

# TTPP: Temporal Transformer with Progressive Prediction for Efficient Action Anticipation

Wen Wang, Xiaojiang Peng, Yanzhou Su, Yu Qiao, Jian Cheng

**Abstract**—Video action anticipation aims to predict future action categories from observed frames. Current state-of-the-art approaches mainly resort to recurrent neural networks to encode history information into hidden states, and predict future actions from the hidden representations. It is well known that the recurrent pipeline is inefficient in capturing long-term information which may limit its performance in prediction task. To address this problem, this paper proposes a simple yet efficient Temporal Transformer with Progressive Prediction (TTPP) framework, which repurposes a Transformer-style architecture to aggregate observed features, and then leverages a light-weight network to progressively predict future features and actions. Specifically, predicted features along with predicted probabilities are accumulated into the inputs of subsequent prediction. We evaluate our approach on three action datasets, namely TVSeries, THUMOS-14, and TV-Human-Interaction. Additionally we also conduct a comprehensive study for several popular aggregation and prediction strategies. Extensive results show that TTPP not only outperforms the state-of-the-art methods but also more efficient.

**Index Terms**—Action anticipation, Transformer, Encoder-Decoder.

## I. INTRODUCTION

Human action anticipation, also aka early action prediction, aiming to predict future unseen actions, is one of the main topics in video understanding with wide applications in security, visual surveillance and human-computer interaction, etc. In contrast to the well-studied action recognition, which infers the action label after observing the entire action execution, action anticipation is to early predict human actions without observing the future action execution. It is a very challenging task because the input videos are temporally incomplete with wide variety of irrelevant background, and decisions must be made based on such incomplete action executions. In all, action anticipation needs to overcome all the difficulties of action recognition and capture sufficient historical and contextual information to make future predictions in untrimmed video streams.

Generally, most of the action anticipation approaches can be divided into two key phases, namely observed information aggregation and future prediction, as shown in Figure 1. Early

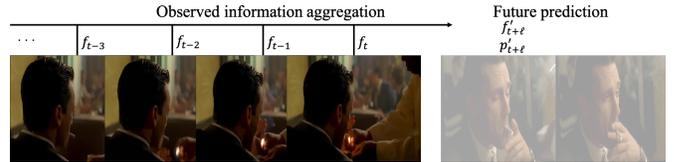


Fig. 1: The summarized generic flowchart for action anticipation, which mainly consists of observed information aggregation and future prediction.

works of action anticipation focus on trimmed action videos, and mainly make efforts on extracting discriminative features from partial videos, *i.e.* observed information aggregation, for early action prediction [1], [28], [30], [42], [57]. In deep learning era, recent works turn to predict future actions in practical untrimmed video streams [11], [13], [56], [62], and mainly repurpose sequential models from the natural language processing (NLP) domain, like long short-term memory (LSTM) [15] and gated recurrent neural networks [6]. For instance, Gao *et al.* [13] propose a Reinforced Encoder-Decoder network, which utilizes an encoder-decoder LSTM network [36], [59] to aggregate historical features and predict future features or actions. Xu *et al.* [62] propose a LSTM-based temporal recurrent network to predict future features for both online action detection and action anticipation. Though the encoder-decoder recurrent networks can be easily transformed from NLP domain to temporal action anticipation, their inherent sequential nature precludes parallelization within training examples and limits the memory power for longer sequence length [52]. Moreover, they are known to have limited improvements in other action understanding tasks [?], [60].

In this paper, we address the two issues of action anticipation via a simple yet efficient *Temporal Transformer with Progressive Prediction* (TTPP) framework. TTPP repurposes a Transformer-style module to aggregate observed information and leverages a light-weight network to progressively predict future features and actions. Specifically, TTPP contains a Temporal Transformer Module (TTM) and an elaborately-designed Progressive Prediction Module (PPM). Given historical and current features, the TTM aggregates the historical features based solely on self-attention mechanisms with the current feature as query, which is inspired by the Transformer in machine translation [52]. The aggregated historical feature along with the current feature are then fed into the PPM. The PPM is comprised of an initial prediction block and a shared-

Wen Wang, Yanzhou Su and Jian Cheng are with the School of Information and Communication Engineering, University of Electronic Science and Technology of China, Chengdu, Sichuan, China, 611731.

Xiaojiang Peng and Yu Qiao are with ShenZhen Key Lab of Computer Vision and Pattern Recognition, SIAT-SenseTime Joint Lab, Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences; SIAT Branch, Shenzhen Institute of Artificial Intelligence and Robotics for Society.

