first example, our ASST successfully localized the sentence “a girl with a blue shirt and black backpack speaks into the camera” and “there is a mountain”. For the third text query, “woman sharing a landscape view”, our ASST failed to generate an accurate prediction, but the predicted clip has an overlap with the ground-truth annotation. The woman is not shown in the last segment, and our model failed to predict this segment as positive.

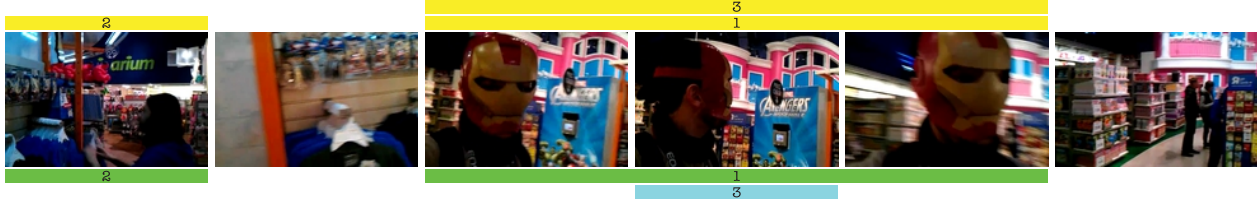
In the second example, our ASST successfully detected “we cant see the man’s head in these frames” and “person with camera tries to hand something to the guy in white lab coat”. Our ASST failed to localize “camera tilts then corrects”. This text query is describing the camera (observer). Events related to the observer needed to be inferred from object and camera motion, which is difficult to be directly modeled.

Charades-STA. We perform experiments on the Charades-STA dataset for natural language description localization in videos. Table 2 shows our results and the baseline provided by [9].

Our ASST achieves 37.04% in Rank@1, *IoU* = 0.5 and 18.04% in Rank@1, *IoU* = 0.7, which significantly outperforms CTRL by 15.62% and 10.89%, respectively. In Rank@5s, our ASST outperforms CTRL by 9.01% in *IoU* = 0.5 and 11.37% in *IoU* = 0.7. We also applied our model with Two-Stream visual feature as we used for the DiDeMo dataset. With Two-Stream visual feature, the performances of our model boosted to 42.72% in Rank@1, *IoU* = 0.5 and 24.06% in Rank@1, *IoU* = 0.7. This demonstrates that Two-Stream feature has a better ability to capture visual clues than C3D feature.

Figure 4 shows some of our localization results with C3D visual feature on the Charades-STA dataset. In the first example, our model successfully localized “a person sitting at a desk eating some food” and “person drinking from a coffee cup”. But for text query “a person is eating at their desk”, our model predicted a long sequence with eating and drinking together, which means our model is not good enough to clearly distinguish subtle actions. In the second example, our ASST successfully localized “person drinks from a glass”. For text query 2 “a smiling person runs into their garage holding a phone”, our model yields a longer prediction than the annotation. Our model failed to detect the action “person open the door”. This action happens at the very beginning of the video and is very short. The confidence of our predicted clip is very small, which means our model did not detect any positive clip. Better visual feature and more detailed modeling can be used to overcome these issues.

Text query 1: we can see someone in a mask
Text query 2: clerk hangs a blue shirt
Text query 3: a scary mask



Text query 1: the animal is eating the food
Text query 2: animal is closest to camera
Text query 3: the animal approaches the cameraman aggressively

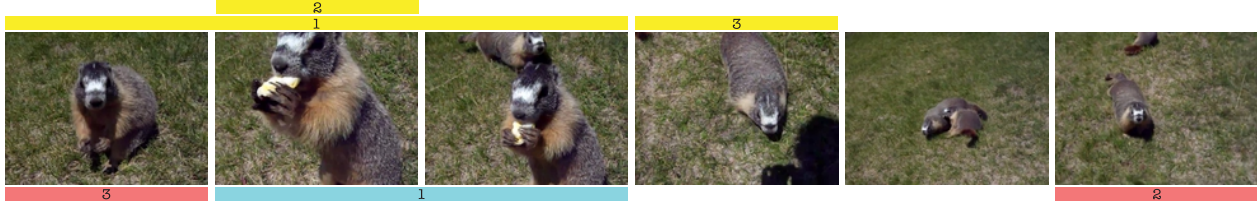


Figure 4. Visualization of our ASST’s experimental results of Rank@1. First two examples are from the DiDeMo dataset, and the last two are from the Charades-STA dataset. The yellow bars above videos indicate ground truths. The bars under videos indicate our predictions. A green bar represents our prediction is correct ($IoU \geq 0.7$ for Charades-STA). A blue bar represents our prediction overlaps with ground-truth annotation, but not accurate ($0.5 \leq IoU < 0.7$ for Charades-STA). A red bar represents our prediction does not overlap with ground truth ($IoU < 0.5$ for Charades-STA). Best viewed in color.

