# BasicTAD: an Astounding RGB-Only Baseline for Temporal Action Detection

Min Yang[a,1], Guo Chen[a,1], Yin-Dong Zheng[a], Tong Lu[a] and Limin Wang [a,*]

[a]*The State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, 210023, China*

## ARTICLE INFO

## ABSTRACT

Temporal action detection (TAD) is extensively studied in the video understanding community by generally following the object detection pipeline in images. However, complex designs are not uncommon in TAD, such as two-stream feature extraction, multi-stage training, complex temporal modeling, and global context fusion. In this paper, we do not aim to introduce any novel technique for TAD. Instead, we study a simple, straightforward, yet must-known baseline given the current status of complex design and low detection efficiency in TAD. In our simple baseline (BasicTAD), we decompose the TAD pipeline into several essential components: data sampling, backbone design, neck construction, and detection head. We extensively investigate the existing techniques in each component for this baseline and, more importantly, perform end-to-end training over the entire pipeline thanks to the simplicity of design. As a result, this simple BasicTAD yields an astounding and real-time RGB-Only baseline very close to the state-of-the-art methods with two-stream inputs. In addition, we further improve the BasicTAD by preserving more temporal and spatial information in network representation (termed as PlusTAD). Empirical results demonstrate that our PlusTAD is very efficient and significantly outperforms the previous methods on the datasets of THUMOS14 and FineAction. Meanwhile, we also perform in-depth visualization and error analysis on our proposed method and try to provide more insights into the TAD problem. Our approach can serve as a strong baseline for future TAD research. The code and model are released at https://github.com/MCG-NJU/BasicTAD.

## 1. Introduction

Video understanding is a fundamental and challenging problem in computer vision research. Temporal action detection (TAD) (Jiang et al., 2014; Heilbron et al., 2015; Liu et al., 2022b), which aims to localize the temporal interval of each action instance in an untrimmed video and recognize its action class, is particularly crucial for long-term video understanding. Due to the high complexity of video and the lack of a clear definition of action boundaries (Moltisanti et al., 2017), past efforts on TAD often employ a relatively sophisticated paradigm to solve this problem. For example, they typically use two-stream inputs for feature extraction in advance (Qing et al., 2021; Yang et al., 2020; Tan et al., 2021; Liu and Wang, 2020; Lin et al., 2021; Wu et al., 2021), a multi-stage training strategy for different components (Lin et al., 2018; Tan et al., 2021), temporal reasoning with complex models like graph neural networks (Xu et al., 2020; Chen et al., 2022; Zeng et al., 2019), and yielding the detection results with the global classification scores (Lin et al., 2019b; Tan et al., 2021; Lin et al., 2018) (e.g., UntrimmedNet (Wang et al., 2017)). Such complex designs need additional extraction of optical flow, action detectors trained on temporal features rather than raw video frames, complex components for multi-stage training and additional classification results. Unlike object detectors in images, these complex designs hinder the existing TAD methods from being a simple and neat detection framework that could be easily trained in an end-to-end manner and efficiently deployed in real-world applications.

Given the current status of complex designs and low detection efficiency in the TAD method, there is a need to step back and reconsider the basic principle of designing an efficient and practical TAD framework. In order to cater to the need for scalability, ease of use, and deployment into real applications for real needs, **simplicity** is the most important requirement for building an efficient and practical TAD method. Inspired by the great success of object detectors (Lin et al., 2017c; Tian et al., 2019), we argue that an optimal action detector should strictly follow a similar *modular design*, where each component has its function, and different components can be easily integrated seamlessly. In addition, **efficient training** is another important property that needs to be considered. Multi-stage training often brings extra computational and storage costs while not being conducive to unleashing the power of deep learning. We argue that an ideal action detector should be end-to-end trainable where the entire pipeline trains and infers directly against raw video frames. Finally, being **free of pre-processing** is another desired property for action detector design, particularly important for the deployment in real applications. We argue that optical flow extraction is the bottleneck for many TAD methods as it often requires a high computational cost to calculate these inputs explicitly. An RGB-only TAD approach should be more favorable and practical for video understanding applications.

---

*Corresponding author (lmwang@nju.edu.cn).
ORCID(s):
[1]M. Yang and G. Chen equally contribute to this work.

Therefore, it is urgent to re-design a simple baseline for Temporal Anomaly Detection (TAD) to meet the abovementioned requirements. In this paper, we do not aim to introduce any new techniques for building a TAD framework. Rather, we investigate a simpler yet crucial baseline for the TAD task. Unlike the complex designs used in existing TAD frameworks, we carefully design a modular temporal detection framework that enables us to conduct in-depth studies on different components to determine the optimal settings, including data augmentation, backbone, neck, and detection head. While the design for standard object detectors has been highly mature and robust, better practice for building a temporal detector in videos is still needed. For each module design, we choose the simplest and most basic options and extensively investigate them to discover the optimal configurations. Additionally, we explore different data sampling and augmentation techniques during training and testing to develop an effective TAD method, which has been largely overlooked in previous methods. Based on our modular design and extensive empirical investigation, we establish an RGB-only baseline for TAD called BasicTAD. Our BasicTAD achieves competitive performance compared to the state-of-the-art methods with two-stream inputs. This performance indicates that the basic design of TAD requires reconsideration, and the search for the optimal basic design is foundational work that must be taken into account.

Encouraged by the outstanding performance of Basic-TAD, we strictly follow its basic principle of designing an ideal TAD method and then further improve it minimally but meaningfully. Our core idea is to preserve complete temporal information in our backbone and spatial information in our neck under enhanced data augmentation. The implementation is quite simple to achieve these objectives by removing the temporal downsampling operations in the backbone and exchanging the spatial and temporal pooling operations in the neck. Such a simple design allows us to further explore larger temporal and spatial inputs for better detection results. These small changes would significantly improve the performance of BasicTAD, increasing the performance from 50.5% to 59.6% mAP on the THUMOS14 dataset. The resulted TAD method is denoted by **PlusTAD**, and we ascribe its good performance to our core idea of keeping rich information and careful implementation. We believe our work is timely and will draw the whole TAD community to reconsider the design of TAD methods, given the current status of complex design and low efficiency. The basic comparisons of our BasicTAD and PlusTAD with previous end-to-end TAD methods are summarized in Table 1. In summary, our contributions are threefold:

- We reconsider the TAD pipeline and present a simple modular detection framework. For each component in our modular framework, we perform extensive studies on the existing basic choices. Through extensive empirical studies, we have developed good practice to build an astounding RGB-Only baseline method for TAD, called BasicTAD.

- Encouraged by the outstanding performance of Basic-TAD, we further improve it with minimal changes to fully unleash the power of deep networks. Our core idea is to preserve richer information during our backbone feature extraction and neck compression. This idea could be easily implemented with high efficiency and the resulting PlusTAD significantly improve the TAD performance on the THUMOS14 dataset.

- The extensive experiments on THUMOS14 and Fine-Action demonstrate that PlusTAD outperforms previous state-of-the-art methods by a large margin. In particular, we obtain an average mAP of 59.6% on the THUMOS14 only with RGB input. In addition, we perform in-depth ablation studies and error analysis to provide more insights for the future TAD pipeline design.

## 2. Related Work

### 2.1. Action Recognition

Action recognition is an important task in video understanding. Current deep learning-based action recognition methods can be mainly divided into two types. The first one is the CNN-based method, which includes the specific video architectures of 2D CNN, 3D CNN, and (2+1)D CNN. 2D CNN methods (Simonyan and Zisserman, 2014; Wang et al., 2016, 2021b) take RGB frames and optical flow as input to capture appearance and motion information, respectively. 3D CNN methods (Tran et al., 2015, 2017; Carreira and Zisserman, 2017; Wang et al., 2018a,b; Feichtenhofer et al., 2019) capture spatiotemporal information between frames by performing 3D convolution on stacked video frames. (2+1)D CNN methods (Qiu et al., 2017; Xie et al., 2018; Tran et al., 2018; Lin et al., 2019a; Liu et al., 2020; Li et al., 2020; Liu et al., 2021d) model spatiotemporal features by decoupling 3D convolution into 2D convolution and 1D convolutions or efficient temporal modules for reducing the computational complexity of the network.

The second type of neural network is the Transformer (Vaswani et al., 2017) architecture, which successfully uses a global self-attention mechanism to address the limitation of CNN in an insufficient receptive field. The great success of image transformers (Dosovitskiy et al., 2021; Touvron et al., 2021; Liu et al., 2021b) has led to the investigation of video transformers (Bertasius et al., 2021; Liu et al., 2021c; Neimark et al., 2021; Arnab et al., 2021; Tong et al., 2022; Wang et al., 2023) for action recognition in videos. However, compared with CNN-based methods, the quadratic complexity of self-attention operations of Transformer-based architectures has led to high training costs and memory consumption. These research efforts on backbone design are orthogonal to our study on the TAD task. Any video backbone could be compatible with our BasicTAD and PlusTAD designs. In this paper, we explore the commonly-used action recognition backbone (Tran et al., 2015; Carreira and Zisserman, 2017; Wang et al., 2018b; Feichtenhofer et al., 2019) and mainly

**Table 1**
**Comparison with previous end-to-end TAD methods only with RGB input on THUMOS14 (Jiang et al., 2014) dataset.**
We categorize components and settings based on their order in the whole pipeline: i) Data Stream: modal, resolution in temporal and spatial; ii) Network: The backbone with $\beta$ times temporal downsampling ($\times\beta$) for feature extraction, the Neck for feature aggregation, and the Head for detecting temporal action segments. The temporal downsample module (TDM), and temporal feature pyramid network (TFPN) in Neck are two different methods for generating multi-scale features in the subsequent discussion. SP-NECK means our improved neck module. iii) Performance: The speed metric is FPS. mAP represents the average mAP from mAP@0.3 to mAP@0.7.

| Method | R-C3D (Xu et al., 2017) | AFSD (Lin et al., 2021) | DaoTAD (Wang et al., 2021a) | **BasicTAD** (Ours) | **PlusTAD** (Ours) |
|---|---|---|---|---|---|
| | | | Data Stream | | |
| Modality | RGB | RGB | RGB | RGB | RGB |
| Resolution | 768 25 FPS 112×112 | 256 10 FPS 112×112 | 768 25 FPS 128×128 | 768 24 FPS 128×128 | 96 3 FPS short-128 |
| | | | Network | | |
| Backbone | C3D(×8) | I3D(×8) | R50-I3D(×8) | SlowOnly(×8) | **TP-SlowOnly**(×1) |
| Neck | TDM | TFPN | TDM+TFPN | TDM+TFPN/TDM | **SP-NECK** |
| Head | R-C3D | AFSD | RetinaNet | **Anchor-based/Anchor-free** | **Anchor-based/Anchor-free** |
| | | | Performance | | |
| mAP | < 36.4 | 43.6 | 50.0 | 50.2/50.5 | 54.9/54.5 |
| Speed | 1030 | 4057 | 6668 | 5454/2702 | 17454/8377 |

choose SlowOnly (Feichtenhofer et al., 2019) to generate the spatiotemporal features due to its good trade-off between accuracy and efficiency.

## 2.2. One-stage Temporal Action Detection

One-stage TAD methods aim to detect the boundaries and categories of action segments in a single shot. The existing one-stage TAD methods can be divided into anchor-based ones (Wang et al., 2021a; Lin et al., 2017b; Buch et al., 2017; Long et al., 2019; Yang et al., 2020) and anchor-free ones (Yang et al., 2020; Zhang et al., 2022). Most existing methods are anchor-based. For example, Lin et al. (2017b) presented the first one-stage TAD method using convolutional networks. (Long et al., 2019) proposed to use Gaussian kernels to optimize the scale of each anchor dynamically. Meanwhile, Wang et al. (2021a) explored the pipeline of RetinaNet (Lin et al., 2017a) in the TAD task with an RGB-only stream. Some works explore the application of anchor-free methods. For instance, Yang et al. (2020) explored the combination of anchor-based and anchor-free methods. Zhang et al. (2022) proposed to use a local transformer encoder as a neck to enhance the video features for TAD.

Our work shares the advantage of simplicity with these one-stage TAD methods by focusing on designing an end-to-end TAD baseline. However, our BasicTAD presents a modular design for easy systematic study over the entire TAD pipeline. Based on this pipeline, we perform a more extensive study on the entire pipeline's components and figure out a simpler yet must-known TAD baseline (BasicTAD). In addition, we further empower our BasicTAD by keeping complete temporal and spatial information to yield our final TAD method of PlusTAD.

