

# Modeling Temporal Concept Receptive Field Dynamically for Untrimmed Video Analysis

Zhaobo Qi<sup>1,2</sup>, Shuhui Wang<sup>2,\*</sup>, Chi Su<sup>3</sup>, Li Su<sup>1,\*</sup>, Weigang Zhang<sup>4</sup>, Qingming Huang<sup>1,2,5</sup>

<sup>1</sup> University of Chinese Academy of Sciences, Beijing, China

<sup>2</sup> Key Lab of Intell. Info. Process., Inst. of Comput. Tech., CAS, Beijing, China

<sup>3</sup> Kingsoft Cloud, Beijing, China

<sup>4</sup> Harbin Inst. of Tech, Weihai, China

<sup>5</sup> Peng Cheng Laboratory, Shenzhen, China

zhaobo.qi@vipl.ict.ac.cn, wangshuhui@ict.ac.cn, suchi@kingsoft.com, suliu@ucas.ac.cn, wgzhang@hit.edu.cn, qmhuang@ucas.ac.cn

## ABSTRACT

Event analysis in untrimmed videos has attracted increasing attention due to the application of cutting-edge techniques such as CNN. As a well studied property for CNN-based models, the receptive field is a measurement for measuring the spatial range covered by a single feature response, which is crucial in improving the image categorization accuracy. In video domain, video event semantics are actually described by complex interaction among different concepts, while their behaviors vary drastically from one video to another, leading to the difficulty in concept-based analytics for accurate event categorization. To model the concept behavior, we study temporal concept receptive field of concept-based event representation, which encodes the temporal occurrence pattern of different mid-level concepts. Accordingly, we introduce temporal dynamic convolution (TDC) to give stronger flexibility to concept-based event analytics. TDC can adjust the temporal concept receptive field size dynamically according to different inputs. Notably, a set of coefficients are learned to fuse the results of multiple convolutions with different kernel widths that provide various temporal concept receptive field sizes. Different coefficients can generate appropriate and accurate temporal concept receptive field size according to input videos and highlight crucial concepts. Based on TDC, we propose the temporal dynamic concept modeling network (TDCMN) to learn an accurate and complete concept representation for efficient untrimmed video analysis. Experiment results on FCVID and ActivityNet show that TDCMN demonstrates adaptive event recognition ability conditioned on different inputs, and improve the event recognition performance of Concept-based methods by a large margin. Code is available at <https://github.com/qzhh/TDCMN>.

## CCS CONCEPTS

• **Computing methodologies** → **Activity recognition and understanding**; *Knowledge representation and reasoning*.

\*Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

MM '20, October 12–16, 2020, Seattle, WA, USA

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-7988-5/20/10...\$15.00

<https://doi.org/10.1145/3394171.3413618>

## KEYWORDS

Temporal Concept Receptive Field, Event Recognition

### ACM Reference Format:

Zhaobo Qi<sup>1,2</sup>, Shuhui Wang<sup>2,\*</sup>, Chi Su<sup>3</sup>, Li Su<sup>1,\*</sup>, Weigang Zhang<sup>4</sup>, Qingming Huang<sup>1,2,5</sup>. 2020. Modeling Temporal Concept Receptive Field Dynamically for Untrimmed Video Analysis. In *Proceedings of the 28th ACM International Conference on Multimedia (MM '20)*, October 12–16, 2020, Seattle, WA, USA. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3394171.3413618>

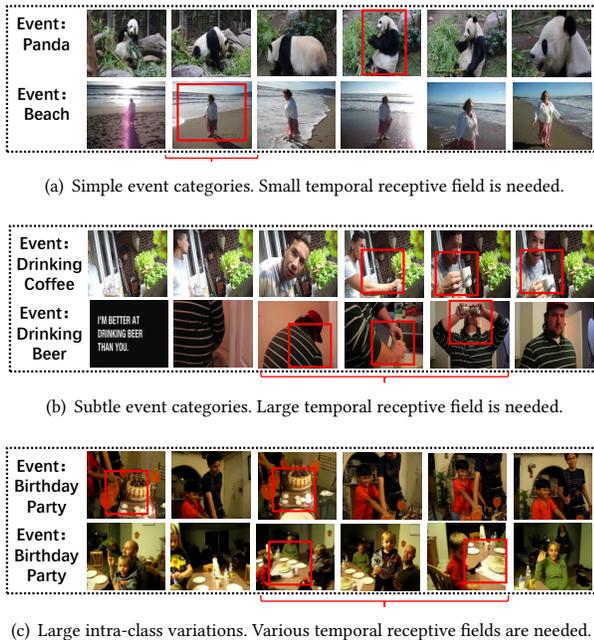
## 1 INTRODUCTION

The receptive field is defined as the region in the input space that a particular CNN feature is looking at [16]. It has been studied extensively to assist us to construct more expressive neural networks in recent years for image recognition [32], semantic segmentation [43], object detection [24], etc. Most methods have been developed to deal with the spatial feature for image-based tasks. However, in video domain, this issue has been less studied in literature.

On the other hand, video analysis has attracted considerable attention due to its extensive applications. Neural-network-based methods have been proposed to promote the development of this field, such as [4, 12, 13, 20, 33, 34]. These methods extract feature representations directly from frames or optical flow, which is recognized as Appearance-based methods. The construction of these models inherits from image-based analysis models and considers the characteristics of video data itself.

Moreover, in order to produce more interpretable recognition results, some researches have focused on Concept-based event recognition methods [1, 5–7, 11, 36, 39, 40]. These models harness simple aggregation method to get video-level concept representations. They ignore the temporal characteristics of concept representations, which limits the feature representation ability. Essentially, the temporal concept receptive field has not been fully exploited, which results in a significant performance gap between Concept-based and Appearance-based event recognition methods.

The various temporal concept receptive field sizes related to the event categories to be classified is another key issue. For some simple event categories (e.g., ‘panda’ and ‘beach’, as shown in Figure 1(a)), crucial concepts existing in short time duration might be sufficient to recognize them. Therefore, the actually suitable temporal receptive field size is small. For some event categories that are subtle events (e.g. ‘drinking coffee’ and ‘drinking beer’, as shown in Figure 1(b)), concepts in a long time duration should be captured to differentiate them, which needs a large temporal receptive field size. This also holds for samples within the same



**Figure 1: The difficulty of making recognition decisions relates to the event category to be classified. For different cases, the needed temporal concept receptive field size for capturing crucial concepts (marked as red bounding box) and differentiating them are varied.**

event category due to large intra-event variations. For example, the concepts ‘cakes’ useful for recognizing the event ‘birthday party’ (as shown in Figure 1(c)) may occur at different time stamps, such as the beginning, the end, or throughout the video. Hence, multiple temporal receptive field sizes for different input videos tend to be more in demand. In a nutshell, a model that can adjust its temporal concept receptive field adaptatively based on different inputs will be more appropriate for flexible and accurate event recognition.