This work was done when Wen Wang was intern at Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences.

Corresponding author: Xiaojiang Peng (xj.peng@siat.ac.cn)

parameter progressive prediction block, each of which is built with two fully-connected (FC) layers, a ReLU activation [38] and a layer normalization (LN) [2]. With the output feature of TTM, the initial prediction block of PPM predicts the immediately following clip feature and action probabilities. The progressive prediction block accumulates the former predictions and the output of TTM, and further predicts a few subsequent future features and actions. The whole TTPP model can be jointly trained in an end-to-end manner with supervision from ground-truth future features and action labels. Compared to previous encoder-decoder methods, the benefits of our TTPP are two-fold. First, the temporal Transformer is more efficient than recurrent methods in capturing historical context by self-attention. Second, the progressive prediction module with skip connections to aggregated historical features can efficiently deliver temporal information and help long-term anticipation. We evaluate our approach on three widely-used action anticipation datasets, namely TVSeries [7], THUMOS-14 [20], and TV-Human-Interaction [39]. Additionally we also conduct a comprehensive study for several popular aggregation and prediction strategies, including temporal convolution, LSTM and single-shot prediction, etc. Extensive results show that TTPP is more efficient than the state-of-the-art methods in both training and inference, and outperforms them with a large margin.

The main contributions of this work can be concluded as follows.

- We propose a simple yet efficient TTPP framework for action anticipation, which leverages a Transformer-style architecture to aggregate information and a light-weight module to predict future actions.
- We elaborately design a progressive prediction module for predicting future features and actions, and achieve the state-of-the-art performance on TVSeries, THUMOS-14, and TV-Human-Interaction.
- We conduct a comprehensive study for several popular aggregation and prediction strategies, including aggregation methods of temporal convolution, Encoder-LSTM, and prediction methods of Decoder-LSTM, single-shot prediction, etc.

The rest of this paper is organized as follows: We first review some related works in Section II. Section III describes the proposed framework with TTM to aggregate observed features and PPM to progressively predict future actions. Afterwards, we show our experimental results on several datasets in section IV and conclude the paper in Section V.

## II. RELATED WORK

**Action recognition.** Action recognition is an important branch of video related research areas and has been extensively studied in the past decades. The existing methods are mainly developed for extracting discriminative action features from temporally complete action videos. These methods can be roughly categorized into hand-crafted feature based approaches and deep learning based approaches. Early methods such as Improved Dense Trajectory (IDT) mainly adopt hand-crafted features, such as HOF [32], HOG [32] and MBH

[55]. Recent studies demonstrate that action features can be learned by deep learning methods such as convolutional neural networks (CNN) and recurrent neural networks (RNN). Two-stream network [46], [57] learns appearance and motion features based on RGB frame and optical flow field separately. RNNs, such as long short-term memory (LSTM) [15] and gated recurrent unit (GRU) [6], have been used to model long-term temporal correlations and motion information in videos, and generate video representation for action classification. A CNN+LSTM model, which uses a CNN to extract frame features and a LSTM to integrate features over time, is also used to recognize activities in videos [9]. C3D [10] network simultaneously captures appearance and motion features using a series of 3D convolutional layers. Recently, I3D [4] networks use two stream CNNs with inflated 3D convolutions on both dense RGB and optical flow sequences to achieve state of the art performance on Kinetics dataset [24].

**Action anticipation.** Many works have been proposed to exploit the partially observed videos for early action prediction or future action anticipation. Recently, Hoai *et al.* [17] propose a max-margin framework with structured SVMs to solve this problem. Ryoo *et al.* [42] develop an early action prediction system by observing some evidences from the temporal accumulated features. Lan *et al.* [31] design a coarse-to-ne hierarchical representation to capture the discriminative human movement at different levels, and use a max-margin framework for final prediction. Cao *et al.* [64] formulate the action prediction problem into a probabilistic framework, which aims to maximize the posterior of activity given observed frames. In their work, the likelihood is computed by feature reconstruction error using sparse coding. However, it suffers from high computational complexity as the inference is performed on the entire training data. Carl *et al.* [54] present a framework that uses large-scale unlabeled data to predict a rich visual representation in the future, and apply it towards anticipating both actions and objects. Kong *et al.* [29] propose a combined CNN and LSTM along with a memory module in order to record “hard-to-predict” samples, they benchmark their results on UCF101 [49] and Sports-1M [23] datasets. Gao *et al.* [13] propose a Reinforced Encoder-Decoder (RED) network for action anticipation, which uses reinforcement learning to encourage the model to make the correct anticipations as early as possible. Ke *et al.* [27] propose an attended temporal feature, which uses multi-scale temporal convolutions to process the time-conditioned observation. In this work, we focus only on recent results on anticipation of action labels, more details can be found in [13] and [62].