Method	Rank@1	Rank@5	$mIoU$
RGB	23.83	68.25	38.23
Flow	28.37	74.43	44.41
Two-Stream	30.91	76.46	47.30
RGB + Flow	29.95	77.15	45.83
Fusion	31.89	80.17	47.66

Table 3. Ablations study of input visual modalities on the DiDeMo dataset.

Method	Rank@1	Rank@5	$mIoU$
No language	21.10	69.35	33.31
Final representation layer	26.58	71.48	42.40
Last dilation layer	29.40	73.83	44.68
First dilation layer	30.12	74.45	46.97
Multiple feedings	30.91	76.46	47.30

Table 5. Ablation study of modeling cross-modal local relevance on the DiDeMo dataset.

Method	Rank@1	Rank@5	$mIoU$
Ours	30.91	76.46	47.30
Ours w/ TEF	30.53	77.34	47.14

Table 7. Ablation study of temporal endpoint feature.

4.4. Ablation studies

We then perform ablation studies from the following three aspects: input visual modality, the importance of cross-modal local relevance and temporal endpoint fea-

Method	R@1	R@1	R@5	R@5
	$IoU=0.5$	$IoU=0.7$	$IoU=0.5$	$IoU=0.7$
C3D	37.04	18.04	68.12	38.28
RGB	31.80	14.46	66.69	33.28
Flow	40.48	21.72	70.32	41.94
Two-Stream	42.72	24.06	71.32	43.98

Table 4. Ablation study of input visual modalities on the Charades-STA dataset.

Method	R@1	R@1	R@5	R@5
	$IoU=0.5$	$IoU=0.7$	$IoU=0.5$	$IoU=0.7$
No language	17.23	7.50	51.13	22.20
Final rep. layer	28.12	13.09	66.94	33.84
Last dilation layer	34.14	17.34	64.54	35.62
First dilation layer	35.99	17.72	62.39	33.63
Multiple feedings	37.04	18.04	68.12	38.28

Table 6. Ablation study of modeling cross-modal local relevance on the Charades-STA dataset.

ture. Ablations are performed on the validation set for the DiDeMo dataset and the test set for the Charades-STA dataset.

4.4.1 Input visual modality

We first perform ablation study on how the input visual modality influences our ASST’s performance. We evaluate our model using RGB, optical flow and concatenation of Two-Stream features as the input visual feature.

DiDeMo. The results of our model on the DiDeMo dataset

are shown in Table 3. Among these three models, RGB model achieves inferior performance. The optical flow model outperforms RGB model significantly by more than 4.54% in Rank@1. By concatenating both RGB and optical flow input feature, Two-Stream model further improves the optical flow model by 2.54% in Rank@1. The results show that both RGB and optical flow modality are important. Optical flow is still a good visual feature for video information modeling on tasks other than action recognition and detection.

By feeding different input video data, our ASST learns relevance between natural language descriptions and video content from different aspects. Directly fusing RGB and optical flow models by a weight of 1 : 2, the fused model achieves a higher performance in Rank@5, but lower performance in Rank@1 and $mIoU$ to single Two-Stream model. We then fuse all three models by a weight of 1 : 2 : 2.3. The performance further improves around 2% in all three metrics. The fusion weights are selected by cross-validation. This fusion model is also our final model used for comparing with baselines.

Charades-STA. The results of our model on the Charades-STA dataset are shown in Table 4. Same as the study on the DiDeMo dataset, RGB model achieves the worst performance among all models. C3D model outperforms RGB model over all metrics with better exploited spatio-temporal clues. Optical flow captures much better motion information than RGB and C3D model, outperforms RGB model and C3D model by 8.68% and 3.44% on Rank@1, $IoU = 0.5$, respectively. Concatenating RGB and optical flow as our input visual feature further improves Rank@1, $IoU = 0.5$ to 42.72%.