Considering the above-mentioned issues, we explore the temporal concept receptive field of concept-based event recognition models in this paper. As shown in Figure 2, we propose the temporal dynamic convolution (TDC) to give stronger flexibility to concept-based event recognition methods. TDC can adjust the temporal concept receptive field size dynamically based on different inputs. The goal of TDC is to learn a set of coefficients to fuse the results of multiple convolutions with different kernel widths that can provide various temporal concept receptive field sizes. Based on the generated coefficients, we can construct appropriate and accurate temporal concept receptive field size for corresponding videos, which can assist us to highlight crucial concepts for efficient video recognition.

Based on the temporal dynamic convolution, we propose temporal dynamic concept modeling network (TDCMN), a general Concept-based event recognition model augmented with the temporal dynamic convolution. Considering the existence of unique temporal pattern of each concept and the relationship between

different types of concepts, our TDCMN contains two key elements, an intra-domain temporal dynamic concept modeling network (InTDCM) and a cross-domain temporal dynamic concept modeling network (CrTDCM). In InTDCM, we apply TDC on each type of concept representation separately. In CrTDCM, we construct two modeling pipelines, one of which is to apply TDC on the concatenation of all types of concept representations, and the other is to extend TDC to cross-domain temporal dynamic convolution. Our TDCMN can learn an accurate and complete concept representation for efficient untrimmed video analysis.

We conduct extensive experiments on two large-scale and challenging video datasets FCVID and ActivityNet. Experiment results show that TDCMN can achieve adaptive inference conditioned on input videos by adjusting its temporal concept receptive field size. Our TDCMN can improve the performance of Concept-based event recognition methods by a large margin. Besides, it can also obtain higher performance compared with some Appearance-based event recognition methods.

The contributions of this paper are three-fold.

- We propose to analysis video event from the temporal concept receptive field. Accordingly, we propose the dynamic temporal convolution, which can adjust temporal concept receptive field based on input adaptatively.
- We propose the temporal dynamic concept modeling network (TDCMN), a general concept-based event recognition model augmented with the temporal dynamic convolution. TDCMN can learn an accurate and complete concept representation for efficient untrimmed video analysis.
- TDCMN achieves adaptive inference conditioned on input videos by adjusting its temporal concept receptive field size. It can improve the performance of Concept-based methods by a large margin and obtain higher performance than some Appearance-based methods on two challenging datasets.

## 2 RELATED WORK

We review related works from two aspects: video event recognition methods, dynamic receptive field.

### 2.1 Video Event Recognition Methods

**Appearance-based methods.** With the rapid development of deep convolution neural networks, plenty of works have been proposed to learn vision and motion features via deep models for video event recognition, *e.g.*, 2D-CNN-based methods [20], two-stream-based methods [13, 34] and 3D-CNN-based methods (C3D [33], I3D [4] and SlowFast [12]). These methods can achieve high performance, but lack of interpretability. Therefore, we focus on more interpretable event recognition models in this paper, and briefly review its recent process below.

**Concept-based methods.** For concept-based methods, the basic is the selection of the concepts. Some define event-driven concepts [8, 42], while others focus on using manually chosen concepts or concept libraries [1, 5–7, 11, 36, 39, 40]. Ye *et al.* [42] build a large scale event-specific concept library EventNet that covers as many real-world events and concepts as possible. In this paper, we consider three key elements when we select concepts, which can ensure the versatility and practicability of our method. First, the

concept detectors can be obtained directly without training from scratch. Second, different types of concept detectors must be used. Third, the concept categories of each type of concept detector must be varied and common.

As for event analysis methods, Chang *et al.* [7] propose a semantic pooling approach which learns the relation between concepts and events through skip-gram model, and the concept relation is used to prioritize the video shot representations. Xu *et al.* [39] build up a multiple feature learning framework for complex event detection. Wu *et al.* [36] introduce Object-Scene semantic Fusion (OSF) network for large-scale video understanding. Compared to these methods, we analysis concept representation from a new perspective, which is the temporal concept receptive field. We propose a temporal dynamic convolution (TDC), which can adjust the temporal concept receptive field adaptively based on different input videos and guarantee more flexible and effective concept modeling.

## 2.2 Dynamic Receptive Field

The receptive field is defined as the region in the input space that a particular CNN feature is looking at (i.e. be affected by)[16]. It is crucial for constructing expressive or light-weight convolution neural networks.

Some prior works increase the receptive field of neural networks by constructing deeper models [15, 31] with downsampling such as pooling operation or dilated convolution [43]. These works can accumulate some fixed-sized receptive fields at the expense of high computational burden to increase the feature expression. Some prior works [2, 10, 25, 45] focus on directly learning the convolution kernel parameters or the location offsets of convolution kernels. These works usually require a great number of parameters and are difficult to extend to multiple-layers neural networks.

Some prior works [9, 23, 41] focus on predicting the coefficients to combine multiple static convolution kernels. These methods can make models more expressive or reduce redundant calculations in convolution neural networks. Especially, Chen *et al.* [9] present a dynamic convolution to aggregate multiple parallel convolution kernels dynamically based on their attention, which can increase model expressive without increasing the network depth or width for constructing light-weight conventional neural networks. They use multiple convolution filters with the same kernel width. Instead, we utilize multiple convolution filters with different kernel widths in our work, which can provide multi-scale information. Li *et al.* [23] propose the SKNet to fuse the results of multiple convolution filters with different kernel sizes, which allows the neural network to adjust its receptive field size based input information adaptively. Inspired by this idea, we develop the dynamic convolution for video event analysis. In contrast to these prior methods, there are several differences. We utilize the dynamic convolution to analysis temporal information for video task. Instead, they use it to process spatial information for image-based tasks. Instead of injecting the dynamic convolution on all layers of deep neural networks, our proposed dynamic convolution is only utilized as a small but crucial part of our model, which can not increase significant computational burden. Moreover, in terms of calculating the fusion coefficients, we take into account the relationships within and between different convolution filters, which can make our model more flexible and

expressive. Instead, prior works only consider the relation between different convolution filters.

## 3 PROPOSED METHODOLOGY

In this section, we will start with the general framework of concept-based event recognition model in Section 3.1, which is the basic model of the following analysis. Then, we will introduce the temporal dynamic convolution (TDC) in Section 3.2. Finally, we propose the temporal dynamic concept modeling network (TDCMN) in Section 3.3.

### 3.1 Concept-based Event Recognition Model

The general framework of concept-based event recognition model is described below. First, each video  $V$  is evenly divided into  $N$  clips  $\{C_1, C_2, \dots, C_N\}$  and  $M$  frames are randomly sampled from each clip. Second, different types of concept detectors  $\{D_1, D_2, \dots, D_O\}$  are applied on all video clips to capture the initial concept representations. Specially, for each concept detector  $D_i$ , we feed all video clips into it and obtain the initial concept representation  $X_i = [x_{i,1}, x_{i,2}, \dots, x_{i,N}] \in \mathbb{R}^{L_i \times N}$ , where  $L_i$  is the concept category number of  $D_i$ . We will show how to select and use these different types of concept detectors in detail in Section 4.2. After that, the aggregation method such as max pooling is applied on  $X_i$  to obtain the aggregated concept representation. Next, we obtain the video level concept representation  $X \in \mathbb{R}^{1 \times L}$  by concatenating all types of aggregated concept representations, where  $L = \sum_{i=1}^O L_i$  is the total number of concept category of all concept detectors. Finally, a simple MLP network is used to predict the event categories. For the sake of simplicity, we use two types of concept detectors  $\{D_1, D_2\}$  and obtain two types of initial concept representations  $X_1 \in \mathbb{R}^{L_1 \times N}$  and  $X_2 \in \mathbb{R}^{L_2 \times N}$  in the following sections.