**Online action detection.** Online action detection is usually solved as an online per-frame labelling task on streaming videos, which requires correctly classifying every frame without accessing future frames. De Geest *et al.* [7] first introduce the problem by introducing a realistic dataset, *i.e.* TVSeries, and benchmarked the existing models. They have shown that a simple LSTM approach is not sufficient for online action detection, and even worse than the traditional pipeline of improved trajectories, Fisher vectors and SVM. Their later work [8] introduces a two-stream feedback network, where one stream processes the input and the other one models

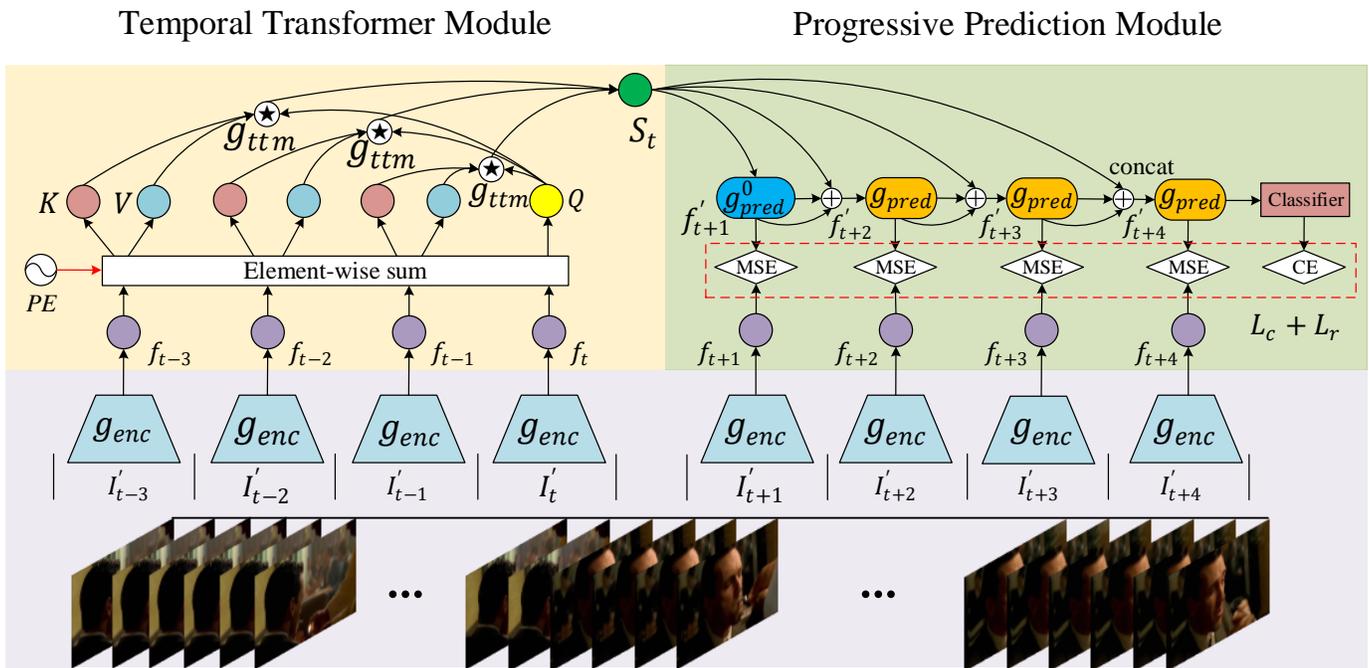


Fig. 2: The flowchart of our TTPP method for action anticipation. Given a continuous video sequence, an encoder network is first used to map each video clip to clip features, and then a Temporal Transformer Module is proposed to aggregate observed clip features, and finally a Progressive Prediction Module is designed for future action anticipation. Note that the future information in the red dashed box is only used in training stage and the classifier is performed on each prediction.

the temporal relations. Gao *et al.* [13] propose a Reinforced Encoder-Decoder network for action anticipation and treat online action detection as a special case of their framework. Xu *et al.* [62] propose the Temporal Recurrent Network (TRN) to model the temporal context by simultaneously performing online action detection and anticipation. Besides, Shou *et al.* [45] address the online detection of action start (ODAS) by encouraging a classification network to learn the representation of action start windows.

