

Weakly-Supervised Action Localization by Hierarchically-structured Latent Attention Modeling

Guiqin Wang^{1†} Peng Zhao¹ Cong Zhao^{3,4} Shusen Yang^{3,4} Jie Cheng² Luziwei Leng²
Jianxing Liao² Qinghai Guo^{2*}

¹ School of Computer Science and Technology, Xi'an Jiaotong University

² ACS Lab, Huawei Technologies

³ School of Mathematics and Statistics, Xi'an Jiaotong University

⁴ National Engineering Laboratory for Big Data Analytics, Xi'an Jiaotong University

Abstract

Weakly-supervised action localization aims to recognize and localize action instances in untrimmed videos with only video-level labels. Most existing models rely on multiple instance learning (MIL), where the predictions of unlabeled instances are supervised by classifying labeled bags. The MIL-based methods are relatively well studied with co-gent performance achieved on classification but not on localization. Generally, they locate temporal regions by the video-level classification but overlook the temporal variations of feature semantics. To address this problem, we propose a novel attention-based hierarchically-structured latent model to learn the temporal variations of feature semantics. Specifically, our model entails two components, the first is an unsupervised change-points detection module that detects change-points by learning the latent representations of video features in a temporal hierarchy based on their rates of change, and the second is an attention-based classification model that selects the change-points of the foreground as the boundaries. To evaluate the effectiveness of our model, we conduct extensive experiments on two benchmark datasets, THUMOS-14 and ActivityNet-v1.3. The experiments show that our method outperforms current state-of-the-art methods, and even achieves comparable performance with fully-supervised methods.

1. Introduction

Action localization is one of the most challenging tasks in video analytics and understanding [16, 43, 50, 17]. The goal is to predict the accurate start and end time stamps of different human actions. Owing to its wide application (e.g., surveillance [40, 42], video summarization [27], high-

light detection [13]), action localization has drawn lots of attention in the community. To tackle this problem, many methods try to solve it in a fully-supervised manner [5, 48, 49], but they rely on massive time-consuming annotations. To alleviate this issue, researchers pay more attention to weakly-supervised action localization (WSAL) [47, 16, 17, 43, 50, 22, 44, 11, 23], which explores a more efficient learning strategy with only video-level categorical labels. Nevertheless, recent works [22, 11, 50, 23] mostly rely on the multiple instance learning (MIL) framework [56]: obtaining a video-level prediction via aggregation and optimization under the video-level supervision. While sig-

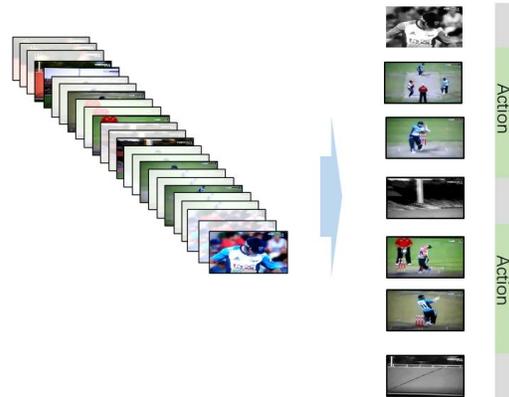


Figure 1. Visualization of the change-point detection component and the co-occurrence component (attention-based classification module) decoupled by our method from a snippet representation. The change-point component helps to detect change-points of temporal variations, which include the change-points of foreground in a video (i.e. the highlighted frames in the left part). Collaborating with the attention module, the points of the foreground are chosen as boundaries of action (i.e. the highlighted frames in the right part).

[†]Work is done during the internship at ACS Lab, Huawei Technologies.

^{*}Corresponding author (guoqinghai@huawei.com).

nificant improvement has been made in prior MIL-based work, there is still a huge performance gap between the weakly-supervised and fully-supervised settings. In consideration of this issue, diverse solutions have been proposed in the literature. For instance, [47, 16, 17] try to erase the most discriminative parts for learning action completeness, [43, 50, 12] learn with pseudo labels generated by manual thresholds and iterative refinement, and [33, 37, 54] formulate the WSAL problem as a video recognition problem and introduce an attention mechanism to construct video-level features, then apply an action classifier to recognize videos.

All the above approaches largely rely on the video-level classification model, which aims at learning the effective classification functions to identify action instances from bags of action instances and non-action frames, but overlook the significance of feature variations. In fact, features usually contain intense semantic information [5, 34], mainly stemming from the temporal and spatial variations in video actions. The variation of features is useful for correcting the wrong action region and adjusting the imprecise boundary of the temporal proposal. Existing solutions often neglect such semantics and thus largely suffer from deviated action boundaries and inaccurate detection.

To better learn the semantics in a given video sample, the model should be able to encode the temporal variations of different time factors. Intuitively, these variations in different timescales disentangle different video fragments, and detection on such variations automatically leads to the detection of change-points, which provide the candidates for action boundaries.

Derived from the above idea, we propose a novel Attention-based Hierarchically-structured Latent Model (AHLM) to model the spatial and temporal features for WSAL task. Specifically, we detect the action boundaries as the change-points of a generative model, where those change-points are determined at the time points with inaccurate generation. Such generative model, is trained by learning the hierarchical representations of the feature semantics in the latent space based on the video inputs. By using an attention-based classification model to select the change-points of the foreground, AHLM localizes the exact action boundaries, see Figure 1 for an illustration.

To our best knowledge, we are the first to consider the temporal variation of feature semantics and study the change-point detection mechanism in WSAL. We design an AHLM that prominently boosts WSAL performance. Our main contributions are summarized as follows:

- To leverage the temporal variations of feature semantics for WSAL, we propose a hierarchically-structured generative latent model that explores spatiotemporal representations and leverages the temporal feature semantics to detect the change-points of videos.
- We build a new framework, AHLM, which firstly proposes the use of an unsupervised change-points detector in a latent space to complement weakly supervised learning, with a novel hierarchical generative model-based change-point detector for complex datasets.
- Based on extensive experiments, we demonstrate that, on two popular action detection datasets, our novel AHLM provides considerable performance gains. On THUMOS14 especially, our method achieves an average mAP of 47.2% when IOU is from 0.1 to 0.7, which is the new state-of-the-art (SotA). On ActivityNet v1.3, our method also achieves the new SotA, with an average mAP of 25.9% when IOU is from 0.5 to 0.95.