### 3.2 Temporal Dynamic Convolution

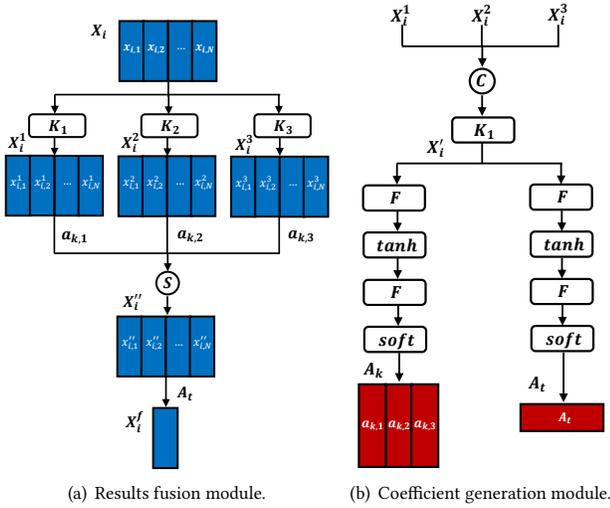
We describe the temporal dynamic convolution (TDC) in this section, and its framework is shown in Figure 2. Our main goal is to learn a set of coefficients to fuse the results of multiple convolutions with different kernel widths that provide plenty of temporal receptive field sizes. Therefore, TDC can adjust the temporal concept receptive field size dynamically based on the different input concept representations.

Especially, given a type of concept representation  $X_i \in \mathbb{R}^{L_i \times N}$ , and three 1-d convolution filters with different kernel widths  $K_1 \in \mathbb{R}^1, K_2 \in \mathbb{R}^3, K_3 \in \mathbb{R}^5$  that provide different temporal receptive fields, we separately perform three 1-d convolutions on  $X_i$  and obtain the feature representations  $X_i^1 \in \mathbb{R}^{L_i \times N}, X_i^2 \in \mathbb{R}^{L_i \times N}, X_i^3 \in \mathbb{R}^{L_i \times N}$  through,

$$X_i^j = Conv(X_i, K_j), \quad j \in \{1, 2, 3\} \quad (1)$$

where  $Conv$  represents a 1-d convolution operation,  $K_j$  is the parameters of the 1-d convolution filter ( $j$  is the convolution kernel index). In fact, we can utilize different numbers and sizes of convolution filters. We will give detail analysis in Section 4.3.

Based on the results of these three convolutions, we employ a coefficient generation module to produce a set of coefficients and a results fusion module to fuse the convolution results.



**Figure 2: Temporal dynamic convolution.**  $K_i$  represents a 1d convolution operation.  $S$ ,  $C$ ,  $F$ ,  $\tanh$  represent summation operation, concatenation operation, fully connected layer, tanh function and softmax function, separately.

**3.2.1 Coefficient Generation Module.** The coefficient generation module is used to generate a set of coefficients based on the multiple convolution results, and its framework is shown in Figure 2(b).

In order to generate the coefficients more flexible, we concatenate the results of all convolutions and apply a 1d convolution with kernel width 1 to reduce its channel number and obtain the feature representation  $X'_i \in \mathbb{R}^{L_i \times N}$  through,

$$X'_i = \text{Conv}(\text{concat}(X_i^1, X_i^2, X_i^3), K_1) \quad (2)$$

where  $\text{Conv}$  and  $\text{concat}$  represent a 1-d convolution operation and concatenation operation, separately.

Based on  $X'_i$ , we will generate two sets of coefficients, which indicate the importance of different representations from two perspectives. They will be used to merge the multiple convolution results dynamically. It can assist us to capture appropriate and accurate temporal concept receptive field size for corresponding videos and highlight crucial concepts for efficient video recognition.

For one thing, we take into account the relations of all channels within the same convolution results and the relations of the same channel among all convolution results to generate a channel coefficient matrix  $A_k = [a_{k,1}, a_{k,2}, a_{k,3}] \in \mathbb{R}^{L_i \times 3}$  through,

$$A_k = \text{softmax}(\tanh(X'_i W_{k,1}) W_{k,2}) \quad (3)$$

where  $W_{k,1} \in \mathbb{R}^{N \times n}$ ,  $W_{k,2} \in \mathbb{R}^{n \times 3}$  are learnable parameters.  $A_k$  will be used to highlight the importance of different channels of each convolution results. For another, we consider all the convolution results and generate a time coefficient matrix  $A_t \in \mathbb{R}^{1 \times N}$  through,

$$A_t = \text{softmax}(W_{t,2} \tanh(W_{t,1} X'_i)) \quad (4)$$

where  $W_{t,1} \in \mathbb{R}^{1 \times L_i}$ ,  $W_{t,2} \in \mathbb{R}^{1 \times l}$  are learnable parameters.  $A_t$  will be used to fuse the feature representation  $X''_i$  through temporal dimension to capture the video level concept representation.

**3.2.2 Results Fusion Module.** The framework of results fusion models is shown in Figure 2(a). Given the obtained convolution results  $\{X_i^1, X_i^2, X_i^3\}$  and generated coefficients matrix  $\{A_k, A_t\}$ , we first apply  $A_k$  on  $\{X_i^1, X_i^2, X_i^3\}$  and sum the results through channel dimension to obtain the feature representation  $X''_i \in \mathbb{R}^{L_i \times N}$  through,

$$\begin{aligned} \hat{x}_{i,t}^j &= a_{k,j,t} \cdot x_{i,t}^j; \quad j \in \{1, 2, 3\}, t \in \{1, 2, \dots, L_i\} \\ X''_i &= \sum_{j=1}^3 \hat{X}_i^j \end{aligned} \quad (5)$$

where  $a_{k,j} = [a_{k,j,1}; a_{k,j,2}; \dots; a_{k,j,L_i}]$ ,  $X_i^j = [x_{i,1}^j; x_{i,2}^j; \dots; x_{i,L_i}^j]$  and  $\hat{X}_i^j = [\hat{x}_{i,1}^j; \hat{x}_{i,2}^j; \dots; \hat{x}_{i,L_i}^j] \in \mathbb{R}^{L_i \times N}$ .  $\cdot$  is scalar-multiplication operation. After that, we capture the video level concept representation  $X_i^f \in \mathbb{R}^{1 \times L_i}$  by matrix multiplication through,

$$X_i^f = A_t X''_i{}^\top \quad (6)$$

### 3.3 Temporal Dynamic Concept Modeling Network

Based on the introduced temporal dynamic convolution, we propose the temporal dynamic concept modeling network (TDCMN) for concept-based event recognition in this section.

