# Distilling Vision-Language Pre-training to Collaborate with Weakly-Supervised Temporal Action Localization

Chen Ju[1*], Kunhao Zheng[1*], Jinxiang Liu[1], Peisen Zhao[2], Ya Zhang[1✉],
Jianlong Chang[2], Yanfeng Wang[1], Qi Tian[2]

[1]Cooperative Medianet Innovation Center, Shanghai Jiao Tong University  [2]Huawei Cloud & AI

{ju_chen, dyekuu, jinxliu, ya_zhang, wangyanfeng}@sjtu.edu.cn, {pszhao93, jianlong.chang, tian.qi1}@huawei.com

## Abstract

*Weakly-supervised temporal action localization (WTAL) learns to detect and classify action instances with only category labels. Most methods widely adopt the off-the-shelf Classification-Based Pre-training (CBP) to generate video features for action localization. However, the different optimization objectives between classification and localization, make temporally localized results suffer from the serious incomplete issue. To tackle this issue without additional annotations, this paper considers to distill free action knowledge from Vision-Language Pre-training (VLP), since we surprisingly observe that the localization results of vanilla VLP have an over-complete issue, which is just complementary to the CBP results. To fuse such complementarity, we propose a novel distillation-collaboration framework with two branches acting as CBP and VLP respectively. The framework is optimized through a dual-branch alternate training strategy. Specifically, during the B step, we distill the confident background pseudo-labels from the CBP branch; while during the F step, the confident foreground pseudo-labels are distilled from the VLP branch. And as a result, the dual-branch complementarity is effectively fused to promote a strong alliance. Extensive experiments and ablation studies on THUMOS14 and ActivityNet1.2 reveal that our method significantly outperforms state-of-the-art methods.*

## 1. Introduction

Temporal action localization (TAL), which aims to localize and classify action instances from untrimmed long videos, has been recognized as an indispensable component of video understanding [11, 66, 85]. To avoid laborious temporal boundary annotations, the weakly-supervised setting (WTAL) [27, 56, 58, 72], *i.e.* only video-level category labels are available, has gained increasing attentions.

To date in the literature, almost all WTAL methods rely on Classification-Based Pre-training (CBP) for video fea-
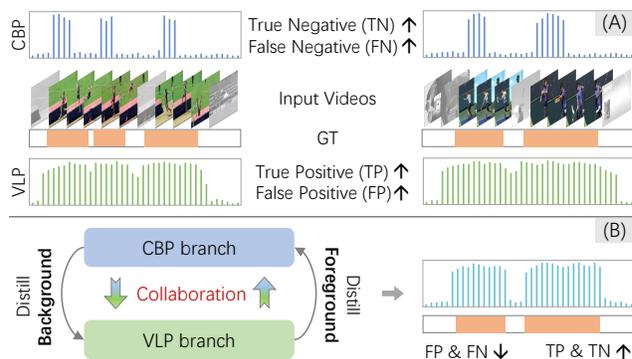


Figure 1. **(A) Complementarity.** Most works use Classification-Based Pre-training (CBP) for localization, causing high TN yet serious FN. Vanilla Vision-Language Pre-training (VLP) confuses action and background, leading to high TP yet serious FP. **(B) Our distillation-collaboration framework** distills foreground from the VLP branch while background from the CBP branch, and promotes mutual collaboration, bringing satisfactory results.

ture extraction [3, 67]. A popular pipeline is to train an action classifier with CBP features, then threshold the frame-level classification probabilities for final localization results. As demonstrated in Figure 1 (A), these CBP methods suffer from the serious **incompleteness issue**, *i.e.* only detecting sparse discriminative action frames and incurring *high false negatives*. The main reason is that the optimization objective of classification pre-training, is to find several discriminative frames for action recognition, which is far from the objective of complete localization. As a result, features from CBP are inevitably biased towards partial discriminative frames. To solve the incompleteness issue, many efforts [29, 39, 50, 89] have been attempted, but most of them are trapped in a 'performance-cost dilemma', namely, solely digging from barren category labels to keep costs low. Lacking location labels fundamentally limits the performance, leaving a huge gap from strong supervision.

To jump out of the dilemma, this paper raises one novel

question: *is there free action knowledge available, to help complete detection results while maintain cheap annotation overhead at the same time?* We naturally turn our sights to the prevalent Vision-Language Pre-training (VLP) [20, 60]. VLP has demonstrated great success to learn joint visual-textual representations from large-scale web data. As language covers rich information about objects, human-object interactions, and object-object relationships, these learned representations could provide powerful human-object co-occurrence priors: valuable gifts for action localization.

We here take one step towards positively answering the question, *i.e.* fill the research gap of distilling action priors from VLP, namely CLIP [60], to solve the incomplete issue for WTAL. As illustrated in Figure 1 (A), we first naively evaluate the temporal localization performance of VLP by frame-wise classification. But the results are far from satisfactory, suffering from the serious **over-complete issue**, *i.e.* confusing multiple action instances into a whole, causing *high false positives*. We conjecture the main reasons are: (1) due to data and computational burden, almost all VLPs are trained using image-text pairs. Hence, VLP lacks sufficient temporal knowledge and relies more on human-object co-occurrence for localization, making it struggle to distinguish the actions with visually similar background contexts; (2) some background contexts have similar (confusing) textual semantics to actions, such as run-up *vs.* running.

Although simply steering VLP for WTAL is infeasible, we fortunately observe the **complementary property** between CBP and VLP paradigms: the former localizes high true negatives but serious false negatives, while the latter has high true positives but serious false positives. To leverage the complementarity, as shown in Figure 1 (B), we design a novel distillation-collaboration framework that uses two branches to play the roles of CBP and VLP, respectively. The design rationale is to distill background knowledge from the CBP branch, while foreground knowledge from the VLP branch, for strong alliances. Specifically, we first warm up the CBP branch using only category supervision to initialize confident background frames, and then optimize the framework via an alternating strategy. *During B step*, we distill background pseudo-labels for the VLP branch to solve the over-complete issue, hence obtaining high-quality foreground pseudo-labels. *During F step*, we leverage high-quality pseudo-labels for the CBP branch to tackle the incomplete issue. Besides, in each step, we introduce both confident knowledge distillation and representation contrastive learning for pseudo-label denoising, effectively fusing complementarity for better results.

On two standard benchmarks: THUMOS14 and ActivityNet1.2, our method improves the average performance by 3.5% and 2.7% over state-of-the-art methods. We also conduct extensive ablation studies to reveal the effectiveness of each component, both quantitatively and qualitatively.

To sum up, our contributions lie in three folds:

• We pioneer the first exploration in distilling free action knowledge from off-the-shelf VLP to facilitate WTAL;

• We design a novel distillation-collaboration framework that encourages the CBP branch and VLP branch to complement each other, by an alternating optimization strategy;

• We conduct extensive experiments and ablation studies to reveal the significance of distilling VLP and our superior performance on two public benchmarks.

## 2. Related Work

**Vision-Language Pre-training (VLP)** aims to learn cross-modal representations [8, 51, 76] from large-scale web data. Comparing to video, image requires fewer costs for annotation and computation, hence almost all VLPs are image-based, *e.g.* [1, 20, 60, 74, 84, 86]. Recently, several studies adopted VLP to provide free visual-semantic knowledge for downstream image tasks, such as detection [14, 87], segmentation [61, 95], human-object interaction [25, 32], synthesis [42], and generation [7, 70, 71]. In terms of the video domain, [38, 46, 73] equipped VLP with temporal transformers for action recognition. [21, 53] introduced prompt learning for efficient retrieval or detection. However, these studies focus more on open-vocabulary scenarios or short video understanding. On the contrary, this paper makes the first exploration to steer VLP for long video temporal localization, under the weakly-supervised setting.