2. Related work

2.1. Weakly-supervised Action Localization

Due to the precise annotation of each action instance in fully-supervised, the Weakly-supervised Action Localization(WSAL) is proposed to reduce the expensive annotation costs. During training, the WSAL methods [5, 34, 43, 50] require only video-level categorical labels. These methods can be grouped into two categories, namely top-down and bottom-up methods. In the top-down pipeline, the video-level classification model is learned first, and then frames with high classification activation values are selected as action locations. [32] and [29] forced foreground features from the same class to be similar, otherwise dissimilar. Unlike the top-down scheme, the bottom-up methods directly produce the attention for each frame from data, and train a classification model with the features weighted by attention. Based on this paradigm, [30] further added a regularization term to encourage the sparsity of action. [45] proposed a method to suppress dominance of the most salient action frames and retrieve less salient ones. [37] proposed a model to learn the class-agnostic frame-wise probability conditioned on the frame attention using conditional Variational Auto-Encoder (VAE). Nevertheless, all of the aforementioned methods overlook the significance of temporal variations with respect to features. Unlike these methods, we focus on modeling the temporal variations of the feature semantics, and utilize an unsupervised change-point detection method to localize the action boundaries.

2.2. Unsupervised Action Analysis

Unsupervised learning targets learning effective feature representations from unlabeled data. [35] proposed a temporally-weighted hierarchical clustering algorithm to represent actions in the video. [7] estimated the similarities across smoothed frames through the difference of actions and external discrepancy across actions. These methods, however, mainly focused on the variation of the frame(*e.g.*, the similarity of frame, the temporal variation of frame).

Similarly, [1] utilized spatial-temporal dynamics of events to learn the visual structure of events in the latent space. However, [1] mainly modeled simple datasets(*i.e.*, Breakfast Actions, 50 Salads, and INRIA Instructional Videos), as the proposed method lacks the capability to represent heterogeneous information in an entangled latent space. Unlike the above methods, we utilize the hierarchically-structured VAE and subjective timescaled transition model to learn spatiotemporal semantics on the multi-scaled latent space, which expands the method’s applicability to more complex datasets(*i.e.*, Thumos14 and Activitynet1.3).

2.3. Hierarchical Generative Model

The hierarchical generative model has experienced a fast development in recent years [36, 39, 46], due to the fact that incorporating hierarchy into latent models improves the expressiveness of spatiotemporal representation. [39] implemented a stable fully-convolutional and hierarchical VAE with the use of separate deterministic bottom-up and top-down information channels. [36] was designed to model temporal data using a hierarchy of latent variables that update over fixed time intervals. Nevertheless, these methods were designed to focus on temporal variations while barely considering temporal semantics. Quite recently, VPR [46] proposed a subjective timescale-based hierarchical structure to model the different temporal dynamics, however, they handled only local information of simple data(*e.g.*, bouncing balls). Consecutive videos, when actions change, have more complex global semantic information leading to detecting more precise boundaries. In our work, we propose to simultaneously model the temporal dynamics and the varied distributions of global semantics which enables the hierarchical generative model to the WSAL tasks.

3. Method

Suppose we have a set of training videos and the corresponding video-level labels. Specifically, let us denote for an untrimmed training video, its ground-truth label as $y \in \mathbb{R}^C$, where C is the number of action categories. Note that y could be a multi-hot vector if more than one action is presented in the video and is normalized with the l_1 -normalization. The goal of temporal action localization is to generate a set of action segments $S = \{(s_i, e_i, c_i, q_i)\}_{i=1}^I$ for a testing video with the number of I segments, where s_i, e_i are the start and end time point of the i -th segment and c_i, q_i are the corresponding class prediction and confidence score.

Our method follows the bottom-up pipeline for WSAL, where we detect change-points $P = (p_i)_{i \in T'}$ directly from data. Here p_i is the change-point frame of the video, which is detected based on the transitions of the observable features, and T' is the set of times when a changed-point occurs. Then we leverage an attention model to optimize

change-points of the video to obtain the refined boundaries.

3.1. Framework Overview

In the localization problem, the target is to predict the boundaries of the action instance, which is essentially equivalent to solving the boundary problem of semantic representations for those instances. To this end, our proposed Attention-based Hierarchically-structured Latent Model(AHLM) divides this problem into two different aspects for boundaries localization, the detection of feature semantic change-points(DFC), and the extraction of foreground change-points(EFC). The second term **EFC** is optimized by discriminative capacity for action classification, which is the main optimization target in previous works [21, 33, 12]. In contrast, the first term **DFC** forces the representation of spatial features to be accurately predicted from the temporal changes, which requires the capability of feature disentanglement. In particular, to learn disentangled spatiotemporal representations, we exploit an action-based hierarchical VAE model to encode the input as hierarchical latent spaces, then construct GRU-based transition models which learn to predict feature changes in the latent spaces optimally.

As Figure 2 shows, AHLM is a two-branch network, including a change-point detection module and an attention-based classification branch. Given the concatenated features $X \in \mathbb{R}^{T \times 2D}$ with T snippets for a video, the change-point detection module predicts the frame-level change-points $P = [p_1, p_2, \dots, p_{T'}]$, conveying the class-agnostic change-point boundaries. In the attention-based classification model, input features X are used to predict an attention-based snippet-level class activation map M by a classification head with background class [33, 37], which learns the frame attention by optimizing the video-level recognition task. Note that, $M = \{m_{c,i}\}_{i=1}^{C+1, T} \in \mathbb{R}^{C+1, T}$, where $m_{c,i}$ is the activation of c -th class for i -th snippet (the $C + 1$ -th class means the background category). By suppressing the background change-points P via the attention-based class activation map M , we can get foreground-based change-points as boundaries of action instances.

With predictions generated as the description above, we further explore the hierarchical-structured detection model (*i.e.* Change-point Detection Module) and attention-based classification module, to facilitate the learning of action classification and boundaries localization.

