# DORi: Discovering Object Relationship for Moment Localization of a Natural-Language Query in Video

Cristian Rodriguez-Opazo[1,2]     Edison Marrese-Taylor[3]     Basura Fernando[4]
Hongdong Li[1,2]     Stephen Gould[1,2]

[1] Australian National University,     [2] Australian Centre for Robotic Vision (ACRV)

{cristian.rodriguez, hongdong.li, stephen.gould}@anu.edu.au

[3] Graduate School of Engineering, The University of Tokyo     [4] A*AI, A*STAR Singapore

emarrese@weblab.t-utokyo.ac.jp fernando.basura@scei.a-star.edu.sg

## Abstract

*This paper studies the task of temporal moment localization in a long untrimmed video using natural language query. Given a query sentence, the goal is to determine the start and end of the relevant segment within the video. Our key innovation is to learn a video feature embedding through a language-conditioned message-passing algorithm suitable for temporal moment localization which captures the relationships between humans, objects and activities in the video. These relationships are obtained by a spatial sub-graph that contextualizes the scene representation using detected objects and human features conditioned in the language query. Moreover, a temporal sub-graph captures the activities within the video through time. Our method is evaluated on three standard benchmark datasets, and we also introduce YouCookII as a new benchmark for this task. Experiments show our method outperforms state-of-the-art methods on these datasets, confirming the effectiveness of our approach.*

## 1. Introduction

Video analysis using natural language has been drawing increasing attention from the computer vision and natural language communities over the past few years, acknowledging the importance of these two modalities to understand video content. While promising results have been achieved on the tasks of video captioning and video question answering, much work still needs to be done to help identify and trim informative video segments in longer videos and align them with relevant textual descriptions. For this reason, tasks such as automatically recognizing *when* an activity is happening in a video have recently become crucial for video analysis.



Figure 1: An illustration of temporal localization of a natural language query in an untrimmed video. Given a query and a video the task is to identify the temporal start and end of the sentence in the video.

Moreover, as the amount of video data continues to grow, searching for specific visual events in large video collections has become increasingly relevant for search engines. This search engine requirement has helped draw increased attention to the task of activity detection in recent years. This task is especially important, considering that manually annotating videos is laborious and error-prone, even for a small number of videos. In this sense, it is clear that search engines have to retrieve videos not only based on video metadata but that they must also consider the videos' content in order to localize a given query accurately. Applications in practical areas such as video surveillance and robotics [33] have also helped bring interest in this task.

Regarding the issue mentioned above, in this paper we are specifically interested in the task of temporal sentence localization, in which given an untrimmed video and a natural language query, the goal is to identify the start and end points of the video segment (i.e., moment) that best corresponds to the given query. This task can be seen as a generalization of temporal action localization task [39, 31, 9, 4, 11, 50].

Many of the existing approaches to the localization problem in vision-and-language, either spatial or temporal, have

1

focused on creating a good multi-modal embedding space and generating proposals based on the given query. In these *propose and rank* approaches, candidate regions are first generated by a particular method and then fed to a classifier to get the probabilities of containing target classes, effectively ranking them. Most recently, approaches that do not rely on proposals to tackle this task have also been proposed, including the work of Ghosh et al. [14] and Rodriguez et al. [40].

Evidence shows that solving grounded language tasks such as ours often requires reasoning about relationships between objects in the context of the task [20]. For example, the work of Sigurdsson et al. [44] showed that the performance in action recognition tasks improves by a large margin if we have a perfect object recognition oracle. Moreover, we note that the majority of the queries that are used for this task are related to human actions. Our primary motivation is to reason about the relationship between humans and objects with the activity that they are performing. One can 'read a book' or 'look at the mobile.' A good way to know what the person is doing is to make use of object clues.

In light of this, in this paper we propose a mechanism to obtain contextualized activity representations based on a language-conditioned message passing algorithm. As activities are usually the result of the composition of several actions or interactions between a subject and objects [24], our algorithm incorporates both spatial and temporal dependencies. Therefore, modeling the relationship between subjects and objects in a scene and how these change over time, supporting the temporal moment localization task.

We conduct experiments on four challenging datasets, Charades-STA [10], ActivityNet [2, 27], TACoS [41] and YouCookII [58, 57], demonstrating the effectiveness of our proposed method and obtaining state-of-the-art performance. Our results highlight the importance of our message-passing algorithm in modeling the relationship between human and object and their interaction to understand the activity, ultimately validating our proposed approach. Our approach is the first to incorporate a language-conditioned message-passing algorithm to obtain contextualized activity representations using the objects and subjects to the best of our knowledge.

## 2. Related Work

Our work is related to the temporal action localization task, which aims to recognize and determine the temporal boundaries of action instances in videos. There is extensive previous work on this task, ranging from models that train existing video feature extractors with a localization loss [43], to systems that generally rely on temporal action proposal, as well as more sophisticated models that perform contextual modelling, capturing objects and their interactions [16, 15].

Since action localization is restricted to a pre-defined list of options, Gao et al. [10] and Hendricks et al. [19] introduced a variation of the task called language-driven temporal moment localization, where the goal is to determine the start and end time of the temporal video segment that best corresponds to a given natural language query. Early approaches for this task, including Liu et al. [32] and Ge et al. [12], were mainly based on generating proposals or candidate clips which could later be ranked. More recently, Chen et al. [5], Chen and Jiang [6], and Xu et al. [50], have worked on reducing the number of proposals by producing query-guided or query-dependent approaches.

Despite their ability to provide coarse control over the video snippets, proposal-based methods suffer from the computationally expensive candidate proposal matching, which has led to the development of methods that can directly output the temporal coordinates of the segment. In this context, Yuan et al. [52] first proposed to use a co-attention-based model, and soon after Ghosh et al. [14] focused directly on predicting the start and end frames using regressions. More recently, Rodriguez et al. [40] used dynamic filters and modeled label uncertainty to further improve performance, while Mun et al. [34] and Zeng et al. [54] proposed more sophisticated modality matching strategies. Compared to these works, although our approach is also proposal-free, we differ in the sense that we aim at incorporating specific spatial information that is useful for the localization problem.

Our work is also related to context modeling in action recognition. In this context, structural-RNN [22] models a spatio-temporal graph using an RNN mixture that is differentiable, with applications on human motion modeling and human activity detection. While we build on top of a concept similar to this, we inject the language component into the spatio-temporal graph and focus on the task of temporal moment localization of a natural language query. Our method adds the language into the pipeline using an attention mechanism that captures the objects' and subjects' interactions at the language level.

Context modeling has also been recently utilized in other computer vision tasks, such as referring expression comprehension [51] and VQA [20]. In the latter, the authors proposed a Language-Conditioned Graph Network (LCGN) where each node represents an object and is described by a context-aware representation from related objects through iterative message-passing conditioned on the textual input. Our work is fundamentally different from this as our task requires us to model the temporal component in our graph. Moreover, LCGN emphasizes the role of edge representations in the graph, whereas our approach is node-centric as connections between two given node types share the same edges.
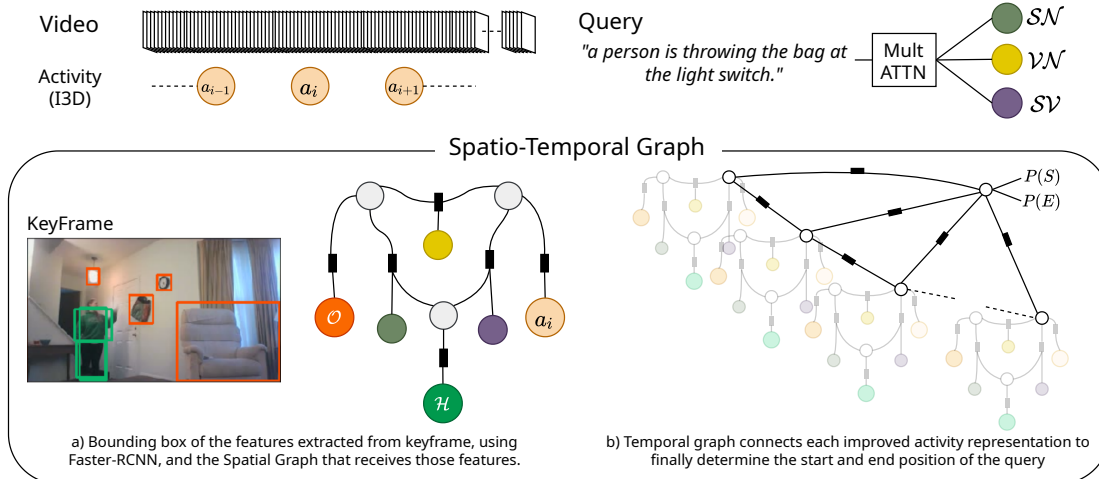
Figure 2: For each activity feature $a_i$, we create a Spatial graph to find the relationship between object and human nodes conditioned in the query, and thus improve the activity representation to be used by the Temporal graph.

Furthermore, Zeng et al. [53] used graph convolutions to obtain contextualized representations for action localization, while Zhang et al. [55] utilized a graph-structured network to model temporal relationships among different moments and thus obtain contextualized moment representations. Their approach is different from ours as they rely on proposals to perform the task. More recent approaches, such as SLTFNet [23], rely on attention instead of message-passing to deal with the spatio-temporal nature of the moment localization task.

Finally, Zhang et al. [56] have recently proposed a novel task that requires not only to perform temporal language-driven moment localization but also to locate the objects mentioned in the query spatially. Their approach is similar to ours in the sense that it also utilizes a spatio-temporal graph. However, the textual clues are incorporated after the graph construction rather than being an explicit part of it.

## 3. Graph-Based Temporal Moment Localization

In temporal moment localization, the objective is to find the temporal location of a natural language query $Q$ in an untrimmed video $V$. The video consists of a sequence of frames $V = [v_t \mid t = 1, \ldots, \ell]$ and the query is a sequence of words $Q = [w_j \mid j = 1, \ldots, m]$ that describes a short moment in the video. We denote the starting and ending times of the moment described by query $Q$ as $t^s$ and $t^e$, respectively.