Our TDCMN consists of an intra-domain temporal dynamic concept modeling network (InTDCM) and a cross-domain temporal dynamic concept modeling network (CrTDCM), which utilizes TDC to capture the accurate and complete concept representation within and between different types of concepts. As described in Section 3.1, we can obtain the initial concept representations  $X_1 \in \mathbb{R}^{L_1 \times N}$  and  $X_2 \in \mathbb{R}^{L_2 \times N}$ . We feed  $X_1$  and  $X_2$  into the InTDCM and CrTDCM to capture intra-domain dynamic concept representation  $X^{in}$  and cross-domain dynamic concept representation  $X^{cr}$ . We will give the detail of InTDCM and CrTDCM in the following subsections.

**3.3.1 Intra-domain Temporal Dynamic Concept Modeling Network.** In InTDCM, given  $\{X_1, X_2\}$ , we fed them into TDC separately and obtain the initial intra-domain dynamic concept representation  $\{X_1^{in}, X_2^{in}\}$ . Finally, we simply concatenate them and obtain the final intra-domain dynamic concept representation  $X^{in} \in \mathbb{R}^{1 \times L}$ .

**3.3.2 Cross-domain Temporal Dynamic Concept Modeling Network.** For CrTDCM, we exploit two methods to get the final cross-domain dynamic concept representation ( $X_{si}^{cr}$  or  $X_{co}^{cr}$ ) for event recognition. The simplest and most direct way is to concatenate the initial concept representations  $\{X_1, X_2\}$  and feed it to the TDC. As described in Section 3.3.1, we will obtain the final cross-domain dynamic concept representation  $X_{si}^{cr}$ . We term this method as CrTDCM<sub>si</sub>.

Furthermore, we construct a more efficient module to capture the cross-domain dynamic concept representation, and its framework is shown in Figure 3. Specially, we feed  $X_1$  and  $X_2$  into three 1-d convolution layers with different kernel widths and obtain two types of concept representations  $\{X_1^1, X_1^2, X_1^3\} \in \mathbb{R}^{L_1 \times N}$  and  $\{X_2^1, X_2^2, X_2^3\} \in \mathbb{R}^{L_2 \times N}$ . Then we produce three sets of coefficients to fuse the convolution results dynamically and obtain the final cross-domain dynamic concept representation.

First, we fuse the concept representation of each type through element-wise summation and then concatenate them to obtain the

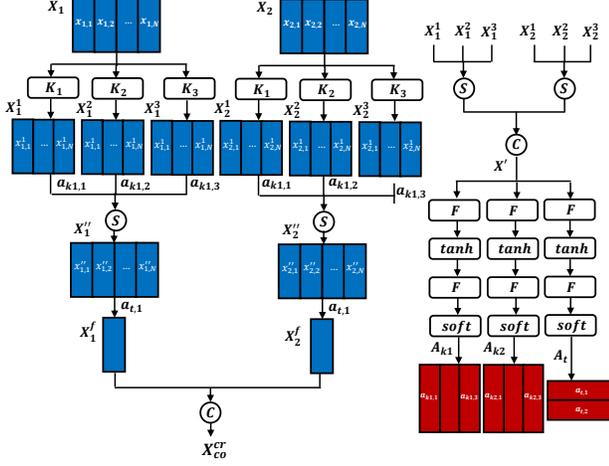


Figure 3: Cross-domain temporal dynamic convolution.

feature representation  $X' \in \mathbb{R}^{L \times N}$  for generating coefficients,

$$X' = \text{concat}(\text{sum}(X_1^1, X_1^2, X_1^3), \text{sum}(X_2^1, X_2^2, X_2^3)) \quad (7)$$

where  $\text{sum}$  and  $\text{concat}$  represent element-wise summation and concatenation operations, separately.

And then we generate three sets of coefficients to fuse multiple convolution results based on  $X'$ . On the one hand, we generate two channel coefficient matrix  $A_{k1} \in \mathbb{R}^{L_1 \times 3}$  and  $A_{k2} \in \mathbb{R}^{L_2 \times 3}$  through,

$$\begin{aligned} A_{k1} &= \text{softmax}(\text{tanh}(X' W_{k1,1}) W_{k1,2}) \\ A_{k2} &= \text{softmax}(\text{tanh}(X' W_{k2,1}) W_{k2,2}) \end{aligned} \quad (8)$$

where  $W_{k1,1} \in \mathbb{R}^{N \times n}$ ,  $W_{k1,2} \in \mathbb{R}^{n \times 3}$ ,  $W_{k2,1} \in \mathbb{R}^{N \times n}$ ,  $W_{k2,2} \in \mathbb{R}^{n \times 3}$  are learnable parameters. We take into account the relations within and between different types of concept representations when we calculate  $A_{k1}$  and  $A_{k2}$ . They will be used to indicate the importance of different channels for each convolution results in different types of concepts, separately. On the other hand, we also generate a time coefficient matrix  $A_t \in \mathbb{R}^{2 \times L_i}$  through,

$$A_t = \text{softmax}(W_{t,2} \text{tanh}(W_{t,1} X')) \quad (9)$$

where  $W_{t,1} \in \mathbb{R}^{L \times L}$ ,  $W_{t,2} \in \mathbb{R}^{2 \times l}$  are learnable parameters.  $A_t$  will be used to fuse the concept representation through temporal dimensionation to get the video level concept representation.

Next, we apply  $A_{k1}$  and  $A_{k2}$  on  $\{X_1^1, X_1^2, X_1^3\}$  and  $\{X_2^1, X_2^2, X_2^3\}$  to obtain the feature representation  $X_1'' \in \mathbb{R}^{L_1 \times N}$  and  $X_2'' \in \mathbb{R}^{L_2 \times N}$ , separately,

$$\begin{aligned} \hat{x}_{1,t}^j &= a_{k1,j,t} \cdot x_{1,t}^j; \quad j \in \{1, 2, 3\}, t \in \{1, 2, \dots, L_1\} \\ \hat{x}_{2,t}^j &= a_{k2,j,t} \cdot x_{2,t}^j; \quad j \in \{1, 2, 3\}, t \in \{1, 2, \dots, L_2\} \\ X_1'' &= \sum_{j=1}^3 \hat{X}_1^j; \quad X_2'' = \sum_{j=1}^3 \hat{X}_2^j \end{aligned} \quad (10)$$

where  $X_1^j = [x_{1,1}^j; x_{1,2}^j; \dots; x_{1,L_1}^j]$ ,  $X_2^j = [x_{2,1}^j; x_{2,2}^j; \dots; x_{2,L_2}^j]$ ,  $A_{k1} = [a_{k1,1}, a_{k1,2}, a_{k1,3}]$ ,  $A_{k2} = [a_{k2,1}, a_{k2,2}, a_{k2,3}]$ ,  $\hat{X}_1^j = [\hat{x}_{1,1}^j; \hat{x}_{1,2}^j; \dots; \hat{x}_{1,L_1}^j]$ ,  $\hat{X}_2^j = [\hat{x}_{2,1}^j; \hat{x}_{2,2}^j; \dots; \hat{x}_{2,L_2}^j]$ ,  $\cdot$  represents scalar-multiplication. Finally,