**Strongly-supervised Temporal Action Localization** has achieved great progress [34, 41, 59, 88], given precise action boundaries and categories. There are two popular pipelines: the top-down framework [4, 12, 36, 63, 65, 69, 75, 77, 78, 96] pre-defines massive anchors based on the action distribution prior, and uses fixed-length sliding windows to generate initial proposals, then regresses to refine boundaries; the bottom-up framework [2, 6, 33, 35, 37, 52, 68, 78, 80, 92, 93] trains frame-wise boundary detectors for extreme frames (start, end, center), then groups extreme frames or estimates action lengths to produce proposals. In addition, some papers [9, 44, 81] proposed various fusion strategies to complement these frameworks. Several studies [40, 88] are devoted to better post-processing. Nevertheless, all the above methods demand precise boundary annotations, which are time-consuming and expensive in reality.

**Weakly-supervised Temporal Action Localization** significantly alleviates annotation costs, training with only category labels. The pivotal component is CAS [19, 56, 58, 72] obtained from Classification-Based Pre-training (CBP). But due to the gap of classification and localization, CBP suffers from the serious incomplete issue: only detect discriminative action fragments or even background. To solve this issue, [23, 50, 94] introduced the erasing strategy. [39, 47, 54] produced multiple CAS in parallel for complementar-
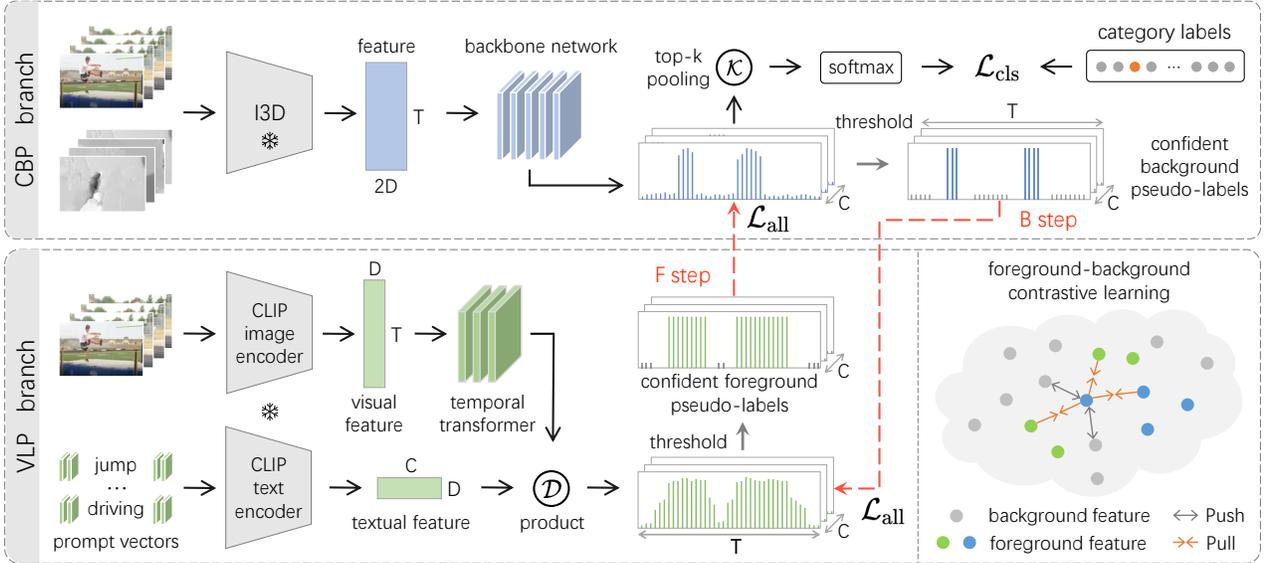
Figure 2. **Distillation-Collaboration Framework.** It covers two parallel branches, named CBP and VLP, and is optimized by an alternating strategy. We warm up the CBP branch in advance. In B step, we freeze both CLIP encoders, and distill confident background pseudo-labels from the CBP branch, to train prompt vectors and temporal Transformer in the VLP branch. In F step, confident foreground pseudo-labels are distilled for the CBP branch. We utilize both knowledge distillation loss and contrastive loss during dual-branch collaboration.

ity. [28, 29, 57, 62] proposed background modeling and context separation. [10, 15, 17, 90] enhanced features by intra- and inter-video modeling. To iteratively refine results, some recent papers [48, 82, 89] introduced the self-training strategy, [45, 64] adopted outer-inner contrastive learning. While encouraging, most of them are trapped in the 'performance-cost dilemma', *i.e.* solely digging from barren category labels to keep costs low. But the lack of annotations leads to a huge performance gap between strong and weak supervisions. To this end, [22, 24, 26, 49, 55, 79] explored the idea of adding instance-number or single-frame annotations for further improvements. As a fresh departure from existing work, to pursue better performance without additional annotation costs, this paper intends to distill free action knowledge from off-the-shelf VLPs, to assist WTAL.

## 3. Method

### 3.1. Notations and Preliminaries

**Task Formulation.** Given $N$ untrimmed videos $\{v_i\}_{i=1}^N$, and their video-level category labels $\{\mathbf{y}_i \in \mathbb{R}^C\}_{i=1}^N$, where $C$ means the total number of action categories, WTAL intends to detect and classify all action instances, in terms of a set of quadruples $\{(s, e, c, p)\}$, where $s$, $e$, $c$, $p$ represent the start time, the end time, the action category and detection score of the action proposal, respectively. Note that each video may contain multiple action instances.

**Motivation.** There exists significant complementarity between the localization results of Classification-Based Pre-training and Vision-Language Pre-training, as concluded in Table 3. The former suffers from incomplete results (serious false negatives), but has good true negatives; the latter suffers from over-complete results (serious false positives), but has good true positives. Such investigations motivate us to collaborate the complementary results for strong alliances, through a distillation-collaboration framework.

**Framework Overview.** As demonstrated in Figure 2, the distillation-collaboration framework consists of two parallel branches, named CBP branch and VLP branch respectively, and is optimized by an alternating strategy. The CBP branch is first warmed-up to produce rich background information with only classification supervision. During B step, we distill confident background pseudo-labels from CBP branch, for VLP branch to tackle the over-complete issue, and thus localize high-quality foreground and background information. During F step, we distill superior pseudo-labels from the well-trained VLP branch, for CBP branch to tackle the incomplete issue. Through such dual-branch collaboration, we effectively fuse the complementary results.

### 3.2. Foreground and Background Distillation

**The CBP Branch** is utilized to identify a large number of background frames as well as several discriminative action frames, by exploiting Classification-Based Pre-training.

Following literature [29, 39], we adopt CB pre-training, *i.e.* the I3D architecture pre-trained on Kinetics [3], to extract RGB and Flow features, and then concatenate them to form the two-stream features $\mathbf{F}_{i3d} \in \mathbb{R}^{T \times 2D}$, where $T$

and $D$ refer to the temporal length and feature dimension. Next, feeding with $\mathbf{F}_{\text{i3d}}$, the CBP branch uses the backbone network for feature fine-tuning and localization, and finally outputs the frame-level action probabilities $\mathbf{P}^{\text{cb}} \in \mathbb{R}^{T \times C}$.

To achieve action classification, we adopt multiple instance learning, *i.e.* for the output $\mathbf{P}^{\text{cb}}$ from the backbone network, we aggregate (pool) the top-$k$ frames' scores as video-level category scores $\widehat{\mathbf{y}} \in \mathbb{R}^C$, then supervise it via the binary cross-entropy loss, which is formulated as:

$$\mathcal{L}_{\text{cls}} = \sum_{c=1}^{C} -y_c \log \widehat{y}_c, \quad \widehat{y}_c = \sigma(\frac{1}{k}\sum \mathcal{K}(\mathbf{P}^{\text{cb}})), \quad (1)$$

where $\mathcal{K}$ denotes the top-$k$ score set in the temporal domain, and $\sigma$ refers to the softmax function.