We propose a model that explicitly captures the relationship between objects and humans, as well as the activities performed in a video, using a spatio-temporal graph. Concretely, we utilize a language-conditioned message-passing algorithm, which allows us to obtain contextualized activity representations for better moment localization. Let

$\mathcal{G} = (\mathcal{V}, \mathcal{E}_S \cup \mathcal{E}_T)$ represent our spatio-temporal graph, where $\mathcal{V}$, $\mathcal{E}_S$ and $\mathcal{E}_T$ are the set of nodes, spatial edges and temporal edges, respectively, as can be seen in Figure 2.

We factorize our spatio-temporal graph into spatial and temporal sub-graphs, denoted by $\mathcal{G}_S = (\mathcal{V}, \mathcal{E}_S)$ and $\mathcal{G}_T = (\mathcal{V}, \mathcal{E}_T)$, respectively.

The spatial graph is designed to improve the *activity representations* by exploiting the relationships between objects and humans in a given scene conditioned on an attended language representation for each of this relationships. As we know [45], actions and moments are characterized by complex interactions between humans as well as human-object interactions. Our spatial sub-graph is designed to exploit these spatial relationships specifically. It is iteratively applied through the video (Figure 2.a.)

On the other hand, the temporal sub-graph is designed to model the relationships between the improved activity representations at different times to more efficiently localize the start and end points of the query in the video (Figure 2.b.)

### 3.1. Spatial graph

Consider the query presented in Figure 2, "*a person is throwing the bag at the light switch*". It describes what action (*verb*) is performed by a *subject* and what *objects* are involved in that action. Our spatial graph is designed to capture the relationships between these visual entities conditioned on the linguistic entities. As such, we decompose our graph into six semantically meaningful nodes, three for representing visual information and three for representing linguistic information.

### 3.1.1 Linguistic Nodes

We create language nodes to capture essential information in the query related to the visual input: the *subject-verb* relationship node $\mathcal{SV}$ (person-throwing), the *subject-object* relation node $\mathcal{SN}$ (person-bag/light switch) and the *verb-object* relation node $\mathcal{VN}$ (throwing-bag/light switch).

To obtain representations for each one of the linguistic nodes, we start by encoding each of the words $w_j$ for $j = 1, \ldots, m$ in the query $Q$ using a function $F_w : w \mapsto h$, which maps each word to a semantic embedding vector $h_j \in \mathbb{R}^{d_w}$, where $d_w$ defines the hidden dimension of the word embedding. Specifically, we use GLoVe embeddings [37] to obtain the vector representations for each word.

We then initialize three-headed multi-head attention module [48] using an aggregated, fixed-length query vector $q$. Concretely, we construct this vector using a bi-directional GRU [8] over the word embeddings and mean pooling, which allows us to more accurately capture the global meaning of the input query by first contextualizing each word representation. We project each of the word embeddings using a linear mapping to obtain the key $k$ components of multi-head attention. In the case of the values $v$ we use the contextualized word representations from the GRU. Each head attends these contextualized vectors and returns a re-weighted combination of them, using softmax$(qk^\top)v$ and aimed at understanding a specific relation between the visual nodes at the linguistic level.

### 3.1.2 Visual Nodes

As mentioned above, our spatial graph contains three semantically meaningful nodes that represent visual information, specifically an activity node $\mathcal{A}$, a human node $\mathcal{H}$ and an object node $\mathcal{O}$. This setting allows us to share factors for semantically similar observations taken from the video [22], which provides several advantages. First, the model can deal with more observations of objects and humans without increasing the number of parameters that need to be learnt. Second, we alleviate the problem of having jittered object detections in videos, specially due to objects appearing and disappearing across frames.

To capture the relationships between activity, human and object observations, we densely connect these nodes within a single video frame. Such relationships are commonly parameterized by factor graphs that convey how a function over the graph factorizes into simpler functions [29]. Similarly, we learn a non-linear mapping function for each of the semantically alike observations that are associated with the same semantic node. In this sense, each semantic node, human $\mathcal{H}$, object $\mathcal{O}$ and activity $\mathcal{A}$, is considered to be a latent representation of the corresponding observation. Let us take as an example the case of the object node $\mathcal{O}$, where we observe a *table* in the video, represented by a feature vector $x$, obtained directly from the object detector. In this case, we use a function $\Psi_{\mathcal{O}} \doteq \tanh(W_{\mathcal{O}} x + b_{\mathcal{O}})$. Similar mapping functions (with different parameters), namely $\Psi_{\mathcal{H}}$ and $\Psi_{\mathcal{A}}$ are defined for the other semantic nodes.

**Activity node**: We use a video encoder that generates a video representation summarizing spatio-temporal patterns directly from the raw input frames. Concretely, let $F_V$ be our video encoding function that maps a video into a sequence of vectors $[a_i \in \mathbb{R}^{d_v} \mid i = 1, \ldots, t]$. These features capture high-level visual semantics in the video. Note that length of the video, $\ell = |V|$, and the number of output features, $t = |F_V(V)|$, are different due to temporal striding. Specifically, in this work we model $F_V$ using I3D [3]. This method inflates the 2D filters of a well-known convolutional neural network, e.g., Inception [47, 21] or ResNet [18] for image classification to obtain 3D filters.

**Human and object nodes**: Activity representations are obtained using small clips of frames. This means that there may be a set of many frames from where to extract spatial information that is semantically relevant for each node. Utilizing every frame is computationally expensive and given the piece-wise smooth nature of video, this could also prove to be redundant. As such, in this work we propose to utilize key-frames associated to each activity representation to extract observations for human and object nodes. Since many frames in a video are blurry due to various reasons, e.g., the natural movement of the objects and the camera motion, we select the sharpest key-frame in the subset of frames. Here we use the Laplace variance algorithm [36], which is a well-known approach for measuring the sharpness of an image.

While our method is agnostic to the choice of object detector, in this work we use Faster RCNN [38, 1] for the detection and spatial representation of the objects in all key-frames. Our Faster RCNN detector is trained on the Visual Genome [28] dataset, which consists of 1,600 object categories. These categories are manually assigned to either the human and object nodes depending on the type of object. The human node receives the set of features $H = \{h_1, ..., h_K\}$ corresponding to the categories associated to human body parts, clothes and subjects, while the object node receives the set of features that are not associated to human labels with that $O = \{o_1, ..., o_J\}$. This label-based categorization is based on a manual analysis of the label names supported by the Faster RCNN detector. In this way, when taking the predicted labels for each object we can use our categorization to re-label them as human or object and thus assign each instance to their corresponding visual node.

### 3.1.3 Language-conditioned message-passing

We argue that the setting of the scene contains important clues to improve the representation of a given activity. Examples of these clues are human clothes, objects that are present in the scene as well as their appearance. To the best of our knowledge, previous work on moment localization has not utilized this information. Therefore, we propose to obtain an activity representation suitable for the moment localization task, by capturing object, human and activity relationships. Concretely, we use a mean-field like approximation of the message-passing algorithm to capture such relationships. The messages sent between nodes are conditioned on the natural language query. We propose to use this approximation instead of the standard message-passing algorithm due to high demand on memory and compute, specially to process all the key-frames in a given video. The messages are iteratively sent a total of $N$ times, which is a hyperparameter of our model. In the equations below, index $n = 1, ..., N$ denotes the iteration step for each of the nodes. Notice that in the rest of this subsection, we drop the temporal index $i$ in the activity feature $a$ since the message-passing is done for each of the activity features independently.

First, we capture the relationship between the visual observations of the nodes human $\mathcal{H}$, object $\mathcal{O}$ and activity $\mathcal{A}$ with the corresponding language nodes $\mathcal{SN}, \mathcal{SV}$ and $\mathcal{VN}$ that connect the semantic meaning of the visual nodes, using a linear mapping function $f$ specific for each node. For instance, in the case of the object observations, the mapping functions $f$ have the following shape.

$$f_{\mathcal{SN},\mathcal{O}}(\mathcal{SN}, o^{j,n}) = W_{sno}[\mathcal{SN}; o^{j,n}] + b_{sno} = \Phi_{\mathcal{SN},\mathcal{O}}^{j,n} \quad (1)$$

$$f_{\mathcal{VN},\mathcal{O}}(\mathcal{VN}, o^{j,n}) = W_{vno}[\mathcal{VN}; o^{j,n}] + b_{vno} = \Phi_{\mathcal{VN},\mathcal{O}}^{j,n} \quad (2)$$

where $j$ is the j-th object observation in the object node $\mathcal{O}$. Similarly, we have specific mapping functions $f_{\mathcal{SV},\mathcal{A}}(\mathcal{SV}, a^n) = \Phi_{\mathcal{SV},\mathcal{A}}^n$, $f_{\mathcal{VN},\mathcal{A}}(\mathcal{VN}, a^n) = \Phi_{\mathcal{VN},\mathcal{A}}^n$, $f_{\mathcal{SN},\mathcal{H}}(\mathcal{SN}, h^{k,n}) = \Phi_{\mathcal{SN},\mathcal{A}}^{k,n}$, and $f_{\mathcal{SV},\mathcal{H}}(\mathcal{SV}, h^{k,n}) = \Phi_{\mathcal{SV},\mathcal{A}}^{k,n}$ for the activity and human observations, where $k$ is the k-th human observation.