3.2. Change-point Detection Module

The cornerstone of supervised learning is to fully leverage the given annotations, especially for the weakly-supervised learning that has limited information. Previous works mainly develop their framework on MIL-paradigm for video-level learning, ignoring the temporal variation of features. Our framework is constructed under the follow-

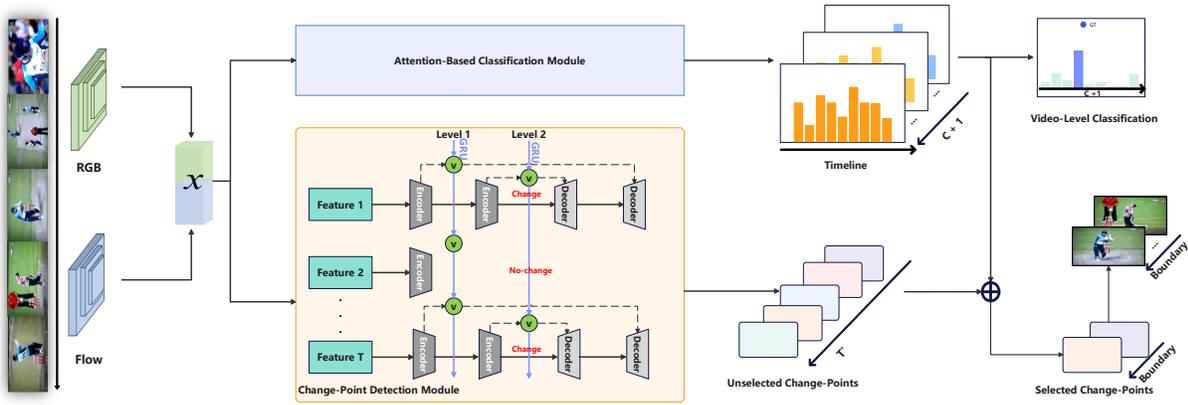


Figure 2. The overall pipeline of the proposed framework. It consists of three parts: Feature Embedding, Change-point Detection Module(DFC), and Attention-based classification Module(EFC). First, the feature embedding stage extracts original snippet features through the I3D network. Subsequently, the DFC is trained to represent spatiotemporal information and change points of feature semantics, supervised by feature distribution and feature reconstruction. Meanwhile, the EFC is trained to distinguish foreground, background and context. For inference, the DFC produces change-points and the EFC selects the change-points of the foreground as boundaries.

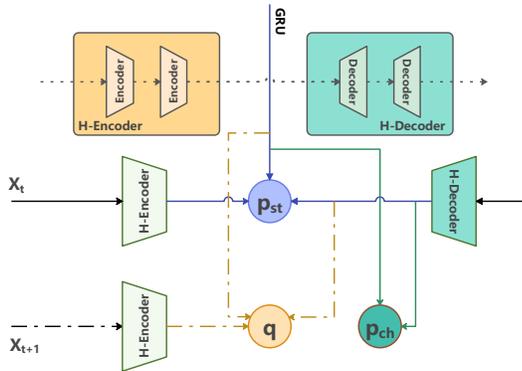


Figure 3. Change-point Detection Module. The black dotted lines indicate the architecture of the Change-point Detection Module. The colored lines indicate the mechanism of the Change-point Detection Module, where dotted lines indicate the next state x_{t+1} , and solid lines indicate the current state x_t . The p_{st} , p_{ch} and q represent the distribution of, current state x_t , predicted result, and next state x_{t+1} in the latent space respectively.

ing cognitive phenomenon: *if one can build a perfect world model, then the change-points will occur when this model does not achieve an ideal prediction along the temporal domain.* Indeed, a perfect prediction should clearly reflect the information at both video-level and feature-level. Such feature-level learning, mainly unsupervised, relies on an effective expression of spatiotemporal information. Recent hierarchical generative models [36, 46] show strong capability regarding this aspect, which make it possible to construct an event boundary detection module based on the pre-

dictions. Inspired by such observations, we develop our change-point detection module based on a hierarchical generative model. We build a 2-level generative model through a variational autoencoder(VAE) structure, combined with a transition model to learn the temporal variations of the video. The first level aims to learn the latent representation for each time point, and the second level further projects the encoded information to higher latent space when observed a change-point based on DFC. Figure 3 indicates the mechanism and architecture of the DFC.

As Figure 3 shows, we use H-Encoder(two 1024-d fully-connected networks) to encode the feature x through $f = f_{H-enc}(x)$, and use H-Decoder(two 1024-d fully-connected networks) to decode the feature u based on the encoder output f , namely, $u = f_{H-dec}(f)$. We use a recurrent GRU model [6](256-d) to learn the temporal transition through $d_{t+1} = f_{tran}(v_t, d_t)$, where t indicates temporal dimension, and v_t is the latent random variable conditioned on the deterministic variables f_t , u_t and d_t . For a given sequence of input observations $\{x_1, x_2, \dots, x_T, x_{T+1}\}$, modeled by the latent variables $v_{1:T}^{1,2}$, where superscripts 1, 2 indicate different levels and subscript 1 : T indicates time sequence, the generative model can be written as the following factorized distribution:

$$p(x_{1:T}, v_{1:T}^1, v_{\tau_1:T_2}^2) = \left[\prod_{t=1}^T p(x_t | v_t^1, v_t^2) \right] \cdot \left[\prod_{t=1}^T p(v_t^1 | v_{1:t-1}^1, v_t^2) \right] \cdot \left[\prod_{t=\tau_1}^{T_2} p(v_t^2 | v_{t'<t}^2) \right], \quad (1)$$

where T_2 denotes the number of change-points detected by level 2, and we define the distribution of the initial state

of v_t^2 as $p(v_t^2) = \mathcal{N}(0, 1)$ a Gaussian prior. Note that $p(v_t^1|v_{1:t-1}^1, v_t^2)$ is a prior distribution of the latent state v_t^1 conditioned on all the past states $v_{1:t-1}^1$ in level 1 and the possible past upper level state v_t^2 , and $p(v_t^2|v_{t'<t}^2)$ is the prior distribution of the latent state v_t^2 conditioned on all the past states $v_{t'<t}^2$, which are all the possible states in the same level 2 before time t .