4.4.2 Modeling cross-modal local relevance

We perform ablation study to evaluate the necessity of early fusing natural language description and video content. We use plain video subnet as our baseline. This model contains no language information, and predicts clip localization directly from video content. To evaluate the effectiveness of cross-modal deep fusion and the importance of modeling cross-modal local relevance, we feed the language feature only once into video subnet via attention module. We choose three feeding positions: final representation layer, the last dilation layer and the first dilation layer.

DiDeMo. The results on the DiDeMo dataset are shown in Table 5. We use Two-Stream feature as our input visual feature. In this table, the model that feeds language information on the final representation layer achieves the worst performance among all models with one-time language feature feeding. By moving language feature feeding towards early stages, the performances increase. Feeding language feature on the first dilation layer achieves better performances

than feeding on the last dilation layer. Both models significantly outperform the model that feeds language information on the final representation layer.

Feeding language information multiple times achieves the best performance. The results demonstrate the importance of modeling video-language local relevance. However, feeding language feature on the final representation layer makes the model unable to build joint video-language hierarchies.

Charades-STA. The results on the DiDeMo dataset are shown in Table 5. We use C3D feature as our input visual feature. We observed similar behavior of feeding position as on the DiDeMo dataset. Moving language feature feeding towards early stages improves localization accuracy on Rank@1s.

4.4.3 Temporal endpoint feature

Finally, we evaluate the effect of temporal endpoint feature (TEF) proposed by Hendricks *et al.* [2]. We concatenate time coordinates to the final representation layer for each frame. The results are shown in Table 7.

TEF improves MCN significantly by 38.68% relatively on Rank@1 [2]. But TEF does not bring significant improvement to our ASST. One possible reason could be that during the process of visual information encoding, every element in the final representation sequence has a large receptive field. The temporal position information is encoded implicitly.

5. Conclusion and future work

In this paper, we proposed an effective attentive sequence to sequence translator for localizing clips by natural language descriptions. We demonstrated the effectiveness of modeling vision-language information jointly. Our standalone video subnet is also an effective model for video temporal modeling. Currently, our ASST only models temporal information of videos. Rich details on video frames are ignored. For future work, we will take detailed spatial information into consideration.

References

- [1] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al. Tensorflow: A system for large-scale machine learning. In *OSDI*, 2016. 6
- [2] L. Anne Hendricks, O. Wang, E. Shechtman, J. Sivic, T. Darrell, and B. Russell. Localizing moments in video with natural language. In *ICCV*, 2017. 1, 2, 3, 6, 7, 9
- [3] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *CVPR*, 2015. 1