For clarity, we explain the message-passing algorithm again using the object node as an example. The object node $\mathcal{O}$ receives messages from the human $\mathcal{H}$ and the activity $\mathcal{A}$ nodes. The message from the human node is constructed using a linear mapping function that receives as an input the concatenation of the object-query relationship $\Phi_{\mathcal{SN},\mathcal{O}}^{j,n}$ and the aggregation of all the human-query relationships $\sum_k \Phi_{\mathcal{SN},\mathcal{H}}^{k,n}$. A similar process is done for the message received from the activity observations, as can be seen in Equation 4, below.

$$\Psi_{\mathcal{H},\mathcal{SN},\mathcal{O}}^{j,n} = f_{\mathcal{H},\mathcal{SN},\mathcal{O}}(\Phi_{\mathcal{SN},\mathcal{O}}^{j,n}, \sum_{k=1}^{K} \Phi_{\mathcal{SN},\mathcal{H}}^{k,n}) \quad (3)$$

$$\Psi_{\mathcal{A},\mathcal{VN},\mathcal{O}}^{j,n} = f_{\mathcal{A},\mathcal{VN},\mathcal{O}}(\Phi_{\mathcal{VN},\mathcal{O}}^{j,n}, \Phi_{\mathcal{VN},\mathcal{A}}^n) \quad (4)$$

$$o^{j,n+1} = \sigma(m_o(\Psi_{\mathcal{H},\mathcal{SN},\mathcal{O}}^{j,n} \odot \Psi_{\mathcal{A},\mathcal{VN},\mathcal{O}}^{j,n}) \odot o^{j,0}) \quad (5)$$

Finally, the new representation of the object observation is computed using Equation 5, where $o^{j,0}$ is the initial object representation, $\sigma$ is an activation function, $\odot$ is the Hadamard product and $m_o$ is a linear function with a bias that constructs the message for the object $o^j$. A similar process is applied for each observation, as can be seen in Equations 6 to 9 below, where we create the message for each edge and Equations 10 to 11 and show how these messages are later used to contextualize the features. Note that the parameters learnt for each specific case are shared. For instance, parameters for $f_{\mathcal{A},\mathcal{SV},\mathcal{H}}$ and $f_{\mathcal{H},\mathcal{SV},\mathcal{A}}$ are the same.

$$\Psi_{\mathcal{H},\mathcal{SV},\mathcal{A}}^n = f_{\mathcal{H},\mathcal{SV},\mathcal{A}}(\Phi_{\mathcal{SV},\mathcal{A}}^n, \sum_{k=1}^{K} \Phi_{\mathcal{SV},\mathcal{H}}^{k,n}) \quad (6)$$

$$\Psi_{\mathcal{O},\mathcal{VN},\mathcal{A}}^n = f_{\mathcal{O},\mathcal{VN},\mathcal{A}}(\Phi_{\mathcal{VN},\mathcal{A}}^n, \sum_{j=1}^{J} \Phi_{\mathcal{VN},\mathcal{O}}^{j,n}) \quad (7)$$

$$\Psi_{\mathcal{O},\mathcal{SN},\mathcal{H}}^{k,n} = f_{\mathcal{O},\mathcal{SN},\mathcal{H}}(\Phi_{\mathcal{SN},\mathcal{H}}^{k,n}, \sum_{j=1}^{J} \Phi_{\mathcal{SN},\mathcal{O}}^{j,n}) \quad (8)$$

$$\Psi_{\mathcal{A},\mathcal{SV},\mathcal{H}}^{k,n} = f_{\mathcal{A},\mathcal{SV},\mathcal{H}}(\Phi_{\mathcal{SV},\mathcal{H}}^{k,n}, \Phi_{\mathcal{SV},\mathcal{A}}^n) \quad (9)$$

$$a^{n+1} = \sigma(m_a(\Psi_{\mathcal{H},\mathcal{SV},\mathcal{A}}^n \odot \Psi_{\mathcal{O},\mathcal{VN},\mathcal{A}}^n) \odot a^0) \quad (10)$$

$$h^{k,n+1} = \sigma(m_h(\Psi_{\mathcal{O},\mathcal{SN},\mathcal{H}}^{k,n} \odot \Psi_{\mathcal{A},\mathcal{SV},\mathcal{H}}^{k,n}) \odot h^{k,0}) \quad (11)$$

### 3.2. Temporal graph

The temporal graph is responsible for predicting the starting and ending points of the moment in the video. It uses the previously computed activity representations $a^{i,N}$ for $i = 1, \ldots, t$ where $N$ is the final iteration in the message passing. The temporal graph is implemented using a 2-layer bi-directional GRU [7] which receives as input the improved activity representation, and it is designed to contextualize the temporal relationship between the activity features. To obtain a probability distribution for the start and end predicted positions, we utilize two different fully connected layers to produce scores associated to the probabilities of each output of the GRU being the start/end of the location. Then, we take the softmax of these scores and thus obtain vectors $\hat{\tau}^s, \hat{\tau}^e \in \mathbb{R}^T$ containing a categorical probability distribution. Even though we do not constrain the starting and ending points to follow the right order in time, this does not result in any difficulties in practice.

### 4. Training

Our method is trained end-to-end on a dataset consisting of annotated tuples $(V, Q, t^s, t^e)$. Note that each video $V$ may include more than one moment and may therefore appear in multiple tuples. We treat each training sample independently. Given a new video and sentence tuple $(V_r, Q_r)$, our model predicts the most likely temporal localization of the moment described by $Q_r$ in terms of its start and end positions, $t_r^{s\star}$ and $t_r^{e\star}$, in the video. We use the

Kullback-Leibler divergence and an spatial loss proposed by Rodriguez et al. [40]. We explain this in more detail in the supplemental material. Given the predicted/ground truth starting/ending times of the moment, we use the following loss function during training:

$$L_{\text{KL}} = D_{\text{KL}}(\hat{\boldsymbol{\tau}}^s \parallel \boldsymbol{\tau}^s) + D_{\text{KL}}(\hat{\boldsymbol{\tau}}^e \parallel \boldsymbol{\tau}^e) \qquad (12)$$

where $D_{\text{KL}}$ is the Kullback-Leibler divergence. Moreover, inspired by Rodriguez et al. [40], we use a spatial loss that aims to create activity features that are good at identifying where the action is occurring. This loss, equation 13, receives as input $\mathbf{y} = \text{softmax}(g(\mathbf{a}))$ where $\mathbf{a}$ is the matrix that results by concatenating the improved activity representations over time, and $g$ is a linear mapping that gives us a score for each activity representation. We apply a softmax function over these and our loss penalizes if this normalized score is large for those features associated to positions that lie outside the temporal location of the query.

$$L_{\text{spatial}} = -\sum_{i=1}^{t}(1 - \delta_{\tau^s \leq i \leq \tau^e})\log(1 - y^i) \qquad (13)$$

where $\delta$ is the Kronecker delta. The final loss for training our method is the sum of the two individual losses defined previously setting $\mathcal{L} = L_{\text{KL}} + L_{\text{spatial}}$. During inference, we predict the starting and ending positions using the most likely locations given by the estimated distributions, using $\hat{\tau}^s = \text{argmax}(\hat{\boldsymbol{\tau}}^s)$ and $\hat{\tau}^e = \text{argmax}(\hat{\boldsymbol{\tau}}^e)$. Since values correspond to positions in the feature domain of the video, so we convert them back to time positions.

## 5. Experiments and Results

To evaluate our proposed approach we work with three widely utilized and challenging datasets, namely Charades-STA [10], ActivityNet Caption [2, 27] and TACoS [41]. In addition to these, we also consider the YouCookII dataset [58, 57]. This decision is motivated by its activity-centric nature as YoucookII is built upon instructional videos making it an excellent candidate to evaluate our proposals.

**Charades-STA**: Collected from the Charades dataset [46] by adding sentences that describe actions in the videos, it consists on a total of 13,898 pairs of queries and temporal locations. We use the predefined train and test splits [10]. Videos are 31 seconds long on average, with 2.4 moments on average, each being 8.2 seconds long on average.

**ActivityNet Captions**: Introduced by [27] this dataset, which was originally constructed for dense video captioning, consists of 20k YouTube videos with an average length of 120 seconds. The videos contain 3.65 temporally sentence descriptions on average, where the average length of the descriptions is 13.48 words. Following previous work, we report the performance of our approach on the two existing validation sets combined [34].

**MPII TACoS**: Built on top of the MPII Compositive dataset, it consists of videos of cooking activities with detailed temporally-aligned text descriptions. There are 18,818 pairs of sentence and video clips in total, with the average video length being 5 minutes. We use the same splits as [10], consisting of 50% for training, 25% for validation and 25% for testing.

**YouCookII**: consists on 2,000 long untrimmed videos from 89 cooking recipes obtained from YouTube by [58]. Each step for cooking these dishes was annotated with temporal boundaries and aligned with the corresponding section of the recipe. Similarly to TACoS, the average video length is 5.26 minutes.

### 5.1. Implementation Details

We first pre-process the videos by extracting features of size 1024 using I3D feaures with average pooling, taking as input the raw frames of dimension $256 \times 256$, at 25fps. We use the pre-trained model trained on Kinetics for TACoS, ActivityNet and YouCookII released by [3]. For Charades-STA, we use the pre-trained model trained on Charades. We extract the top 15 objects detected in terms of confidence for each of the key-frames using Faster-RCNN.

All of our models are trained in an end-to-end fashion using ADAM [26] with a learning rate of $10^{-4}$ and weight decay $10^{-3}$. As mentioned earlier, our temporal graph is modeled using a two-layer BiGRU. We use a hidden size of 256 and to prevent over-fitting we add a dropout of 0.5 between the two layers.

### 5.2. Evaluation

We evaluate our model using two widely used metrics proposed by [10]. Firstly, we measure recall at various thresholds of the temporal Intersection over Union ($R@\alpha$) obtaining the percentage of predictions that have tIoU with ground truth larger than certain $\alpha$, with threshold values 0.3, 0.5 and 0.7. In addition to that, we also compute and report mean or averaged tIoU (mIoU).

### 5.3. Ablation Study

To show the effectiveness of our proposals we perform several ablation studies each aimed at assessing the contribution of different components of our model. All of our ablative experiments are based on a segment of the training split of the Charades-STA dataset. As mIoU provides a more comprehensive evaluation of the performance of our model we utilized this metric to select the best model configuration for the rest of the experiments in this paper.