**Attention for video understanding.** The attention mechanism which directly models long-term interactions with self-attention has led to state-of-the-art models for action understanding tasks, such as video-based and skeleton-based action recognition [14], [34], [35], [44], [48]. Our work is related to the recent Video Action Transformer Network [14], which uses the Transformer architecture as the “head” of a detection framework. Specifically, it uses the ROI-pooled I3D feature of a target region as query and aggregates contextual information from the spatial-temporal feature maps of an input video clip. Our work differs from it in the following aspects: (1) The problem is different from spatial-temporal action detection. To the best of our knowledge, we are the first to use Transformer architecture for action anticipation. (2) We have task-specific considerations. For instance, our Transformer unit takes the current frame feature as query and the historical frame features as memory. (3) We elaborately design a light-weight progressive prediction module for efficient action anticipation.

### III. OUR APPROACH

In this section, we present our temporal Transformer with progressive prediction for the action anticipation task. We propose two module, temporal Transformer module (TTM) to aggregate observed information and progressive prediction module (PPM) to anticipate future actions.

#### A. Problem Formulation

The action anticipation task aims to predict the action class  $y$  for each frame in the future from an observed action video  $V$ . More formally, let  $V_1^L = [I_1, I_2, \dots, I_L]$  be a video with  $L$  frames. Given the first  $t$  frames  $V_1^t = [I_1, I_2, \dots, I_t]$ , the task is to predict the actions happening from frames  $t+1$  to  $L$ . That is, we aim to assign action labels  $y_{t+1}^L = [y_{t+1}, y_{t+2}, \dots, y_L]$  to each of the unobserved frames.

#### B. Overall Framework

Two crucial issues of action anticipation are i) how to aggregate observed information and ii) how to predict future actions. We address these two issues with a simple yet efficient framework, termed as *Temporal Transformer with Progressive Prediction* (TTPP). As illustrated in Figure 2, a long video is first segmented into multiple non-overlapped chunks  $[I'_1, I'_2, \dots, I'_t]$  with each clip containing an equal number of consecutive frames. Then, a network, *i.e.*  $g_{enc}$ , maps each video chunk into a representation  $f_t = g_{enc}(I'_t)$ . More details on video pre-processing and feature extraction are presented in Section IV-C. Subsequently, a Temporal Transformer Module (TTM), *i.e.*  $g_{ttm}$ , temporally aggregates  $t$  consecutive

chunk representations into a historical representation  $S_t = g_{ttm}(f_1, f_2, \dots, f_t)$ . Finally, a Progressive Prediction Module (PPM) progressively predicts future features and actions. The PPM is comprised of an *initial prediction block*, i.e.  $g_{pred}^0$ , and a *shared-parameter progressive prediction block*, i.e.  $g_{pred}$ .  $g_{pred}^0$  takes  $S_t$  as input and predicts the immediately following clip feature and action probability.  $g_{pred}$  accumulates the former predictions and  $S_t$ , and further predicts a few subsequent future features and actions.

### C. Temporal Transformer Module (TTM)

**Transformer revisit.** Transformer was originally proposed to replace traditional recurrent models for machine translation [52]. The core idea of Transformer is to model correlation between contextual signals by an attention mechanism. Specifically, it aims to encode the input sequence to a higher-level representation by modeling the relationship between queries ( $Q$ ) and memory (keys ( $K$ ) and values ( $V$ )) with,

$$Attention(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_m}}\right)V, \quad (1)$$

where  $Q \in R^{L_q \times d_m}$ ,  $K \in R^{L_k \times d_m}$  and  $V \in R^{L_v \times d_v}$ . This architecture becomes ‘‘self-attention’’ with  $Q = K = V = \{f_1, f_2, \dots, f_T\}$  which is also known as the non-local networks [58]. A self-attention module maps the sequence to a higher-level representation like RNNs but without recurrence.