Follow the general VAE [20], we define the corresponding variational evidence lower bound(ELBO) is as:

$$\mathcal{L}_{ELBO} = \sum_{t=1}^T \mathbb{E}_{q(v_t^{1,2})} [\log p(x_t|v_t^{1,2})] - \sum_{n=1}^2 \sum_{t=1}^T \mathbb{E}_{q(v_t^{>n}, v_{<t}^n)} [D_{KL}(q_\phi(v_t^n|x_t, v_t^{>n}, v_{<t}^n) || p_\theta(v_t^n|v_t^{>n}, v_{<t}^n))], \quad (2)$$

where p_θ is the prior model with θ representing the parameters defined by x, u and d , q_ϕ is the posterior model and $v_t^{>n}$ represents the latent states at time t in the level higher than n .

The first term in Equ. (2) represents the likelihood of the reconstructed x_t given the latent variables $v_t^{1,2}$, which measures the reconstruction loss. The second term is the KL-divergence of the prior distribution $p(v)$ and the posterior distribution $q(v)$. The loss function of DFC defines by

$$\mathcal{L}_{DFC} = -\mathcal{L}_{ELBO}. \quad (3)$$

Our model is trained in a way that the second level updates the latent state v_t^2 in a subjective time scale, while the determined time points correspond to the changes in the observable features over time.

The key component of the DFC relies on a Bayesian inference mechanism under the static assumption on the level 2, and the changes are detected when the updated posterior violates such an assumption. Specifically, as Figure 3 shows, given the current feature inputs from the H-Encoder f_t and H-Decoder u_t , we construct the static assumption and change assumption respectively (note we omit the level index for better readability in the following since we only consider the detection in level 2). Under the static assumption, the prior is calculated as

$$p_{st} = p_\theta(v_{t+1}|f_t, d_t, u_t), \quad (4)$$

while under the change assumption, the prior is calculated as

$$p_{ch} = p_\theta(v_{t+1}|f_t, d_{t+1}, u_t), \quad (5)$$

where we trigger the transition model to predict a next temporal state d_{t+1} to produce a new prior.

The above can be seen as the model's belief over the observable features f_t at the latest time step t under static and change assumptions, respectively. Given a new input f_{t+1} , the updated posterior is computed by

$$q = q_\phi(v_{t+1}|f_{t+1}, d_{t+1}, u_t). \quad (6)$$

We then use the KL-divergence, $D_{st} = D_{KL}(q||p_{st})$ and $D_{ch} = D_{KL}(q||p_{ch})$, to measure how much the features have changed compared to the last time step under different assumptions. A change-point boundary is considered to be detected when the static assumption based update is less accurate than the change assumption based update. Particularly, we define such boundary condition as

$$D_{KL}(q||p_{st}) > \beta D_{KL}(q||p_{ch}), \quad (7)$$

where β is an empirical hyperparameter, $\beta \in [0.15, 0.9]$. Satisfying this criterion indicates that the model's prediction produced a belief state more consistent with the change assumption, suggesting that it contains a change-point in the features. In other words, Equ. (7) compares the difference between the predicted result and observation.

The GRUs, used in the transition model, nevertheless, have been observed to suffer from the state saturation problems in very long sequences [3, 10], and this issue is further aggregated for highly heterogeneous spatial-temporal features as in videos. To counter this problem, we propose to utilize network resetting in our GRU model for d_t . That is, after we detect a change-point frame x_t , the network is reset to take the next observation x_{t+1} as an initial input, together with the initialization of network parameters. Furthermore, we leverage a dynamic β in the boundary condition in Equ. (7) to preserve stable detection by the following rule:

$$\beta(t+1) = \begin{cases} \beta(t) + \alpha & (\text{change-point}); \\ \beta(t) - \alpha & (\text{no-change-point}), \end{cases} \quad (8)$$

where α is a hyperparameter that is set as $\alpha = 0.15$, and (no-)change-point means that we detect a (no-)change-point at the frame x_t .

3.3. Attention-based Classification Module

The attention-based classification module learns the attention of features for optimizing change-points by distinguishing foreground and background.

For the classification module, we follow previous work [33, 17, 37], applying the cross-entropy loss function between the predicted video-level label and the ground truth label to classify different action classes in a video.

$$\mathcal{L}_{clf} = \sum_{c=1}^{C+1} -y_c(x) \log(p_c(x)), \quad (9)$$

where $y_c(x)$ is the ground truth video action probability distribution and $p_c(x)$ is the predicted video-level action probability distribution.

For the attention module, we utilize three-branch class activation sequences to represent the foreground, context,

and background individually. In specific, we set the video-level instance label

$$y = (y(1) = 0, \dots, y(n) = 1, \dots, y(C) = 0, y(C + 1) = 0), \quad (10)$$

which represents the ground truth of the video in the n -th category and $C + 1$ is the background label index. We then optimize the following attention-based *foreground* model:

$$\mathcal{L}_{fg} = \sum_{c=1}^{C+1} -y(x) \log(p_c(x)). \quad (11)$$

With the foreground attention weighting, background and action context snippets have been suppressed, as shown in temporal class activation sequences (CAS) [55, 33].

Similarly, we can set the video-level background label $y = (y(n) = 0, y(C + 1) = 1)$ and context label $y = (y(n) = 1, y(C + 1) = 1)$ to optimize the attention-based *background* model \mathcal{L}_{bg} and *context* model \mathcal{L}_{ct} . Following analogical arguments as in the foreground attention branch, we implement contexts and background branches by using CAS.

After obtaining three attention-based classification loss \mathcal{L}_{fg} , \mathcal{L}_{bg} , and \mathcal{L}_{ct} , we compose the overall loss \mathcal{L}_{EFC} for extraction of foreground change-points as:

$$\mathcal{L}_{EFC} = \mathcal{L}_{fg} + \mathcal{L}_{bg} + \mathcal{L}_{ct}. \quad (12)$$

For inference, based on the attention-based classification module, we choose the change-points of the foreground as action boundaries. In addition, we add the longest common sub-sequence(LCS) contrasting to optimize the boundaries from the DFC module. In specific, for two adjacent change-points A and B , we construct two snippets $l_{AC} = \{l_1, l_2, \dots\}$ and $l_{BC} = \{m_1, m_2, \dots\}$ by connecting A and B with a third change-point C . We then calculate the cosine similarity of l_i and m_j ($i, j = 1, 2, \dots$) and compare the results with a given threshold 0.65 to get the LCS of l_{AC} and l_{BC} . Based on the length of LCS, we delete the redundant points from adjacent change-points.