Since feed-forwarding through our proposed spatial graph is an iterative process, we first studied the impact on performance of the number of iterations ($N$) utilized in the message-passing algorithm of our full model. As shown in Table 1, we experimented setting $N$ to a minimum value

Table 1: Performance when using a different number of iterations ($N$) for the message-passing algorithm, on a subsection of the training split of Charades-STA.

| $N$ | R@0.3 | R@0.5 | R@0.7 | R@0.9 | mIoU |
|---|---|---|---|---|---|
| 0 | 47.46 | 22.88 | 14.38 | 6.00 | 33.67 |
| 1 | 73.21 | 55.32 | 36.02 | 11.48 | 52.11 |
| 2 | 79.01 | 67.16 | 48.71 | 17.97 | 59.30 |
| 3 | **79.25** | **68.41** | **50.56** | **19.14** | **60.29** |
| 4 | 70.99 | 60.31 | 44.16 | 17.32 | 54.01 |

Table 2: Results of our ablation studies, performed on a section of the training split of Charades-STA.

| Model | R@0.3 | R@0.5 | R@0.7 | R@0.9 | mIoU |
|---|---|---|---|---|---|
| (1) No Graph | 44.32 | 13.46 | 7.66 | 2.50 | 31.09 |
| (2) No Node Types | 74.78 | 61.24 | 43.35 | 14.59 | 55.18 |
| (3) No $\mathcal{H}$ Node | 75.46 | 60.60 | 43.31 | 15.51 | 55.32 |
| (4) No $\mathcal{O}$ Node | 75.66 | 61.28 | 44.08 | 15.39 | 56.13 |
| (5) No LA | 76.79 | 66.32 | 49.92 | 20.87 | 58.93 |
| (6) No $L_{\text{spatial}}$ | 76.79 | 66.60 | **52.54** | **23.57** | 59.95 |
| Full Model | **79.25** | **68.41** | 50.56 | 19.14 | **60.29** |

of 0 (where nodes are not updated at all) up to a maximum number of 4 iterations. As expected, performance tends to improve with larger values of $N$, with a saturation point at $N = 3$. Based on these results, all of our models in the rest of this paper are trained utilizing three iterations.

Table 2 summarizes the results of our ablation studies, which include: (1) Concatenating the mean-pooling of the features extracted by Faster RCNN directly with the activity representation, therefore eliminating the human and object nodes (No Graph) to assess the relevance of our graph in using the spatial information. (2) Evaluating the importance of distinguishing between human versus object features by testing how our model performs when assigning all the detected features to one spatial node (No Node Types). In (3) and (4) we remove the use of human (No $\mathcal{H}$) and object (No $\mathcal{O}$) spatial information, respectively. (5) Assessing the contribution of the linguistic nodes (No LA) by modifying our graph so that it only contains a single textual node connected to the rest of the graph in a way analogous to our full model. (6) Testing the importance of the spatial loss $L_{\text{spatial}}$ which encourages our model to focus on the features within the segment of interest. As can be seen, the importance of each one of our studied components is validated as ablations always result in consistent performance drops in terms of both mIoU as well as tIoU at the majority of $\alpha$ thresholds. For details about our ablated models please check Sections B.1 and B.2 in the Supplementary Material.

### 5.4. Comparison with the state-of-the-art

We start by presenting our results on the YouCookII dataset, which we introduce as a baseline for this task since

Table 3: Performance comparison on YouCookII for different tIoU $\alpha$ levels.

| Method | R@0.3 | R@0.5 | R@0.7 | mIoU |
|---|---|---|---|---|
| Random | 4.84 | 1.72 | 0.60 | - |
| TMLGA | 33.48 | 20.65 | 10.94 | 23.07 |
| DORi | **43.73** | **29.93** | **17.61** | **30.43** |

so far it has not been considered by previous work. On this dataset we consider a random baseline that simply selects an arbitrary video segment as the moment for each example, and also used the official implementation of TMLGA [40] released by its authors as an additional baseline since it is a direct alternative to our approach also being proposal-free. Table 3 summarizes our obtained results, where it can be seen that DORi is able to outperform both the random baseline and TMLGA by a large margin, specially on the lower $\alpha$ bands.

Regarding benchmark datasets, we compare the performance of our proposed approach against several prior work selected from the literature. We consider a broad selection of models based on different approaches, specifically proposal-based techniques including CTRL [10], SAP [6], MAN [55] and CBP [49], as well as TripNet [17], a method based on reinforcement learning. In addition to that, we also compare our approach to more recent methods that do not rely on proposals, including ABLR [52], ExCL [14], TMLGA [40] and LGVTI [34], as well as our random baseline.

Table 4 summarizes our results on Charades-STA, ActivityNet Captions and TACoS, while also comparing the obtained performance to relevant prior work. It is possible to see that our method is able to outperform previous work by a consistent margin, specially for the $\alpha = 0.7$ band and also in terms of the mean tIoU (mIoU). Comparing results across these datasets, we also see that the performance of all models drops substantially on ActivityNet Captions and TACoS, compared to Charades-STA. We think this is mainly due to the nature of the dataset, which contains a considerable amount of video moments that only span a few seconds.

Finally, we also study the effect of using different a pre-trained model to obtain activity representations in our proposed approach. Concretely, we test the performance of our model using VGG-16 features instead of I3D on Charades-STA. Table 5 summarizes our obtained results and compares them to prior work also utilizing these features. As can be seen, although VGG features provide lower performance than I3D in our experiments, therefore experimentally validating our choice, our model is still able to outperform existing approaches also using these features by a large margin, showing the superiority of our proposed approach.

Table 4: Performance comparison of our approach with existing methods for different tIoU $\alpha$ levels. Values are reported on the validation split of Charades-STA and ActivityNet Captions, and test splits for the TACoS datasets. † Results for ABLR are as reported by [6].

| Method | Charades-STA | | | | ActivityNet | | | | TACoS | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | R@0.3 | R@0.5 | R@0.7 | mIoU | R@0.3 | R@0.5 | R@0.7 | mIoU | R@0.3 | R@0.5 | R@0.7 | mIoU |
| Random | - | 8.51 | 3.03 | - | 5.60 | 2.50 | 0.80 | | 1.81 | 0.83 | | - |
| CTRL | - | 21.42 | 7.15 | - | 28.70 | 14.00 | - | 20.54 | 18.90 | 13.30 | - | - |
| ABLR † | - | 24.36 | 9.00 | - | 55.67 | 36.79 | - | 36.99 | 18.90 | 9.30 | - | - |
| TripNet | 51.33 | 36.61 | 14.50 | - | 48.42 | 32.19 | 13.93 | - | 23.95 | 19.17 | 9.52 | - |
| CBP | 50.19 | 36.80 | 18.87 | 35.74 | 54.30 | 35.76 | 17.80 | 36.85 | 27.31 | 24.79 | 19.10 | 21.59 |
| MAN | - | 46.53 | 22.72 | - | - | - | - | - | - | - | - | - |
| EXCL | 65.10 | 44.10 | 22.60 | - | - | - | - | - | **44.20** | 28.00 | 14.60 | - |
| TMLGA | 67.53 | 52.02 | 33.74 | 48.22 | 51.28 | 33.04 | 19.26 | 37.78 | 24.54 | 21.65 | 16.46 | 22.06 |
| LGVTI | **72.96** | 59.46 | 35.48 | 51.38 | **58.52** | **41.51** | 23.07 | 41.13 | - | - | - | - |
| DORi | 72.72 | **59.65** | **40.56** | **53.28** | 57.89 | 41.49 | **26.41** | **42.78** | 31.80 | **28.69** | **24.91** | **26.42** |

Query: *"cover the dish with mashed potatoes"*



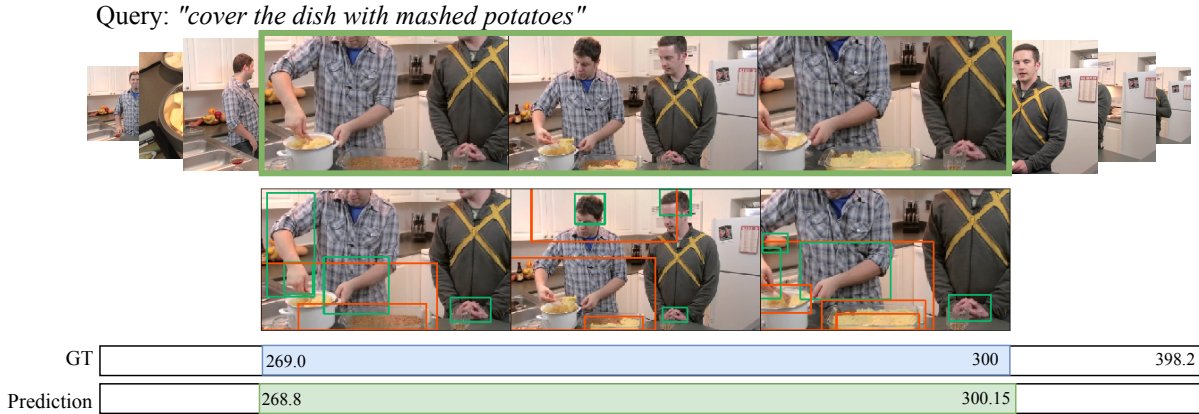| GT | | 269.0 | | 300 | | 398.2 |
| Prediction | | 268.8 | | 300.15 | | |

Figure 3: Visualization of a success case of our method in the YouCookII dataset. The second row shows the observations associated to the Human node (green) and Object node (orange).