*Remark.* Under CB Pre-training and category-only supervision, $\mathbf{P}^{\text{cb}}$ is well known for focusing on sparse discriminative action frames, *i.e.* high true negatives but serious false negatives, thus can provide rich background information.

**The VLP Branch** is designed to mine free action knowledge from VL Pre-training, *e.g.* CLIP [60]. Since image-text pre-training lacks sufficient temporal priors, its vanilla localization results have a serious over-complete issue, *i.e.* high false positives. To tackle this issue, we propose to fine-tune CLIP using extensive background samples.

Instead of linear probing, we use efficient prompt learning [21, 30] for fine-tuning: *we freeze the CLIP backbone, only optimize several prompt vectors and temporal layers*.

Concretely, for the visual stream, we first split the video into consecutive frames, and then utilize the CLIP image encoder to extract frame-level features $\mathbf{F}_{\text{vis}} \in \mathbb{R}^{T \times D}$. For temporal relationship construction, we strengthen $\mathbf{F}_{\text{vis}}$ into $\mathbf{F}_{\text{vid}} \in \mathbb{R}^{T \times D}$ through simple temporal transformer layers $\Phi_{\text{temp}}(\cdot)$. While for the textual stream, we first prepend and append several learnable prompt vectors $\Phi_{\text{prmp}}(\cdot)$ to category names, then feed them into the CLIP text encoder, to obtain textual features $\mathbf{F}_{\text{txt}} \in \mathbb{R}^{C \times D}$. Formally,

$$\mathbf{F}_{\text{vid}} = \Phi_{\text{temp}}(\mathbf{F}_{\text{vis}}), \quad \mathbf{F}_{\text{txt}} = \Phi_{\text{txt}}(\Phi_{\text{prmp}}(C_{\text{name}})), \quad (2)$$

where $C_{\text{name}}$ refers to action category names, and $\Phi_{\text{txt}}(\cdot)$ is the CLIP textual encoder. Thereafter, the frame-level localization results $\mathbf{P}^{\text{vl}}$ for this branch can be calculated as:

$$\mathbf{P}^{\text{vl}} = \sigma(\mathbf{F}_{\text{vid}} \cdot \mathbf{F}_{\text{txt}}^{\mathsf{T}}) \in \mathbb{R}^{T \times C}. \quad (3)$$

*Remark.* For the VLP branch, we only optimize lightweight model parameters for false-positive suppression, naturally bringing two main benefits: (1) the frozen CLIP backbone preserves the action prior knowledge in pre-training, thus maintaining high true-positive results; (2) it matches the demand for less supervision data under weakly-supervised settings, and also saves the memory footprint.

**Confident Pseudo-labels.** Since both $\mathbf{P}^{\text{vl}}$ and $\mathbf{P}^{\text{cb}}$ contain somewhat noise, to make reliable use of complementary information, we distill confident location pseudo-labels of foreground and background respectively. That is, for the CBP branch, we distill extensive background pseudo-labels from the output $\mathbf{P}^{\text{cb}}$; while for the VLP branch, we distill sufficient foreground pseudo-labels from the output $\mathbf{P}^{\text{vl}}$.

For both branches, we leverage double thresholds $\delta_h$ and $\delta_l$ ($\delta_h > \delta_l$), to convert localization results $\mathbf{P}$ into *ternary* pseudo-labels $\mathbf{H} \in \mathbb{R}^{T \times C}$, which are formally written as:

$$h_{t,c} = \begin{cases} 1 & \text{if } p_t > \delta_h \text{ and } p_c = y_c \\ 0 & \text{if } p_t < \delta_l \text{ or } p_c \neq y_c \\ -1 & \text{otherwise} \end{cases} \quad (4)$$

where the subscripts $c$ and $t$ refer to the indices of category and frame. More specifically, for either branch, we regard frames with scores more than $\delta_h$ and the correct action category as the foreground; while frames with scores less than $\delta_l$ or with the wrong action category are treated as the background; the remaining frames are considered uncertain. As a result, the pseudo-labels $\mathbf{H}^{\text{cb}}$ cover vast confident background frames, as well as trivial foreground frames; while the pseudo-labels $\mathbf{H}^{\text{vl}}$ contain dense confident foreground frames, and partial background frames. Note that, for both branches, we generate positive and negative frames to avoid trivial solutions, and facilitate contrastive learning for feature enhancement, as detailed in the following section.

### 3.3. Dual-Branch Collaborative Optimization

In this section, we encourage two branches to collaborate with each other, such that forming a strong alliance of their complementary localization results. To reduce serious noises in pseudo-labels, we introduce an alternate training strategy for dual-branch collaborative optimization.

The design rationale is to distill **B**ackground knowledge from the CBP branch for **B** step, while distill **F**oreground knowledge from the VLP branch in **F** step. To be specific, we warm up the CBP branch in advance, using only category supervision, to initialize reliable background frames. During B Step, we freeze the well-trained CBP branch, and then generate confident background pseudo-labels $\mathbf{H}^{\text{cb}}$ to supervise the VLP branch. As a result, these false-positive confusions from vanilla CLIP pre-training are greatly tackled, and the resultant pseudo-labels contain a large number of confident foreground frames and background frames. During F Step, the high-quality pseudo-labels $\mathbf{H}^{\text{vl}}$ are distilled from the frozen VLP branch, to guide the CBP branch for the false-negative suppression. Under such an alternating strategy, these two branches not only complement each other, but also correct each other, thus jointly contributing to more precise and complete action localization.

During each step, to supervise either branch, we adopt both the knowledge distillation loss $\mathcal{L}_{\text{kd}}$ and foreground-

background contrastive loss $\mathcal{L}_{\text{fb}}$. The total optimization loss can be written with a balancing ratio $\lambda$, as follows:

$$\mathcal{L}_{\text{all}} = \mathcal{L}_{\text{kd}}(\mathbf{H}', \mathbf{P}) + \lambda \mathcal{L}_{\text{fb}}(\Psi^+, \Psi^-). \quad (5)$$

Here, $\mathcal{L}_{\text{kd}}$ regularizes one branch to output similar detection results with pseudo-labels from the other branch. Given that there exist some noises in pseudo-labels, *i.e.* uncertain frames, we only make supervision on confident frames.

$$\mathcal{L}_{\text{kd}}(\mathbf{H}', \mathbf{P}) = \frac{1}{O} \sum_{c=1}^{C} \sum_{t=1}^{O} D_{\text{KL}}(h'_{t,c} \,||\, p_{t,c}), \quad (6)$$

where $D_{\text{KL}}(p(x) \,||\, q(x))$ refers to the Kullback-Leibler divergence of distribution $p(x)$ from distribution $q(x)$, $O$ is the total number of confident frames, and $\mathbf{H}'$ denotes the pseudo-labels from the other branch. Note that, the pseudo-labels contain two types of confident frames: foreground and background, which could help to avoid falling into the trivial solutions under one single type of labels.

Moreover, in untrimmed long videos, some background contexts could appear visually similar with the action (foreground). We further introduce contrastive learning to pull foreground features and push background features. Specifically, we treat confident foreground frames from the same action category as the positive set $\Psi^+$, while all confident background frames as the negative set $\Psi^-$, then foreground-background contrastive loss is formulated as:

$$\mathcal{L}_{\text{fb}}(\Psi_i^+, \Psi_i^-) = \sum_i -\log \frac{\sum_{m \in \Psi_i^+} \exp(\mathbf{f}_i \cdot \mathbf{f}_m / \tau)}{\sum_{j \in *} \exp(\mathbf{f}_i \cdot \mathbf{f}_j / \tau)}, \quad (7)$$

where $\mathbf{f} \in \mathbb{R}^D$ refers to the frame feature, $\tau$ means the temperature hyper-parameter for scaling, and $*$ means the union of $\Psi_i^+$ and $\Psi_i^-$. Take the VLP branch as an example, sufficient background pseudo-labels contain considerable hard negative samples, to help distinguish features of foreground and contexts. In addition, the enhanced features of uncertain frames also become more discriminative, further facilitating more complete temporal action localization.