## 2.3. Multi-stage Temporal Action Detection

Multi-stage TAD methods often involve multiple stages to generate and refine action detection results. These methods might focus on different aspects to obtain better detection results. Most of them (Lin et al., 2018, 2019b; Su et al., 2021; Xu et al., 2020; Bai et al., 2020; Chen et al., 2022) focus on improving the quality of generated proposals, while a few others (Liu et al., 2021a; Zhao et al., 2020) try to improve the quality of classification results. (Lin et al., 2018, 2019b; Su et al., 2021; Xu et al., 2020; Bai et al., 2020; Chen et al., 2022; Gao et al., 2017, 2020; Chao et al., 2018; Xu et al., 2017) generated candidate action proposals at first, which is called temporal action proposal generation, and then further classified them into action categories possibly with the global classification results (e.g., Untrimmed-Net (Wang et al., 2017)). For proposal generation, Lin et al. (2018, 2019b); Su et al. (2021); Xu et al. (2020); Bai et al. (2020); Chen et al. (2022) were boundary-based methods that predict each frame's start and end confidence and then match start and end frames to generate the proposals with confidence evaluation. (Gao et al., 2017, 2020; Chao et al., 2018; Xu et al., 2017) generated proposals based on pre-defined sliding window anchors and trained a classifier to filter anchors. (Lin et al., 2021) designed a saliency-based refinement module to refine the detection results generated by a one-stage detector. Wu et al. (2021); Liu et al. (2021a); Tan et al. (2021) adopted different query-based dynamic networks along with multi-stage refinement modules to generate a direct sparse action proposal, effectively removing the post-processing steps of NMS. To improve classification results, Zhao et al. (2020) proposed regressing another completeness score to complement the classification score. Liu et al. (2021a) designed a post-processing technique to refine the confidence score based on actionness regression.

Unlike these multi-stage TAD methods, our BasicTAD aims to provide a one-stage and end-to-end TAD baseline method without requiring multi-stage processing. This simpler design would allow us to focus on an extensive investigation of basic yet overlooked designs of the TAD pipeline. As a result, we obtain a much simpler yet must-known TAD method, which obtains significant improvement over these complex multi-stage TAD methods on the THUMOS14 benchmark.

## 2.4. Training Strategies for Temporal Action Detection

There are various training strategies for optimizing TAD frameworks. The entire training process of TAD is usually split into multiple independent steps to reduce the optimization difficulty. Mainstream training strategies can be divided into two types. The first type consists of two steps (Lin et al., 2018, 2019b; Su et al., 2021; Gao et al., 2020; Liu et al., 2021a; Tan et al., 2021). First, the backbone networks are pre-trained on the TAD or action recognition datasets and used to perform feature extraction in a sliding window manner. Then, a separate TAD head network is trained without fine-tuning over the backbone networks to generate action segment proposals or directly predict the action segment boundaries and categories. This multi-step training paradigm would fail to unleash the potential of end-to-end representation learning on TAD.

The second type is to train the whole TAD pipeline (i.e., backbone and detection head) from RGB frames or both RGB frames and optical flow in an end-to-end manner. Recently, several works (Lin et al., 2017b; Wu et al., 2021; Lin et al., 2021; Wang et al., 2021a; Liu et al., 2022a), adopted this training strategy. Among them, the concurrent work (Liu et al., 2022a) compared the difference between head-only training and end-to-end learning and explored different backbones and detection heads. It also tried to balance detection results and computing overhead. However, it ignored the detailed investigation of some basic components in the entire TAD pipeline, such as the data augmentation and the neck design. Meanwhile, our result is better than its performance on the THUMOS14 dataset thanks to our simpler design and more detailed empirical study.

Our work shares the same advantage of end-to-end training with these methods. Compared with these works, our work conducts more detailed studies on the overall pipeline of end-to-end TAD and includes in-depth investigation for each component. Our solution is simpler yet more effective, achieving much better performance than these methods. BasicTAD could be easily extended to PlusTAD and has a fast inference speed, meeting the requirement of real-time TAD. We hope our extensive study could encourage future research to focus on designing an end-to-end TAD pipeline.

## 3. Methodology

In this section, we reconsider the TAD pipeline design by focusing on simplicity, efficient training, and eliminating complex pre-processing steps. We present a modular TAD framework consisting of four key components: data sampling (augmentation), backbone design, neck construction, and detection head. These components can be seamlessly integrated to yield a simple and efficient TAD framework. For each component, we explore the basic options to determine the optimal configuration, resulting in the **BasicTAD** framework. Additionally, we make minimal modifications to the BasicTAD framework by proposing three new design principles and creating an improved version of **PlusTAD**. Thanks to the straightforward design of both BasicTAD and PlusTAD, both frameworks enjoy end-to-end training and fully unleash the power of representation learning for the TAD task.

### 3.1. BasicTAD

To establish a straightforward and universal pipeline for facilitating the analysis and development of TAD methods, we break down the TAD pipeline into four fundamental components according to their functions: data augmentation, backbone design, neck construction, and one-stage detection head design. In general, these components are in analogy with the design in the common object detectors (Tian et al., 2019; Lin et al., 2017a). An overview of the modular TAD framework is illustrated in Fig 1, and we will delve into each component in detail in the following subsections.

#### 3.1.1. Data Sampling and Augmentation

As the TAD task involves temporal localization and prediction, the detection precision greatly relies on the robustness of the extracted feature sequence. The temporal sampling method is a crucial factor that affects feature extraction. The most basic sampling method is dense sampling with a fixed sampling rate. However, due to the limited memory of GPUs, it is not feasible to load an untrimmed video in its entirety. Therefore, we employ a sliding window strategy to divide the untrimmed video into overlapping clips. Within each video clip, we utilize the simple dense sampling method to downsample the original frame sequence with a fixed frame rate per second (FPS). This hyper-parameter FPS must be tuned to balance the detection accuracy and memory consumption afforded by the common GPU devices. We carefully tune these basic sampling options in our BasicTAD pipeline to demonstrate their influence on the final detection performance, which has been largely ignored in the existing TAD methods.

Another vital factor in feature extraction is data augmentation. Drawing on the success of multiple image-level data augmentation techniques in action recognition, we incorporate a similar data augmentation strategy in the first step of BasicTAD. Random cropping and horizontal flipping have been commonly utilized in existing end-to-end work, and we use the two methods by default. Furthermore, since video data may contain objects with various poses and scenes with different brightness levels, we adopt more image-based data augmentation techniques in our BasicTAD, including photo distortion and random rotation. In subsequent experiments, we employ all four data augmentation methods mentioned
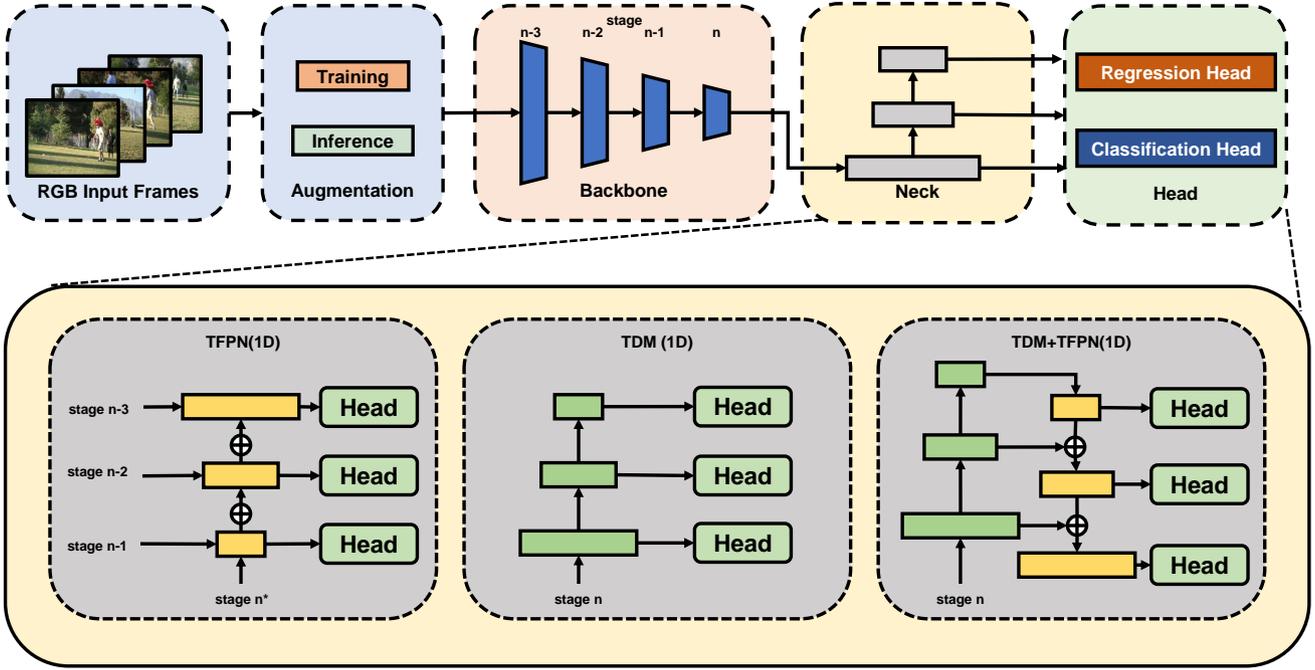
**Figure 1: BasicTAD Pipeline.** Our BasicTAD exhibits a modular design framework for the TAD task, composed of data sampling (augmentation), backbone, neck, and detection head. (a) The backbone is a spatial-temporal network to extract temporal features. (b) The neck module of BasicTAD offers three different ways of construction to leverage the extracted features. In the neck module, $n$ means the number of stages in the backbone, and stage $n$ is the last stage of the backbone. Stage $n*$ contains extra pooling operations on stage $n$ to construct a temporal feature with low temporal resolution. Particularly, $\oplus$ represents adding two features. As the temporal dimension of features varies across different scales, we interpolate the high-level features along the temporal dimension to align them with the low-level features before addition. (c) We adopt a typical one-stage network as the detection head, implemented by anchor-based and anchor-free methods.

above. It is worth noting that although these augmentation techniques are commonplace in images, their effects on TAD tasks have yet to be explored in previous works.

### 3.1.2. Backbone

The second component in BasicTAD is the backbone which is responsible for extracting spatiotemporal features from videos using a multi-layer network. However, the temporal modeling of early backbones only learns a simple consensus of RGB scene changes, making it suffering to capture temporal dynamics. Therefore, optical flow is naturally introduced in the previous TAD methods to explicitly supplement the short-form motion information. Nowadays, the latest designs of 3D backbones like SlowOnly can implicitly model such temporal dynamics of RGB input by learning to attain good performance on Something-something V2 (Materzynska et al., 2020). Hence, we believe that RGB input is enough, and we adopt an off-the-shelf 3D action recognition backbone to encode temporal dynamics for RGB-only input.

To balance the computing overhead and detection precision, we choose widely used SlowOnly (Feichtenhofer et al., 2019) as the backbone of BasicTAD by default. For the early TAD work, 8× spatial-temporal downsampling operation is applied in the backbones such as C3D (Tran et al., 2015) and I3D (Carreira and Zisserman, 2017) to reduce the temporal

resolution for saving computational overhead and increasing the receptive field to capture actions at larger scales. Since SlowOnly does not contain 8× temporal downsampling, we insert three 2× temporal downsampling operations to align the settings and fairly compare with these methods. We will perform ablation studies to explore the impact of the temporal downsampling operation. We also perform comparative studies over the different choices of video backbones on the final TAD performance.

### 3.1.3. Neck

The neck module, located after the backbone and before the detection head, plays a critical role in the TAD task by facilitating the alignment between the video features and downstream detection tasks. Its objective is to construct multi-resolution representations that can flexibly handle the vast variety of temporal durations of action instances. In this section, we introduce three neck modules and consider them candidates for the neck in BasicTAD.

After obtaining the extracted spatiotemporal features $F_{3d} \in R^{C \times T \times H \times W}$, we squeeze the H and W dimensions using spatial average pooling to get the temporal feature $F_{1d} \in R^{C \times T}$. We adopt multi-scale temporal features for our subsequent one-stage detection heads, helping the one-stage methods (anchor-based or anchor-free) better detect action segments of many scales. We use two basic feature pyramid

networks, *i.e.*, Temporal Feature Pyramid Network (TFPN) and Temporal Down Network (TDM) to build necks for creating or enhancing multi-scale temporal features from the backbone. TFPN enhances multi-scale temporal features by performing up-sampling and lateral addition on given multi-scale temporal features. Formally, the operation of each layer in TFPN can be represented as follows:

$$F_i^{'} = \text{Conv}(F_i + \text{Upsample}(F_{i+1})), \tag{1}$$

where $F_i$ represents the features in the $i$-th layer of the temporal feature pyramids, and $F_i^{'}$ denotes the new features that aggregate both high-level features $F_{i+1}$ and current scale features $F_i$. TDM utilizes multiple down-sampling operations to create multi-scale temporal features from a single-scale feature. Given a single-scale temporal feature sequence as the first-layer feature of feature pyramids, each layer of TDM can be expressed as follows:

$$F_{i+1} = \text{Downsample}(F_i) \tag{2}$$

where Downsample represents a max-pool or convolution layer, aggregating the temporal semantic information of the features and halving the temporal dimension. Based on both basic networks, as shown in Fig 1, we construct three necks: Lateral-TFPN, Post-TDM, and Post-TDM-TFPN.

**Lateral-TFPN.** One of the most intuitive approaches to building a TFPN is by leveraging multi-scale spatiotemporal features extracted from different stages of the backbone. This approach facilitates the integration of high-level semantic information with low-level features, much like what is accomplished by FPN in object detection. Although shallow layers offer a greater temporal resolution, the spatial features at each temporal location may need to capture more rich temporal semantic information, which could hinder the effectiveness of the TAD task.

**Post-TDM.** TDM generates multi-scale temporal features on pre-extracted temporal features. While the last layer of backbone features provides rich semantic information at each temporal location, increasing the number of feature layers may prevent the learned high-level semantics from being fed back to low-level features. TDM, on the other hand, utilizes features from the backbone directly and is, therefore, easier to optimize compared to TFPN.