Table 5: Performance of our method with VGG-16 features, We compared to relevant prior work that uses the same type of features.

| Model | R@0.3 | R@0.5 | R@0.7 | R@0.9 | mIoU |
|---|---|---|---|---|---|
| SAP | - | 27.42 | 13.36 | - | - |
| MAN | - | 41.24 | 20.54 | - | - |
| DORi | 61.83 | 43.47 | 26.37 | 7.63 | 42.52 |

## 5.5. Qualitative results

Figure 9 presents a success case of our method on YouCookII dataset. The visualization is presenting a sub-sample of the key-frames inside of the prediction with their corresponding spatial observations, with green observations associated with the human node $\mathcal{H}$ and orange to the object node $\mathcal{O}$. Moreover, each visualization is presenting the ground-truth localization and predicted localization of the given query. As shown in Figure 9, given the query "cover

the dish with mashed potatoes", our method could localize the moment at a tIoU of 98.88%. The most relevant features extracted by Faster-RCNN to localize the query are *'arm', 'bowl', 'cake', 'hand', 'kitchen', 'man', 'mug', 'spoon', 'stove', 'tray'*. Additional qualitative examples of success and failure cases of our method are included in Section D of the Supplementary Material.

## 6. Conclusion

We have presented a novel approach to temporal moment localization in video. Our approach consists of a spatial-temporal graph for capturing the relationships between detected humans, objects and activities over time. Conditioned on a natural language query, we proposed a message-passing algorithm that propagates information across the graph to ultimately infer the arbitrarily long segment in the video most likely described by the query. Using our approach we are able to achieve state-of-the-art results on several benchmark datasets.

# References

[1] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6077–6086, 2018.

[2] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 961–970, 2015.

[3] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*, 2017.

[4] Yu-Wei Chao, Sudheendra Vijayanarasimhan, Bryan Seybold, David A. Ross, Jia Deng, and Rahul Sukthankar. Rethinking the faster R-CNN architecture for temporal action localization. *CVPR*, 2018.

[5] Jingyuan Chen, Xinpeng Chen, Lin Ma, Zequn Jie, and Tat-Seng Chua. Temporally grounding natural sentence in video. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 162–171, Brussels, Belgium, 2018. Association for Computational Linguistics.

[6] Shaoxiang Chen and Yu-Gang Jiang. Semantic proposal for activity localizaiton in videos via sentence query. *AAAI*, 2019.

[7] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar, Oct. 2014. Association for Computational Linguistics.

[8] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.

[9] Victor Escorcia, Fabian Caba Heilbron, Juan Carlos Niebles, and Bernard Ghanem. DAPs: Deep Action Proposals for Action Understanding. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016*, Lecture Notes in Computer Science, pages 768–784. Springer International Publishing, 2016.

[10] Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. Tall: Temporal activity localization via language query. In *ICCV*, 2017.

[11] Jiyang Gao, Zhenheng Yang, Chen Sun, Kan Chen, and Ram Nevatia. TURN TAP: temporal unit regression network for temporal action proposals. *ICCV*, 2017.

[12] Runzhou Ge, Jiyang Gao, Kan Chen, and Ram Nevatia. Mac: Mining activity concepts for language-based temporal localization. In *WACV*, 2019.

[13] Cornelia Gerhardt, Maximiliane Frobenius, and Susanne Ley. *Culinary Linguistics*. John Benjamins Publishing, 2013.

[14] Soham Ghosh, Anuva Agarwal, Zarana Parekh, and Alexander Hauptmann. Excl: Extractive clip localization using natural language descriptions. *arXiv preprint arXiv:1904.02755*, 2019.

[15] Rohit Girdhar, Joao Carreira, Carl Doersch, and Andrew Zisserman. Video Action Transformer Network. In *CVPR*, pages 244–253, 2019.

[16] Chunhui Gu, Chen Sun, David A. Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sudheendra Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, Cordelia Schmid, and Jitendra Malik. AVA: A Video Dataset of Spatio-Temporally Localized Atomic Visual Actions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6047–6056, 2018.

[17] Meera Hahn, Asim Kadav, James M Rehg, and Hans Peter Graf. Tripping through time: Efficient localization of activities in videos. *arXiv preprint arXiv:1904.09936*, 2019.

[18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.

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

[20] Ronghang Hu, Anna Rohrbach, Trevor Darrell, and Kate Saenko. Language-Conditioned Graph Networks for Relational Reasoning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 10294–10303, 2019.

[21] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.

[22] Ashesh Jain, Amir R. Zamir, Silvio Savarese, and Ashutosh Saxena. Structural-RNN: Deep Learning on Spatio-Temporal Graphs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5308–5317, 2016.

[23] Bin Jiang, Xin Huang, Chao Yang, and Junsong Yuan. SLTFNet: A spatial and language-temporal tensor fusion network for video moment retrieval. *Information Processing & Management*, 56(6):102104, Nov. 2019.

[24] Yu-Gang Jiang, Subhabrata Bhattacharya, Shih-Fu Chang, and Mubarak Shah. High-level event recognition in unconstrained videos. *International journal of multimedia information retrieval*, 2(2):73–101, 2013.

[25] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2014. cite arxiv:1412.6980Comment: Published as a conference paper at the 3rd International Conference for Learning Representations, San Diego, 2015.

[26] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, 2014.

[27] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *ICCV*, 2017.

[28] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, Michael Bernstein, and Li Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. 2016.

[29] Frank R Kschischang, Brendan J Frey, and H-A Loeliger.

Factor graphs and the sum-product algorithm. *IEEE Transactions on information theory*, 47(2):498–519, 2001.

[30] Jing Lin, Chris Mellish, and Ehud Reiter. Style Variation in Cooking Recipes. page 5.

[31] Tianwei Lin, Xu Zhao, and Zheng Shou. Single Shot Temporal Action Detection. In *Proceedings of the 25th ACM International Conference on Multimedia*, MM '17, pages 988–996, New York, NY, USA, 2017. ACM. event-place: Mountain View, California, USA.

[32] Meng Liu, Xiang Wang, Liqiang Nie, Xiangnan He, Baoquan Chen, and Tat-Seng Chua. Attentive moment retrieval in videos. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pages 15–24. ACM, 2018.

[33] Yuxuan Liu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Imitation from observation: Learning to imitate behaviors from raw video via context translation. 2019.

[34] Jonghwan Mun, Minsu Cho, and Bohyung Han. Local-Global Video-Text Interactions for Temporal Grounding. *arXiv:2004.07514 [cs]*, Apr. 2020.

[35] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library, 2019. cite arxiv:1912.01703Comment: 12 pages, 3 figures, NeurIPS 2019.

[36] José Luis Pech-Pacheco, Gabriel Cristóbal, Jesús Chamorro-Martinez, and Joaquín Fernández-Valdivia. Diatom autofocusing in brightfield microscopy: a comparative study. In *Proceedings 15th International Conference on Pattern Recognition. ICPR-2000*, volume 3, pages 314–317. IEEE, 2000.

[37] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *EMNLP*, 2014.

[38] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NIPS*, 2015.

[39] Alexander Richard, Hilde Kuehne, Ahsan Iqbal, and Juergen Gall. Neuralnetwork-viterbi: A framework for weakly supervised video learning. In *IEEE Conf. on Computer Vision and Pattern Recognition*, volume 2, 2018.

[40] Cristian Rodriguez-Opazo, Edison Marrese-Taylor, Fatemeh Sadat Saleh, Hongdong Li, and Stephen Gould. Proposal-free temporal moment localization of a natural-language query in video using guided attention. *WACV*, 2020.

[41] Anna Rohrbach, Marcus Rohrbach, Wei Qiu, Annemarie Friedrich, Manfred Pinkal, and Bernt Schiele. Coherent multi-sentence video description with variable level of detail. In Xiaoyi Jiang, Joachim Hornegger, and Reinhard Koch, editors, *Pattern Recognition*, 2014.

[42] Marcus Rohrbach, Michaela Regneri, Mykhaylo Andriluka, Sikandar Amin, Manfred Pinkal, and Bernt Schiele. Script data for attribute-based recognition of composite activities. In *European Conference on Computer Vision*, pages 144–

157. Springer, 2012.

[43] Zheng Shou, Dongang Wang, and Shih-Fu Chang. Temporal action localization in untrimmed videos via multi-stage cnns. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.

[44] Gunnar A Sigurdsson, Olga Russakovsky, and Abhinav Gupta. What actions are needed for understanding human actions in videos? In *ICCV*, 2017.

[45] Gunnar A Sigurdsson, Olga Russakovsky, and Abhinav Gupta. What actions are needed for understanding human actions in videos? In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2137–2146, 2017.

[46] Gunnar A. Sigurdsson, Gül Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav Gupta. Hollywood in homes: Crowdsourcing data collection for activity understanding. In *European Conference on Computer Vision*, 2016.

[47] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015.

[48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.

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

[50] Huijuan Xu, Kun He, L Sigal, S Sclaroff, and K Saenko. Multilevel language and vision integration for text-to-clip retrieval. In *AAAI*, 2019.

[51] Licheng Yu, Zhe Lin, Xiaohui Shen, Jimei Yang, Xin Lu, Mohit Bansal, and Tamara L. Berg. MAttNet: Modular Attention Network for Referring Expression Comprehension. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1307–1315, 2018.

[52] Yitian Yuan, Tao Mei, and Wenwu Zhu. To find where you talk: Temporal sentence localization in video with attention based location regression. *AAAI*, 2019.

[53] Runhao Zeng, Wenbing Huang, Mingkui Tan, Yu Rong, Peilin Zhao, Junzhou Huang, and Chuang Gan. Graph Convolutional Networks for Temporal Action Localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7094–7103, 2019.

[54] Runhao Zeng, Haoming Xu, Wenbing Huang, Peihao Chen, Mingkui Tan, and Chuang Gan. Dense Regression Network for Video Grounding. *arXiv:2004.03545 [cs]*, Apr. 2020.

[55] Da Zhang, Xiyang Dai, Xin Wang, Yuan-Fang Wang, and Larry S Davis. Man: Moment alignment network for natural language moment retrieval via iterative graph adjustment. *CVPR*, 2019.