*Discussion.* Comparing to various fusion strategies, our alternating strategy produces more precise and complete results (see Table 4 for details). This strategy plays a similar role as the multi-view co-training, where CBP and VLP branches can be deemed as two distinctive views, thus being quite robust to pseudo-label noises. For B step, we prompt the VLP branch for high-quality pseudo-labels, where some frames with conflicting predictions are still treated as uncertain. In F step, we use feature contrastive loss to make their results more discriminative, *i.e.* further denoising.

### 3.4. Inference

At testing time, we leverage the results from the CBP branch for post-processing, as vision-language pre-training cannot handle Optical Flow, which is essential for WTAL. Given an input video, we first obtain video-level category probabilities and frame-level localization scores. For action classification, we select the classes with probability greater than $\theta_{cls}$; and for localization, we threshold detection scores with $\theta_{loc}$, concatenate consecutive snippets as action proposals, and eliminate redundant proposals with soft non-maximum suppression (NMS). Each proposal is scored with the detection maximum in the proposal interval.

## 4. Experiments

### 4.1. Implementation

**Datasets.** **THUMOS14** owns 413 untrimmed videos from 20 categories, and each video contains an average of 15 instances. As conventions, we train on 200 validation videos, and evaluate on 213 testing videos. Despite its small scale, this dataset is challenging, since video lengths vary widely and actions occur frequently. **ActivityNet1.2** covers 9682 videos of 100 categories, dividing into 4619 training videos, 2383 validation videos, and 2480 testing videos. Almost all videos contain one single category, and action regions take up more than half of the duration in most videos. We train on the training set and evaluate on the validation set.

**Metrics.** To evaluate localization performance, we follow the standard protocol to use mean Average Precision (mAP) at different intersections over union (IoU) thresholds. Note that a proposal is regarded as positive only if both the category prediction is correct and IoU exceeds the set threshold. To clearly evaluate the quality of pseudo-labels, we also report mean Intersection over Union (mIoU) averaged over the foreground categories and the background category.

**Details.** To handle the large variety in video durations, we randomly sample $T$ consecutive snippets for each video. $T$ is set to 1000 on THUMOS14, and 400 on ActivityNet1.2. We utilize the TV-L1 algorithm to extract optical flow from RGB data. For the CBP branch, we use Transformer architectures (multi-head self-attention, layer norm, and MLPs) as the backbone network. For the VLP branch, we use a 2-layer temporal Transformer, prepend and append 16 prompt vectors to textual inputs, both are initialized by $\mathcal{N}(0, 0.01)$. Both CLIP image encoder and text encoder are adopted from ViT-B/16. The framework is optimized by Adam with the learning rate of $10^{-4}$. All hyper-parameters are set by grid search: pseudo-label thresholds $\delta_h = 0.3$, $\delta_l = 0.1$, inference thresholds $\theta_{cls} = 0.85$, $\theta_{loc} = 0.45$, the balancing ratio $\lambda = 0.05$, the temperature $\tau = 0.07$.

### 4.2. Comparison with State-of-the-art Methods

Here, we make comprehensive comparisons with current state-of-the-art methods across multiple IoU thresholds.

**Performance.** The comparisons on THUMOS14 are provided in Table 1. Here, we separate two levels of supervi-

| Supervision | Method | Feature | mAP@IoU | | | | | | | AVG (0.1-0.5) | AVG (0.3-0.7) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | | |
| Strong | DBS [13] | I3D | 56.7 | 54.7 | 50.6 | 43.1 | 34.3 | 24.4 | 14.7 | 47.9 | 33.4 |
| | BUMR [92] | | - | - | 53.9 | 50.7 | 45.4 | 38.0 | 28.5 | - | 43.3 |
| | GCM [88] | | 72.5 | 70.9 | 66.5 | 60.8 | 51.9 | - | - | **64.5** | - |
| | RCL [75] | | - | - | 70.1 | 62.3 | 52.9 | 42.7 | 30.7 | - | 51.7 |
| | ActionFormer [91] | | - | - | 82.1 | 77.8 | 71.0 | 59.4 | 43.9 | - | **66.8** |
| Weak | HAM-Net [19] | I3D | 65.4 | 59.0 | 50.3 | 41.1 | 31.0 | 20.7 | 11.4 | 49.4 | 30.9 |
| | UM [29] | | 63.0 | 56.6 | 49.0 | 40.9 | 30.4 | 21.0 | 10.4 | 48.0 | 30.3 |
| | TS-PCA [43] | | 67.6 | 61.1 | 53.4 | 43.4 | 34.3 | 24.7 | 13.7 | 52.0 | 33.9 |
| | FTCL [10] | | 69.6 | 63.4 | 55.2 | 45.2 | 35.6 | 23.7 | 12.2 | 53.8 | 34.4 |
| | ACGNet [83] | | 68.1 | 62.6 | 53.1 | 44.6 | 34.7 | 22.6 | 12.0 | 52.6 | 33.4 |
| | DCC [31] | | 69.0 | 63.8 | 55.9 | 45.9 | 35.7 | 24.3 | 13.7 | 54.1 | 35.1 |
| | ASM-Loc [15] | | 71.2 | 65.5 | 57.1 | 46.8 | 36.6 | 25.2 | 13.4 | 55.4 | 35.8 |
| | CO2-Net [16] | | 70.1 | 63.6 | 54.5 | 45.7 | 38.3 | 26.4 | 13.4 | 54.4 | 35.7 |
| | RSKP [18] | | 71.3 | 65.3 | 55.8 | 47.5 | 38.2 | 25.4 | 12.5 | 55.6 | 35.9 |
| | DELU [5] | | 71.5 | 66.2 | 56.5 | 47.7 | 40.5 | 27.2 | 15.3 | 56.5 | 37.4 |
| | CO2-Net† | I3D + CLIP | 66.5 | 59.4 | 50.7 | 41.7 | 34.2 | 22.4 | 12.0 | 50.5 | 32.2 |
| | CO2-Net‡ | | 68.8 | 62.0 | 51.7 | 42.2 | 35.4 | 22.3 | 11.7 | 51.9 | 32.7 |
| | DELU† | | 68.5 | 61.2 | 52.1 | 43.1 | 35.0 | 23.1 | 12.7 | 52.0 | 33.2 |
| | DELU‡ | | 70.5 | 64.5 | 55.2 | 45.7 | 38.5 | 25.7 | 13.8 | 54.9 | 35.8 |
| | **Ours** | | **73.5** | **68.8** | **61.5** | **53.8** | **42.0** | **29.4** | **16.8** | **60.0** | **40.8** |

Table 1. **Comparison with state-of-the-art methods on THUMOS14.** For fair comparisons, we reproduce the results of SOTA methods: CO2-Net [16] and DELU [5], by inputting both I3D [3] and CLIP [60] features. † and ‡ refer to averaging or concatenating these two features. AVG(0.1-0.5) and AVG(0.3-0.7) are the average mAP from IoU 0.1 to 0.5 and from IoU 0.3 to 0.7. Our framework significantly surpasses all weakly-supervised competitors using identical features, and is even comparable to early strongly-supervised methods.