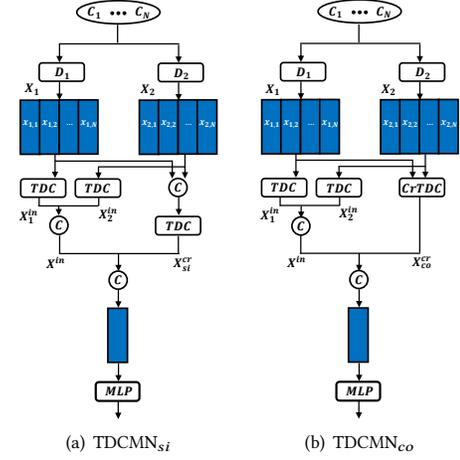


Figure 4: Temporal dynamic concept modeling network.

we capture the final cross-domain dynamic concept representation  $X_{co}^{cr} \in \mathbb{R}^{1 \times L}$  through,

$$X_{co}^{cr} = \text{concat}(a_{t,1} X_1''^\top, a_{t,2} X_2''^\top) \quad (11)$$

where  $A_t = [a_{t,1}, a_{t,2}]$  and  $\text{concat}$  represents concatenation operation. We term this method as CrTDCM<sub>co</sub>.

**3.3.3 Temporal Dynamic Concept Modeling Network.** As discussed in Section 3.3.2, we have constructed two cross-domain temporal dynamic concept modeling networks CrTDCM<sub>si</sub> and CrTDCM<sub>co</sub>, therefore we have two temporal dynamic concept modeling networks, which we call them as TDCMN<sub>si</sub> and TDCMN<sub>co</sub>, and the framework are shown in Figure 4(a) and Figure 4(b). For each model, we concatenate the obtained intra-domain and cross-domain dynamic concept representation and use an MLP network for final event recognition.

## 4 EXPERIMENTS

### 4.1 Experimental Setup

**Evaluation Datasets:** We use two large scale datasets for video event recognition: Fudan-Columbia Video Dataset (FCVID)[18] and ActivityNet Dataset (ActivityNet)[3] in our experiments. The FCVID dataset contains 91,223 videos annotated with 239 event categories, including social events (e.g., ‘tailgate party’) and procedural events (e.g., ‘making a cake’). Some videos can not be downloaded, and there are some corrupted videos that can not be used to extract video frames. We filter out these videos and end up with a training set of 45,416 videos and a testing set of 45,437 videos. The ActivityNet dataset is a large-scale untrimmed video dataset. We use the latest released version of the dataset (v1.3), which contains about 20K videos from 200 activity categories. Note that the annotations of the testing set have not been publicly available, we only use the training and validation set, and report results on the validation set.

**Evaluation Metric:** We compute average precision (AP) for each event class and report the mean Average Precision (mAP) over all event categories.

**Table 1: Ablation study about different types of concept detectors on FCVID.**

Model	mAP (%)
Baseline <sub>so</sub>	72.9
Baseline <sub>soa</sub>	75.3
TDCMN <sub>si,so</sub>	80.6
TDCMN <sub>si,soa</sub>	81.9
TDCMN <sub>co,so</sub>	81.2
TDCMN <sub>co,soa</sub>	82.2

**Implementation Details:** We implement the proposed models based on the Pytorch framework. All proposed models are trained by SGD optimizer. The momentum and weight decay are set to 0.9 and 0.0005, respectively. We set the  $N$  and  $M$  as 16 and 8, respectively. For FCVID dataset, the batch size is set as 12. The initial learning rate is set to 0.5 and drops by 0.1 every 32 epochs. The training procedure stops after 40 epochs. For ActivityNet, the batch size is set as 12. The initial learning rate is set to 0.5 and drops by 0.1 every 40 epochs. The training procedure stops after 48 epochs.

## 4.2 The Selection and Use of Concept Detectors

To guarantee the practicability of our method, we take into account several crucial elements when we pick concept detectors. First, in order to be practical, the concept detectors must be readily available. In other words, we can obtain them directly without training from scratch. Second, to guarantee the representation power of our model, different types of concept detectors must be employed. Third, the category of concepts included in each detector must be varied and universal, which assures the versatility of our model. Therefore, they can represent more event classes, and our model can be applied to more video event datasets.

Based on the above considerations, we select the scene, object and action concept detectors, as these are the essential concepts for recognizing events. All these concept detectors are pre-trained on standard datasets and are readily available. We only use them to obtain initial concept representations. All layers of each detector except the final classification layer are fixed when we train our model. We will show how to use them below.

**Scene Concept Detectors**  $D_s$ . ResNet-50 [15] based  $D_s$  pre-trained on Places365-Standard [47] dataset is utilized.  $D_s$  contains 365 scene classes. The scene concept representation is extracted from the output of the last classification layer. Specially, an event video is evenly divided into  $N$  clips and  $M$  frames are randomly sampled from each clip. For each clip  $C_n$ , the concept representation  $X'_{s,n} \in \mathbb{R}^{365 \times M}$  is acquired by inputting  $C_n$  into  $D_s$ . And then, the clip-level concept representation  $X_{s,t} \in \mathbb{R}^{365}$  is obtained by imposing maximum pooling on  $X'_{s,n}$  over all frames. Lastly, the video-level initial scene concept representation  $X_s \in \mathbb{R}^{365 \times N}$  is obtained by concatenating the representations of all the video clips.

**Object Concept Detectors**  $D_o$ . We use ResNet-50 based object concept detector  $D_o$  pre-trained on ImageNet [28] as  $D_o$ . Similarly, we obtain the video-level initial object concept representation  $X_o \in \mathbb{R}^{1000 \times N}$ , where 1000 is the number of object concepts in  $D_o$ .

**Table 2: Ablation study about modeling intra-domain and cross-domain concept representation on FCVID.**

Model	InTDCM	CrTDCM <sub>si</sub>	CrTDCM <sub>co</sub>	mAP (%)
Baseline <sub>so</sub>				72.9
+InTDCM	✓			77.9
+CrTDCM <sub>si</sub>		✓		78.4
TDCMN <sub>si,so</sub>	✓	✓		<b>80.6</b>
+CrTDCM <sub>co</sub>			✓	77.6
TDCMN <sub>co,so</sub>	✓		✓	<b>81.2</b>

**Action Concept Detectors**  $D_a$ . We use the I3D-based [4]  $D_a$  pre-trained on kinetics [21] dataset. There are 400 action classes in  $D_a$ . We also employ the output of the last classification layer as the action concept representation. We input each video clip  $C_n$  to  $D_a$  and get the clip-level action concept representation  $X_{a,n} \in \mathbb{R}^{400}$ . Finally, we concatenate the representation of all clips and obtain the video-level initial action concept representation  $X_a \in \mathbb{R}^{400 \times N}$ .

## 4.3 Ablation Studies

In this section, we construct abundant experiments to testify the efficiency of each choice of our model.

**4.3.1 The number of concept detectors.** The video comprises a variety of concepts, which work together to describe the event of this video. In this subsection, we conduct six experiments on FCVID to show the efficiency of different types of concepts. The details setting are shown below, and the results are shown in Table 1.