[56] Zhu Zhang, Zhou Zhao, Yang Zhao, Qi Wang, Huasheng Liu, and Lianli Gao. Where Does It Exist: Spatio-Temporal Video Grounding for Multi-Form Sentences. *arXiv:2001.06891 [cs]*, Jan. 2020.

[57] Luowei Zhou, Nathan Louis, and Jason J Corso. Weakly-supervised video object grounding from text by loss weighting and object interaction. In *British Machine Vision Confer-*

*ence*, 2018.

[58] Luowei Zhou, Chenliang Xu, and Jason J Corso. Towards automatic learning of procedures from web instructional videos. In *AAAI Conference on Artificial Intelligence*, pages 7590–7598, 2018.

# Appendices

## A. Datasets

**Charades-STA**: built upon the Charades dataset [46] which provides time-based annotations using a pre-defined set of activity classes, and general video descriptions. In [10], the sentences describing the video were semi-automatically decomposed into smaller chunks and aligned with the activity classes, which were later verified by human annotators. As a result of this process, the original class-based activity annotations are effectively associated to their natural language descriptions, totalling 13,898 pairs. We use the predefined train and test sets, containing 12,408 and 3,720 moment-query pairs respectively. Videos are 31 seconds long on average, with 2.4 moments on average, each being 8.2 seconds long on average. In our ablation studies we randomly split the train set in 80% for training and 20% for evaluation of the experiments.

**ActivityNet Captions**: Introduced by [27], this dataset which was originally constructed for dense video captioning, consists of 20k YouTube videos with an average length of 120 seconds. The videos contain 3.65 temporally localized time intervals and sentence descriptions on average, where the average length of the descriptions is 13.48 words. Following the previous methods, we report the performance of our algorithm on the combined two validation set.

**MPII TACoS**: built on top of the MPII Compositive dataset [42], it consists of videos of cooking activities with detailed temporally-aligned text descriptions. There are 18,818 pairs of sentence and video clips in total with the average video length being 5 minutes. A significant feature of this dataset is that due to the atomic nature of many of the descriptions —e.g. "takes out the knife" and "chops the onion"— the associated video moments only span over a few seconds, with 8.4% of them being less than 1.6 seconds long. This makes this dataset specially challenging for our task, as the relative brevity of the moments allows for a smaller margin of error. When it comes to splits, we use the same as in [10], consisting of 50% for training, 25% for validation and 25% for testing.

**YouCookII**: consists on 2,000 long untrimmed videos from 89 cooking recipes obtained from YouTube by [58]. Each step for cooking these dishes was annotated with temporal boundaries and aligned with the corresponding section of the recipe. Recipes are written following the usual style of the domain [30, 13], which includes very specific instruction-like statements with a wide degree of detail. The videos on this dataset are taped by individual persons at their houses while following the recipes using movable cameras. Similarly to TACoS, the average video length is 5.26 minutes. In terms of relevant moment segments, each video has 7.73 moments on average, with each segment being 19.63 seconds long on average. Videos have a minimum of 3 and a maximum of 16 moments.

Table 6 below summarizes the details of the exact sizes of the train/validation/test splits for each dataset.

| Dataset | Train | Validation | Test |
|---|---|---|---|
| Charades-STA | 12,408 | 3,720 | - |
| ActivityNet Captions | 37,414 | 17,502 | - |
| YouCookII | 10,337 | 3,492 | - |
| TACoS | 10,146 | 4,589 | 4,083 |

Table 6: Exact sizes of the train/validation/test splits for each dataset. Test sets for Charades-STA, ActivityNet and YouCookII are withheld, therefore, the common practice is to report results on the validation set instead.

## B. Ablated models

In the following sub-sections we give details about the ablated models presented in Section 4.3.

### B.1. No Node Type

This ablation experiment is intended to show the importance of considering the Faster-RCNN features related to human labels as a different source of information. The experiment consists of assigning the same 15 object features extracted for each of the keyframes only to the Object node $\mathcal{O}$. In this way limit the ability of the network to only be able to find relations between objects and activity representations, but without reducing the total amount of data that is available to it. We consider this experiment is very relevant as it shows that the additional information provided by the objects detected is not the only reason to explain the performance improvements, but rather the way in which this data is used is more relevant. In fact enabling the model to obtain state-of-the-art performance in different and challenging benchmarks.

### B.2. No Language Attention

In this case we replace the set of linguistic nodes by a single query node $\mathcal{Q}$. It receives a high-dimensional representation (denoted by $q$) of the natural language query $Q$, as can be seen in Figure 4. This high-dimensional representation is constructed using a function $F_Q : Q \mapsto q$ that first

maps each word $w_j$ for $j = 1, \ldots, m$ in the query to a semantic embedding vector $h_j \in \mathbb{R}^{d_w}$, where $d_w$ defines the hidden dimension of the word embedding. Representations for each word are then aggregated using mean pooling to get a semantically rich representation of the whole query.
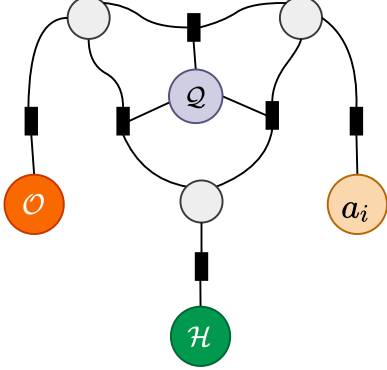


Figure 4: Spatial graph with a single query node $\mathcal{Q}$

Although the query node is generic, in this work we use a bi-directional GRU [7] on top of GLoVe word embeddings, which are pre-trained on a large collection of documents, for computing the $h_j$. Therefore, our query function $F_Q$ is parameterized by both GLoVe embedding and the GRU.

Again we capture capture the relationship between this high-dimensional representation of the query and any observation of the nodes human $\mathcal{H}$, object $\mathcal{O}$ and activity $\mathcal{A}$, using a linear mapping function $f$ specific for each node, as follows:

$$\Phi_{\mathcal{Q},\mathcal{A}}^{n} = f_{\mathcal{Q},\mathcal{A}}(q, a^n) \tag{14}$$

$$\Phi_{\mathcal{Q},\mathcal{O}}^{j,n} = f_{\mathcal{Q},\mathcal{O}}(q, o^{j,n}) \tag{15}$$

$$\Phi_{\mathcal{Q},\mathcal{H}}^{k,n} = f_{\mathcal{Q},\mathcal{H}}(q, h^{k,n}) \tag{16}$$

where functions $f_{\mathcal{Q},\mathcal{A}}, f_{\mathcal{Q},\mathcal{O}}, f_{\mathcal{Q},\mathcal{H}}$ are simple linear projections. For example, in the case of the object observations we have $f_{\mathcal{Q},\mathcal{O}}(q, o^{j,n}) = W_{qo}[q; o^{j,n}] + b_{qo}$, where the subindex $qo$ denotes the dependency of the parameters of the linear function which are specific for each relation. To compute the messages that are passed between the nodes, we utilize the following functions:

$$\Psi_{\mathcal{H},\mathcal{Q},\mathcal{O}}^{j,n} = f_{\mathcal{H},\mathcal{Q},\mathcal{O}}(\Phi_{\mathcal{Q},\mathcal{O}}^{j,n}, \sum_{k=1}^{K} \Phi_{\mathcal{Q},\mathcal{H}}^{k,n}) \tag{17}$$

$$\Psi_{\mathcal{A},\mathcal{Q},\mathcal{O}}^{j,n} = f_{\mathcal{A},\mathcal{Q},\mathcal{O}}(\Phi_{\mathcal{Q},\mathcal{O}}^{j,n}, \Phi_{\mathcal{Q},\mathcal{A}}^{n}) \tag{18}$$

$$\Psi_{\mathcal{H},\mathcal{Q},\mathcal{A}}^{n} = f_{\mathcal{H},\mathcal{Q},\mathcal{A}}(\Phi_{\mathcal{Q},\mathcal{A}}^{n}, \sum_{k=1}^{K} \Phi_{\mathcal{Q},\mathcal{H}}^{k,n}) \tag{19}$$

$$\Psi_{\mathcal{O},\mathcal{Q},\mathcal{A}}^{n} = f_{\mathcal{O},\mathcal{Q},\mathcal{A}}(\Phi_{\mathcal{Q},\mathcal{A}}^{n}, \sum_{j=1}^{J} \Phi_{\mathcal{Q},\mathcal{O}}^{j,n}) \tag{20}$$

$$\Psi_{\mathcal{O},\mathcal{Q},\mathcal{H}}^{k,n} = f_{\mathcal{O},\mathcal{Q},\mathcal{H}}(\Phi_{\mathcal{Q},\mathcal{H}}^{k,n}, \sum_{j=1}^{J} \Phi_{\mathcal{Q},\mathcal{O}}^{j,n}) \tag{21}$$

$$\Psi_{\mathcal{A},\mathcal{Q},\mathcal{H}}^{k,n} = f_{\mathcal{A},\mathcal{Q},\mathcal{H}}(\Phi_{\mathcal{Q},\mathcal{H}}^{k,n}, \Phi_{\mathcal{Q},\mathcal{A}}^{n}) \tag{22}$$

where again $f_{\mathcal{H},\mathcal{Q},\mathcal{O}}, f_{\mathcal{A},\mathcal{Q},\mathcal{O}}, f_{\mathcal{H},\mathcal{Q},\mathcal{A}}, f_{\mathcal{O},\mathcal{Q},\mathcal{A}}, f_{\mathcal{O},\mathcal{Q},\mathcal{H}}$ and $f_{\mathcal{A},\mathcal{Q},\mathcal{H}}$ are linear mappings, each receiving as input a con-

catenations of the corresponding features capturing. Finally, we update the representation of the human, action and object nodes based on the following formulas.

$$o^{j,n+1} = \sigma(m_o(\Psi_{\mathcal{H},\mathcal{Q},\mathcal{O}}^{j,n} \odot \Psi_{\mathcal{A},\mathcal{Q},\mathcal{O}}^{j,n}) \odot o^{j,0}) \tag{23}$$

$$a^{n+1} = \sigma(m_a(\Psi_{\mathcal{H},\mathcal{Q},\mathcal{A}}^{n} \odot \Psi_{\mathcal{O},\mathcal{Q},\mathcal{A}}^{n}) \odot a^{0}) \tag{24}$$

$$h^{k,n+1} = \sigma(m_h(\Psi_{\mathcal{O},\mathcal{Q},\mathcal{H}}^{k,n} \odot \Psi_{\mathcal{A},\mathcal{Q},\mathcal{H}}^{k,n}) \odot h^{k,0}) \tag{25}$$

where $\odot$ is the element-wise product and $m_o, m_a, m_h$ are again linear functions.