| Method | Feature | mAP@IoU | | | AVG |
|---|---|---|---|---|---|
| | | 0.5 | 0.75 | 0.95 | |
| CleanNet [45] | I3D | 37.1 | 20.3 | 5.0 | 21.6 |
| CMCS [39] | | 36.8 | 22.0 | 5.6 | 22.4 |
| TSCN [89] | | 37.6 | 23.7 | 5.7 | 23.6 |
| BasNet [27] | | 38.5 | 24.2 | 5.6 | 24.3 |
| DGAM [62] | | 41.0 | 23.5 | 5.3 | 24.4 |
| UM [29] | | 41.2 | 25.6 | 6.0 | 25.9 |
| ACGNet [83] | | 41.8 | 26.0 | 5.9 | 26.1 |
| D2-Net [54] | | 42.3 | 25.5 | 5.8 | 26.0 |
| CO2-Net [16] | | 43.3 | 26.3 | 5.2 | 26.4 |
| DELU [5] | | 44.2 | 26.7 | 5.4 | 26.9 |
| CO2-Net† | I3D + CLIP | 44.3 | 26.6 | 5.4 | 26.9 |
| CO2-Net‡ | | 44.7 | 26.9 | 5.8 | 27.4 |
| DELU† | | 44.9 | 26.9 | 5.6 | 27.2 |
| DELU‡ | | 45.6 | 27.5 | 5.8 | 27.8 |
| **Ours** | | **48.3** | **29.3** | **6.1** | **29.6** |

Table 2. **Comparison with state-of-the-art methods on ActivityNet1.2.** For fair comparisons, we reproduce CO2-Net [16] and DELU [5], by inputting I3D [3] and CLIP [60] features. † and ‡ refer to averaging or concatenating these features. AVG is the average mAP at the thresholds 0.5:0.05:0.95. Our method significantly surpasses all competitors, especially at loose IoU thresholds.

sion: strong and weak, for better quantification. Generally speaking, our framework achieves new state-of-the-art on all IoU regimes. Comparing with recent methods, the gain of the average mAP (0.3-0.7) even reaches 4-5%, further narrowing the performance gap between weak and strong supervisions. Moreover, our method achieves considerable gains on strict and loose evaluations. For example, when comparing to DELU [5], the gains are 3.5% average mAP (0.1-0.5) and 3.4% average mAP (0.3-0.7), indicating that our results are complete and precise. Furthermore, despite being weakly-supervised settings, at several low IoU thresholds, our framework even performs comparably with some earlier strongly-supervised methods [12, 13, 93].

Table 2 shows the comparison results on ActivityNet1.2. On all IoU thresholds, our designed framework surpasses existing methods by a large margin. In terms of the average mAP, the performance improvement can reach 2.7%, taking the state-of-the-art to a new level. However, due to the lack of precise location annotations, the gain decreases as the IoU threshold becomes stricter, e.g. 4.1% @IoU 0.5 vs. 0.7% @IoU 0.95, when comparing to DELU [5].

**Source of Gain.** Comparing to existing methods, we leverage Vision-Language Pre-training for free knowledge. To make fair comparisons, we also input both I3D and CLIP features into two SOTA methods [5, 16], by one early fusion mode (average or concatenate these features). In general, simply adding CLIP features gives only slight or even negative gains on both datasets, which is due to the bad over-

| Setting | THUMOS | | ActivityNet | |
|---|---|---|---|---|
| | Fore | Back | Fore | Back |
| CB Pre-training | 56.0 | 88.7 | 54.7 | 88.3 |
| VL Pre-training | 75.6 | 34.1 | 72.1 | 45.0 |
| Ours | 72.4 | 80.0 | 69.5 | 71.7 |

Table 3. **Complementarity of pre-training.** CB Pre-training has good background mIoU but inferior foreground mIoU. VL Pre-training has good foreground mIoU but inferior background mIoU. Our method achieves both high foreground and background mIoU.

| Model | $\mathcal{L}_{kd}$ | $\mathcal{L}_{fb}$ | mAP@IoU | | | AVG |
|---|---|---|---|---|---|---|
| | | | 0.3 | 0.5 | 0.7 | (0.3-0.7) |
| A1 | ✓ | | 57.1 | 37.0 | 12.6 | 36.0 |
| A2 | | ✓ | 51.7 | 31.2 | 9.4 | 30.9 |
| A3 | ✓ | ✓ | **61.5** | **42.0** | **16.8** | **40.8** |

Table 5. **Contribution of optimization losses on THUMOS14.** The single knowledge distillation loss $\mathcal{L}_{kd}$ has brought gratifying localization results, and the additional foreground-background contrastive loss $\mathcal{L}_{fb}$ further boosts the performance to the best.

| Fusion Strategy | mAP@IoU | | | AVG |
|---|---|---|---|---|
| | 0.3 | 0.5 | 0.7 | (0.3-0.7) |
| Only F step  (AVG) | 40.0 | 16.1 | 3.3 | 19.1 |
| Only F step  (Weight) | 52.8 | 25.5 | 6.1 | 27.7 |
| Only B step  (AVG) | 56.7 | 33.6 | 11.3 | 34.1 |
| Only B step  (Weight) | 58.5 | 38.8 | 14.8 | 37.8 |
| Alternating | **61.5** | **42.0** | **16.8** | **40.8** |

Table 4. **Comparison of optimization strategies.** "Only F step": extract pseudo-labels from vanilla VLP to train the CBP branch. "Only B step": get pseudo-labels from the warm-up CBP branch to train the VLP branch. We combine two branches by averaging and weighting. Our alternating strategy shows clear superiority.

| Method | mAP@IoU | | | AVG |
|---|---|---|---|---|
| | 0.3 | 0.5 | 0.7 | (0.3-0.7) |
| UM [29] | 49.0 | 30.4 | 10.4 | 30.3 |
| UM + Ours | **53.9** | **34.7** | **13.6** | **34.3** |
| CO2-Net [16] | 54.5 | 38.3 | 13.4 | 35.7 |
| CO2-Net + Ours | **56.2** | **39.7** | **15.9** | **37.5** |

Table 6. **Framework generalization on THUMOS14.** Our proposed framework is generalized, *i.e.* existing methods can serve as our CBP branch, to achieve further improvements.

complete issue of CLIP (detailed in Table 3). For THUMOS with complex and frequent actions, over-complete worsens results. While in ActivityNet, most videos only contain one action that takes up half of the video duration, thus benefiting a bit from over-complete. Nevertheless, with identical features, our method significantly outperforms existing competitors, proving the effectiveness of our framework.

### 4.3. Ablation Study and Comparison

Here, we evaluate the contributions of each component and framework designs, to further dissect our framework.

**Complementarity of pre-training.** As the detection performance is mainly determined by pseudo-labels, here we show the quality of pseudo-labels, in terms of foreground mIoU and background mIoU. Table 3 provides the comprehensive results for various settings on both benchmarks.

For common Classification-Based Pre-training, its localization results suffer from the incomplete issue. In detail, the background mIoU is impressive, *i.e.* high true negatives, while the foreground mIoU is poor, *i.e.* serious false negatives. The main reason is that features pre-trained on action classification datasets only highlight sparse discriminative frames. While the results of Vision-Language Pre-training are just the opposite: suffering from the over-complete issue. The foreground mIoU is considerable, *i.e.* high true positives, but the background mIoU is terrible, *i.e.* serious false positives. The main reason is that VLP using image-

text pairs lacks temporal priors. The above results strongly prove the complementarity between these two pre-training.

Moreover, on both datasets, our method achieves high foreground and background mIoU, that is, more precise and complete localization. This mainly benefits from extensive background supervision provided by the CBP branch, and extensive foreground supervision distilled from the VLP branch. Besides, comparing to VL Pre-training, our results achieve immediate gains on the background mIoU, while only slight drops on the foreground mIoU. This reveals that the lightweight trainable parameters, *i.e.* prompts + Transformer, indeed retain the action prior knowledge in VLP, and also significantly suppress false-positive results.