**Post-TDM-TFPN.** Constructing TDM and TFPN after the backbone can simultaneously address the abovementioned problems. TFPN integrates high-level semantic information into low-level features, which could improve the network's capacity to aggregate temporal context while ensuring adequate spatial semantics. However, the introduction of additional learnable parameters could impede network optimization.

In the subsequent experiments, we study and analyze the different choices of the three necks and provide our final combination in BasicTAD.

### 3.1.4. Head

The detection head is the final component in our modular TAD framework. It is responsible for completing the detection task by generating the temporal interval of the action instance and its corresponding label. Typically, this component includes sub-networks designed for classification and regression tasks, respectively. Moreover, the specific sample assignment is critical for training these detection heads effectively. Both aspects complement each other to produce the best possible results in the final detection performance.

For simplicity and end-to-end training, we adopt anchor-based and anchor-free mechanisms as the basic detection head of our BasicTAD, by following the basic design principle in object detection (Ren et al., 2015; Lin et al., 2017a; Tian et al., 2019). Since these two detection methods have advantages and disadvantages, there is no clear conclusion on which is better. We briefly introduce two detection methods in our TAD pipeline and provide specific implementation details. Both methods share the same sub-networks composed of four temporal convolutional layers followed by a normalization and activation layer. Both methods share two classification and regression sub-networks composed of four temporal convolutional layers followed by a normalization and activation layer. The regression and classification branches for both methods can be formulated as follows:

$$F_i = \text{Conv}(F_{i-1}), \tag{3}$$

where $i \in [1, 4]$ is the output features of the $i$-th convolution layer and $F_0$ represents input features of head. In the last layer, the output of the classification branch is represented as $F_4^{cls} \in R^{N_a \times N_c \times T}$, where $N_a$ represents the number of anchors ($N_a = 1$ for anchor-free method) and $N_c$ represents the number of categories. The output of the regression branch denotes $F_4^{reg} \in R^{N_a \times 2 \times T}$, where 2 represents the variables related to boundaries. Since the sub-network design of both methods is identical, we will primarily focus on introducing the two methods from the perspective of the sample assignment mechanism.

**Anchor-based Method.** Anchor-based methods generate temporal proposals by assigning dense and multi-scale intervals with pre-defined lengths to uniformly distributed temporal locations in the input video. We use translation-invariant anchors, which have an increasing temporal size from the bottom to the top of the feature pyramid network. At each level, we add anchors of 5 sizes $\{2^0, 2^{1/5}, 2^{2/5}, 2^{3/5}, 2^{4/5}\}$ of the original set of default anchors for dense scale coverage. Anchors are assigned to ground-truth action segments using the temporal Intersection-over-Union (tIoU) threshold of 0.6 and a background with a tIoU lower than 0.4. Other anchors overlapping [0.4, 0.6) will be eliminated during training. We obtain the predicted action boundaries by optimizing the relative offset between anchors and ground truths. Dense anchor matching is required at each temporal position of the feature for anchor-based methods, which need rich semantic context information.

**Anchor-free Method.** Anchor-free methods directly regress the offsets to action boundaries at each temporal location and then use these offsets to generate temporal proposals. We first compute the regression offsets for each location on all feature levels. Any location which falls into any ground-truth box will be set as a positive sample. The others are negative samples. Each level is responsible for a range of motion detection. If an action proposal at one temporal location is beyond this range, this location will be ignored. We define $m_i$ as the maximum range border that the feature level $i$ needs to regress. In this work, $m_2, m_3, m_4, m_5, m_6, m_7$ are set as $-1, 5, 10, 20, 40$ and $\infty$, respectively. Even with multi-level prediction, if a location is still assigned to more than one ground-truth box, we choose the ground-truth box with minimal area as its target. Subsequently, we employ $e^{s_i x}$ with a trainable scalar $s_i$ to automatically adjust the base of the exponential function for feature level $i$. This approach enables us to obtain the predicted actions' boundaries accurately. Compared with anchor-based methods, anchor-free methods have a more flexible matching mechanism.

**Post Processing.** We also apply post-processing to suppress redundant predictions to yield the final detection results for both anchor-free and anchor-based methods. Specifically, we choose two suppression algorithms: Non-Maximum Suppression (NMS) (Neubeck and Gool, 2006) and Non-Maximum Weighting (NMW) (Ning et al., 2017) for both anchor-based and anchor-free methods. NMS is a crucial technique that ensures the algorithm produces only one detection per object. It selects the proposal with the highest confidence in each iteration, chooses all remaining proposals with a high overlap rate with the current proposal, and then suppresses them. After that, it saves the currently selected proposal and proceeds to the next iteration, excluding the selected proposal. This procedure continues until all proposals have been processed. NMW is an improved version of NMS since proposals with the highest confidence scores may not be accurately positioned, while other well-located proposals may exist. NMW uses the confidence score and Intersection over Union (IoU) to calculate a weighted average of all proposal coordinates of the same type. We will conduct ablation studies on them in the later section.

The above detection methods and post-processing mechanisms are our optional basic components in head. Validation and more in-depth studies on these components will be shown in ablation studies.

## 3.2. From BasicTAD to PlusTAD

Based on our modular TAD framework, we perform comprehensive studies of the basic options and come up with a simple yet effective TAD baseline, termed BasicTAD. Furthermore, based on the BasicTAD, we introduce three customized improvements to fully unleash the power of this simple and end-to-end detection pipeline, and the upgraded framework is called **PlusTAD**. Specifically, we propose two network structure designs, Temporal Preservation and Spatial Preservation, to improve the quality of the temporal features of networks' backbone and neck parts. We further study and adopt stronger data augmentations for training and multi-view ensemble methods in the testing phase.

### 3.2.1. Temporal Preservation for Backbone

The existing TAD methods use backbones that downsample temporal and spatial information and require dense frames sampled at high FPS as input.

Intuitively, denser inputs generally lead to better detection performance. To alleviate the huge computational overhead imposed by dense inputs, the existing practice is to perform temporal downsampling on the inputs in the early stages of the backbone. However, our subsequent experiments show that this approach does not always yield gains. Due to insufficient discriminative power between adjacent frames in shallow layers, it is hard for the early stages of the backbone to learn how to extract useful temporal signals from these subtle frame changes. Therefore, we propose **Temporal Preservation** (termed as "TP" for short) design that feeds frames into the model from the beginning at an explicitly sampled frame rate (sparse or dense) without downsampling in the backbone. For long actions in datasets like THUMOS14, sparse input without temporal downsampling in the backbone is enough to predict actions. Though sparse input may not be suitable for fine-grained action detection, maintaining the temporal dimension and adopting dense input can encourage the backbone to capture their subtle changes. TP design is effective and helpful for both scenarios above.

### 3.2.2. Spatial Preservation for Neck

Many existing TAD works (Lin et al., 2021; Wang et al., 2021a) squeeze the spatial dimension of features before feeding them into the neck module. We argue that the squeezing operation prematurely drops the spatial dimension, making the feature pyramid constructed by the neck module unable to capture spatial-sensitive multi-scale features.

In our further experiments, we find it is worth preserving spatial dimension. So we propose a design named **Spatial Preservation** (termed as "SP" for short) that postpones the spatial squeezing until after the neck and replaces the 1D operators with 3D operators. It introduces the local spatial context and enhances the robustness of multi-scale spatiotemporal information. We discuss its performance improvement in the experimental section.

### 3.2.3. Enhanced Data Augmentation

In Section 3.1.1, we equip BasicTAD with various image data augmentations in the training phase. To further improve the performance of our PlusTAD pipeline, we explore more data augmentation methods. Through ablation experiments, we investigate the effects of more different spatial resolutions on TAD performance in the training stage.

Furthermore, considering the data augmentation in the testing phase mainly works on the temporal and spatial levels, we adopt various test augmentation methods to enhance the robustness of the predictions. In detail, we employ two spatial-level data augmentation methods for each temporal

window, namely "ThreeCrop" and "Flip", individually or simultaneously. We also use the reverse sliding window, termed "Backward", to increase the density of temporal windows and combine the predictions of each temporal window for post-processing.

## 3.3. Training and Testing
### 3.3.1. Training

In the training phase, we randomly sample a fix-sized temporal window of consecutive frames in each untrimmed video per iteration and feed them into our model. We train BasicTAD and PlusTAD with both heads as a multi-task loss function $L$, including a classification loss $L_{cls}$, a regression loss $L_{reg}$.

$$L = L_{cls} + \alpha L_{reg}, \qquad (4)$$

where $\alpha$ is a hyper-parameter to balance these two terms and set them to 1. classification loss: We use focal loss (Lin et al., 2017c) as classification loss for its ability to solve the problem of positive and negative sample imbalance. The formal expression of focal loss is as follows:

$$L_{cls} = -(1 - p_t)^\gamma log(p_t) \qquad (5)$$

where $p_t$ represents the probability of each action category after softmax, and $\gamma$ represents the modulation factor that focuses on hard samples. Meanwhile, Distance-IoU (DIoU) loss (Zheng et al., 2020) is adopted as regression loss for faster convergence during training and more accurate boundary regression. The formal expression of DIoU loss is as follows:

$$L_{reg} = 1 - \left| \frac{B \cap B^{gt}}{B \cup B^{gt}} \right| + \frac{\rho^2(b, b^{gt})}{c^2} \qquad (6)$$

where $B$ represents the boundary of the proposed action, $B^{gt}$ represents the boundary of ground truth action, $\rho^2(b, b^{gt})$ represents the square of the Euclidean distance between the center points of $B$ and $B^{gt}$ and $c$ represents the shortest length that simultaneously covers $B$ and $B^{gt}$.

### 3.3.2. Testing

We use sliding windows with overlap to sample fixed-number frames and feed them into BasicTAD and PlusTAD in the testing stage. Our model outputs $\{(s_i, e_i, p_i)\}_{i=1}^m$ as the predicted action set, where $i$ and $m$ is the $i$-th action and total number of predicted action. $s_i$, $e_i$ represent the start and end time of the $i$-th action $i$ and $p_i$ is its score. We adopt NMW (Ning et al., 2017) for both heads to remove the redundant action segments.

## 4. Experiments
## 4.1. Datasets

We perform extensive experiments on three datasets, summarized in Fig 2 and Table 2, to demonstrate the effectiveness of our BasicTAD and PlusTAD. **THUMOS14** (Jiang et al., 2014) is a commonly-used dataset in TAD, containing

**Table 2**
**Summary of three datasets.** Video number, category number and other information about THUMOS14, FineAction and ActivityNet-v1.3 are listed in this table.

| Dataset | Video | Category | Instance | Duration | Type |
|---|---|---|---|---|---|
| THUMOS14 | 413 | 20 | 6,316 | 4.3 s | sports |
| FineAction | 16,732 | 106 | 103,324 | 7.1 s | daily events |
| ActivityNet-v1.3 | 19,994 | 200 | 23,064 | 49.2 s | daily events |



**Figure 2: Samples of three TAD datasets.** (a): A few examples in THUMOS14. All actions are sports-related. (b): A few examples in FineAction. These samples' categories are shown in the form of "top-level category"-"bottom-level category". Its action scenes are more varied than THUMOS14. (c): A few examples in ActivityNet-v1.3. It has plenty of action scenes. Unlike FineAction, actions in ActivityNet-v1.3 are mostly long-term actions.

200 validation videos and 212 test videos with labeled temporal annotations from 20 action categories in sports. **FineAction** (Liu et al., 2022b) is a newly collected large-scale fine-grained TAD dataset containing 57,752 training instances from 8,440 videos and 24,236 validation instances from 4,174 videos and 21,336 testing instances from 4,118 videos. It contains 106 action categories within a new taxonomy of three-level granularity. This taxonomy consists of 4 top-level categories, 17 middle-level categories, and 106 bottom-level categories. The 4 top-level categories are "Household Activities", "Personal Care", "Socializing, Relaxing" and "Sports, Exercise". With richer action scenes and finer-grained action categories, FineAction will be a challenging data set in the TAD task. **ActivityNet-v1.3** (Heilbron et al., 2015) is a large-scale dataset containing 10,024 training videos, 4,926 validation videos, and 5,044 test videos belonging to 200 activities covering sports, household, and working actions. Different from THUMOS14 and FineAction, which consist mostly of short actions, most videos in ActivityNet-v1.3 contain only one long-term action instance, and the action frames exceed 64% of all frames.

## 4.2. Evaluation Metric

Following previous work, we report the **mean average precision (mAP)** with tIoU thresholds [0.3 : 0.1 : 0.7] on the test set of THUMOS14. And on the validation set

of FineAction and ActivityNet-v1.3, the thresholds are set as $[0.5, 0.05, 0.95]$. We use "Avg" to represent the average mAP on THUMOS14, FineAction, and ActivityNet-v1.3.