**Temporal Transformer.** To efficiently aggregate observed information, our TTPP framework resorts to a Transformer-style architecture, termed as Temporal Transformer Module (TTM). The TTM takes as input the video chunk features and maps them into a query feature and memory features. For online action anticipation, considering that the last observed feature  $f_t$  would be the most relevant one to the future actions, we use  $f_t$  as the query of TTM. The memory of TTM is intuitively set as the historical features  $[f_1, f_2, \dots, f_{t-1}]$ . Formally, the query and memory are as follows,

$$Q = f_t, K = V = [f_1, f_2, \dots, f_{t-1}]. \quad (2)$$

Since temporal information is lost in the attention operation, we add the positional encoding [52] into the input representations. Given sequence feature  $f_{in} = [f_1, f_2, \dots, f_T] \in R^{T \times d_m}$ , the  $i$ -th value of the positional vector in temporal position  $pos$  is defined as,

$$PE_{(pos,i)} = \begin{cases} \sin(pos/10000^{i/d_m}) & \text{if } i \text{ is even} \\ \cos(pos/10000^{i/d_m}) & \text{otherwise.} \end{cases} \quad (3)$$

The original feature vector  $f_{pos}$  is then updated by  $f_{pos} = f_{pos} + PE_{(pos,:)}$  which provides information about temporal position of each clip feature.

To model complicated action videos, our TTM further leverages the multi-head attention mechanism as follows,

$$A_t = \text{MultiHead}(Q, K, V) = \text{Concat}(h_1, \dots, h_n)W^o, \quad (4)$$

where  $h_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$ ,

where  $n$  is the number of attention heads, and  $W_i^Q \in R^{d_m \times d_q}$ ,  $W_i^K \in R^{d_m \times d_k}$ ,  $W_i^V \in R^{d_m \times d_v}$  are parameters for the  $i$ -th attention head which are used for linear projection, and  $W^o \in$

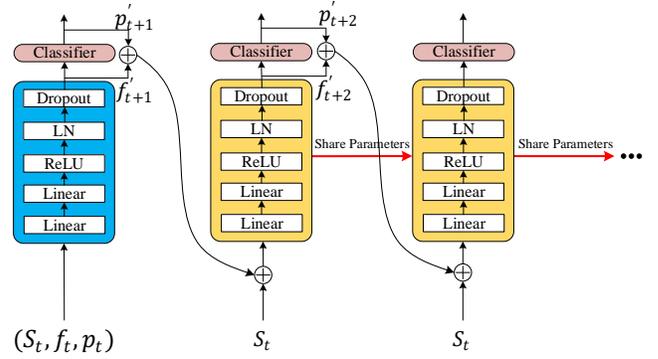


Fig. 3: The illustration of our PPM. It consists of an *initial prediction block* highlighted in blue and a *shared-parameter progressive prediction block* highlighted in yellow, where each block is built with two fully-connected layers, followed by a ReLU activation. In addition, we use layer normalization and dropout to improve regularization.

$R^{n_{d_v} \times d_m}$  is the projection matrix to reduce the dimension of the concatenated attention vector. For each head, we use  $d_k = d_q = d_v = \frac{d_m}{n}$ . Considering the importance of  $f_t$  for anticipation, we view  $A_t$  as an extra information and add it to the original query feature via a shortcut connection. The final output feature of TTM is  $S_t = A_t + f_t$  with dimension  $d_m$ .

### D. Progressive Prediction Module (PPM)

Partially inspired by WaveNet [51], we design a Progressive Prediction Module (PPM) to better exploit the aggregated historical knowledge for future prediction. As illustrated in Figure 3, the PPM is comprised of an *initial prediction block* and a *shared-parameter progressive prediction block*, where each block is built with two fully-connected (FC) layers, a ReLU activation [38] and a layer normalization (LN) [2].

Assume we predict  $l$  steps in the future from time  $t + 1$  to  $t + l$ . At the first time step  $t + 1$ , the *initial prediction block* takes as input the aggregated historical representation  $S_t \in R^{d_m}$  and predicts the feature  $f'_{t+1} \in R^{d_m}$  and action probability  $p'_{t+1} \in R^C$ . Formally, this block is as follows,

$$p_t = \text{Softmax}(W_c f_t) \quad (5)$$

$$f'_{t+1} = g_{pred}^0(S_t \oplus f_t \oplus p_t), \quad (6)$$

$$p'_{t+1} = \text{Softmax}(W_c f'_{t+1}), \quad (7)$$

where  $W_c$  is the multi-class ( $C$  action classes) action classifier. At other time step  $t + i$  ( $i > 1$ ), the previously predicted embedding  $f'_{t+i-1}$  and action probability  $p'_{t+i-1}$  are first concatenated with  $S_t$  in channel-wise, and then fed into the *progressive prediction block*. Formally, this block is defined as follows,

$$f'_{t+i} = g_{pred}(S_t \oplus f'_{t+i-1} \oplus p'_{t+i-1}), \quad (8)$$

$$p'_{t+i} = \text{Softmax}(W_c f'_{t+i}), \quad (9)$$

where ‘ $\oplus$ ’ denotes concatenate operation. Due to the concatenation, the input dimension of the *progressive prediction block* is  $2d_m + C$ . For both blocks, we use two fully-connected (FC) layers with the first FC reducing the input dimension to  $\frac{d_m}{2}$  and the second FC generating output vector of dimension  $d_m$ . It is worth noting that different steps in the *progressive prediction block* share parameters. Thus, the whole PPM is a light-weight network.