**Comparison of optimization strategy.** To evaluate the efficacy of our alternating strategy, Table 4 makes comparison with another two solutions. (1) 'Only F step': extract foreground pseudo-labels from vanilla VLP to train the CBP branch; (2) 'Only B step': distill background pseudo-labels from the warm-up CBP branch to train the VLP branch. For each solution, we combine the results from two branches via averaging and weighting, respectively.

'Only F step' performs the worst, mainly suffering from heavy noise in vanilla VLP pseudo-labels (over-complete). By fine-tuning vanilla VLP with many background frames, 'Only B step' gets great improvements, but still lacks full fusion of complementarity. Moreover, the weighted operation could suppress noise somewhat, and thus causes better results than the average operation. By comparison, our alternating strategy shows great advantages over competitors, proving the *non-trivial* nature of fusing complementary in-
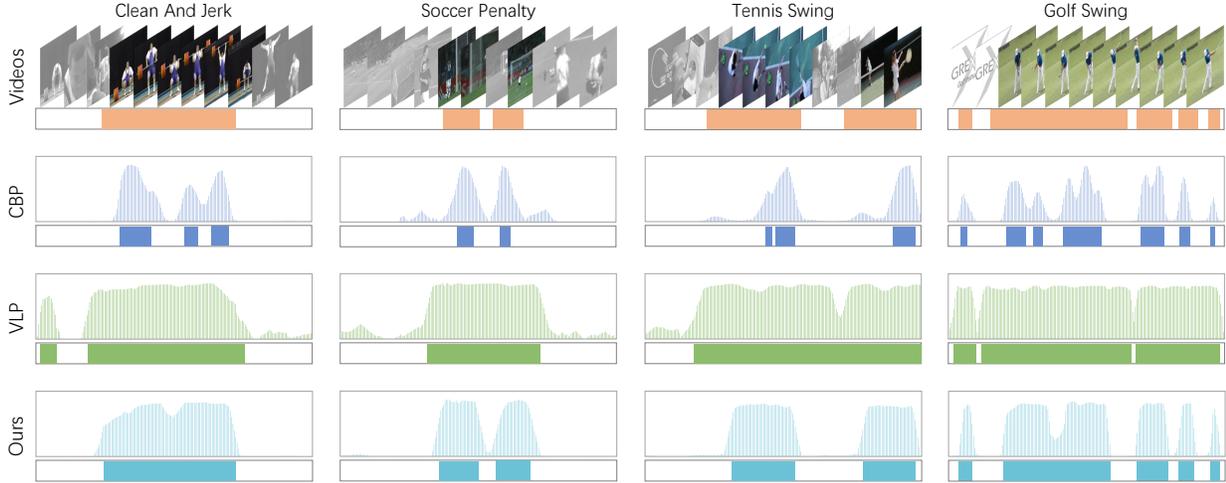
Figure 3. **Qualitative comparisons.** The first two rows are videos and ground-truth action intervals. The last six rows are frame-level action probabilities, and localization results of Classification-Based Pre-training (CBP), Vision-Language Pre-training (VLP), and our framework, respectively. CBP suffers from the incomplete issue, while VLP has the over-complete issue. Our framework distills foreground knowledge from VLP and background from CBP, for the strong collaborative alliance, thus bringing more complete and precise results.

formation, and robust denoising to pseudo-labels.

**Contribution of various losses.** To train our method, both knowledge distillation loss $\mathcal{L}_{kd}$ and foreground-background contrastive loss $\mathcal{L}_{fb}$ are leveraged. In Table 5, we analyze the effectiveness of optimization losses on THUMOS14. The single $\mathcal{L}_{kd}$ (model A1) already brings gratifying performance. This reveals that the confident pseudo-labels of foreground and background distilled from two branches, are well fused by the alternate optimization. On the other hand, with only $\mathcal{L}_{fb}$, model A2 also performs barely satisfactory results. This is because $\mathcal{L}_{fb}$ differentiates foreground features and background features, thus eliminating extensive uncertain frames, and also making the localization task easier. Overall, these two losses jointly contribute to the best performance, indicating that both are essential.

**Framework generalization.** Our distillation-collaboration framework is generalized, which means that existing WTAL methods can be employed as the CBP branch. Table 6 takes two typical methods: CO2-Net [16] and UM [29], as examples. Our framework further improves their performance by up to 2-4% average mAP, showing a good generalization to other methods and other backbone designs.

### 4.4. Qualitative Results

To intuitively demonstrate the superiority, Figure 3 visualizes detection results from various types of videos.

In general, Classification-Based Pre-training highlights only several discriminative action frames (the incomplete issue), which is more prominent for videos covering low-frequency actions. On the contrary, Vision-Language Pre-training tends to over-activate the action foreground to the background (the over-complete issue), which is especially

evident in videos with high-frequency actions. We design the distillation-collaboration framework to fuse the complementarity from these two pre-training. In B step, extensive confident background information is distilled from the well-trained CBP branch, to supervise the VLP branch for false-positive suppression. In F step, sufficient confident foreground locations are distilled from the VLP branch, to guide the CBP branch for false-negative elimination. Hence, our method establishes the strong alliance by collaborative optimization. The detection results are more precise and more complete, regardless of dense or sparse actions.

## 5. Conclusion

This work proposes the novel distillation-collaboration framework to distill free knowledge from Vision-Language Pre-training, for weakly-supervised temporal action localization. Our core insight is that existing VLP often localizes over-complete actions, which is just complementary to the incomplete results of conventional Classification-Based Pre-training. And to form strong alliances, we optimize the framework containing complementary dual branches by an alternating strategy: distill confident background pseudo-labels from the CBP branch, and the confident foreground pseudo-labels from the VLP branch, for collaborative training. Extensive experiments show the significance of distilling VLP and our superior performance. Thorough ablations are studied both quantitatively and qualitatively.

## 6. Limitations and Future Work

For the CBP branch, we freeze the I3D architecture pre-trained on Kinetics [3], to extract RGB and Flow features.

Such one frozen extractor could save computing resources, but may somewhat limit the performance.

For the VLP branch, we leverage the CLIP [60], which is pre-trained with 400M image-text pairs collected from web, thus could potentially bias towards web data.

As the future work, we expect more computing resources available, to further optimize our distillation-collaboration framework into the end-to-end training setups, also rendering asynchronous online training.

# References

[1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. In *Advances in Neural Information Processing Systems*, 2022. 2

[2] Yueran Bai, Yingying Wang, Yunhai Tong, Yang Yang, Qiyue Liu, and Junhui Liu. Boundary content graph neural network for temporal action proposal generation. In *Proceedings of the European Conference on Computer Vision*, 2020. 2

[3] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 1, 3, 6, 8

[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. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 2

[5] Mengyuan Chen, Junyu Gao, Shicai Yang, and Changsheng Xu. Dual-evidential learning for weakly-supervised temporal action localization. In *Proceedings of the European Conference on Computer Vision*, 2022. 6

[6] Feng Cheng and Gedas Bertasius. Tallformer: Temporal action localization with long-memory transformer. In *Proceedings of the European Conference on Computer Vision*, 2022. 2

[7] Katherine Crowson, Stella Biderman, Daniel Kornis, Dashiell Stander, Eric Hallahan, Louis Castricato, and Edward Raff. Vqgan-clip: Open domain image generation and editing with natural language guidance. In *Proceedings of the European Conference on Computer Vision*, 2022. 2

[8] Andrea Frome, Greg S Corrado, Jon Shlens, Samy Bengio, Jeff Dean, Marc'Aurelio Ranzato, and Tomas Mikolov. Devise: A deep visual-semantic embedding model. In *Advances in Neural Information Processing Systems*, 2013. 2

[9] Jiyang Gao, Kan Chen, and Ram Nevatia. Ctap: Complementary temporal action proposal generation. In *Proceedings of the European Conference on Computer Vision*, 2018. 2