## C. Language Attention

In the following Figures 5 and 6, we present a set of samples of the multihead attention to the query sentence on the Charades-STA dataset.
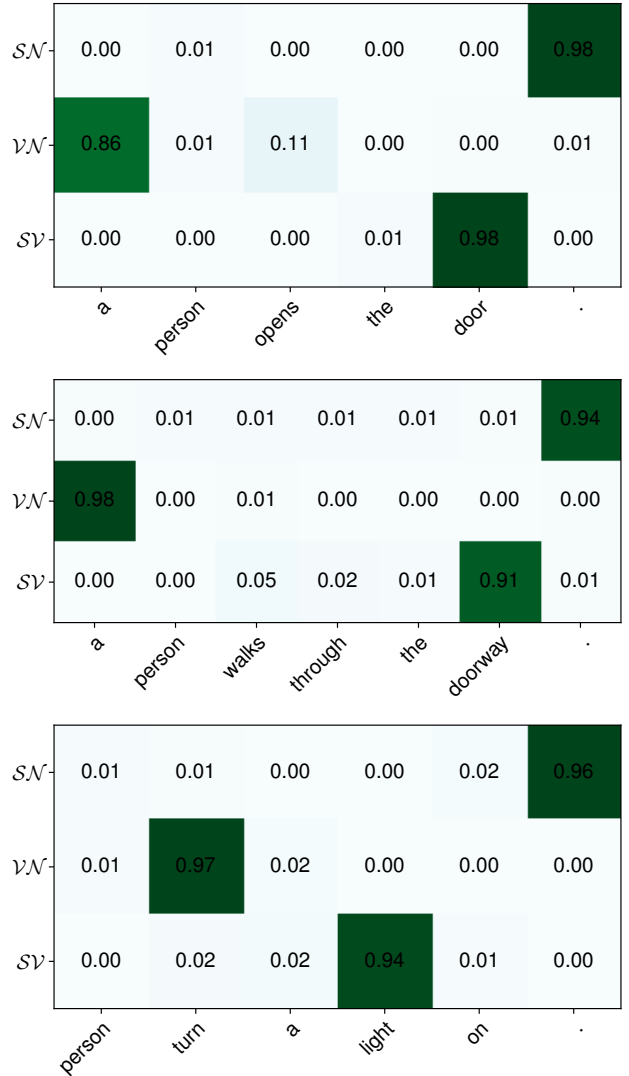

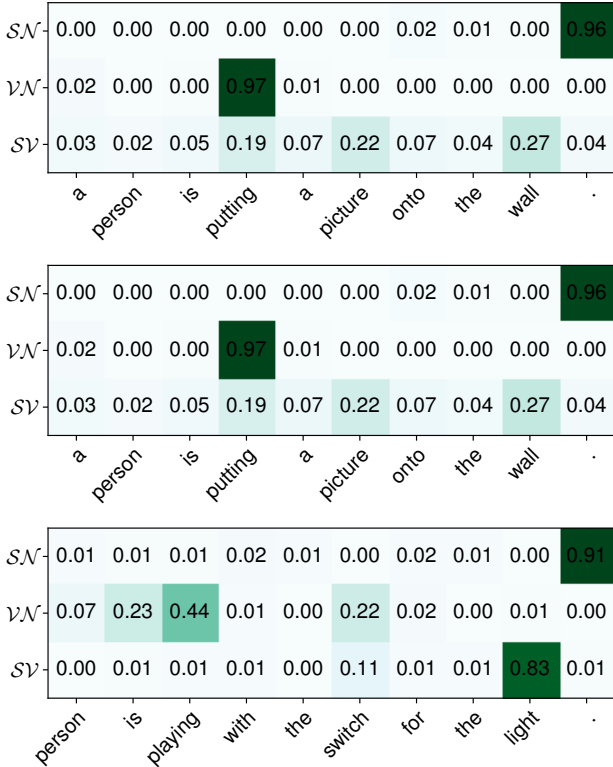
Figure 5: Linguistic nodes attentions on Charades-STA.

| | a | person | is | putting | a | picture | onto | the | wall | . |
|---|---|---|---|---|---|---|---|---|---|---|
| $\mathcal{SN}$ | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.02 | 0.01 | 0.00 | 0.96 |
| $\mathcal{VN}$ | 0.02 | 0.00 | 0.00 | 0.97 | 0.01 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| $\mathcal{SV}$ | 0.03 | 0.02 | 0.05 | 0.19 | 0.07 | 0.22 | 0.07 | 0.04 | 0.27 | 0.04 |

| | a | person | is | putting | a | picture | onto | the | wall | . |
|---|---|---|---|---|---|---|---|---|---|---|
| $\mathcal{SN}$ | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.02 | 0.01 | 0.00 | 0.96 |
| $\mathcal{VN}$ | 0.02 | 0.00 | 0.00 | 0.97 | 0.01 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| $\mathcal{SV}$ | 0.03 | 0.02 | 0.05 | 0.19 | 0.07 | 0.22 | 0.07 | 0.04 | 0.27 | 0.04 |

| | person | is | playing | with | the | switch | for | the | light | . |
|---|---|---|---|---|---|---|---|---|---|---|
| $\mathcal{SN}$ | 0.01 | 0.01 | 0.01 | 0.02 | 0.01 | 0.00 | 0.02 | 0.01 | 0.00 | 0.91 |
| $\mathcal{VN}$ | 0.07 | 0.23 | 0.44 | 0.01 | 0.00 | 0.22 | 0.02 | 0.00 | 0.01 | 0.00 |
| $\mathcal{SV}$ | 0.00 | 0.01 | 0.01 | 0.01 | 0.00 | 0.11 | 0.01 | 0.01 | 0.83 | 0.01 |

Figure 6: Linguistic nodes attentions on Charades-STA.

# D. Examples

In the following Figures, we present success and failure cases of our method on Charades-STA, YouCookII and TACoS dataset. Each visualization is showing a subsample of the keyframes inside of the prediction with their corresponding spatial observations. In green observations associated with the human node $\mathcal{H}$ and orange for the object node $\mathcal{O}$. Moreover, each visualization is presenting the ground-truth and predicted localization in seconds of the given query.

## D.1. Charades-STA

Success cases of our algorithm on the Charades-STA dataset can be seen in Figure 7. In Figure 7a, given the query "a person cooks a sandwich on a panini maker" our method could localize the moment at a tIoU of 99.56%. The label of the features extracted by Faster-RCNN to localize the query are *'bottle'*, *'counter'*, *'door'*, *'drawer'*, *'faucet'*, *'floor'*, *'glasses'*, *'hair'*, *'jacket'*, *'jeans'*, *'kitchen'*, *'microwave'*, *'pants'*, *'shelf'*, *'shirt'*, *'sink'*, *'stove'*, *'sweater'*, **'toaster'**, *'wall'*, *'window'*, *'woman'*.

In the case of Figure 7b, given the query "the person closes a cupboard door." our method could localize the moment at a tIoU of 97.88%. The features extracted by

Faster RCNN for this query are *'arm'*, *'building'*, **'cabinet'**, *'counter'*, *'door'*, *'faucet'*, *'hair'*, *'hand'*, *'head'*, *'jacket'*, *'kitchen'*, *'man'*, *'microwave'*, *'refrigerator'*, *'shirt'*, *'sink'*, *'sleeve'*, *'stove'*, *'sweater'*, *'wall'*, *'window'*, *'woman'*.

Failure cases of our method are presented in Figure 8. In the first example, given a query "a person opens a door goes into a room." our method could detect correct spatial features, such as 'door' and 'knob', and the correct span of the query, according to our qualitative evaluation. However, in this case, the annotation for the query is localized incorrectly in the video. It refers to the last part of the video, where a person is using a laptop, as can be seen at the right of Figure 8a. In Fig. 8b we can see our method localizing the query "person walks over to the refrigerator open it up", however, the annotation is not considering that the moment is performed two times in the video.

## D.2. YouCookII

Although videos in YouCookII are much longer than videos in Charades-STA, our method still can get good localization performance. In Figure 9a given the query "spread the sauce onto the dough" our method localize the query at a tIoU of 98.57%. The label of the feature extracted by Faster-RCNN on this case are *'bacon'*, *'bird'*, *'board'*, *'bottle'*, *'bowl'*, *'cabinet'*, *'cake'*, *'cherry'*, *'chocolate'*, *'cookie'*, *'counter'*, *'cutting board'*, *'dessert'*, *'door'*, *'drawer'*, *'finger'*, *'floor'*, *'fork'*, *'fruit'*, *'glass'*, *'grape'*, *'ground'*, *'hand'*, *'handle'*, *'jeans'*, *'ketchup'*, *'knife'*, *'meat'*, *'olive'*, *'pancakes'*, *'pepperoni'*, *'person'*, *'phone'*, **'pizza'**, *'plant'*, *'plate'*, **'sauce'**, *'saucer'*, *'shirt'*, *'sleeve'*, *'spoon'*, *'table'*, *'towel'*, *'tree'*, *'wall'*.