## 4.3. Implementation Details

We perform experiments using BasicTAD and PlusTAD on THUMOS14 (Jiang et al., 2014). Though THUMOS14 is a scene-biased dataset that can not reflect the importance of temporal dynamics in classification, it can provide temporal dynamics for the regression of action boundaries. This means that the conclusions obtained from the ablation experiments on THUMOS14 are universal to a certain extent, so we choose THUMOS14 to do ablation studies. We sample temporal windows of 32 seconds for both configurations, covering over 99.7% action instances. We use 768 frames at 24FPS for BasicTAD and 96 frames at 3FPS for PlusTAD. In the training phase, we resize the **original size** (the size of the original frames) to $128 \times 128$ for BasicTAD and short-128 (the short side of the frame is set to 128) for PlusTAD. We randomly sample a temporal window in each untrimmed video per iteration. We set the **crop size** (the size of the cropped images) to $112 \times 112$. The photo distortion and random spatial rotation are consistent with (Wang et al., 2021a; Liu et al., 2016). In the testing phase, we sample uniform sliding windows whose sampling stride between adjacent sliding windows is set to 25% of the window length. And the center crop of $112 \times 112$ is used. If there is no special emphasis, the above configuration will be the default configuration for subsequent experiments based on THUMOS14.

We use PlusTAD to perform further experiments on FineAction (Liu et al., 2022b) and ActivityNet-v1.3 (Heilbron et al., 2015). On FineAction, we sample RGB frames at 2FPS. In the training phase, we set the original size to short-256 and the crop size to $224 \times 224$. In the testing phase, we use the center crop of $224 \times 224$. On ActivityNet-v1.3, we follow the settings in AFSD (Lin et al., 2021) that we sample frames at different fps and ensure the number of frames in each video is the same. We set the original size to $224 \times 224$.

We use SlowOnly (Feichtenhofer et al., 2019) pretrained on Kinetics (Kay et al., 2017) as our backbone by default, where all batch normalization layers are frozen. We train the model using SGD with a momentum of 0.9 and weight decay of 0.0001. The batch size is set as 16. The learning rate schedule is annealing down from 0.01 to 0.0001 every 1200 iterations on THUMOS14, 20 epochs on FineAction, and ActivityNet-v1.3 using the cosine decay rule.

## 4.4. Ablation Studies on BasicTAD

We first begin the ablation study on the design of each component in our proposed modular TAD framework. In the previous section, we introduced the basic and simple options for each step in our modular framework. In this section, we will perform extensive studies over these basic designs through in-depth and step-by-step ablation experiments. Our goal is to discover a simple yet must-known TAD baseline.

### 4.4.1. Study on different backbones

The first and most important choice in our modular TAD framework is the backbone design. The previous works (Wang et al., 2021a; Xu et al., 2017; Liu and Wang, 2020; Lin et al., 2021; Wu et al., 2021) with end-to-end training manner usually input many frames and use backbones with 8× temporal downsampling, such as C3D (Tran et al., 2015), I3D (Carreira and Zisserman, 2017), and R50-I3D (Wang et al., 2018b). This method can increase the utilization of the data source as much as possible. Meanwhile, using 8× spatiotemporal downsampling can reduce the computational overhead. However, these backbones are relatively outdated and inferior to the recent SlowOnly backbone (Feichtenhofer et al., 2019). To fully unleash the power of our proposed modular TAD framework with end-to-end training, we also try SlowOnly backbone in our TAD framework. For a fair comparison with previous backbones, we align with the settings of these previous works, and 8× temporal downsampling is adopted. Hence, we insert a 2× downsampling layer before res-2, res-3, and res-5(res-x denotes the x-th res-stage in SlowOnly) of the SlowOnly backbone, respectively.

Inputting **768 frames at 24FPS** under end-to-end training, we compare the performance of BasicTAD, which uses C3D, I3D, R50-I3D, and SlowOnly with 8× downsampling in Table 3. These four backbone encoders are all 3D CNN methods that capture spatiotemporal information between frames by performing 3D convolution. As shown in Table 3, BasicTAD with C3D does not perform well due to its limited representation power caused by relatively shallow networks. Meanwhile, I3D performs worse than SlowOnly with 8× downsampling. This can imply that the ResNet50 (He et al., 2016) network outperforms the Inception-v1 (Ioffe and Szegedy, 2015) network in the TAD task. R50-I3D is slightly weaker than SlowOnly with 8× downsampling because they both are based on ResNet. R50-I3D performs 8× downsampling at the beginning of the backbone, while SlowOnly with 8× downsampling delays the timing of backbone downsampling. This difference will contribute to the slight performance difference. This downsampling location will be analyzed in the next section.

### 4.4.2. Study on the gain of end-to-end training

In this section, we study the gain of end-to-end training BasicTAD with different backbones. Table 4 shows that abandoning the end-to-end training strategy leads to significantly worse detection results for the anchor-based method across all four backbones. While this method reduces training costs, it fails to fully leverage the modeling capabilities of the backbone, given the gap between TAD and action recognition tasks. This underscores the importance of a trainable backbone in the TAD pipeline.

### 4.4.3. Study on Downsampling Locations in backbone

Following the previous section, we further study downsampling locations in the SlowOnly backbone. To construct a structure with 8× temporal downsampling, the backbone

**Table 3**
**Performance of BasicTAD on different backbones on THUMOS14 (Jiang et al., 2014). We replace SlowOnly with C3D (Tran et al., 2015), I3D (Carreira and Zisserman, 2017) and R50-I3D (Wang et al., 2018b) for our BasicTAD to perform experiments, where I3D only retains RGB branch.**

| Method | Backbone | mAP@0.3 | mAP@0.4 | mAP@0.5 | mAP@0.6 | mAP@0.7 | Avg |
|---|---|---|---|---|---|---|---|
| Anchor-based | C3D | 54.4 | 50.8 | 45.7 | 37.4 | 26.1 | 42.9 |
| | I3D | 59.5 | 56.0 | 51.4 | 41.8 | 28.3 | 47.4 |
| | R50-I3D | 62.8 | 59.5 | 53.8 | 43.6 | 30.1 | 50.0 |
| | SlowOnly (×8) | **63.1** | **59.5** | **54.3** | **43.6** | **30.5** | **50.2** |
| Anchor-free | C3D | 56.8 | 50.1 | 44.7 | 35.2 | 24.3 | 42.2 |
| | I3D | 61.7 | 56.9 | 49.3 | 38.7 | 26.1 | 46.6 |
| | R50-I3D | **63.7** | 58.5 | 51.6 | 41.0 | 30.8 | 49.1 |
| | SlowOnly (×8) | 63.2 | **59.7** | **52.6** | **44.5** | **32.7** | **50.5** |

**Table 4**
**Study on the effectiveness of end-to-end training based on the anchor-based BasicTAD with different backbones on THUMOS14 (Jiang et al., 2014). "e2e" is short for end-to-end training. We freeze all layers in the backbone to construct a non-end-to-end training manner.**

| Backbone | e2e | mAP@0.3 | mAP@0.4 | mAP@0.5 | mAP@0.6 | mAP@0.7 | Avg |
|---|---|---|---|---|---|---|---|
| C3D | ✔ | **54.4** | **50.8** | **45.7** | **37.4** | **26.1** | **42.9** |
| | ✘ | 32.0 | 27.2 | 22.0 | 14.4 | 7.6 | 20.6 |
| I3D | ✔ | **59.5** | **56.0** | **51.4** | **41.8** | **28.3** | **47.4** |
| | ✘ | 35.2 | 29.7 | 23.1 | 15.9 | 8.1 | 22.4 |
| R50-I3D | ✔ | **62.8** | **59.5** | **53.8** | **43.6** | **30.1** | **50.0** |
| | ✘ | 36.5 | 31.8 | 26.7 | 19.6 | 12.3 | 25.4 |
| SlowOnly (×8) | ✔ | **63.1** | **59.5** | **54.3** | **43.6** | **30.5** | **50.2** |
| | ✘ | 37.3 | 32.4 | 26.6 | 19.8 | 13.0 | 25.8 |

of BasicTAD need to choose three different places to operate 2× downsampling. Since R50-based SlowOnly has five stages, we can choose three of four intervals between these res-stages. In this sense, we can downsample the spatiotemporal features before res-2, res-3, res-4, and res-5 (res-$x$ denotes the $x$-th res-stage in SlowOnly). Inserting downsampling layers at different locations can lead to differences in model performance due to different temporal receptive fields and aggregations. Meanwhile, their corresponding computing costs will be different as well.

Following the input in the previous section, we compare the performance of BasicTAD with different downsampling locations. The results are shown in Table 5. In the backbone's first half stages (res-2 and res-3), the features are highly redundant in the temporal dimension. If we put the temporal downsampling layers at the first three stages, it can save the computational cost most but also achieves the worst detection mAP. However, if we simply change a single downsampling location from res-4 to res-5, it leads to a large performance improvement of around 5% mAP.

**Table 5**
**Comparison between different locations of downsampling in BasicTAD on THUMOS14 (Jiang et al., 2014). We choose three of four intervals between these res-layers to insert the max-pool downsampling operation. In this table, choosing an interval before a res-layer will be indicated by a tick.**

| Method | res-2 | res-3 | res-4 | res-5 | FLOPs | mAP@0.3 | mAP@0.4 | mAP@0.5 | mAP@0.6 | mAP@0.7 | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Anchor-based | ✓ | ✓ | ✓ | | 227.7G | 55.9 | 52.5 | 47.3 | 40.3 | 29.2 | 45.1 |
| | ✓ | ✓ | | ✓ | 280.2G | **63.1** | **59.5** | **54.3** | **43.6** | **30.5** | **50.2** |
| | ✓ | | ✓ | ✓ | 329.9G | 55.3 | 51.5 | 46.8 | 38.0 | 26.9 | 43.7 |
| | | ✓ | ✓ | ✓ | 395.4G | 56.5 | 53.1 | 48.2 | 40.1 | 29.0 | 45.4 |
| Anchor-free | ✓ | ✓ | ✓ | | 245.8G | 59.2 | 55.0 | 47.5 | 37.2 | 27.3 | 45.3 |
| | ✓ | ✓ | | ✓ | 298.4G | 63.2 | **59.7** | **52.6** | **44.5** | **32.7** | **50.5** |
| | ✓ | | ✓ | ✓ | 348.1G | 63.8 | 58.2 | 50.4 | 41.8 | 30.1 | 48.9 |
| | | ✓ | ✓ | ✓ | 413.5G | **64.6** | 58.8 | 51.4 | 42.6 | 31.2 | 49.7 |

**Table 6**
**Comparison between three kinds of the neck on THUMOS14 (**Jiang et al., 2014**).** We test three different kinds of neck and compare the results on THUMOS14. We denote the down-sampling operator of Lateral-TFPN as "-" because it uses the multi-scale feature from the backbone.

| Method | Neck | Operator | mAP@0.3 | mAP@0.4 | mAP@0.5 | mAP@0.6 | mAP@0.7 | Avg |
|---|---|---|---|---|---|---|---|---|
| | Lateral-TFPN | - | 59.1 | 55.4 | 50.0 | 40.9 | 27.9 | 46.7 |
| | Post-TDM | Conv | 56.7 | 53.1 | 48.0 | 38.8 | 27.0 | 44.7 |
| Anchor-based | Post-TDM | Maxpool | 54.1 | 52.0 | 46.3 | 37.9 | 27.8 | 43.4 |
| | Post-TDM-TFPN | Conv | **63.1** | **59.5** | **54.3** | **43.6** | 30.5 | **50.2** |
| | Post-TDM-TFPN | Maxpool | 58.7 | 55.2 | 50.7 | 42.0 | **30.6** | 47.5 |
| | Lateral-TFPN | - | 56.1 | 50.2 | 43.6 | 34.1 | 24.5 | 41.7 |
| | Post-TDM | Conv | **63.9** | 58.7 | 50.2 | 40.0 | 29.1 | 48.4 |
| Anchor-free | Post-TDM | Maxpool | 63.2 | **59.7** | **52.6** | **44.5** | **32.7** | **50.5** |
| | Post-TDM-TFPN | Conv | 62.3 | 56.5 | 49.1 | 38.3 | 24.6 | 46.2 |
| | Post-TDM-TFPN | Maxpool | 61.8 | 56.6 | 49.2 | 39.9 | 28.7 | 47.2 |

Meanwhile, we notice that other configurations of down-sampling locations achieve a weaker performance than the previous version yet with higher computational costs. From the above results, it can be found that the temporal features are highly redundant in the first half stages, and the latter half stages are key stages of feature encoding, so temporal downsampling cannot be easily performed there. Thus, we keep the temporal downsampling locations before res-2, res-3, and res-5 by default.

### 4.4.4. Study on Neck Design

After the ablation studies on the backbone design, we now turn to the exploration of the neck design. Based on the best result in Table 5, we set the down-sampling locations at res-2, res-3, and res-5 in this ablation study. Specifically, we compare three kinds of neck modules for anchor-based and anchor-free TAD methods, namely Lateral-TFPN, Post-TDM, and Post-TDM-TFPN, as introduced above. These neck modules with different down-sampling operators yield multi-resolution representations for action detection. The results are listed in Table 6. As shown in the table, Lateral-TFPN achieves the worst performance because the feature maps lack enough high-level semantic information in the early stages of the backbone. In order to keep high-level semantic information in the neck module, it is necessary to build the neck representation right behind the backbone.