- [4] Y.-W. Chao, S. Vijayanarasimhan, B. Seybold, D. A. Ross, J. Deng, and R. Sukthankar. Rethinking the faster r-cnn architecture for temporal action localization. In *CVPR*. 12
- [5] J. Dai, Y. Li, K. He, and J. Sun. R-fcn: Object detection via region-based fully convolutional networks. In *NIPS*, 2016. 5
- [6] X. Dai, B. Singh, G. Zhang, L. S. Davis, and Y. Qiu Chen. Temporal context network for activity localization in videos. In *ICCV*, 2017. 12
- [7] Y. N. Dauphin, A. Fan, M. Auli, and D. Grangier. Language modeling with gated convolutional networks. In *ICML*, 2017. 5
- [8] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, T. Mikolov, et al. Devise: A deep visual-semantic embedding model. In *NIPS*, 2013. 2
- [9] J. Gao, Z. Yang, K. Chen, C. Sun, and R. Nevatia. Turn tap: Temporal unit regression network for temporal action proposals. In *ICCV*, 2017. 1, 2, 3, 6, 7, 12
- [10] J. Gao, Z. Yang, and R. Nevatia. Cascaded boundary regression for temporal action detection. 2017. 12
- [11] J. Gehring, M. Auli, D. Grangier, D. Yarats, and Y. N. Dauphin. Convolutional sequence to sequence learning. In *ICML*, 2017. 5
- [12] R. Girshick. Fast r-cnn. In *ICCV*, 2015. 5
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 4
- [14] R. Hu, H. Xu, M. Rohrbach, J. Feng, K. Saenko, and T. Darrell. Natural language object retrieval. In *CVPR*, 2016. 3
- [15] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama, and K. Murphy. Speed/accuracy trade-offs for modern convolutional object detectors. In *CVPR*, 2017. 5
- [16] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 6, 7
- [17] Y.-G. Jiang, J. Liu, A. Roshan Zamir, G. Toderici, I. Laptev, M. Shah, and R. Sukthankar. THUMOS challenge: Action recognition with a large number of classes. <http://crcv.ucf.edu/THUMOS14/>, 2014. 1, 2, 12
- [18] N. Kalchbrenner, L. Espeholt, K. Simonyan, A. v. d. Oord, A. Graves, and K. Kavukcuoglu. Neural machine translation in linear time. *arXiv preprint arXiv:1610.10099*, 2016. 3
- [19] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. In *CVPR*. 7
- [20] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 7
- [21] C. Lea, M. D. Flynn, R. Vidal, A. Reiter, and G. D. Hager. Temporal convolutional networks for action segmentation and detection. In *CVPR*, 2017. 3
- [22] T. Lin, X. Zhao, and Z. Shou. Single shot temporal action detection. In *ACM MM*, 2017. 1, 3, 12
- [23] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 4
- [24] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollar. Focal loss for dense object detection. In *ICCV*, 2017. 4, 5
- [25] J. Liu, L. Wang, M.-H. Yang, et al. Referring expression generation and comprehension via attributes. In *CVPR*. 3
- [26] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot multibox detector. In *ECCV*, 2016. 3, 5
- [27] J. Mao, J. Huang, A. Toshev, O. Camburu, A. L. Yuille, and K. Murphy. Generation and comprehension of unambiguous object descriptions. In *CVPR*, 2016. 3, 5
- [28] A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016. 3
- [29] C. Pantofaru, C. Sun, C. Gu, C. Schmid, D. Ross, G. Toderici, J. Malik, R. Sukthankar, S. Vijayanarasimhan, S. Ricco, et al. Ava: A video dataset of spatio-temporally localized atomic visual actions. 2017. 2
- [30] J. Pennington, R. Socher, and C. Manning. Glove: Global vectors for word representation. In *EMNLP*, 2014. 6
- [31] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 5
- [32] A. Richard and J. Gall. Temporal action detection using a statistical language model. In *CVPR*, 2016. 12
- [33] A. Rohrbach, M. Rohrbach, R. Hu, T. Darrell, and B. Schiele. Grounding of textual phrases in images by reconstruction. In *ECCV*, 2016. 3
- [34] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015. 7
- [35] Z. Shou, J. Chan, A. Zareian, K. Miyazawa, and S.-F. Chang. Cdc: Convolutional-de-convolutional networks for precise temporal action localization in untrimmed videos. In *CVPR*, 2017. 1, 3, 12
- [36] Z. Shou, D. Wang, and S.-F. Chang. Temporal action localization in untrimmed videos via multi-stage cnns. In *CVPR*, 2016. 12
- [37] A. Shrivastava, R. Sukthankar, J. Malik, and A. Gupta. Beyond skip connections: Top-down modulation for object detection. *arXiv preprint arXiv:1612.06851*, 2016. 4
- [38] G. A. Sigurdsson, G. Varol, X. Wang, A. Farhadi, I. Laptev, and A. Gupta. Hollywood in homes: Crowdsourcing data collection for activity understanding. In *ECCV*, 2016. 1, 6
- [39] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. In *NIPS*, 2014. 1, 3, 6
- [40] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 7
- [41] K. Soomro, A. Roshan Zamir, and M. Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. In *CRCV-TR-12-01*, 2012. 7, 12
- [42] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features with 3d convolutional networks. In *ICCV*, 2015. 1, 3, 6, 7
- [43] S. Venugopalan, M. Rohrbach, J. Donahue, R. Mooney, T. Darrell, and K. Saenko. Sequence to sequence – video to text. In *ICCV*, 2015. 3