[10] Junyu Gao, Mengyuan Chen, and Changsheng Xu. Fine-grained temporal contrastive learning for weakly-supervised temporal action localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 3, 6

[11] Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. Tall: Temporal activity localization via language query. In *Proceedings of the International Conference on Computer Vision*, 2017. 1

[12] Jiyang Gao, Zhenheng Yang, Kan Chen, Chen Sun, and Ram Nevatia. Turn tap: Temporal unit regression network for temporal action proposals. In *Proceedings of the International Conference on Computer Vision*, 2017. 2, 6

[13] Zhanning Gao, Le Wang, Qilin Zhang, Zhenxing Niu, Nanning Zheng, and Gang Hua. Video imprint segmentation for temporal action detection in untrimmed videos. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019. 6

[14] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. In *Proceedings of the International Conference on Learning Representations*, 2021. 2

[15] 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 Conference on Computer Vision and Pattern Recognition*, 2022. 3, 6

[16] Fa-Ting Hong, Jia-Chang Feng, Dan Xu, Ying Shan, and Wei-Shi Zheng. Cross-modal consensus network for weakly supervised temporal action localization. In *Proceedings of ACM International Conference on Multimedia*, 2021. 6, 7, 8

[17] Linjiang Huang, Liang Wang, and Hongsheng Li. Foreground-action consistency network for weakly supervised temporal action localization. In *Proceedings of the International Conference on Computer Vision*, 2021. 3

[18] Linjiang Huang, Liang Wang, and Hongsheng Li. Weakly supervised temporal action localization via representative snippet knowledge propagation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 6

[19] Ashraful 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*, 2021. 2, 6

[20] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V. Le, Yunhsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *Proceedings of the International Conference on Machine Learning*, 2021. 2

[21] Chen Ju, Tengda Han, Kunhao Zheng, Ya Zhang, and Weidi Xie. Prompting visual-language models for efficient video understanding. In *Proceedings of the European Conference on Computer Vision*, 2022. 2, 4

[22] Chen Ju, Peisen Zhao, Siheng Chen, Ya Zhang, Yanfeng Wang, and Qi Tian. Divide and conquer for single-frame temporal action localization. In *Proceedings of the International Conference on Computer Vision*, 2021. 3

[23] Chen Ju, Peisen Zhao, Siheng Chen, Ya Zhang, Xiaoyun Zhang, and Qi Tian. Adaptive mutual supervision for weakly-supervised temporal action localization. *IEEE Transactions on Multimedia*, 2022. 2

[24] Chen Ju, Peisen Zhao, Ya Zhang, Yanfeng Wang, and Qi Tian. Point-level temporal action localization: Bridging fully-supervised proposals to weakly-supervised losses. *arXiv preprint arXiv:2012.08236*, 2020. 3

[25] Apoorv Khandelwal, Luca Weihs, Roozbeh Mottaghi, and Aniruddha Kembhavi. Simple but effective: Clip embeddings for embodied ai. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 2

[26] Pilhyeon Lee and Hyeran Byun. Learning action completeness from points for weakly-supervised temporal action localization. In *Proceedings of the International Conference on Computer Vision*, 2021. 3

[27] 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*, 2020. 1, 6

[28] 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*, 2020. 3

[29] 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*, 2021. 1, 3, 6, 7, 8

[30] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the Conference on Empirical Methods in Natural Language Processinng*, 2021. 4

[31] 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 Conference on Computer Vision and Pattern Recognition*, 2022. 6

[32] Yue Liao, Aixi Zhang, Miao Lu, Yongliang Wang, Xiaobo Li, and Si Liu. Gen-vlkt: Simplify association and enhance interaction understanding for hoi detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 2

[33] Chuming Lin, Jian Li, Yabiao Wang, Ying Tai, Donghao Luo, Zhipeng Cui, Chengjie Wang, Jilin Li, Feiyue Huang, and Rongrong Ji. Fast learning of temporal action proposal via dense boundary generator. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020. 2

[34] 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 Conference on Computer Vision and Pattern Recognition*, 2021. 2

[35] Tianwei Lin, Xiao Liu, Xin Li, Errui Ding, and Shilei Wen. Bmn: Boundary-matching network for temporal action proposal generation. In *Proceedings of the International Conference on Computer Vision*, 2019. 2

[36] Tianwei Lin, Xu Zhao, and Zheng Shou. Single shot temporal action detection. In *Proceedings of ACM International Conference on Multimedia*, 2017. 2

[37] Tianwei Lin, Xu Zhao, Haisheng Su, Chongjing Wang, and Ming Yang. Bsn: Boundary sensitive network for temporal action proposal generation. In *Proceedings of the European Conference on Computer Vision*, 2018. 2

[38] Ziyi Lin, Shijie Geng, Renrui Zhang, Peng Gao, Gerard de Melo, Xiaogang Wang, Jifeng Dai, Yu Qiao, and Hongsheng Li. Frozen clip models are efficient video learners. In *Proceedings of the European Conference on Computer Vision*, 2022. 2

[39] Daochang Liu, Tingting Jiang, and Yizhou Wang. Completeness modeling and context separation for weakly supervised temporal action localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 1, 2, 3, 6

[40] Qinying Liu and Zilei Wang. Progressive boundary refinement network for temporal action detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020. 2

[41] Xiaolong Liu, Song Bai, and Xiang Bai. An empirical study of end-to-end temporal action detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 2

[42] Xihui Liu, Dong Huk Park, Samaneh Azadi, Gong Zhang, Arman Chopikyan, Yuxiao Hu, Humphrey Shi, Anna Rohrbach, and Trevor Darrell. More control for free! image synthesis with semantic diffusion guidance. *arXiv preprint arXiv:2112.05744*, 2021. 2

[43] Yuan Liu, Jingyuan Chen, Zhenfang Chen, Bing Deng, Jianqiang Huang, and Hanwang Zhang. The blessings of unlabeled background in untrimmed videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021. 6

[44] Yuan Liu, Lin Ma, Yifeng Zhang, Wei Liu, and Shih-Fu Chang. Multi-granularity generator for temporal action proposal. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 2

[45] Ziyi Liu, Le Wang, Qilin Zhang, Zhanning Gao, Zhenxing Niu, Nanning Zheng, and Gang Hua. Weakly supervised temporal action localization through contrast based evaluation networks. In *Proceedings of the International Conference on Computer Vision*, 2019. 3, 6

[46] Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li. Clip4clip: An empirical study of clip for end to end video clip retrieval. *Neurocomputing*, 2021. 2

[47] Wang Luo, Tianzhu Zhang, Wenfei Yang, Jingen Liu, Tao Mei, Feng Wu, and Yongdong Zhang. Action unit memory network for weakly supervised temporal action localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021. 2

[48] Zhekun Luo, Devin Guillory, Baifeng Shi, Wei Ke, Fang Wan, Trevor Darrell, and Huijuan Xu. Weakly-supervised action localization with expectation-maximization multi-instance learning. In *Proceedings of the European Conference on Computer Vision*, 2020. 3

[49] Fan Ma, Linchao Zhu, Yi Yang, Shengxin Zha, Gourab Kundu, Matt Feiszli, and Zheng Shou. Sf-net: Single-frame supervision for temporal action localization. In *Proceedings of the European Conference on Computer Vision*, 2020. 3

[50] 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*, 2020. 1, 2

[51] Yasuhide Mori, Hironobu Takahashi, and Ryuichi Oka. Image-to-word transformation based on dividing and vector quantizing images with words. In *Proceedings of ACM International Conference on Multimedia*, 1999. 2

[52] Sauradip Nag, Xiatian Zhu, Yi-Zhe Song, and Tao Xiang. Temporal action detection with global segmentation mask learning. In *Proceedings of the European Conference on Computer Vision*, 2022. 2