For Post-TDM and Post-TDM-TFPN, we explore two basic operators in the downsampling procedure, namely convolution, and max-pool. As demonstrated in Table 6, Post-TDM with the max-pool operator and Post-TDM-TFPN with the convolution operator is more suitable for the anchor-free method and the anchor-based method, respectively. This difference is due to the different detection and training mechanisms of the anchor-free method and the anchor-based method. The anchor-based TAD method adopts the global sample assignment mechanism based on tIoU. It depends on the global multi-scale context aggregated by Post-TDM-TFPN with convolution operations. Instead, the anchor-free TAD method uses a sample assignment mechanism based on

the multi-scale center points, so it is subject to local high-frequent boundary information in the temporal dimension to directly regress the locations of action boundaries.

Based on the above results and analysis, we apply Post-TDM with the max-pool operator and Post-TDM-TFPN with the convolution operator to the anchor-free method and the anchor-based method in BasicTAD.

### 4.4.5. Study on Head Design

To ensure the simplicity principle in BasicTAD, we introduce basic one-stage anchor-based and anchor-free methods in our study. In the previous ablation experiments on backbone and neck design for BasicTAD, we conducted experiments with both detection heads. The results are shown in Table 3, Table 5, and Table 6. From these results, we can also provide some comparative analysis for head design.

As shown in Table 5, the anchor-based TAD method has slightly lower FLOPs than the anchor-free TAD method due to its smaller number of channels in the detection head. Although the max-pool operator does not contain any learnable parameters, the anchor-free method does not reduce the dimension of the last layer features from SlowOnly, which is 2048. It leads to more FLOPs compared with the anchor-based method. For detection accuracy, the anchor-free TAD method achieves a slightly better mAP than the anchor-based TAD method.

In our design, the anchor-based and anchor-free BasicTAD are both dense detectors. As dense prediction methods, anchor-based and anchor-free heads rely on post-processing to suppress redundant predictions to yield the final detection results. In our BasicTAD, we generally follow the box suppression methods in object detection. Specifically, various classical suppression algorithms, *i.e.*, Non-Maximum Suppression (NMS) (Neubeck and Gool, 2006), Non-Maximum Weighting (NMW) (Ning et al., 2017) are introduced to suppress bounding boxes in post-processing. We perform a comparative study on them to determine the best option for anchor-free and anchor-based TAD methods.

The experiment results on the choice of the post-processing algorithms are reported in Table 7. For the anchor-free method, the detection result is not sensitive to the choice of

**Table 7**
Comparison of Post-Processing strateg on THUMOS14 (Jiang et al., 2014). We take both Non-Maximum Suppression (NMS) and Non-Maximum Weighting (NMW) into account for both anchor-free and anchor-based methods. NMW performs better than NMS for its weighted average mechanism for redundant accurate boundaries.

| Method | Post Processing | mAP@0.3 | mAP@0.4 | mAP@0.5 | mAP@0.6 | mAP@0.7 | Avg |
|---|---|---|---|---|---|---|---|
| Anchor-based | NMS | 61.4 | 57.8 | 51.8 | 42.1 | **30.6** | 48.7 |
| | NMW | **63.1** | **59.5** | **54.3** | **43.6** | 30.5 | **50.2** |
| Anchor-free | NMS | 63.2 | 59.7 | 52.6 | **44.5** | 32.7 | 50.5 |
| | NMW | **64.7** | **60.0** | **52.8** | 43.7 | **33.0** | **50.8** |

**Table 8**
Ablation study on improvements of PlusTAD on THUMOS14 (Jiang et al., 2014). Temporal Preservation is for backbone design (indicated as "TP"), and Spatial Preservation is for neck design (indicated as "SP").

| Method | TP | SP | FPS | Frames | mAP@0.3 | mAP@0.4 | mAP@0.5 | mAP@0.6 | mAP@0.7 | Avg | FLOPs |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Anchor-based | | | 24 | 768 | 63.1 | 59.5 | 54.3 | 43.6 | 30.5 | 50.2 | 280.2G |
| | ✓ | | 3 | 96 | 63.9 | 59.4 | 53.7 | 44.3 | 31.0 | 50.5 | 133.3G |
| | ✓ | ✓ | 3 | 96 | 63.2 | 59.7 | 54.4 | 45.6 | 32.2 | 51.0 | 136.4G |
| | ✓ | | 6 | 192 | **65.4** | 61.8 | 56.4 | 47.9 | 33.8 | 53.1 | 266.0G |
| | ✓ | ✓ | 6 | 192 | 65.1 | **62.0** | **57.1** | **48.9** | **33.9** | **53.4** | 272.9G |
| Anchor-free | | | 24 | 768 | 64.7 | 60.0 | 52.8 | 43.7 | 33.0 | 50.8 | 298.4G |
| | ✓ | | 3 | 96 | 65.7 | 60.7 | 52.8 | 41.7 | 29.3 | 50.0 | 151.4G |
| | ✓ | ✓ | 3 | 96 | 67.3 | 61.6 | 54.4 | 41.2 | 29.5 | 50.8 | 151.5G |
| | ✓ | | 6 | 192 | 67.5 | 62.6 | 55.4 | **45.7** | **33.7** | 53.0 | 284.2G |
| | ✓ | ✓ | 6 | 192 | **69.9** | **64.3** | **56.8** | 44.7 | 32.1 | **53.6** | 284.2G |

the post-processing algorithm, and using NMW is slightly better. However, the results of the anchor-based method are quite different, where the choice of the post-processing algorithm significantly affects the detection results. One possible reason is that in TAD tasks, annotations tend to be very sparse, and the anchor-based method generates too many redundant predictions that we need an advanced suppression algorithm. NMW performs better on anchor-based and anchor-free methods, so we use NMW in the subsequent experiments.

### 4.5. Ablation Studies on PlusTAD

We have explored several basic settings of each module in our modular pipeline and come up with a very simple TAD baseline method, termed BasicTAD. These experiments are all based on the basic design described in Section 3.1. In this subsection, we further perform ablation experiments to validate the improvements proposed over the BasicTAD design. Specifically, we provide empirical results for the three improvements, namely temporal preservation of backbone design, spatial preservation of neck, and more training augmentations. With these three improvements, we obtain another more powerful TAD baseline method, termed as PlusTAD.

#### 4.5.1. Temporal Preservation in Backbone Design

Temporal sampling is an important but less studied factor influencing the design of the TAD method. We have presented a Temporal Preservation (TP) principle in our PlusTAD design as we need to keep as much temporal semantic information as possible for efficient localization of each action instance. We conduct the following experiments to study the effectiveness of TP for the backbone design. In detail, we directly remove 8× temporal downsampling in the backbone design. Thus the temporal size of feature maps keeps the same with the input frames for all stages in ResNet. In order to keep the final temporal feature the same as the BasicTAD, we use **96 frames at 3FPS** to replace **768 frames at 24FPS** in BasicTAD. We change the sampling FPS to keep the temporal window size the same.

As is shown in Table 8, compared with BasicTAD (the first row in the table), using a video clip of 96 frames at 3FPS achieves a similar performance for both anchor-based and anchor-free TAD methods. However, the computation overhead is only half of BasicTAD, indicating that our proposed TP principle in backbone design is a good practice for increasing TAD running speed without mAP loss. To fully unleash the modeling power of our TP-equipped backbone in TAD, we further increase the sampling FPS and sampling frame number to provide richer information. When using video clips of 192 frames at 6FPS, the average mAP of the anchor-based and the anchor-free methods is improved by 2.9 and 2.2 over the BasicTAD, respectively. But, their overall FLOPs are still lower than that of BasicTAD. This performance demonstrates that preserving the temporal resolution in the backbone design makes it more effective and efficient in capturing temporal information. Therefore, in our PlusTAD, we will employ the TP principle in our design by default.

**Table 9**
Comparison of the different spatial sizes in the training stage on THUMOS14 (Jiang et al., 2014). Original size means input frame size, and the crop size indicates the cropped patch resolution as the network input.

| Method | Original size | Crop size | mAP@0.3 | mAP@0.4 | mAP@0.5 | mAP@0.6 | mAP@0.7 | Avg | FLOPs |
|---|---|---|---|---|---|---|---|---|---|
| Anchor-based | $128 \times 128$ | $112 \times 112$ | 63.2 | 59.7 | 54.4 | 45.6 | 32.2 | 51.0 | 136.4G |
| | $171 \times 128$ | $112 \times 112$ | 66.8 | 63.3 | 57.8 | 48.8 | 32.9 | 53.9 | 136.4G |
| | short-128 | $112 \times 112$ | 68.4 | 65.0 | 58.6 | 49.2 | 33.5 | 54.9 | 136.4G |
| | short-180 | $160 \times 160$ | **71.2** | **66.8** | **61.2** | **50.1** | **36.3** | **57.1** | 262.1G |
| Anchor-free | $128 \times 128$ | $112 \times 112$ | 67.3 | 61.6 | 54.4 | 41.3 | 29.5 | 50.8 | 151.5G |
| | $171 \times 128$ | $112 \times 112$ | 68.1 | 64.3 | 55.2 | 45.3 | 31.9 | 53.0 | 151.5G |
| | short-128 | $112 \times 112$ | 70.4 | 65.5 | 57.6 | 46.0 | 33.2 | 54.5 | 151.5G |
| | short-180 | $160 \times 160$ | **72.5** | **66.8** | **59.1** | **48.4** | **35.0** | **56.4** | 275.9G |

**Table 10**
Ablation study results on testing augmentation on THUMOS14 (Jiang et al., 2014). For ThreeCrop and Flip augmentation, we fuse the features after the neck. For backward augmentation, we fuse the detection results in the post-processing phase.

| Method | Test Aug | Crop Size | mAP@0.3 | mAP@0.4 | mAP@0.5 | mAP@0.6 | mAP@0.7 | Avg |
|---|---|---|---|---|---|---|---|---|
| Anchor-based | CenterCrop | $160 \times 160$ | 71.2 | 66.8 | 61.2 | 50.1 | 36.3 | 57.1 |
| | CenterCrop | $180 \times 180$ | 71.6 | 67.3 | 61.6 | 50.8 | 34.7 | 57.2 |
| | CenterCrop+Backward | $180 \times 180$ | 71.1 | 67.4 | 61.1 | 51.0 | 34.5 | 57.0 |
| | ThreeCrop | $180 \times 180$ | 71.7 | **67.9** | 62.0 | 50.7 | **35.6** | **57.6** |
| | CenterCrop+Flip | $180 \times 180$ | **71.9** | 67.7 | **62.1** | **51.0** | 35.2 | **57.6** |
| Anchor-free | CenterCrop | $160 \times 160$ | 72.5 | 66.8 | 59.1 | 48.4 | 35.0 | 56.4 |
| | CenterCrop | $180 \times 180$ | 72.9 | 66.3 | 59.5 | 48.2 | 35.1 | 56.4 |
| | CenterCrop+Backward | $180 \times 180$ | **72.9** | 66.7 | **59.9** | 47.6 | 33.9 | 56.2 |
| | ThreeCrop | $180 \times 180$ | 72.3 | **67.6** | 59.0 | 48.5 | **35.9** | **56.7** |
| | CenterCrop+Flip | $180 \times 180$ | 72.4 | 66.8 | 59.7 | **48.9** | 35.0 | 56.6 |

### 4.5.2. Spatial Preservation in Neck Design

Our proposed Spatial Preservation (SP) design aims to build spatial-sensitive temporal multi-scale features in the neck module. We maintain the spatial dimension of features in the neck module and perform the spatial squeezing after completing the multi-scale feature construction.

As shown in Table 8, with the same other settings, using SP for the neck module with 96 and 192 frames can obtain 0.8 and 0.6 mAP improvement in the anchor-free method, and 0.5 and 0.3 mAP improvement in the anchor-based method, respectively. This is due to the fact that additional spatial dimension can make features capture rich structure on the spatial locations of actions occurring at different scales. Meanwhile, we observe that the computational cost of PlusTAD with SP is almost the same as without SP. The performance of SP demonstrates that it is a simple module for improving temporal action detection performance. Therefore, by default, we incorporate the SP design principle in our PlusTAD.

### 4.5.3. Effectiveness of Training Augmentation

In the previous ablation studies, we used the default data augmentation techniques for training and testing. For simplicity, we set the original size to $128 \times 128$ and the crop size to $112 \times 112$. But a larger resolution of original frames can provide more structure information for modeling and more crops for data augmentation. Therefore, we perform a detailed study on the input frame resolution to investigate its effect on the TAD performance. To study its effect on detection performance, we adopt the same settings of using **96 frames at 3FPS** as the input and enable **SP** and **TP** in our PlusTAD.

Specifically, we use two choices to expand the frame resolution. One choice is to use the fixed resolution as "171× 128", and the other is to use "short-128" to keep the original aspect ratio. Among them, "171 × 128" were first adopted in (Xu et al., 2017) for network training, and we follow its setting to increase the input resolution. First, we only change the frame resolution and fix the crop size as $112 \times 112$. As shown in the first three rows in Table 9 for both anchor-based and anchor-free PlusTAD, a larger original size can contribute to a better TAD performance. Notably, it can bring around 4% performance improvement for both anchor-based and anchor-free detection methods. This significant performance improvement indicates that larger input resolution can provide more detailed spatial structure information and generate more diverse training samples.