Figure 9b shows the query "cook the pizza in the oven", which belong to the same video. In this case the label of the features extracted by Faster-RCNN are *'arm'*, *'bar'*, *'board'*, *'building'*, *'cabinet'*, *'car'*, *'ceiling'*, *'cheese'*, *'cord'*, *'counter'*, *'crust'*, *'cucumber'*, *'curtain'*, *'door'*, *'drawer'*, *'fireplace'*, *'floor'*, *'food'*, *'fork'*, *'glass'*, *'grill'*, *'hand'*, *'hotdog'*, *'key'*, *'keyboard'*, *'kitchen'*, *'knife'*, *'knob'*, *'laptop'*, *'leaf'*, *'leaves'*, *'leg'*, *'light'*, *'man'*, *'microwave'*, *'mouse'*, **'oven'**, *'oven door'*, *'person'*, **'pizza'**, *'plate'*, *'pole'*, *'rack'*, *'roof'*, *'room'*, *'salad'*, *'screen'*, *'shadow'*, *'sleeve'*, *'slice'*, *'spinach'*, *'stove'*, *'table'*, *'television'*, *'thumb'*, *'tracks'*, *'train'*, *'tray'*, *'vegetable'*, *'vegetables'*, *'wall'*, *'window'*, *'wood'* and our method could localize the query with a temporal intersection over union of 97.60%.

Failure cases of our method on YouCookII dataset are presented in Figure 10. In these cases, it is possible to see that our approach is able to recognize the activity *add* and *mix* correctly. However, the objects "dressing, ginger and garlic" are not detected by Faster-RCNN, probably given that the object detector has not been trained to deal with some of the kinds of objects present on this dataset. We

think this naturally hinders the disambiguation capabilities of our model, specially in terms of the repetitive actions such as as adding, mixing and pouring, which are often performed throughout recipes like the one depicted in the example.

### D.3. TACoS

Figures 11 and 12 show two examples of success and failure cases on the TaCoS dataset, respectively. It is possible to see the how challenging this dataset is in general, as in the the cases where our approach fails it is in fact difficult even for us to localize the given query.

## E. Experimental Information

Our models are implemented using PyTorch [35] and are trained using the Adam [25] optimizer, with a batch size of 6. Experiments for different datasets were run in two different machines:

- First server machine with an Intel Core i7-6850K CPU with two NVIDIA Titan Xp (Driver 430.40, CUDA 10.1) GPUs, and one NVIDIA Quadro P5000, running ArchLinux

- An additional server machine with an Intel Xeon 4215 CPU, with three NVIDIA RTX8000 (Driver 430.44, CUDA 10.1) GPUs, running Ubuntu 16.04

We used PyTorch version 1.4. Our method has 10.865.155 trainable parameters. In training takes 1.56 hours per epoch in Charades-STA, 4.3 hours per epoch in TACoS, 5.4 hours per epoch in YouCookII and 6.7 hours per epoch in ActivityNet. In average our method takes 0.015 seconds to localize one query.

Query: *"a person cooks a sandwich on a panini maker."*



(a) Example of success 1.

Query: *"the person closes a cupboard door"*



(b) Example of success 2.

Figure 7: Success examples of our method on Charades-STA dataset.
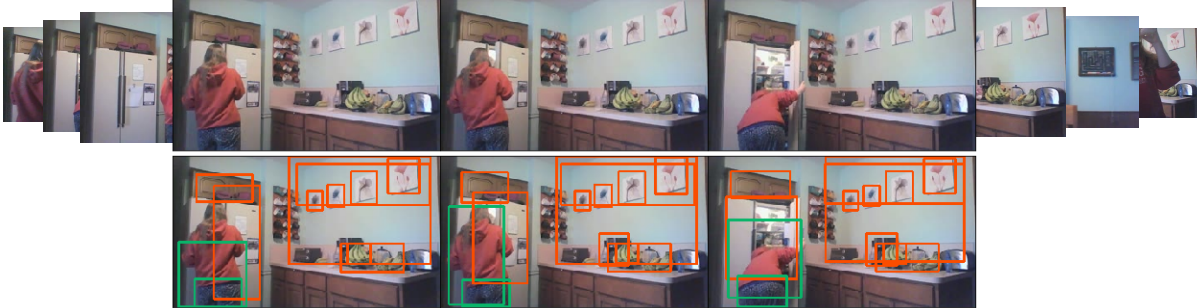
14

Query: *"a person opens a door goes into a room."*



| GT | | 16.8 | 30.4 |
| Prediction | 0.0 | 5.94 | |

(a) Example of failure 1.

Query: *"person walks over to the refrigerator open it up"*
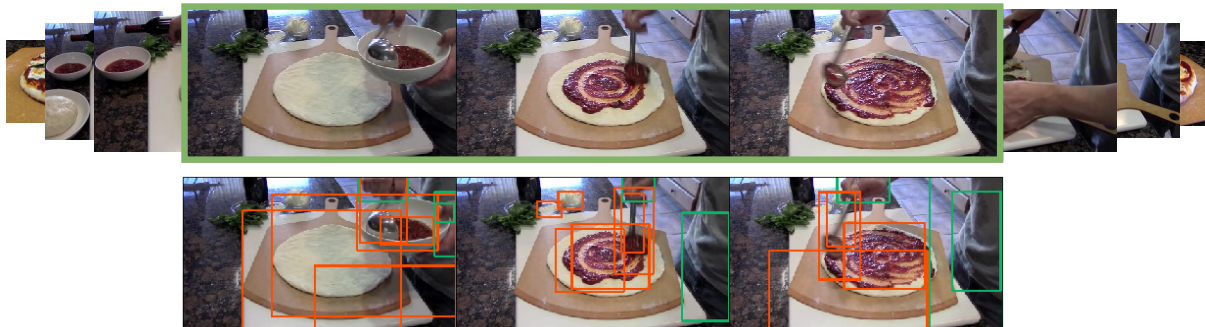


| GT | 0.0 | 4.9 | | 33.8 |
| Prediction | | 12.87 | 18.44 | |

(b) Example of failure 2.

Figure 8: Failure examples of our method on Charades-STA.

Query: *"spread the sauce onto the dough"*



| GT | 131.0 | 158.0 | 366.8 |
|---|---|---|---|
| Prediction | 130.64 | 157.97 | |

(a) Success example 1.

Query: *"cook the pizza in the oven"*



| GT | 257.0 | 288.0 | 366.8 |
|---|---|---|---|
| Prediction | 256.29 | 287.95 | |

(b) Success example 2.

Figure 9: Success examples of our method in the YouCookII dataset.

Query: *"pour the dressing over the salad and mix"*



| GT | 153.0 | 162.0 | | 206.88 |
| Prediction | | 164.84 | 186.86 | |

(a) Failure case 1.

Query: *"add oil ginger and garlic to a pot"*
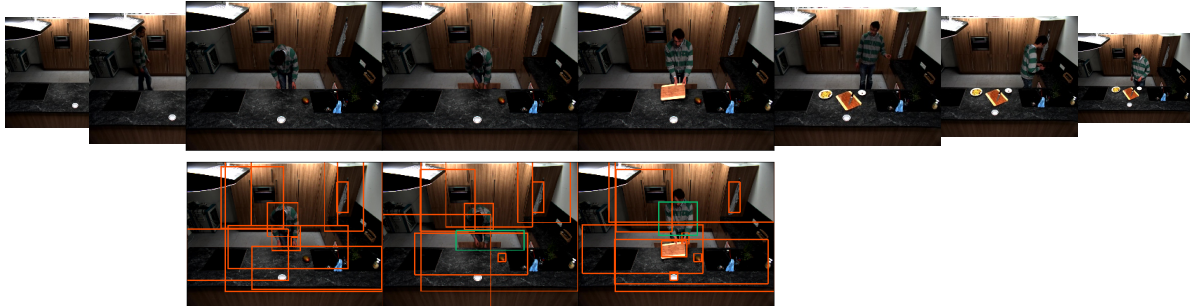


| GT | | 275.0 | 301.0 | 517.2 |
| Prediction | 135.39 | 141.05 | | |

(b) Failure case 2.

Figure 10: Failure cases of our method in the YouCookII dataset.

Query: *"The person gets out a cutting board."*



| GT | | 34.3 | | 42.24 | | | 365 |
| Prediction | | 34.39 | | 42.28 | | | |

(a) Success example 1.

Query: *"The person takes a bottle of oil and an onion from the pantry."*
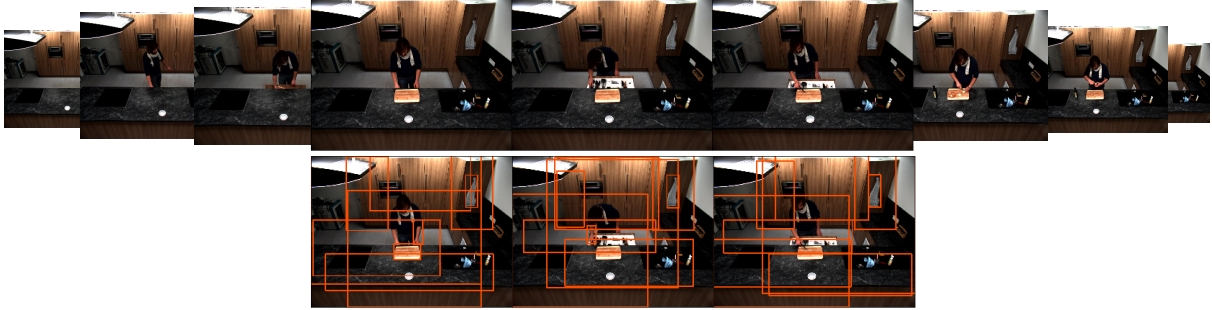


| GT | | 23.53 | | 34.58 | | | 574.5 |
| Prediction | | 23.50 | | 34.39 | | | |

(b) Success example 2.

Figure 11: Success examples of our method in the TACoS dataset.

Query: *"He takes the skin off of the onion."*



(a) Failure case 1.

Query: *"He sliced mango"*



(b) Failure case 2.

Figure 12: Failure cases of our method in the TACoS dataset.