**Training.** Our TTPP framework is trained in an end-to-end manner with supervision on the PPM module. Specifically, we use two types of loss functions, namely a feature reconstruction loss  $L_r$  and a classification loss  $L_c$ .  $L_r$  is the mean squared error loss (MSE) between predicted features and ground-truth features, which is defined as,

$$L_r = \sum_{i=1}^l \|f'_{t+i} - f_{t+i}\|^2. \quad (10)$$

$L_c$  is the sum of cross-entropy loss (CE) on all the prediction steps, which is defined as,

$$L_c = - \sum_{i=1}^l \sum_{j=1}^C y_{(t+i,j)} \log p'_{(t+i,j)}, \quad (11)$$

where  $y_{(t+i,:)}$  is the one-hot ground-truth vector at time  $t+i$ . The total loss is formulated as,

$$L = L_c + \lambda L_r, \quad (12)$$

where  $\lambda$  is a trade-off weight for feature reconstruction loss. We experimentally find the final performance is not sensitive to the value of weight, we set  $\lambda = 1$  for simplicity in our experiments.

#### IV. EXPERIMENTS

The proposed method was evaluated on three datasets, *i.e.* TVSeries [62], THUMOS-14 [20] and TV-Human-Interaction [39]. We choose these datasets because they include videos from diverse perspectives and applications: TVSeries was recorded from television and contains a variety of everyday activities, THUMOS-14 is a popular dataset of sports-related actions, and TV-Huamn-Interaction contains human interaction actions collected from tv shows. In this section, we report experimental results and detailed analysis.

##### A. Datasets

**TVSeries** [7] is originally proposed for online action detection, which consists of 27 episodes of 6 popular TV series, namely *Breaking Bad* (3 episodes), *How I Met Your Mother* (8), *Mad Men* (3), *Modern Family* (6), *Sons of Anarchy* (3), and *Twenty-four* (4). It contains totally 16 hours of video. The dataset is temporally annotated at the frame level with 30 realistic, everyday actions (*e.g.*, *pick up*, *open door*, *drink*, etc.). It is challenging with diverse actions, multiple actors, unconstrained viewpoints, heavy occlusions, and a large proportion of non-action frames.

**THUMOS-14** [20] is a popular benchmark for temporal action detection. It contains over 20 hours of sport videos annotated with 20 actions. The training set (*i.e.* UCF101 [49])

contains only trimmed videos that cannot be used to train temporal action detection models. Following prior works [13], [62], we train our model on the validation set (including  $3K$  action instances in 200 untrimmed videos) and evaluate on the test set (including  $3.3K$  action instances in 213 untrimmed videos).

**TV-Human-Interaction** (TV-HI) [39]. We also evaluate our method on TV-Human-Interaction which is also used in [13]. The dataset contains 300 trimmed video clips extracted from 23 different TV shows. It is annotated with four interaction classes, namely *hand shake*, *high five*, *hug*, and *kissing*. It also contains a *negative* class with 100 videos, that have none of the listed interactions. We use the suggested experimental setup of two train/test splits.

##### B. Evaluation Protocols

For each class on TVSeries, we use the per-frame calibrated average precision (cAP) which is proposed in [7],

$$cAP = \frac{\sum_k cPrec(k) * I(k)}{P}, \quad (13)$$

where calibrated precision  $cPrec = \frac{TP}{TP+FP/w}$ ,  $I(k)$  is an indicator function that is equal to 1 if the cut-off frame  $k$  is a true positive,  $P$  denotes the total number of true positives, and  $w$  is the ratio between negative and positive frames. The mean cAP over all classes is reported for final performance. The advantage of cAP is that it is fair for class imbalance condition. For THUMOS-14, we report per-frame mean Average Precision (mAP) performance. For TV-Human-Interaction, we report classification accuracy (ACC).

##### C. Implementation Details

To make fair comparisons with state-of-the-art methods [7], [13], [62], we follow their experimental settings on each dataset.