Furthermore, we perform another comparative study by increasing the crop size. More concrete, we increase the origin frame size to "short-180" and crop size to "160×160". This setting has a similar ratio of crop size and original size to the original PlusTAD. As shown in Table 9, under this setting, the average mAP obtains an improvement of about 2.0%. This performance improvement demonstrates the effectiveness of increasing input frame resolution and keeping the resolution correspondence between the original

**Table 11**
Ablation Study results on spatial and temporal resolution on THUMOS14 (Jiang et al., 2014). Overhead represents the increase in computational overhead after expanding the temporal or spatial resolution, and the basic one denotes by ×1.

| Method | Original Size | Crop Size | FPS | Frames | mAP@0.3 | mAP@0.4 | mAP@0.5 | mAP@0.6 | mAP@0.7 | Avg | FLOPs |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Anchor-based | short-128 | $112 \times 112$ | 3 | 96 | 68.4 | 65.0 | 58.6 | 49.2 | 33.5 | 54.9 | 136.4G |
| | short-180 | $160 \times 160$ | 3 | 96 | 71.2 | 66.8 | 61.2 | 50.1 | 36.3 | 57.1 | 262.1G |
| | short-128 | $112 \times 112$ | 6 | 192 | 70.2 | 66.3 | 60.6 | 50.3 | 36.2 | 56.7 | 272.9G |
| | short-180 | $160 \times 160$ | 6 | 192 | **72.3** | **68.4** | **62.0** | **52.4** | **37.0** | **58.4** | 519.3G |
| Anchor-free | short-128 | $112 \times 112$ | 3 | 96 | 70.4 | 65.5 | 57.6 | 46.0 | 33.2 | 54.5 | 151.5G |
| | short-180 | $160 \times 160$ | 3 | 96 | 72.5 | 67.2 | 59.1 | 48.2 | 35.3 | 56.4 | 275.9G |
| | short-128 | $112 \times 112$ | 6 | 192 | 71.7 | 66.9 | 59.0 | 49.2 | 35.3 | 56.4 | 284.2G |
| | short-180 | $160 \times 160$ | 6 | 192 | **75.5** | **70.8** | **63.5** | **50.9** | **37.4** | **59.6** | 533.1G |

frame and crop patch. Therefore, we choose the input frame resolution as "short-180" and the crop size as $160 \times 160$ by default in the following studies.

### 4.5.4. Effectiveness of Testing Augmentation

After the ablation on the spatial resolution of training frames and crop patches, we investigate the influence of testing augmentation on the TAD performance. The augmentation methods in the testing phase can improve the robustness of the predictions without additional training costs. We conduct ablation experiments to study the effect of different cropping methods. We attempt to apply four following kinds of testing augmentation. "CenterCrop" is cropping the center area the same size as the training crop from images. "ThreeCrop" is cropping three areas of the same size as the training crop from the frames along the long side. "Flip" is flipping the frames for horizontal augmentation. "Backward" complements a reverse sliding window process to sample windows. Each augmentation method produces a new clip view to be inferred by the network separately and fuses predictions of these views to produce the final TAD results. Specifically, for ThreeCrop and Flip augmentation, we fuse the features after the neck. For backward augmentation, we fuse the detection results in the post-processing phase.

The results of different testing augmentation are shown in Table 10. We find that "ThreeCrop" and "CenterCrop+Flip" achieve very similar results for both anchor-based and anchor-free TAD methods, which can improve the performance of "CenterCrop" by around 0.5%. Unexpectedly, "+Backward" seems to hurt the final precision. A possible explanation is that increasing the number of predicted actions degenerates the distribution of the original prediction. Considering the extra computation cost brought by extra testing views and its small performance improvement, we set the default testing scheme of our PlusTAD as the simple "Center Crop".

### 4.5.5. Study on the spatial and temporal resolutions

From previous studies, we have concluded that larger frame resolution and more input frames can contribute to a higher TAD performance. We further perform ablation studies to investigate how to keep a balance between spatial resolution and temporal frame under a fixed computational budget. Considering video data contains two dimensions, spatial and temporal. Using **96 frames at 3FPS** as a baseline, we can add sources of information along two alternative paths, one of which is increasing the number of frames (higher temporal resolution) and the other is using larger frames (higher spatial resolution). We conduct experiments to explore the gain of different data sources and search for an optimal route to boost performance.

As reported in Table 11, doubling the number of frames and doubling the frame size separately brings similar performance improvement for both detection heads. It is worth noting that, compared with more frames, a larger resolution of frames may deliver more benefits to the anchor-based method. When we double the number of frames and double the frame size, both detection heads are further improved and create a new record of temporal action detection. We obtain the mAP of 59.6% for anchor-free PlusTAD. In summary, we find that temporal resolution and spatial resolution are both important for improving TAD performance, and we can choose a reasonable baseline method under the computational resource available.

### 4.6. Comparison with the State of the Art

In previous subsections, we have performed thorough ablation studies on the basic design in our modular Basic-TAD framework. We also provide extensive investigation on our proposed new design principles to enhance the TAD performance of BasicTAD. In this subsection, we turn to compare our PlusTAD with previous state-of-the-art methods. In this comparison, we direct transfer these optimal configurations discovered from THUMOS14 (Jiang et al., 2014) to the other two large-scale benchmarks: FineAction (Liu et al., 2022b) and ActivityNet v1.3 (Heilbron et al., 2015). Table 12 reports our configurations on the three datasets. We use $\text{PlusTAD}_{R,F}^{S}$ to represent the configuration name on the datasets with the window sampling strategy, such as THUMOS14 and FineAction. R, F, S is the frame rate of the sampling, the number of frames, and the size of the frames. For the datasets based on the video-level sampling strategy, such as ActivityNet v1.3, we delete R of the subscript due to the sampled frames representing the whole video. We

**Table 12**
**Summary of configurations** of the optimal PlusTAD on the datasets of THUMOS14, FineAction, and ActivityNet-v1.3.

| Dataset | Configuration | Frame Representation | Temporal Resolution | Spatial Resolution | Window Size |
|---|---|---|---|---|---|
| THUMOS14 | PlusTAD$_{3,96}^{112}$ | Window | 96 | 112×112 | 32s |
| THUMOS14 | PlusTAD$_{6,192}^{160}$ | Window | 192 | 160×160 | 32s |
| FineAction | PlusTAD$_{2,96}^{224}$ | Window | 96 | 224×224 | 48s |
| ActivityNet-v1.3 | PlusTAD$_{96}^{224}$ | Video | 96 | 224×224 | - |

**Table 13**
**Comparison with state of the art on the THUMOS14.** "RGB-Only" means whether to use other input modalities besides RGB input.

| Type | Method | Backbone | RGB-Only | mAP @0.3 | mAP @0.4 | mAP @0.5 | mAP @0.6 | mAP @0.7 | mAP @Avg | FLOPs |
|---|---|---|---|---|---|---|---|---|---|---|
| Multi-stage | BSN (Lin et al., 2018) | TSN | ✗ | 53.5 | 45.0 | 36.9 | 28.4 | 20.0 | 36.8 | - |
| | MGG (Liu et al., 2019) | TSN | ✗ | 53.9 | 46.8 | 37.4 | 29.5 | 21.3 | 37.8 | - |
| | BMN (Lin et al., 2019b) | TSN | ✗ | 56.0 | 47.4 | 38.8 | 29.7 | 20.5 | 38.5 | - |
| | DBG (Lin et al., 2019b) | TSN | ✗ | 57.8 | 49.4 | 39.8 | 30.2 | 21.7 | 39.8 | - |
| | RTD-Net (Tan et al., 2021) | I3D | ✗ | 58.5 | 53.1 | 45.1 | 36.4 | 25.0 | 43.6 | - |
| | TCANet (Qing et al., 2021) | TSN | ✗ | 60.6 | 53.2 | 44.6 | 36.8 | 26.7 | 44.4 | - |
| | G-TAD (Xu et al., 2020) | TSN | ✗ | 66.4 | 60.4 | 51.6 | 37.6 | 22.9 | 47.8 | - |
| | AFSD (Lin et al., 2021) | I3D | ✗ | 67.3 | 62.4 | 55.5 | 43.7 | 31.1 | 52.0 | 2780.0G |
| | DCAN (Chen et al., 2022) | TSN | ✗ | 68.2 | 62.7 | 54.1 | 43.9 | 32.6 | 52.3 | - |
| | SP-TAD (Wu et al., 2021) | I3D | ✗ | 69.2 | 63.3 | 55.9 | 45.7 | 33.4 | 53.5 | - |
| | TadTR (Liu et al., 2022a) | R50-SlowFast | ✔ | 69.4 | 64.3 | 56.0 | 46.4 | 34.9 | 54.2 | 475.0G |
| One-stage | SSAD (Lin et al., 2017b) | TSN | ✗ | 43.0 | 35.0 | 24.6 | - | - | - | - |
| | DBS (Gao et al., 2019) | TSN | ✗ | 50.6 | 43.1 | 34.3 | 24.4 | 14.7 | 33.4 | - |
| | A2Net (Yang et al., 2020) | I3D | ✗ | 58.6 | 54.1 | 45.5 | 32.5 | 17.2 | 41.6 | - |
| | PBRNet (Liu and Wang, 2020) | I3D | ✗ | 58.5 | 54.6 | 51.3 | 41.8 | 29.5 | 47.1 | - |
| | R-C3D (Xu et al., 2017) | C3D | ✔ | 44.8 | 35.6 | 28.9 | - | - | - | 1360.0G |
| | GTAN (Long et al., 2019) | P3D | ✔ | 57.8 | 47.2 | 38.8 | - | - | - | - |
| | DaoTAD (Wang et al., 2021a) | R50-I3D | ✔ | 62.8 | 59.5 | 53.8 | 43.6 | 30.1 | 50.0 | 206.7G |
| | DaoTAD$_{3,96}^{112}$ (Wang et al., 2021a) | R50-SlowOnly | ✔ | 63.2 | 59.7 | 54.4 | 45.6 | 32.2 | 51.0 | 133.3G |
| | **PlusTAD$_{3,96}^{112}$ (Anchor-based)** | R50-SlowOnly | ✔ | **68.4** | **65.0** | **58.6** | **49.2** | **33.5** | **54.9** | 136.4G |
| | **PlusTAD$_{3,96}^{112}$ (Anchor-free)** | R50-SlowOnly | ✔ | **70.4** | **65.5** | **57.6** | **46.0** | **33.2** | **54.5** | 151.5G |
| | **PlusTAD$_{6,192}^{160}$ (Anchor-based)** | R50-SlowOnly | ✔ | **72.3** | **68.4** | **62.0** | **52.4** | **37.0** | **58.4** | 519.3G |
| | **PlusTAD$_{6,192}^{160}$ (Anchor-free)** | R50-SlowOnly | ✔ | **75.5** | **70.8** | **63.5** | **50.9** | **37.4** | **59.6** | 533.1G |

compare these configurations with state-of-the-art methods on the three datasets.

**THUMOS14.** We compare other state-of-the-art methods in Table 13 on the THUMOS14 dataset. We compute the FLOPs of other end-to-end methods (Lin et al., 2021; Wang et al., 2021a; Liu et al., 2022a; Xu et al., 2017) in the table for a fair comparison with our method. For some two-stage and head-only methods, their FLOPs are unavailable and indicated as "-" in the table. From this table, we see that our PlusTAD significantly outperforms previous methods by a large margin with only RGB as input. In particular, PlusTAD$_{6,192}^{160}$ uses more frames with a larger size and improves the mAP jumps by 4 to 5 points. Meanwhile, our overall FLOPS are still less than these methods, and our simple design allows for very fast deployment. These superior results demonstrate the effectiveness of our PlusTAD thanks to its simplicity and end-to-end training.

**FineAction.** We also compare our proposed PlusTAD with the state-of-the-art methods on FineAction (Liu et al., 2022b). FineAction is a new and large-scale TAD dataset, and few works have reported results on this new benchmark. As shown in Table 14. We compare our PlusTAD$_{2,96}^{224}$ with three representative works (Lin et al., 2019b, 2020; Xu et al., 2020), which are provided by (Liu et al., 2022b). All of them are based on pre-extracted I3D features. We find that our method can obtain a better performance on Average mAP. However, our method performs worse on mAP@0.95. One possible reason is that all three methods above are multi-stage methods that can get refined proposals from the previous detection or proposal generation stage. This leads to better localization results for metrics at high thresholds like mAP@0.95.