- [44] L. Wang, Y. Xiong, Z. Wang, and Y. Qiao. Towards good practices for very deep two-stream convnets. *arXiv preprint arXiv:1507.02159*, 2015. [3](#)
- [45] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *ECCV*, 2016. [1](#), [7](#), [12](#)
- [46] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, et al. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*, 2016. [5](#)
- [47] H. Xu, A. Das, and K. Saenko. R-c3d: Region convolutional 3d network for temporal activity detection. In *ICCV*, 2017. [1](#), [3](#), [12](#)
- [48] F. Yu and V. Koltun. Multi-scale context aggregation by dilated convolutions. In *ICLR*, 2016. [3](#)
- [49] Y. Zhao, Y. Xiong, L. Wang, Z. Wu, X. Tang, and D. Lin. Temporal action detection with structured segment networks. In *ICCV*, 2017. [1](#), [3](#), [12](#)

Appendix

To validate our video subnet’s ability for temporal visual information modeling, we perform experiments for a relevant task, “action detection”, on the THUMOS 14 dataset [17].

THUMOS 14 is a video dataset for action recognition and detection. There are 101 classes for action recognition task and 20 classes for detection task. For action detection task, usually untrimmed videos from validation set and testing set are used as training data and testing data. There are 1,010 videos in the validation set and 1,574 in the testing set respectively. Only 200 and 213 videos in these two sets contain annotations within the 20 action classes. We only use these 413 videos for training and evaluation following most recent literatures. The metrics for THUMOS 14 action detection task are $mAPs$ (mean average precision) on different $IoUs$ (Intersection over Union).

We use the detection model of standalone video subnet as our model. This model contains four dilation layers, six squeezing layers and six expansion layers, which is the same as we used for localizing clips by natural language descriptions task. Language subnet and attention module are removed for action detection task. Clip samplers with 21-way classifiers (20 classes and 1 background class) and temporal coordinate regression are built upon the final representation sequence. For a fair comparison, we use Inception-BN trained on ImageNet training set as our RGB network, and Inception-BN trained on UCF101 [41] split 1 training set [45] as our optical flow network, following [49].

We use Stochastic Gradient Descent (SGD) with batch size 128 to train our model. The learning rate starts with $1e - 2$, and multiplies 0.9 every 1,000 iterations. During training, we sample clips with positive:negative ratio 1 : 5. For inference, we slide our model along temporal domain for every testing video, followed by Non-Maximum-Suppression (NMS) with threshold 0.3 as post processing.

We show our action detection performances on the THUMOS 14 dataset in Table 8. We report our results at IoU 0.3 to 0.7 following Shou *et al.* [35].

Our model is able to outperform state-of-the-art approaches over all metrics. Our method outperforms earlier approaches [49, 10] by a large margin of more than 10% on $mAP@IoU$ 0.5. Comparing to very recent approaches [4], our model still achieves very competitive results of around 1% improvement over all metrics.

The video subnet encodes temporal visual context via multiple convolutional layers, which learns deep contextual temporal structure with multiple granularities, and is more capable of capturing video content at different scales, thereby achieving a better performance. The experimental results for action detection on the THUMOS 14 dataset show that our video subnet is able to capture good temporal visual structures for clip localization. The ability of mod-

Method	0.3	0.4	0.5	0.6	0.7
Richard <i>et al.</i> [32]	30.0	23.2	15.2	—	—
Shou <i>et al.</i> [36]	36.3	28.7	19.0	10.3	5.3
Shou <i>et al.</i> [35]	40.1	29.4	23.3	13.1	7.9
Lin <i>et al.</i> [22]	43.0	35.0	24.6	—	—
Dai <i>et al.</i> [6]	—	33.3	25.6	15.9	9.0
Gao <i>et al.</i> [9]	44.1	34.9	25.6	—	—
Xu <i>et al.</i> [47]	44.8	35.6	28.9	—	—
Zhao <i>et al.</i> [49]	51.9	41.0	29.8	19.6	10.7
Gao <i>et al.</i> [10]	50.1	41.3	31.0	19.1	9.9
Chao <i>et al.</i> [4]	53.2	48.5	42.8	33.8	20.8
Ours	54.9	50.3	43.7	34.4	22.2

Table 8. Action detection results on the THUMOS 14 dataset (in percentage). The IoU threshold used in evaluation varies from 0.3 to 0.7. - indicates the results in corresponding papers are unavailable.

eling temporal visual structures is an essential ingredient of our ASST, enables our ASST to better understand visual information and linguistic information jointly.