**Chunk-level feature extraction.** We extract frames from all videos at 24 Frames Per Second (*FPS*). The video chunk size is set to 6, *i.e.* 0.25 second. We use three different feature extractors as the visual encoder  $g_{enc}$ , VGG-16 [47] network pre-trained on UCF101 [49], two-stream (TS) [61] network<sup>1</sup> pre-trained on ActivityNet-1.3 [3], and inflated 3D ConvNet<sup>2</sup> (I3D) [5] pre-trained on Kinetics [25]. VGG-16 features (4096-D) are extracted at the *fc6* layer for the central frame of each chunk. For the two-stream features in each chunk, the appearance CNN feature is extracted on the central frame which is the output of *Flatten 673 layer* in ResNet-200 [16], and the motion feature is extracted on the 6 optical flow frames of each chunk which is output of *global pool layer* in a pre-trained BN-Inception model [19]. The motion feature and appearance feature are then concatenated into a TS feature (4096-D) for each chunk. Different from prior works [13], [62], we also use recent I3D features. The I3D model is originally trained with 64-frame video snippets, thus may not be a good idea for per-frame action anticipation. Nevertheless,

<sup>1</sup><https://github.com/yjxiong/anet2016-cuhk>

<sup>2</sup><https://github.com/piergiaj/pytorch-i3d>

Method	Inputs	Time predicted into the future (seconds)								Avg
		0.25s	0.5s	0.75s	1.0s	1.25s	1.5s	1.75s	2.0s	
ED [13]	VGG	71.0	70.6	69.9	68.8	68.0	67.4	67.0	66.7	68.7
RED [13]	VGG	71.2	71.0	70.6	70.2	69.2	68.5	67.5	66.8	69.4
Ours	VGG	<b>72.7</b>	<b>72.3</b>	<b>71.9</b>	<b>71.6</b>	<b>71.3</b>	<b>70.9</b>	<b>69.9</b>	<b>69.3</b>	<b>71.3</b>
ED [13]	TS	78.5	78.0	76.3	74.6	73.7	72.7	71.7	71.0	74.5
RED [13]	TS	79.2	78.7	77.1	75.5	74.2	73.0	72.0	71.2	75.1
TRN [62]	TS	79.9	78.4	77.1	75.9	74.9	73.9	73.0	72.3	75.7
Ours	TS	<b>81.2</b>	<b>80.3</b>	<b>79.3</b>	<b>77.6</b>	<b>76.9</b>	<b>76.7</b>	<b>76.0</b>	<b>74.9</b>	<b>77.9</b>

TABLE I: Comparison with the state-of-the-art methods on TVSeries in terms of mean cAP (%).

Method	Inputs	Time predicted into the future (seconds)								Avg
		0.25s	0.5s	0.75s	1.0s	1.25s	1.5s	1.75s	2.0s	
ED [13]	TS	43.8	40.9	38.7	36.8	34.6	33.9	32.5	31.6	36.6
RED [13]	TS	45.3	42.1	39.6	37.5	35.8	34.4	33.2	32.1	37.5
TRN [62]	TS	45.1	42.4	40.7	39.1	37.7	36.4	35.3	34.3	38.9
Ours	TS	<b>45.9</b>	<b>43.7</b>	<b>42.4</b>	<b>41.0</b>	<b>39.9</b>	<b>39.4</b>	<b>37.9</b>	<b>37.3</b>	<b>40.9</b>
Ours	I3D	46.8	45.5	44.6	43.6	41.9	41.1	40.4	38.7	<b>42.8</b>

TABLE II: Comparison with the state-of-the-art methods on THUMOS-14 in terms of mAP (%).

we input the 6 frames of each chunk to I3D and extract the output (1024-D) of the last *global average pooling layer* as I3D-based feature.

**Hyperparameter setting.** We implement our proposed method in PyTorch and perform all experiments on a system with 8 Nvidia TITAN X graphic cards. We use the SGD optimizer with a learning rate 0.001, a momentum of 0.9, and batch size 32. The input sequence length is set to 8 by default, corresponding to 2 seconds. We use single-layer multi-head setting for our TTM, and the number of heads is set to 4 by default.

#### D. Popular Baselines

Here we present several advanced baselines for temporal information aggregation and future prediction.

**Temporal convolution** (*i.e.* Conv1D) aggregates temporal features with 1-D convolution operations in temporal axis. We apply 3 Conv1D layers with kernel size 3 and stride 2 on two-stream features for this baseline.