**ActivityNet-v1.3.** Due to the complex semantics of the action class on ActivityNet-v1.3, we adapt our PlusTAD to generate binary action proposals and obtain the detection

**Table 14**
**Comparison with state of the art on the FineAction dataset.** "RGB-Only" means whether to use other input modalities besides RGB input.

| Method | Backbone | RGB-Only | mAP@0.5 | mAP@0.75 | mAP@0.95 | Avg |
|---|---|---|---|---|---|---|
| BMN (Lin et al., 2019b) | I3D | ✔ | 12.56 | 7.49 | 2.62 | 7.86 |
| BMN (Lin et al., 2019b) | I3D | ✗ | 14.44 | 8.92 | 3.12 | 9.25 |
| DBG (Lin et al., 2020) | I3D | ✔ | 8.57 | 5.01 | 1.93 | 5.31 |
| DBG (Lin et al., 2020) | I3D | ✗ | 10.65 | 6.43 | 2.50 | 6.75 |
| G-TAD (Xu et al., 2020) | I3D | ✔ | 10.88 | 6.52 | 2.19 | 6.87 |
| G-TAD (Xu et al., 2020) | I3D | ✗ | 13.74 | 8.83 | **3.06** | 9.06 |
| **PlusTAD$_{2,96}^{224}$(Anchor-based)** | R50-SlowOnly | ✔ | **24.34** | **10.57** | 0.43 | **12.15** |
| **PlusTAD$_{2,96}^{224}$(Anchor-free)** | R50-SlowOnly | ✔ | 22.37 | 10.36 | 0.83 | 11.66 |

**Table 15**
**Comparison with state of the art on the ActivityNet-1.3 dataset.** "RGB-Only" means whether to use other input modalities besides RGB input.

| Method | Backbone | RGB-Only | mAP@0.5 | mAP@0.75 | mAP@0.95 | Avg |
|---|---|---|---|---|---|---|
| BSN (Lin et al., 2018) | TSN | ✗ | 46.45 | 29.96 | 8.02 | 30.03 |
| P-GCN (Zeng et al., 2019) | - | ✗ | 48.26 | 33.16 | 3.27 | 31.11 |
| BMN (Lin et al., 2019b) | TSN | ✗ | 50.07 | 34.78 | 8.29 | 33.85 |
| G-TAD (Xu et al., 2020) | TSN | ✗ | 50.36 | 34.60 | **9.02** | 34.09 |
| AFSD (Lin et al., 2021) | I3D | ✗ | **52.40** | **35.30** | 6.50 | **34.40** |
| SP-TAD (Wu et al., 2021) | I3D | ✗ | 50.06 | 32.92 | 8.44 | 32.99 |
| BMN (Lin et al., 2019b) | TSN | ✔ | 41.93 | 30.10 | **9.00** | 29.23 |
| G-TAD (Xu et al., 2020) | TSN | ✔ | 45.68 | 31.36 | 7.42 | 30.98 |
| AFSD (Lin et al., 2021) | I3D | ✔ | - | - | - | 32.90 |
| **PlusTAD$_{96}^{224}$(Anchor-based)** | R50-SlowOnly | ✔ | 50.04 | **33.79** | 2.75 | 32.13 |
| **PlusTAD$_{96}^{224}$(Anchor-free)** | R50-SlowOnly | ✔ | **51.20** | 33.41 | 7.57 | **33.12** |

results by applying video-level action classifiers. This is a common setting for many previous state-of-the-art methods on ActivityNet-v1.3. This is mainly due to the annotation sparsity of ActivityNet-v1.3, and thus the detection results can benefit from the video-level classification. As shown in Table 12, due to the sparse action instance distribution of ActivityNet-v1.3, we sparsely sample 96 frames from the entire video to represent the whole video. This sparse sampling would greatly increase the temporal receptive field but also lose detailed temporal information that might be useful for TAD.

As shown in Table 15, our PlusTAD$_{96}^{224}$ performs better than (Lin et al., 2019b; Xu et al., 2020; Lin et al., 2021) when only using RGB frames as input. Meanwhile, our method is a one-stage detector without any refinement sub-network, such as AFSD (Lin et al., 2021), to refine boundaries. However, other head-only methods use pre-extracted features from the entire video to get richer information. Note in this case, they can use more frames (not only 96) for feature extraction. Their detection results can be further improved by introducing an additional optical flow modality. Considering these several factors, our PlusTAD still remains a powerful baseline method due to its simplicity, effectiveness, and efficiency.

### 4.7. Efficiency Analysis

In the above subsections, we have demonstrated that our BasicTAD and PlusTAD can achieve very high TAD mAP on the standard benchmarks. Due to its simplicity in design, our PlusTAD enjoys high efficiency as well. In this subsection, we report the running speed of our PlusTAD and compare it with the other state-of-the-art methods with end-to-end training strategies.

To make the comparison fair and rigorous, we first define the inference speed for TAD tasks. We use "FPS" to represent the number of frames per second processed by the model when processing the video stream. "FPS" here is calculated differently between the image field and TAD. Given an input video with $s$ seconds and its $f$ FPS, the total number of frames to be processed is $s \times f$. If we re-sample frames before feeding them into the network, we can get a new input with $s$ seconds, $f'$ FPS, and $s \times f'$ frames. Suppose we spend $t$ seconds to process re-sampled frames. The speed (FPS) of processing the origin video is $\frac{s \times f}{t}$.

We use 96 frames at 3FPS to evaluate the inference speed of PlusTAD with both heads because the average mAP is close to other SOTA methods under this setting. Considering using the same GPU, as shown in Table 16, PlusTAD with the anchor-based method yields 17454 FPS using V100 GPU, which is 4× faster than (Lin et al., 2021).

**Table 16**
**Comparison of inference speed.** The running speed of previous methods is directly cited from their papers. Note that these methods all use optical flow as inputs, and their running time does not include the optical flow calculation. So, their real running speed is lower than the reported speed. † donate DaoTAD (Wang et al., 2021a) tested in our platform.

| Method | FPS | Frames | Size | GPU | RGB-Only | FPS | mAP |
|---|---|---|---|---|---|---|---|
| SS-TAD (Buch et al., 2017) | - | - | - | TITAN XM | ✗ | < 701 | - |
| R-C3D (Xu et al., 2017) | 25 | 768 | $112 \times 112$ | TITAN XP | ✗ | < 1030 | - |
| PBRNet (Liu and Wang, 2020) | 10 | 256 | $96 \times 96$ | 1080Ti | ✗ | < 1488 | 47.1 |
| AFSD (Lin et al., 2021) | 10 | 256 | $96 \times 96$ | 1080Ti | ✗ | < 3259 | 52.0 |
| AFSD (Lin et al., 2021) | 10 | 256 | $96 \times 96$ | V100 | ✗ | < 4057 | 52.0 |
| DaoTAD (Wang et al., 2021a) | 25 | 768 | $112 \times 112$ | 1080Ti | ✔ | 6668 | 50.0 |
| SP-TAD (Wu et al., 2021) | 10 | 256 | $96 \times 96$ | V100 | ✗ | <5574 | 53.5 |
| TadTR (Liu et al., 2022a) | 10 | 256 | $96 \times 96$ | TITAN XP | ✔ | 5076 | 54.2 |
| DaoTAD† (Wang et al., 2021a) | 25 | 768 | $112 \times 112$ | TITAN XP | ✔ | 7064 | 50.0 |
| DaoTAD† (Wang et al., 2021a) | 25 | 768 | $112 \times 112$ | V100 | ✔ | 8989 | 50.0 |
| **PlusTAD**$_{3,96}^{112}$**(Anchor-based)** | 3 | 96 | $112 \times 112$ | TITAN XP | ✔ | **13715** | **54.9** |
| **PlusTAD**$_{3,96}^{112}$**(Anchor-based)** | 3 | 96 | $112 \times 112$ | V100 | ✔ | **17454** | **54.9** |
| **PlusTAD**$_{3,96}^{112}$**(Anchor-free)** | 3 | 96 | $112 \times 112$ | TITAN XP | ✔ | **7143** | **54.5** |
| **PlusTAD**$_{3,96}^{112}$**(Anchor-free)** | 3 | 96 | $112 \times 112$ | V100 | ✔ | **8377** | **54.5** |

Although PlusTAD with the anchor-free method is much slower because of more channels in the neck and head, it still has a higher FPS and mAP than AFSD (Lin et al., 2021) and TadTR (Liu et al., 2022a).

Furthermore, to eliminate the influence of other hardware, such as CPU, RAM, and disk, we also compare the inference speed of DaoTAD (Wang et al., 2021a) reproduced on our platform. In Table 16, the FPS on 1080Ti reported by DaoTAD (Wang et al., 2021a) is slower by 5.9% than our report. It roughly matches the performance gap between 1080Ti and TITAN XP. From the above analysis, our Plus-TAD still maintains an advantage in inference speed.

These comparisons show that it is necessary to re-design a simple and efficient TAD baseline method, and our proposed PlusTAD can achieve a new astounding baseline of high effectiveness and efficiency.

### 4.8. Feature Activation Analysis

After conducting extensive quantitative studies on the accuracy and efficiency of our proposed BasicTAD and PlusTAD, we present some visualization analysis on THU-MOS14 in this subsection. We start by visualizing the feature maps of the backbone for both anchor-based and anchor-free heads. To enhance the visualization clarity, we perform average pooling to squeeze the feature map spatially and retain only the temporal dimension. As illustrated in Figure 3, we depict the feature activation of anchor-based PlusTAD at the top, the original frames with ground truth in the middle, and the feature activation of anchor-free PlusTAD at the bottom. We provide three examples in total. As evident in Figure 3, the anchor-based method exhibits a preference for activating a large area of target regions. On the other hand, the anchor-free method demonstrates activation at a specific and single location within the target action. Figure 3 shows that the anchor-based method prefers to activate at a large area of target regions. In contrast, the anchor-free

method activates at a specific and single location inside the target action. While these findings cannot be considered definitive conclusions, we believe they can serve as useful guidelines for future research to investigate the application of anchor-based and anchor-free methods in TAD and draw more comprehensive conclusions.

### 4.9. Error Analysis

In addition to the feature map visualization, we provide the error analysis on our detection results as well. To better understand the detection errors, we use the diagnostic tool provided by Alwassel et al. (2018) to analyze the detection results of PlusTAD with 192 frames at 6FPS. We visualize the error analysis of both anchor-based and anchor-free methods to make a direct comparison between them and provide some insightful analysis of both TAD methods.

*False Positive Analysis* The left of Figure 4 shows false positive (FP) profiling. As shown on the left of the figure, top-1G predictions contain the most TP (true positive) predictions. The majority of error type is localization error and background error. When we consider more predictions, the true positive rate will decrease dramatically. In this sense, the vast majority of ground truth actions can be successfully predicted within top-1G predictions. Comparing the anchor-based and anchor-free methods from top-2G to top-10G, we could see that the anchor-free method suffers more from "Background Err". It may be due to limited anchors causing more predictions to fail to match ground truth. Conversely, the anchor-based method has sufficient anchors to reduce "Background Err", but may bring other errors such as "Wrong Label Err".

*False Negative Analysis* False negative (FN) profiling is illustrated in the right of Figure 4. From the visualization, we see that under the measure of Coverage and Length, the

**Figure 3:** Visualization of our anchor-based and anchor-free methods on three different actions "TennisSwing", "SoccerPenalty" and "GolfSwing". For each action, the top feature map represents the anchor-based method's action activation feature map, while the bottom represents the anchor-free method's action activation map. The middle one is the original action frame, and the video order is from left to right. For example, we choose "TennisSwing" with two ground truth actions in the red and orange boxes. Anchor-based and anchor-free methods each have two activation parts of target actions, and we mark them with red and orange boxes.



(a) False Positive Profiles                    (b) False Negative Profiles

**Figure 4:** (a) **False Positive Profiles.** Left: False positive profiles of both methods. Each profile demonstrates the FP error breakdown in the top-10G predictions. Right: Improvement gained from removing all predictions that cause each type of error. The higher the value, the greater the effect on average-mAP$_N$. (b) **False Negative Profiles.** Average false negative rate across algorithms for each characteristic on both methods.

FN rate of the anchor-based method is significantly lower than the anchor-free method. Since short actions of the same category usually have more numbers in THUMOS14, the anchor-based method can benefit from the dense sample matching mechanism and tends to better detect short actions.

While, for the anchor-free method, it lacks the flexibility of detecting very small or large action instances, which might be due to the difficulty of directly regressing the boundaries of these hard action instances. In general, we can conclude that the anchor-based method might be more robust than the

anchor-free method for dealing with extremely hard cases, but not their detection results are not as accurate as the anchor-free method.

# 5. Conclusion and Future Work

In this paper, we have reconsidered the design of the TAD pipeline and presented a simple modular detection framework. Based on this modular design, we perform extensive investigation over the basic options in each component and finally, come up with a simple yet effective TAD baseline, termed BasicTAD. Furthermore, we improve the BasicTAD with minimal changes by following the core design of preserving rich information in the backbone and neck, and the resulted detector is termed PlusTAD. Our end-to-end BasicTAD and PlusTAD are free of pre-processing and can be used in real-time application scenarios. Extensive experiments demonstrate that our PlusTAD significantly outperforms previous state-of-the-art methods on the THU-MOS14 and FineAction, and achieves quite competitive results on the ActivityNet-v1.3 datasets. We also provide in-depth ablation studies on each TAD component in a step-by-step manner and detailed visualization results to figure out the main property and major error of our PlusTAD. We hope our approach can serve as a strong baseline for future TAD research.

In the future, we can enhance the temporal modeling capacity of our BasicTAD and PlusTAD by incorporating long-term memory. This addition will enable us to tackle action instances of longer duration with greater precision and accuracy. In addition, we can expand our TAD framework to include Transformer backbones with global receptive fields to boost its representation power. Doing so can significantly improve the final TAD performance, further enhancing our ability to analyze complex temporal sequences.

# References

Alwassel, H., Heilbron, F.C., Escorcia, V., Ghanem, B., 2018. Diagnosing error in temporal action detectors, in: ECCV, Springer. pp. 264–280.

Arnab, A., Dehghani, M., Heigold, G., Sun, C., Lucic, M., Schmid, C., 2021. Vivit: A video vision transformer, in: ICCV, pp. 6816–6826.