4. Experiment

In this section, we first describe datasets and evaluation metrics. Then, we evaluate our model’s effectiveness followed by the main results and ablation study.

4.1. Datasets and Evaluation Metrics

To validate the effectiveness of our model, we conduct extensive experiments on commonly-used benchmark THUMOS14 [18] and ActivityNet v1.3 [8].

THUMOS14 [18] It contains 101 categories of videos and is composed of four parts: training, validation, testing

and a background set. Each set includes 13320, 1010, 1574 and 2500 videos, respectively. Following the common setting in [18], we used 200 videos in the validation set for training, and 213 videos in the testing set for evaluation.

ActivityNet [8] We evaluate our method on the ActivityNet release 1.3, which contains samples from 200 categories of activities and 19994 videos in total. It includes untrimmed video classification and activity detection tasks. It is divided into training, validation, and test sets with a ratio of 2 : 1 : 1, containing 10024, 4926 and 5044 videos respectively. Following [12, 33], we use the training set to train our model and the validation set for evaluation.

Evaluation Metrics We follow the standard evaluation protocol and report mean Average Precision (mAP) at the different intersections over union (IoU) thresholds. The results are calculated using the benchmark code provided by the ActivityNet official codebase¹.

4.2. Implementation Details

For the feature extraction, we first sample RGB frames at 25 fps for each video and apply the TV-L1 algorithm [48] to generate optical flow frames. Then, we divide each video into non-overlapping snippets with consecutive 16 frames. Thereafter, we perform the I3D networks [2] pre-trained on the Kinetics dataset [2] to obtain the video features. Note that, for a fair comparison, we do not introduce any other feature fine-tuning operations to the pre-trained I3D model.

For the training process on the THUMOS-14, we set the batch size as 16, and apply the Adam optimizer [19] with learning rate 10^{-4} and weight decay 5×10^{-4} . We set the video snippets length $T = 750$. For the training process on the ActivityNet-v1.3 dataset, we set the training video batch size to 64, applying the same Adam optimizer with THUMOS-14, and set the video snippets length $T = 150$.

4.3. Main Results

In Table 1, we compare AHLM with current methods on Thumos14. Selected fully-supervised methods are presented for reference. We observe that AHLM outperforms all the previous WSAL methods and establishes new state of the art on THUMOS-14 with 47.2% average mAP for IoU thresholds 0.1 : 0.7. In particular, our approach outperforms ACM [33], which also utilizes an attention model to distinguish foreground and background but without explicit temporal variations of feature semantics. Even compared with the fully supervised methods, AHLM outperforms BU-TAL and GTAD and achieves comparable results with AFSD and TRA when the IoU threshold is low. The results demonstrate the superior performance of our approach with temporal variations of feature semantics modeling.

We also conduct experiments on ActivityNet-v1.3 and the comparison results are summarized in Table 2. Again,

¹<https://github.com/activitynet/ActivityNet/tree/master/Evaluation>

Type	Model	Publication	THUMOS14							
			0.1	0.2	0.3	0.4	0.5	0.6	0.7	avg.
Fully-supervised	BU-TAL [53]	<i>ECCV20</i>	-	-	53.9	50.7	45.4	38.5	28.0	-
	G-TAD [41]	<i>CVPR20</i>	-	-	54.5	47.6	40.2	30.8	23.4	-
	GCM [49]	<i>TPAMI21</i>	72.5	70.9	66.5	60.8	51.9	-	-	-
	AFSD [24]	<i>CVPR21</i>	72.2	70.8	67.1	62.2	55.5	43.7	31.1	57.6
	TadTR [26]	<i>TIP22</i>	-	-	74.8	69.1	60.1	46.6	32.8	-
	RefactorNet [26]	<i>CVPR22</i>	-	-	70.7	65.4	58.6	47.0	32.1	-
	TRA [54]	<i>TIP22</i>	73.7	72.6	70.0	64.3	57.4	46.2	31.1	59.3
Weakly-supervised	CMCS [25]	<i>CVPR19</i>	57.4	50.8	41.2	32.1	23.1	15.0	7.0	32.4
	WSAL-BM [31]	<i>ICCV19</i>	60.4	56.0	46.6	37.5	26.8	19.6	9.0	36.6
	DGAM [37]	<i>CVPR20</i>	60.0	54.2	46.8	38.2	28.8	19.8	11.4	37.0
	A2CL-PT [28]	<i>ECCV20</i>	61.2	56.1	48.1	39.0	30.1	19.2	10.6	37.8
	HAM-Net [16]	<i>AAAI21</i>	65.9	59.6	52.2	43.1	32.6	21.9	12.5	41.1
	FAC-Net [14]	<i>ICCV21</i>	67.6	62.1	52.6	44.3	33.4	22.5	12.7	42.2
	CoLA [51]	<i>CVPR21</i>	66.2	59.5	51.5	41.9	32.2	22.0	13.1	40.9
	ACM-Net [33]	<i>TIP21</i>	68.9	62.7	55.0	44.6	34.6	21.8	10.8	42.6
	FTCL [9]	<i>CVPR22</i>	69.6	63.4	55.2	45.2	35.6	23.7	12.2	43.6
	ASM-Loc [12]	<i>CVPR22</i>	71.2	65.5	57.1	46.8	36.6	25.2	13.4	45.1
	DCC [23]	<i>CVPR22</i>	69.0	63.8	55.9	45.9	35.7	24.3	13.7	44.0
	RSKP [15]	<i>CVPR22</i>	71.3	65.3	55.8	47.5	38.2	25.4	12.5	45.1
	DELU [4]	<i>ECCV22</i>	71.5	66.2	56.5	47.7	40.5	27.2	15.3	46.4
	StochasticFormer [38]	<i>TIP23</i>	66.5	61.1	52.5	43.9	33.5	22.6	13.2	41.9
	ASCN [38]	<i>TMM23</i>	71.4	65.6	57.0	48.2	39.8	26.8	14.4	46.1
	Ours	-	75.1	68.9	60.2	48.9	38.3	26.8	14.7	47.2