[53] Sauradip Nag, Xiatian Zhu, Yi-Zhe Song, and Tao Xiang. Zero-shot temporal action detection via vision-language prompting. In *Proceedings of the European Conference on Computer Vision*, 2022. 2

[54] Sanath Narayan, Hisham Cholakkal, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. D2-net: Weakly-supervised action localization via discriminative embeddings and denoised activations. In *Proceedings of the International Conference on Computer Vision*, 2021. 2, 6

[55] 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 International Conference on Computer Vision*, 2019. 3

[56] Phuc Nguyen, Ting Liu, Gautam Prasad, and Bohyung Han. Weakly supervised action localization by sparse temporal pooling network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 1, 2

[57] Phuc Xuan Nguyen, Deva Ramanan, and Charless C Fowlkes. Weakly-supervised action localization with background modeling. In *Proceedings of the International Conference on Computer Vision*, 2019. 3

[58] Sujoy Paul, Sourya Roy, and AmitK Roy-Chowdhury. W-talc: Weakly-supervised temporal activity localization and classification. In *Proceedings of the European Conference on Computer Vision*, 2018. 1, 2

[59] Zhiwu Qing, Haisheng Su, Weihao Gan, Dongliang Wang, Wei Wu, Xiang Wang, Yu Qiao, Junjie Yan, Changxin Gao, and Nong Sang. Temporal context aggregation network for temporal action proposal refinement. In *Proceedings of the International Conference on Computer Vision*, 2021. 2

[60] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *Proceedings of the International Conference on Machine Learning*, 2021. 2, 4, 6, 9

[61] Yongming Rao, Wenliang Zhao, Guangyi Chen, Yansong Tang, Zheng Zhu, Guan Huang, Jie Zhou, and Jiwen Lu. Denseclip: Language-guided dense prediction with context-aware prompting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 2

[62] Baifeng Shi, Qi Dai, Yadong Mu, and Jingdong Wang. Weakly-supervised action localization by generative attention modeling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 3, 6

[63] Zheng Shou, Jonathan Chan, Alireza Zareian, Kazuyuki Miyazawa, and Shih-Fu Chang. Cdc: Convolutional-de-convolutional networks for precise temporal action localization in untrimmed videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2

[64] Zheng Shou, Hang Gao, Lei Zhang, Kazuyuki Miyazawa, and Shih-Fu Chang. Autoloc: Weakly-supervised temporal action localization in untrimmed videos. In *Proceedings of the European Conference on Computer Vision*, 2018. 3

[65] Zheng Shou, Dongang Wang, and Shih-Fu Chang. Temporal action localization in untrimmed videos via multi-stage cnns. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 2

[66] Tianmin Shu, Dan Xie, Brandon Rothrock, Sinisa Todorovic, and Song Chun Zhu. Joint inference of groups, events and human roles in aerial videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015. 1

[67] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In *Advances in Neural Information Processing Systems*, 2014. 1

[68] Haisheng Su, Weihao Gan, Wei Wu, Yu Qiao, and Junjie Yan. Bsn++: Complementary boundary regressor with scale-balanced relation modeling for temporal action proposal generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021. 2

[69] Jing Tan, Jiaqi Tang, Limin Wang, and Gangshan Wu. Relaxed transformer decoders for direct action proposal generation. In *Proceedings of the International Conference on Computer Vision*, 2021. 2

[70] Guy Tevet, Brian Gordon, Amir Hertz, Amit H Bermano, and Daniel Cohen-Or. Motionclip: Exposing human motion generation to clip space. In *Proceedings of the European Conference on Computer Vision*, 2022. 2

[71] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Clip-nerf: Text-and-image driven manipulation of neural radiance fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 2

[72] Limin Wang, Yuanjun Xiong, Dahua Lin, and Luc Van Gool. Untrimmednets for weakly supervised action recognition and detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 1, 2

[73] Mengmeng Wang, Jiazheng Xing, and Yong Liu. Actionclip: A new paradigm for video action recognition. *arXiv preprint arXiv:2109.08472*, 2021. 2

[74] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *Proceedings of the International Conference on Machine Learning*, 2022. 2

[75] Qiang Wang, Yanhao Zhang, Yun Zheng, and Pan Pan. Rcl: Recurrent continuous localization for temporal action detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 2, 6

[76] Jason Weston, Samy Bengio, and Nicolas Usunier. WSA-BIE: Scaling up to large vocabulary image annotation. In *Proceedings of the International Joint Conference on Artificial Intelligence*, 2011. 2

[77] Huijuan Xu, Abir Das, and Kate Saenko. R-c3d: Region convolutional 3d network for temporal activity detection. In *Proceedings of the International Conference on Computer Vision*, 2017. 2

[78] 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 Conference on Computer Vision and Pattern Recognition*, 2020. 2

[79] Yunlu Xu, Chengwei Zhang, Zhanzhan Cheng, Jianwen Xie, Yi Niu, Shiliang Pu, and Fei Wu. Segregated temporal assembly recurrent networks for weakly supervised multiple action detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019. 3

[80] Haosen Yang, Wenhao Wu, Lining Wang, Sheng Jin, Boyang Xia, Hongxun Yao, and Hujie Huang. Temporal action proposal generation with background constraint. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022. 2

[81] Le Yang, Houwen Peng, Dingwen Zhang, Jianlong Fu, and Junwei Han. Revisiting anchor mechanisms for temporal action localization. *IEEE Transactions on Image Processing*, 2020. 2

[82] 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 Conference on Computer Vision and Pattern Recognition*, 2021. 3

[83] Zichen Yang, Jie Qin, and Di Huang. Acgnet: Action complement graph network for weakly-supervised temporal action localization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022. 6

[84] Lewei Yao, Runhui Huang, Lu Hou, Guansong Lu, Minzhe Niu, Hang Xu, Xiaodan Liang, Zhenguo Li, Xin Jiang, and Chunjing Xu. Filip: Fine-grained interactive language-image pre-training. In *Proceedings of the International Conference on Learning Representations*, 2022. 2

[85] Ting Yao, Tao Mei, and Yong Rui. Highlight detection with pairwise deep ranking for first-person video summarization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 1

[86] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. In *Advances in Neural Information Processing Systems*, 2022. 2

[87] Alireza Zareian, Kevin Dela Rosa, Derek Hao Hu, and Shih-Fu Chang. Open-vocabulary object detection using captions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021. 2

[88] 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 International Conference on Computer Vision*, 2019. 2, 6

[89] 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*, 2020. 1, 3, 6

[90] 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 Conference on Computer Vision and Pattern Recognition*, 2021. 3

[91] Chenlin Zhang, Jianxin Wu, and Yin Li. Actionformer: Localizing moments of actions with transformers. In *Proceedings of the European Conference on Computer Vision*, 2022. 6

[92] 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*, 2020. 2, 6

[93] Yue Zhao, Yuanjun Xiong, Limin Wang, Zhirong Wu, Xiaoou Tang, and Dahua Lin. Temporal action detection with structured segment networks. In *Proceedings of the International Conference on Computer Vision*, 2017. 2, 6

[94] Jia-Xing Zhong, Nannan Li, Weijie Kong, Tao Zhang, Thomas H Li, and Ge Li. Step-by-step erasion, one-by-one collection: A weakly supervised temporal action detector. In *Proceedings of ACM International Conference on Multimedia*, 2018. 2

[95] Chong Zhou, Chen Change Loy, and Bo Dai. Extract free dense labels from clip. In *Proceedings of the European Conference on Computer Vision*, 2022. 2

[96] Zixin Zhu, Wei Tang, Le Wang, Nanning Zheng, and Gang Hua. Enriching local and global contexts for temporal action localization. In *Proceedings of the International Conference on Computer Vision*, 2021. 2