We only use scene and object concept detectors in Baseline<sub>so</sub>. We concatenate the results of these two detectors and use a fully connected layer to perform event recognition. Baseline<sub>soa</sub> is similar to Baseline<sub>so</sub> except that the action concept detectors are also used. TDCMN<sub>si,so</sub> is TDCMN<sub>si</sub> with scene and object concept detectors. TDCMN<sub>si,soa</sub> is TDCMN<sub>si</sub> with scene, object and action concept detectors. TDCMN<sub>co,so</sub> is TDCMN<sub>co</sub> with use scene and object concept detectors. TDCMN<sub>co,soa</sub> is TDCMN<sub>co</sub> with scene, object and action concept detectors.

By comparing the results of Baseline<sub>so</sub> and Baseline<sub>soa</sub> (or the results of TDCMN<sub>si,so</sub> and TDCMN<sub>si,soa</sub> or the results of TDCMN<sub>co,so</sub> and TDCMN<sub>co,soa</sub>), we can find that the event recognition performance is higher when more concept detectors are used. This is because more concept detectors can give us richer concept representations, which can describe the event more comprehensiveness and accuracy. By comparing the results of our model with baseline methods, we can find that the efficiency of our model is testified no matter how many concept detectors are used. Therefore, we will only use scene and object concept detectors in the following experiments of this section.

**4.3.2 The necessary of modeling intra-domain and cross-domain concept representations.** We propose two temporal dynamic concept modeling networks TDCMN<sub>si</sub> and TDCMN<sub>co</sub>. Each network models both intra-domain and cross-domain concept representations. We construct experiments to demonstrate the necessity

**Table 3: Ablation study about temporal dynamic convolution on FCVID.**

Model	mAP (%)
w/o TDC	76.0
with TDC	<b>77.9</b>

**Table 4: Ablation study about convolution filters with different kernel widths on FCVID.  $K_1, K_2, K_3$  and  $K_4$  represent 1-d convolution filters with kernel width 1, 3, 5 and 7, separately.**

Exp.	$K_1$	$K_2$	$K_3$	$K_4$	mAP (%)
1	✓	✓			76.8
2	✓		✓		76.7
3		✓	✓		77.4
4	✓	✓	✓		<b>77.9</b>
5	✓	✓	✓	✓	77.8

of modeling intra-domain and cross-domain concept representations, and the results are shown in Table 2. We can see that the InTDCM and the CrTDCM<sub>si</sub> (or CrTDCM<sub>co</sub>) both can improve the event recognition performance to some extent compared to Baseline<sub>so</sub>, which indicates the efficiency of each module. Finally, the event recognition performance has further promoted by modeling the intra-domain and cross-domain concept representations at the same time in TDCMN<sub>si,so</sub> (or TDCMN<sub>co,so</sub>). This phenomenon verifies the necessary and complementarity of modeling intra-domain and cross-domain concept representations. Hence, we will only conduct experiments on InTDCM in the following subsections.

**4.3.3 Temporal dynamic convolution.** The core of this paper is the proposed temporal dynamic convolution, we prepare experiments to verify its effectiveness, and the results are shown in Table 3. According to the previous analysis, we construct experiments on InTDCM to see the performance change when we remove the temporal dynamic convolution from it. From Table 3, we can see that the event recognition performance dropped by 1.9% mAP, which proves the effectiveness of the temporal dynamic convolution. The proposed TDC can generate two sets of coefficients, which indicates the importance of different channels within and between each convolution results and the importance of different video clips. Therefore, we can obtain crucial concepts through TDC.

**4.3.4 The number of convolution filters with different kernel widths.** Up to now, we verify the necessity of modeling intra-domain and cross-domain concept representations and the effectiveness of temporal dynamic convolution. Therefore, we construct experiments on InTDCM to see the event recognition performance change when a different number of convolution filters with different kernel widths are used. The results are shown in Table 4. We can find that event recognition performance is increased when we use more convolution filters. It is mainly because that convolution with different kernel widths can provide different sizes of temporal concept receptive field. We can obtain various temporal receptive field sizes for different concepts. Hence, we can find the appropriate temporal concept receptive field size and capture useful concepts

**Table 5: The comparison between our method and other Concept-based methods on FCVID and ActivityNet dataset.**

Method	FCVID	ActivityNet
Early Fusion-NN [36]	75.2	55.9
Fusion-SVM [36]	75.5	55.8
SVM-MKL [22]	74.9	56.3
OSF [36]	76.5	56.8
TDCMN <sub>si,soa</sub>	81.9	84.3
TDCMN <sub>co,soa</sub>	<b>82.2</b>	<b>84.6</b>

**Table 6: The comparison between our method and other Appearance-based methods on FCVID and ActivityNet.**

Method	FCVID	Method	ActivityNet
DMF [29]	72.5	AdaFrame [38]	71.57
DASD [17]	72.8	LiteEval [37]	72.7
M-DBM [30]	74.4	TSN [34]	76.6
rDNN [18]	76.0	KeylessAttention [26]	78.5
GSFMN-all [46]	76.9	P3D [27]	78.9
Pivot CorrNN [19]	77.6	MLSME [44]	83.0
LiteEval [37]	80.0	MARL-based [35]	83.8
AdaFrame [38]	80.2	IMGAUD2VID [14]	84.2
TDCMN <sub>si,soa</sub>	81.9	TDCMN <sub>si,soa</sub>	84.3
TDCMN <sub>co,soa</sub>	<b>82.2</b>	TDCMN <sub>co,soa</sub>	<b>84.6</b>

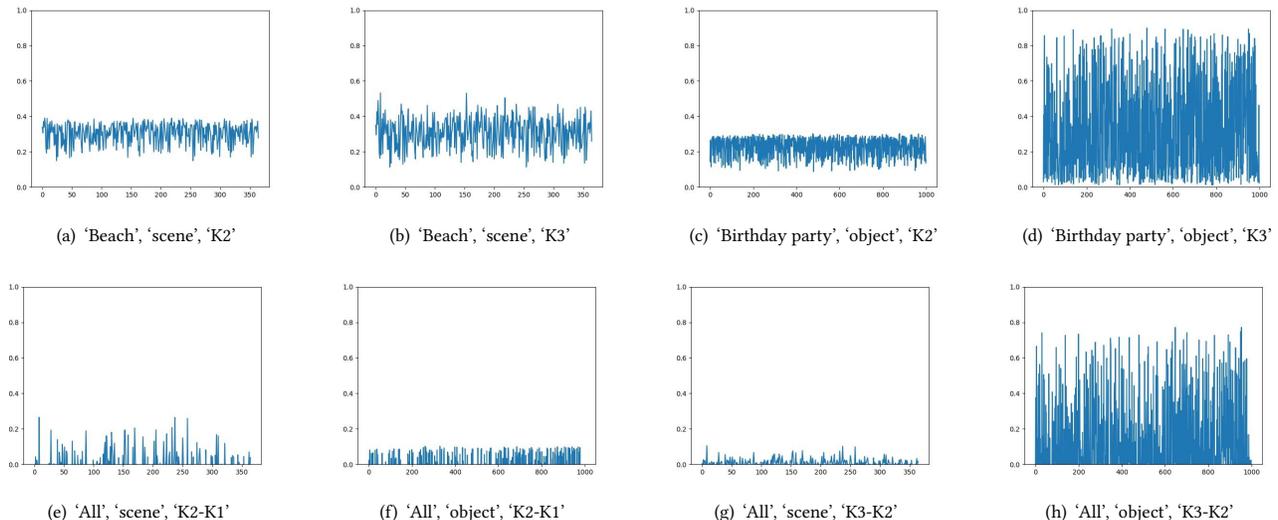
for event recognition. Besides, by comparing the results of Exp.4 and Exp.5 in Table 4, we can find a slight drop in performance when using too many convolution kernels. Therefore, we only use three convolution filters with kernel widths 1, 3, 5 in all experiments.