Table 1. Performance comparison with SotA methods on THUMOS14, measured by mAP at different IoU thresholds.

our AHLM obtains a new state-of-the-art performance of 25.9% average mAP, surpassing the latest works (e.g. FTCL, ASM-Loc). The consistent superior results on both datasets justify the effectiveness of our AHLM. Especially, unlike other methods which only achieve marginal improvement on ActivityNet-v1.3, AHLM maintains the significant improvement as on THUMOS-14. This shows that the performance of AHLM does not rely on the length of video snippets, since the average length of videos in the THUMOS14 dataset is much longer than those in ActivityNet-v1.3. An important reason is that AHLM learns the temporal variations of the feature semantics in the latent space by GRU with dynamic resetting, which is appropriate to detect temporal variations of semantics in different timescales.

4.4. Ablation Study

To demonstrate the reasonableness of our AHLM, we analyze the effect of every submodule and some function operations in this subsection.

Model	Publication	ActivityNet v1.3			
		0.5	0.75	0.95	avg.
STPN [30]	<i>CVPR 2018</i>	29.3	16.9	2.6	16.3
ASSG [52]	<i>MM 2019</i>	32.3	20.1	4.0	18.8
Bas-Net [21]	<i>AAAI 2020</i>	34.5	22.5	4.9	22.2
TS-PCA [43]	<i>CVPR 2021</i>	37.4	23.5	5.9	23.7
FAC-Net [14]	<i>ICCV 2021</i>	37.6	24.2	6.0	24.0
ACM-Net [33]	<i>TIP 2021</i>	37.6	24.7	6.5	24.4
FTCL [9]	<i>CVPR 2022</i>	40.0	24.3	6.4	24.8
ASM-Loc [12]	<i>CVPR 2022</i>	41.0	24.9	6.2	25.1
Ours	-	42.3	24.8	6.9	25.9

Table 2. Performance comparison with state-of-the-art methods on ActivityNet-v1.3 dataset.

4.4.1 Contribution of each design in AHLM

We study the influence of each component in AHLM on overall performance. We start with the basic model that directly optimizes the attention based on foreground loss \mathcal{L}_{fg} . Then we add the background loss \mathcal{L}_{bg} and the context loss \mathcal{L}_{ct} gradually. These three types of loss constitute \mathcal{L}_{EFC} .

The variational generative model loss \mathcal{L}_{ELBO} , which indicates \mathcal{L}_{DFC} , is further included. Note that adding \mathcal{L}_{ELBO} indicates involving the hierarchically-structured modeling, where reconstruction loss of VAE is also simultaneously optimized.

Table 3 summarizes the performances by considering one more factor at each stage on THUMOS14. We first observe that adding the background loss \mathcal{L}_{bg} and the context loss \mathcal{L}_{ct} largely enhances the performance of the foreground-based model. The two losses encourage the sparsity in the foreground attention weights by pushing the background attention weights to be 1 at background snippets, and therefore improve the foreground-background separation. Based on the EFC, our change-point detection module further contributes a significant increase of 2.4% and the performance of AHLM finally reaches 47.0%. Further, a more explicit ablation study by adding DFC on each part of the EFC proves our method’s effectiveness, in particular, our DFC module contributes an increase of 2.1% and 2.5% respectively based on foreground and foreground-background.

\mathcal{L}_{EFC}			\mathcal{L}_{DFC}	THUMOS14				
\mathcal{L}_{fg}	\mathcal{L}_{bg}	\mathcal{L}_{ct}	\mathcal{L}_{ELBO}	0.1	0.3	0.5	0.7	avg.
✓	-	-	-	49.9	32.9	16.6	5.3	26.2
✓	-	-	✓	54.4	36.4	16.8	5.4	28.3
✓	✓	-	-	55.9	41.9	23.0	7.1	32.0
✓	✓	-	✓	57.3	44.2	26.6	10.3	34.5
✓	✓	✓	-	71.2	57.1	36.6	13.4	44.6
✓	✓	✓	✓	75.1	60.2	38.3	14.7	47.0

Table 3. *mAP* at different overlap IoU thresholds performance comparison of each design on THUMOS14.

4.4.2 Effectiveness of Change-point Module

It is obvious that the proposed change-point detection strategy can play a complementary role over existing methods in localizing boundaries of action instances. To see this, we conduct the experiment by directly adding the detected points by our change-point detection module into a MIL-based method [32], which indicates the classification method with CAS [55]. Specifically, for the proposed snippets based on the change-points, we calculate the score using the MIL-based method.

Table 4 shows, compared to the original MIL-based method, our change-point detection module contributes a significant increase of 3.7% and the performance finally reaches 28.5% on average *mAP*.

Following the common setting of WSAL task, as Table 1 and Table 2 show, we chose THUMOS14 and ActivityNet dataset as our benchmark and achieves SOTA performance. Regarding the scalability and generalization, essentially, the effectiveness of our change-point module is related to

	THUMOS14				
	0.1	0.3	0.5	0.7	avg.
MIL-based [32]	46.5	31.2	16.9	4.4	24.8
MIL-based + Ours	54.2	37.1	17.3	5.3	28.5

Table 4. Comparison with MIL-based methods on THUMOS14.

the representation ability of the hierarchical-VAE, which is proved in the literature (*e.g.*, CW-VAE [36], NVAE [39], VPR [46]), hence guarantees the scalability of our change-point module.



Figure 4. Qualitative result visualization on the THUMOS-14 dataset. From the above qualitative results, we can conclude that our proposed AHLM mechanism is greatly beneficial to suppress locate action instances and help us achieve more precise temporal action localization results.

4.4.3 Qualitative Results

Fig 4 shows the visualization comparisons between the attention-based model [33], MIL-based model [32] and our AHLM. From the figure, it can be found that the common errors of current methods are mainly about missed detection of short actions and imprecise localization of the action. Through learning temporal variation of feature semantics, AHLM locates the boundaries of short actions (*e.g.*, examples 1) and more precise boundaries (*e.g.*, examples 2, 3). We see our method depends on an attention-based classification model to filter the change-points from fore- and background. This verifies the importance of improving the quality of separating foreground and background and should be further studied in future work.