Bai, Y., Wang, Y., Tong, Y., Yang, Y., Liu, Q., Liu, J., 2020. Boundary content graph neural network for temporal action proposal generation, in: ECCV, pp. 121–137.

Bertasius, G., Wang, H., Torresani, L., 2021. Is space-time attention all you need for video understanding?, in: Meila, M., Zhang, T. (Eds.), Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event, PMLR. pp. 813–824.

Buch, S., Escorcia, V., Ghanem, B., Fei-Fei, L., Niebles, J.C., 2017. End-to-end, single-stream temporal action detection in untrimmed videos, in: BMVC.

Carreira, J., Zisserman, A., 2017. Quo vadis, action recognition? A new model and the kinetics dataset, in: CVPR, pp. 4724–4733.

Chao, Y., Vijayanarasimhan, S., Seybold, B., Ross, D.A., Deng, J., Sukthankar, R., 2018. Rethinking the faster R-CNN architecture for temporal action localization, in: CVPR, pp. 1130–1139.

Chen, G., Zheng, Y., Wang, L., Lu, T., 2022. DCAN: improving temporal action detection via dual context aggregation, in: AAAI, AAAI Press. pp. 248–257.

Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S.,

Uszkoreit, J., Houlsby, N., 2021. An image is worth 16x16 words: Transformers for image recognition at scale, in: 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021, OpenReview.net.

Feichtenhofer, C., Fan, H., Malik, J., He, K., 2019. Slowfast networks for video recognition, in: ICCV, pp. 6201–6210.

Gao, J., Shi, Z., Wang, G., Li, J., Yuan, Y., Ge, S., Zhou, X., 2020. Accurate temporal action proposal generation with relation-aware pyramid network, in: AAAI, pp. 10810–10817.

Gao, J., Yang, Z., Sun, C., Chen, K., Nevatia, R., 2017. TURN TAP: temporal unit regression network for temporal action proposals, in: ICCV, pp. 3648–3656.

Gao, Z., Wang, L., Zhang, Q., Niu, Z., Zheng, N., Hua, G., 2019. Video imprint segmentation for temporal action detection in untrimmed videos, in: The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019, AAAI Press. pp. 8328–8335.

He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016, IEEE Computer Society. pp. 770–778.

Heilbron, F.C., Escorcia, V., Ghanem, B., Niebles, J.C., 2015. Activitynet: A large-scale video benchmark for human activity understanding, in: CVPR, pp. 961–970.

Ioffe, S., Szegedy, C., 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift, in: Bach, F.R., Blei, D.M. (Eds.), Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015, JMLR.org. pp. 448–456.

Jiang, Y.G., Liu, J., Roshan Zamir, A., Toderici, G., Laptev, I., Shah, M., Sukthankar, R., 2014. THUMOS challenge: Action recognition with a large number of classes. http://crcv.ucf.edu/THUMOS14/.

Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P., Suleyman, M., Zisserman, A., 2017. The kinetics human action video dataset. CoRR abs/1705.06950.

Li, Y., Ji, B., Shi, X., Zhang, J., Kang, B., Wang, L., 2020. TEA: temporal excitation and aggregation for action recognition, in: CVPR, pp. 906–915.

Lin, C., Li, J., Wang, Y., Tai, Y., Luo, D., Cui, Z., Wang, C., Li, J., Huang, F., Ji, R., 2020. Fast learning of temporal action proposal via dense boundary generator, in: Proceedings of the AAAI Conference on Artificial Intelligence, pp. 11499–11506.

Lin, C., Xu, C., Luo, D., Wang, Y., Tai, Y., Wang, C., Li, J., Huang, F., Fu, Y., 2021. Learning salient boundary feature for anchor-free temporal action localization, in: IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021, Computer Vision Foundation / IEEE. pp. 3320–3329.

Lin, J., Gan, C., Han, S., 2019a. TSM: temporal shift module for efficient video understanding, in: ICCV, pp. 7082–7092.

Lin, T., Goyal, P., Girshick, R.B., He, K., Dollár, P., 2017a. Focal loss for dense object detection, in: ICCV, IEEE Computer Society. pp. 2999–3007.

Lin, T., Liu, X., Li, X., Ding, E., Wen, S., 2019b. BMN: boundary-matching network for temporal action proposal generation, in: ICCV, pp. 3888–3897.

Lin, T., Zhao, X., Shou, Z., 2017b. Single shot temporal action detection, in: Liu, Q., Lienhart, R., Wang, H., Chen, S.K., Boll, S., Chen, Y.P., Friedland, G., Li, J., Yan, S. (Eds.), ACM MM, pp. 988–996.

Lin, T., Zhao, X., Su, H., Wang, C., Yang, M., 2018. BSN: boundary sensitive network for temporal action proposal generation, in: Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y. (Eds.), ECCV, pp. 3–21.

Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P., 2017c. Focal loss for dense object detection, in: ICCV.

Liu, Q., Wang, Z., 2020. Progressive boundary refinement network for temporal action detection, in: The Thirty-Fourth AAAI Conference on

Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020, AAAI Press. pp. 11612–11619.

Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S.E., Fu, C., Berg, A.C., 2016. SSD: single shot multibox detector, in: ECCV.

Liu, X., Bai, S., Bai, X., 2022a. An empirical study of end-to-end temporal action detection. CoRR abs/2204.02932.

Liu, X., Wang, Q., Hu, Y., Tang, X., Bai, S., Bai, X., 2021a. End-to-end temporal action detection with transformer. CoRR abs/2106.10271.

Liu, Y., Ma, L., Zhang, Y., Liu, W., Chang, S., 2019. Multi-granularity generator for temporal action proposal, in: CVPR, pp. 3604–3613.

Liu, Y., Wang, L., Wang, Y., Ma, X., Qiao, Y., 2022b. FineAction: A fine-grained video dataset for temporal action localization. IEEE Trans. Image Process. 31, 6937–6950.

Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B., 2021b. Swin transformer: Hierarchical vision transformer using shifted windows, in: 2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021, IEEE. pp. 9992–10002.

Liu, Z., Luo, D., Wang, Y., Wang, L., Tai, Y., Wang, C., Li, J., Huang, F., Lu, T., 2020. Teinet: Towards an efficient architecture for video recognition, in: AAAI, pp. 11669–11676.

Liu, Z., Ning, J., Cao, Y., Wei, Y., Zhang, Z., Lin, S., Hu, H., 2021c. Video swin transformer. CoRR abs/2106.13230.

Liu, Z., Wang, L., Wu, W., Qian, C., Lu, T., 2021d. TAM: temporal adaptive module for video recognition, in: ICCV, pp. 13688–13698.

Long, F., Yao, T., Qiu, Z., Tian, X., Luo, J., Mei, T., 2019. Gaussian temporal awareness networks for action localization, in: CVPR, pp. 344–353.

Materzynska, J., Xiao, T., Herzig, R., Xu, H., Wang, X., Darrell, T., 2020. Something-else: Compositional action recognition with spatial-temporal interaction networks, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020, Computer Vision Foundation / IEEE. pp. 1046–1056.

Moltisanti, D., Wray, M., Mayol-Cuevas, W., Damen, D., 2017. Trespassing the boundaries: Labeling temporal bounds for object interactions in egocentric video, in: Proceedings of the IEEE International Conference on Computer Vision, pp. 2886–2894.

Neimark, D., Bar, O., Zohar, M., Asselmann, D., 2021. Video transformer network. CoRR abs/2102.00719.

Neubeck, A., Gool, L.V., 2006. Efficient non-maximum suppression, in: 18th International Conference on Pattern Recognition (ICPR 2006), 20-24 August 2006, Hong Kong, China, IEEE Computer Society. pp. 850–855.

Ning, C., Zhou, H., Song, Y., Tang, J., 2017. Inception single shot multibox detector for object detection, in: 2017 IEEE International Conference on Multimedia & Expo Workshops, ICME Workshops, Hong Kong, China, July 10-14, 2017, IEEE Computer Society. pp. 549–554.

Qing, Z., Su, H., Gan, W., Wang, D., Wu, W., Wang, X., Qiao, Y., Yan, J., Gao, C., Sang, N., 2021. Temporal context aggregation network for temporal action proposal refinement, in: CVPR, pp. 485–494.

Qiu, Z., Yao, T., Mei, T., 2017. Learning spatio-temporal representation with pseudo-3d residual networks, in: ICCV, pp. 5534–5542.

Ren, S., He, K., Girshick, R.B., Sun, J., 2015. Faster R-CNN: towards real-time object detection with region proposal networks, in: NIPS.

Simonyan, K., Zisserman, A., 2014. Two-stream convolutional networks for action recognition in videos, in: Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N.D., Weinberger, K.Q. (Eds.), NIPS, pp. 568–576.

Su, H., Gan, W., Wu, W., Qiao, Y., Yan, J., 2021. BSN++: complementary boundary regressor with scale-balanced relation modeling for temporal action proposal generation, in: AAAI, pp. 2602–2610.

Tan, J., Tang, J., Wang, L., Wu, G., 2021. Relaxed transformer decoders for direct action proposal generation, in: ICCV, pp. 13526–13535.

Tian, Z., Shen, C., Chen, H., He, T., 2019. FCOS: fully convolutional one-stage object detection, in: ICCV, IEEE. pp. 9626–9635.

Tong, Z., Song, Y., Wang, J., Wang, L., 2022. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. CoRR abs/2203.12602.

Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jégou, H., 2021. Training data-efficient image transformers & distillation through attention, in: Meila, M., Zhang, T. (Eds.), Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event, PMLR. pp. 10347–10357.

Tran, D., Bourdev, L.D., Fergus, R., Torresani, L., Paluri, M., 2015. Learning spatiotemporal features with 3d convolutional networks, in: ICCV, pp. 4489–4497.

Tran, D., Ray, J., Shou, Z., Chang, S., Paluri, M., 2017. Convnet architecture search for spatiotemporal feature learning. CoRR abs/1708.05038.

Tran, D., Wang, H., Torresani, L., Ray, J., LeCun, Y., Paluri, M., 2018. A closer look at spatiotemporal convolutions for action recognition, in: CVPR, pp. 6450–6459.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I., 2017. Attention is all you need, in: Guyon, I., von Luxburg, U., Bengio, S., Wallach, H.M., Fergus, R., Vishwanathan, S.V.N., Garnett, R. (Eds.), Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pp. 5998–6008.

Wang, C., Cai, H., Zou, Y., Xiong, Y., 2021a. RGB stream is enough for temporal action detection. CoRR abs/2107.04362.

Wang, L., Huang, B., Zhao, Z., Tong, Z., He, Y., Wang, Y., Wang, Y., Qiao, Y., 2023. VideoMAE V2: Scaling video masked autoencoders with dual masking, in: CVPR.

Wang, L., Li, W., Li, W., Gool, L.V., 2018a. Appearance-and-relation networks for video classification, in: CVPR, pp. 1430–1439.

Wang, L., Tong, Z., Ji, B., Wu, G., 2021b. TDN: temporal difference networks for efficient action recognition, in: CVPR, pp. 1895–1904.

Wang, L., Xiong, Y., Lin, D., Gool, L.V., 2017. Untrimmednets for weakly supervised action recognition and detection, in: CVPR, pp. 6402–6411.

Wang, L., Xiong, Y., Wang, Z., Qiao, Y., Lin, D., Tang, X., Gool, L.V., 2016. Temporal segment networks: Towards good practices for deep action recognition, in: ECCV, pp. 20–36.

Wang, X., Girshick, R.B., Gupta, A., He, K., 2018b. Non-local neural networks, in: 2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018, Computer Vision Foundation / IEEE Computer Society. pp. 7794–7803.

Wu, J., Sun, P., Chen, S., Yang, J., Qi, Z., Ma, L., Luo, P., 2021. Towards high-quality temporal action detection with sparse proposals. CoRR abs/2109.08847.

Xie, S., Sun, C., Huang, J., Tu, Z., Murphy, K., 2018. Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification, in: Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y. (Eds.), ECCV, pp. 318–335.

Xu, H., Das, A., Saenko, K., 2017. R-C3D: region convolutional 3d network for temporal activity detection, in: ICCV.

Xu, M., Zhao, C., Rojas, D.S., Thabet, A.K., Ghanem, B., 2020. G-TAD: sub-graph localization for temporal action detection, in: CVPR, Computer Vision Foundation / IEEE. pp. 10153–10162.

Yang, L., Peng, H., Zhang, D., Fu, J., Han, J., 2020. Revisiting anchor mechanisms for temporal action localization. IEEE Transactions on Image Processing 29, 8535–8548.

Zeng, R., Huang, W., Tan, M., Rong, Y., Zhao, P., Huang, J., Gan, C., 2019. Graph convolutional networks for temporal action localization, in: ICCV.

Zhang, C., Wu, J., Li, Y., 2022. Actionformer: Localizing moments of actions with transformers. CoRR abs/2202.07925.

Zhao, Y., Xiong, Y., Wang, L., Wu, Z., Tang, X., Lin, D., 2020. Temporal action detection with structured segment networks. Int. J. Comput. Vis. 128, 74–95.

Zheng, Z., Wang, P., Liu, W., Li, J., Ye, R., Ren, D., 2020. Distance-iou loss: Faster and better learning for bounding box regression, in: AAAI, pp. 12993–13000.