## 4.4 Comparison to state of the art

We compare our model with both Appearance-based methods and Concept-based methods on FCVID and ActivityNet datasets. The results are shown in Table 5 and Table 6.

**Concept-based Models.** The comparison results between our method and other Concept-based methods on FCVID and ActivityNet datasets are shown in Table 5. As can be seen from the table, our models can improve the performance of Concept-based event recognition models by a large margin, especially on ActivityNet. Particularly, our model is at least 5.4% mAP higher than OSF on FCVID. Though OSF uses vgg features and resnet are used in our model, OSF utilizes the generic vision feature that directly extracted from videos and more than 20000 concepts. Besides, OSF uses average pooling over the concept representations of all video clips and ignores the temporal characteristic of concept representations, which causes performance reduction.

**Appearance-based Models.** The comparison results between our method and some Appearance-based methods on FCVID and ActivityNet datasets are shown in Table 6. We can find that our model achieves higher event recognition performance compared with some Appearance-based methods on both datasets. Besides, our model is also better than Pivot CorrNN [19], which uses seven types of pre-extracted features to perform event recognition.



**Figure 5: Visualized the distribution of coefficients generated by temporal dynamic convolution in InTDCM module. ‘Beach’, ‘scene’, ‘K2’ means this line graph shows the channel coefficient distribution generated for convolution kernel ‘K2’ of ‘scene’ concept on ‘Beach’ category. ‘All’, ‘scene’, ‘K2-K1’ means this line graph shows the difference of the channel coefficients between convolution kernel ‘K2’ and ‘K1’ of ‘scene’ concept on all event categories. The other subgraphs have similar meanings.**

## 4.5 Visualization

To understand the temporal dynamic convolution more clearly, we visualize the distribution of the coefficients generated by TDC in InTDCM module, and the results are shown in Figure 5. On the one hand, we seek to see the behavior of TDC from the perspective of event categories. Especially, we randomly sample two event categories ‘Beach’ and ‘Birthday party’. Then, we calculate the average of the channel coefficients generated for each convolution filters across all validation videos in each category. From Figure 5(a) and 5(b), we can find that the channel coefficient distribution of convolutions with different kernel widths is similar for scene concept of ‘Beach’ event category. Instead, the channel coefficient distribution is quite different for object concept of ‘Birthday party’ (see Figure 5(c) and 5(d)). This phenomenon is consistent with our motivation. For ‘Beach’, the temporal existence patterns of scene concepts at different time scales are similar. Therefore, the coefficient distribution for different convolution results are similar. In contrast, the temporal existence patterns of object concepts at different time scales are quite different for ‘Birthday party’. Hence, the needed temporal receptive field size is different. Accordingly, the learned coefficients of each convolution result are different.

On the other hand, we also exploit the behavior of TDC from the perspective of concept types. Notably, we calculate the difference of the channel coefficients generated for different convolution filters of each type of concept across all validation videos in all event categories. From Figure 5(e) and 5(f), we can find that the difference between K2 and K1 are similar for scene and object concepts. In contrast, the difference between K3 and K2 are different for scene and object concepts by comparing Figure 5(g) and 5(h). This is mainly because different types of concepts have its unique temporal existence pattern. Our model can generate corresponding temporal

concept receptive field sizes adaptively based on different types of concepts.

## 5 CONCLUSION

In this paper, we explore the temporal receptive field of concept-based event recognition methods for efficiently untrimmed video event analysis. First, we introduce a temporal dynamic convolution (TDC) to give stronger flexibility to concept-based event recognition networks, which can adjust its receptive field size adaptively based on different inputs. Based on TDC, we propose the temporal dynamic concept modeling network (TDCMN) to learn an accurate and complete concept representation for efficient untrimmed video analysis. TDCMN employs TDC to analyze the temporal characteristic of concepts within the same type and between different types. To demonstrate the effectiveness of our model, we apply TDCMN on two challenging video datasets FCVID and ActivityNet. TDCMN can improve event recognition performance by a large margin compared with other concept-based event recognition methods.

## ACKNOWLEDGMENTS

This work was supported in part by the National Key R&D Program of China under Grant 2018AAA0102003 and 2018YFE0118400, in part by National Natural Science Foundation of China: 61672497, 61620106009, 61836002, 61931008, 61976069, 61650202 and U1636214, and in part by Key Research Program of Frontier Sciences, CAS: QYZDJ-SSW-SYS013. We acknowledge Kingsoft Cloud for the helpful discussion and free GPU cloud computing resource support.