5. Conclusion

In this paper, we propose a novel hierarchically-structured attention mechanism to model temporal variation of feature semantics by disentangling the spatial and temporal information on the latent space. Our weakly supervised action localization framework AHLM mainly consists

of feature embedding, change-point detection module and attention-based classification module. We leverage temporal variation of the features to locate the change-points and optimize by attention-based classification model for the WSAL task. Our method outperforms the prior work with a remarkable margin on two popular datasets, achieving the SotA results on both. The results demonstrate that our exploration on temporal variation of feature semantic information effectively improves WSAL performance, which narrows the performance gap between the weakly-supervised and fully-supervised settings. For the future work, we believe the hierarchically-structured latent modeling will be a promising direction for various weakly supervised and unsupervised learning tasks. It is also worth to explore such mechanism in other related tasks.

Acknowledgement. we want to thank Alexey Zakharov and Zafeirios Fountas from Huawei UK, for many helpful discussions and inspirations.

References

- [1] Sathyanarayanan N Aakur and Sudeep Sarkar. A perceptual prediction framework for self supervised event segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1197–1206, 2019.
- [2] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6299–6308, 2017.
- [3] Shuo-Yiin Chang, Bo Li, Gabor Simko, Tara N Sainath, Anshuman Tripathi, Aaron van den Oord, and Oriol Vinyals. Temporal modeling using dilated convolution and gating for voice-activity-detection. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 5549–5553, 2018.
- [4] Mengyuan Chen, Junyu Gao, Shicai Yang, and Changsheng Xu. Dual-evidential learning for weakly-supervised temporal action localization. In *European Conference on Computer Vision*, pages 192–208, 2022.
- [5] Peihao Chen, Chuang Gan, Guangyao Shen, Wenbing Huang, Runhao Zeng, and Mingkui Tan. Relation attention for temporal action localization. *IEEE Transactions on Multimedia*, 22:2723–2733, 2019.
- [6] Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2014.
- [7] Zexing Du, Xue Wang, Guoqing Zhou, and Qing Wang. Fast and unsupervised action boundary detection for action segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3323–3332, 2022.
- [8] Bernard Ghanem Fabian Caba Heilbron, Victor Escorcia and Juan Carlos Nibbles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 961–970, 2015.
- [9] Junyu Gao, Mengyuan Chen, and Changsheng Xu. Fine-grained temporal contrastive learning for weakly-supervised temporal action localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19999–20009, 2022.
- [10] Albert Gu, Caglar Gulcehre, Thomas Paine, Matt Hoffman, and Razvan Pascanu. Improving the gating mechanism of recurrent neural networks. In *Proceedings of the International Conference on Machine Learning*, pages 3800–3809, 2020.
- [11] Bo He, Xitong Yang, Le Kang, Zhiyu Cheng, Xin Zhou, and Abhinav Shrivastava. Asm-loc: Action-aware segment modeling for weakly-supervised temporal action localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13925–13935, 2022.
- [12] Bo He, Xitong Yang, Le Kang, Zhiyu Cheng, Xin Zhou, and Abhinav Shrivastava. Asm-loc: Action-aware segment modeling for weakly-supervised temporal action localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13925–13935, 2022.
- [13] Fa-Ting Hong, Xuanteng Huang, Wei-Hong Li, and Wei-Shi Zheng. Mini-net: Multiple instance ranking network for video highlight detection. In *Proceedings of the European Conference on Computer Vision*, pages 345–360, 2020.
- [14] Linjiang Huang, Liang Wang, and Hongsheng Li. Foreground-action consistency network for weakly supervised temporal action localization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8002–8011, 2021.
- [15] Linjiang Huang, Liang Wang, and Hongsheng Li. Weakly supervised temporal action localization via representative snippet knowledge propagation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3272–3281, 2022.
- [16] Ashraf Islam, Chengjiang Long, and Richard Radke. A hybrid attention mechanism for weakly-supervised temporal action localization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1637–1645, 2021.
- [17] Yuan Ji, Xu Jia, Huchuan Lu, and Xiang Ruan. Weakly-supervised temporal action localization via cross-stream collaborative learning. In *Proceedings of the ACM International Conference on Multimedia*, pages 853–861, 2021.
- [18] 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://cvc.ucf.edu/THUMOS14/>, 2014.
- [19] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International Conference on Learning Representations*, 2015.
- [20] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations*, 2014.
- [21] Pilhyeon Lee, Youngjung Uh, and Hyeran Byun. Background suppression network for weakly-supervised temporal action localization. In *Proceedings of the AAAI conference on artificial intelligence*, pages 11320–11327, 2020.