## REFERENCES

- [1] Subhabrata Bhattacharya, Mahdi M Kalayeh, Rahul Sukthankar, and Mubarak Shah. 2014. Recognition of complex events: Exploiting temporal dynamics between underlying concepts. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2235–2242.
- [2] Egor Burkov and Victor Lempitsky. 2018. Deep Neural Networks with Box Convolutions. (2018), 6211–6221.
- [3] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. 2015. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 961–970.
- [4] Joao Carreira and Andrew Zisserman. 2017. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 6299–6308.
- [5] Xiaojun Chang, Yao-Liang Yu, Yi Yang, and Alexander G Hauptmann. 2015. Searching persuasively: Joint event detection and evidence recounting with limited supervision. In *Proceedings of the 23rd ACM international conference on Multimedia*. ACM, 581–590.
- [6] Xiaojun Chang, Yao-Liang Yu, Yi Yang, and Eric P Xing. 2016. They are not equally reliable: Semantic event search using differentiated concept classifiers. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1884–1893.
- [7] Xiaojun Chang, Yao-Liang Yu, Yi Yang, and Eric P Xing. 2017. Semantic pooling for complex event analysis in untrimmed videos. *IEEE transactions on pattern analysis and machine intelligence* 39, 8 (2017), 1617–1632.
- [8] Jiawei Chen, Yin Cui, Guangnan Ye, Dong Liu, and Shih-Fu Chang. 2014. Event-driven semantic concept discovery by exploiting weakly tagged internet images. In *Proceedings of International Conference on Multimedia Retrieval*. ACM, 1.
- [9] Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu. 2019. Dynamic Convolution: Attention over Convolution Kernels. *arXiv: Computer Vision and Pattern Recognition* (2019).
- [10] Bert De Brabandere, Xu Jia, Tinne Tuytelaars, and Luc Van Gool. 2016. Dynamic filter networks. (2016), 667–675.
- [11] Hehe Fan, Xiaojun Chang, De Cheng, Yi Yang, Dong Xu, and Alexander G Hauptmann. 2017. Complex event detection by identifying reliable shots from untrimmed videos. In *Proceedings of the IEEE International Conference on Computer Vision*. 736–744.
- [12] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. 2019. Slow-fast networks for video recognition. In *Proceedings of the IEEE International Conference on Computer Vision*. 6202–6211.
- [13] Christoph Feichtenhofer, Axel Pinz, and Andrew Zisserman. 2016. Convolutional two-stream network fusion for video action recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1933–1941.
- [14] Ruohan Gao, Taehyun Oh, Kristen Grauman, and Lorenzo Torresani. 2019. Listen to Look: Action Recognition by Previewing Audio. *arXiv: Computer Vision and Pattern Recognition* (2019).
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [16] Dang Ha The Hien. [n.d.]. A Guide to Receptive Field Arithmetic for Convolutional Neural Networks. <https://syncedreview.com/2017/05/11/a-guide-to-receptive-field-arithmetic-for-convolutional-neural-networks/>.
- [17] Yu-Gang Jiang, Qi Dai, Jun Wang, Chong-Wah Ngo, Xiangyang Xue, and Shih-Fu Chang. 2012. Fast semantic diffusion for large-scale context-based image and video annotation. *IEEE Transactions on Image Processing* 21, 6 (2012), 3080–3091.
- [18] Yu-Gang Jiang, Zuxuan Wu, Jun Wang, Xiangyang Xue, and Shih-Fu Chang. 2018. Exploiting feature and class relationships in video categorization with regularized deep neural networks. *IEEE transactions on pattern analysis and machine intelligence* 40, 2 (2018), 352–364.
- [19] Sunghun Kang, Junyeong Kim, Hyunsoo Choi, Sungjin Kim, and Chang D Yoo. 2018. Pivot Correlational Neural Network for Multimodal Video Categorization. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 386–401.
- [20] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. 2014. Large-scale video classification with convolutional neural networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*. 1725–1732.
- [21] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. 2017. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950* (2017).
- [22] Marius Kloft, Ulf Brefeld, Sören Sonnenburg, and Alexander Zien. 2011. Lp-norm multiple kernel learning. *Journal of Machine Learning Research* 12, Mar (2011), 953–997.
- [23] Xiang Li, Wenhai Wang, Xiaolin Hu, and Jian Yang. 2019. Selective Kernel Networks. (2019), 510–519.
- [24] Yanghao Li, Yuntao Chen, Naiyan Wang, and Zhaoxiang Zhang. 2019. Scale-Aware Trident Networks for Object Detection. *arXiv: Computer Vision and Pattern Recognition* (2019).
- [25] Vasileios Lioutas and Yuhong Guo. 2020. Time-aware Large Kernel Convolutions. *arXiv: Learning* (2020).
- [26] Xiang Long, Chuang Gan, Gerard De Melo, Xiao Liu, Yandong Li, Fu Li, and Shilei Wen. 2018. Multimodal keyless attention fusion for video classification. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- [27] Zhaofan Qiu, Ting Yao, and Tao Mei. 2017. Learning spatio-temporal representation with pseudo-3d residual networks. In *proceedings of the IEEE International Conference on Computer Vision*. 5533–5541.
- [28] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. 2015. Imagenet large scale visual recognition challenge. *International journal of computer vision* 115, 3 (2015), 211–252.
- [29] John R Smith, Milind Naphade, and Apostol Natsev. 2003. Multimedia semantic indexing using model vectors. In *2003 International Conference on Multimedia and Expo. ICME'03. Proceedings (Cat. No. 03TH8698)*, Vol. 2. IEEE, II–445.
- [30] Nitish Srivastava and Ruslan R Salakhutdinov. 2012. Multimodal learning with deep boltzmann machines. In *Advances in neural information processing systems*. 2222–2230.
- [31] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. 2016. Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning. (2016), 4278–4284.
- [32] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. 2016. Rethinking the Inception Architecture for Computer Vision. (2016), 2818–2826.
- [33] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. 2015. Learning spatiotemporal features with 3d convolutional networks. In *Proceedings of the IEEE international conference on computer vision*. 4489–4497.
- [34] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. 2016. Temporal segment networks: Towards good practices for deep action recognition. In *European conference on computer vision*. Springer, 20–36.
- [35] Wenhao Wu, Dongliang He, Xiao Tan, Shifeng Chen, and Shilei Wen. 2019. Multi-Agent Reinforcement Learning Based Frame Sampling for Effective Untrimmed Video Recognition. In *Proceedings of the IEEE International Conference on Computer Vision*. 6222–6231.
- [36] Zuxuan Wu, Yanwei Fu, Yu-Gang Jiang, and Leonid Sigal. 2016. Harnessing object and scene semantics for large-scale video understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4489–4497.
- [37] Zuxuan Wu, Caiming Xiong, Yu-Gang Jiang, and Larry S Davis. 2019. LiteEval: A Coarse-to-Fine Framework for Resource Efficient Video Recognition. In *Advances in Neural Information Processing Systems*. 7778–7787.
- [38] Zuxuan Wu, Caiming Xiong, Chih-Yao Ma, Richard Socher, and Larry S Davis. 2019. AdaFrame: Adaptive Frame Selection for Fast Video Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1278–1287.
- [39] Zhongwen Xu, Ivor W Tsang, Yi Yang, Zhigang Ma, and Alexander G Hauptmann. 2014. Event detection using multi-level relevance labels and multiple features. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 97–104.
- [40] Yan Yan, Yi Yang, Haoquan Shen, Deyu Meng, Gaowen Liu, Alex Hauptmann, and Nicu Sebe. 2015. Complex event detection via event oriented dictionary learning. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*.
- [41] Brandon Yang, Gabriel Bender, Quoc V Le, and Jiquan Ngiam. 2019. CondConv: Conditionally Parameterized Convolutions for Efficient Inference. *arXiv: Computer Vision and Pattern Recognition* (2019).
- [42] Guangnan Ye, Yitong Li, Hongliang Xu, Dong Liu, and Shih-Fu Chang. 2015. Eventnet: A large scale structured concept library for complex event detection in video. In *Proceedings of the 23rd ACM international conference on Multimedia*. ACM, 471–480.
- [43] Fisher Yu and Vladlen Koltun. 2016. Multi-Scale Context Aggregation by Dilated Convolutions. (2016).
- [44] Ji Zhang, Kuizhi Mei, Yu Zheng, and Jianping Fan. 2019. Exploiting Mid-Level Semantics for Large-Scale Complex Video Classification. *IEEE Transactions on Multimedia* (2019).
- [45] Linguang Zhang, Maciej Halber, and Szymon Rusinkiewicz. 2019. Accelerating Large-Kernel Convolution Using Summed-Area Tables. *arXiv: Learning* (2019).
- [46] Rui-Wei Zhao, Qi Zhang, Zuxuan Wu, Jianguo Li, and Yu-Gang Jiang. 2019. Visual Content Recognition by Exploiting Semantic Feature Map with Attention and Multi-task Learning. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)* 15, 1s (2019), 6.
- [47] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. 2018. Places: A 10 million image database for scene recognition. *IEEE transactions on pattern analysis and machine intelligence* 40, 6 (2018), 1452–1464.