- [22] Pilhyeon Lee, Jinglu Wang, Yan Lu, and Hyeran Byun. Weakly-supervised temporal action localization by uncertainty modeling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1854–1862, 2021.
- [23] Jingjing Li, Tianyu Yang, Wei Ji, Jue Wang, and Li Cheng. Exploring denoised cross-video contrast for weakly-supervised temporal action localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19914–19924, 2022.
- [24] Chuming Lin, Chengming Xu, Donghao Luo, Yabiao Wang, Ying Tai, Chengjie Wang, Jilin Li, Feiyue Huang, and Yanwei Fu. Learning salient boundary feature for anchor-free temporal action localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3320–3329, 2021.
- [25] Daochang Liu, Tingting Jiang, and Yizhou Wang. Completeness modeling and context separation for weakly supervised temporal action localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1298–1307, 2019.
- [26] Xiaolong Liu, Qimeng Wang, Yao Hu, Xu Tang, Shiwei Zhang, Song Bai, and Xiang Bai. End-to-end temporal action detection with transformer. *IEEE Transactions on Image Processing*, 31:5427–5441, 2022.
- [27] Yu-Fei Ma, Xian-Sheng Hua, Lie Lu, and Hong-Jiang Zhang. A generic framework of user attention model and its application in video summarization. *IEEE Transactions on Multimedia*, 7:907–919, 2005.
- [28] Kyle Min and Jason J Corso. Adversarial background-aware loss for weakly-supervised temporal activity localization. In *Proceedings of the European conference on computer vision*, pages 283–299, 2020.
- [29] Sanath Narayan, Hisham Cholakkal, Fahad Shahbaz Khan, and Ling Shao. 3c-net: Category count and center loss for weakly-supervised action localization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8679–8687, 2019.
- [30] Phuc Nguyen, Ting Liu, Gautam Prasad, and Bohyung Han. Weakly supervised action localization by sparse temporal pooling network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6752–6761, 2018.
- [31] Phuc Xuan Nguyen, Deva Ramanan, and Charles C Fowlkes. Weakly-supervised action localization with background modeling. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5502–5511, 2019.
- [32] Sujoy Paul, Sourya Roy, and Amit K Roy-Chowdhury. W-talc: Weakly-supervised temporal activity localization and classification. In *Proceedings of the European Conference on Computer Vision*, pages 563–579, 2018.
- [33] Sanqing Qu, Guang Chen, Zhijun Li, Lijun Zhang, Fan Lu, and Alois Knoll. Acn-net: Action context modeling network for weakly-supervised temporal action localization. *arXiv preprint arXiv:2104.02967*, 2021.
- [34] Niamul Quader, Md Mafijul Islam Bhuiyan, Juwei Lu, Peng Dai, and Wei Li. Weight excitation: Built-in attention mechanisms in convolutional neural networks. In *Proceedings of the European Conference on Computer Vision*, pages 87–103, 2020.
- [35] Saquib Sarfraz, Naila Murray, Vivek Sharma, Ali Diba, Luc Van Gool, and Rainer Stiefelhagen. Temporally-weighted hierarchical clustering for unsupervised action segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11225–11234, 2021.
- [36] Vaibhav Saxena, Jimmy Ba, and Danijar Hafner. Clockwork variational autoencoders. *Advances in Neural Information Processing Systems*, 34:29246–29257, 2021.
- [37] Baifeng Shi, Qi Dai, Yadong Mu, and Jingdong Wang. Weakly-supervised action localization by generative attention modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1009–1019, 2020.
- [38] Haichao Shi, Xiao-Yu Zhang, and Changsheng Li. Stochasticformer: Stochastic modeling for weakly supervised temporal action localization. *IEEE Transactions on Image Processing*, 32:1379–1389, 2023.
- [39] Arash Vahdat and Jan Kautz. Nvae: A deep hierarchical variational autoencoder. *Advances in Neural Information Processing Systems*, 33:19667–19679, 2020.
- [40] Sarvesh Vishwakarma and Anupam Agrawal. A survey on activity recognition and behavior understanding in video surveillance. *The Visual Computer*, 29:983–1009, 2013.
- [41] Mengmeng Xu, Chen Zhao, David S Rojas, Ali Thabet, and Bernard Ghanem. G-tad: Sub-graph localization for temporal action detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10156–10165, 2020.
- [42] Shusen Yang, Liwen Zhang, Chen Xu, Hanqiao Yu, Jianqing Fan, and Zongben Xu. Massive data clustering by multi-scale psychological observations. *National Science Review*, 9(2):nwab183, 2022.
- [43] Wenfei Yang, Tianzhu Zhang, Xiaoyuan Yu, Tian Qi, Yongdong Zhang, and Feng Wu. Uncertainty guided collaborative training for weakly supervised temporal action detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 53–63, 2021.
- [44] Wenfei Yang, Tianzhu Zhang, Yongdong Zhang, and Feng Wu. Uncertainty guided collaborative training for weakly supervised and unsupervised temporal action localization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [45] Yuan Yuan, Yueming Lyu, Xi Shen, Ivor W Tsang, and Dit-Yan Yeung. Marginalized average attentional network for weakly-supervised learning. In *International Conference on Learning Representations*, 2019.
- [46] Alexey Zakharov, Qinghai Guo, and Zafeirios Fountas. Variational predictive routing with nested subjective timescales. In *International Conference on Learning Representations*, 2022.
- [47] Runhao Zeng, Chuang Gan, Peihao Chen, Wenbing Huang, Qingyao Wu, and Mingkui Tan. Breaking winner-takes-all: Iterative-winners-out networks for weakly supervised temporal action localization. *IEEE Transactions on Image Processing*, 28:5797–5808, 2019.

- [48] 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/CVF International Conference on Computer Vision*, pages 7094–7103, 2019.
- [49] Runhao Zeng, Wenbing Huang, Mingkui Tan, Yu Rong, Peilin Zhao, Junzhou Huang, and Chuang Gan. Graph convolutional module for temporal action localization in videos. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [50] Yuanhao Zhai, Le Wang, Wei Tang, Qilin Zhang, Junsong Yuan, and Gang Hua. Two-stream consensus network for weakly-supervised temporal action localization. In *Proceedings of the European Conference on Computer Vision*, pages 37–54, 2020.
- [51] Can Zhang, Meng Cao, Dongming Yang, Jie Chen, and Yuexian Zou. Cola: Weakly-supervised temporal action localization with snippet contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16010–16019, 2021.
- [52] Chengwei Zhang, Yunlu Xu, Zhanzhan Cheng, Yi Niu, Shiliang Pu, Fei Wu, and Futai Zou. Adversarial seeded sequence growing for weakly-supervised temporal action localization. In *Proceedings of the ACM international conference on multimedia*, pages 738–746, 2019.
- [53] Peisen Zhao, Lingxi Xie, Chen Ju, Ya Zhang, Yanfeng Wang, and Qi Tian. Bottom-up temporal action localization with mutual regularization. In *Proceedings of the European Conference on Computer Vision*, pages 539–555, 2020.
- [54] Yibo Zhao, Hua Zhang, Zan Gao, Weili Guan, Jie Nie, Anan Liu, Meng Wang, and Shengyong Chen. A temporal-aware relation and attention network for temporal action localization. *IEEE Transactions on Image Processing*, 31:4746–4760, 2022.
- [55] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2921–2929, 2016.
- [56] Zhi-Hua Zhou. Multi-instance learning: A survey. *Department of Computer Science & Technology, Nanjing University, Tech. Rep*, 1, 2004.