**LSTM** takes sequence features as input and recurrently updates its hidden states over time. The **Encoder-LSTM** summarizes historical information into the final hidden state for information aggregation. The **Decoder-LSTM** recurrently decodes information into hidden states as predicted features. We use a single-layer LSTM architecture with 4096 hidden units for this baseline.

**Single-shot prediction** (SSP). We implement a single-shot prediction method similar to [13], [54]. With the aggregated historical feature, this method uses two FC layers to anticipate the single future feature at  $T_a$ , where  $T_a \in \{t + 1, t + 2, \dots, t + l\}$ . This prediction method is equal to our PPM without the progressive process.

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

We compare our proposed TTPP method to several state-of-the-art methods on TVSeries, THUMOS-14, and TV-HI. The

Method	Vondrick <i>et al.</i> [54]	RED-VGG [13]	RED-TS [13]	Ours-TS
ACC (%)	43.6	47.5	50.2	<b>53.5</b>

TABLE III: Anticipation results on TV-Human-Interaction at  $T_a = 1s$  in terms of ACC (%).

Method	Time predicted into the future (seconds)								Avg
	0.25s	0.5s	0.75s	1.0s	1.25s	1.5s	1.75s	2.0s	
Conv1D <sub>-</sub>	77.9	77.1	75.9	74.3	73.6	72.7	71.9	71.0	74.3
Conv1D	79.4	77.9	76.6	75.4	73.9	73.6	72.8	72.4	75.3
LSTM <sub>-</sub>	79.3	78.2	76.9	74.6	73.0	72.3	71.6	69.8	74.5
LSTM	78.9	77.9	76.3	75.5	74.3	73.8	72.8	72.0	75.2
TTM <sub>-</sub>	79.1	78.6	77.9	77.0	76.4	75.6	75.1	74.1	76.7
TTM (Ours)	81.2	80.3	79.3	77.6	76.9	76.7	76.0	74.9	<b>77.9</b>

(a) TVSeries

Method	Time predicted into the future (seconds)								Avg
	0.25s	0.5s	0.75s	1.0s	1.25s	1.5s	1.75s	2.0s	
Conv1D	41.8	40.9	39.6	38.1	37.2	36.7	35.9	35.5	38.2
LSTM	41.6	40.5	39.1	37.9	35.6	34.9	34.3	33.3	37.0
TTM (Ours)	45.9	43.7	42.4	41.0	39.9	39.4	37.9	37.3	<b>40.9</b>

(b) THUMOS-14

TABLE IV: Evaluation of temporal *aggregation* methods with two-stream features on TVSeries and THUMOS-14. For fair comparison, PPM is used for future prediction. “<sub>-</sub>” denotes that the aggregated feature is directly used for prediction without the shortcut connection to current feature.

results are presented in Table I, Table II, and Table III, respectively. Our method consistently outperforms all these methods in all the predicted steps. With two-stream features, our TTPP achieves 77.9% (mean cAP), 40.9% (mAP), and 53.5% (ACC) on these datasets, which outperforms these recent advanced methods by 2.2%, 2.0%, and 3.3%, respectively.

On both TVSeries and THUMOS-14, the improvements over other methods are more evident on long-term predictions. For instance, with two-stream features, our TTPP outperforms ED (Encoder-Decoder LSTM) [13] by 2.1% at  $T_a = 0.25s$  and 5.7% at  $T_a = 2.0s$  on THUMOS-14, and these numbers are 2.7% and 3.9% on TVSeries. With VGG features, our method improves the Reinforcement ED by 2.6% in average cAP on TVSeries. Since the VGG and TS features are relatively old, we also test the I3D features, which updates a new state-of-the-art on THUMOS-14 with 42.8% in average mAP over time.

#### F. Ablation Study of TTM and PPM

To further investigate the effectiveness of our proposed TTPP, we conduct extensive evaluations for both TTM and PPM by comparing them to recent temporal aggregation and prediction methods, respectively.

For temporal aggregation, we compare our TTM to Conv1D and Encoder-LSTM on both THUMOS-14 and TVSeries with the PPM as prediction phase. Since we use a shortcut connection in TTM to highlight the current frame information, we also evaluate the benefits of this design for all the aggregation methods. The results are shown in Table IV. Several observations can be concluded as follows. First, the proposed TTM is superior to both Conv1D and LSTM regardless

Method (A-P)*	Time predicted into the future (seconds)								Avg
	0.25s	0.5s	0.75s	1.0s	1.25s	1.5s	1.75s	2.0s	
LSTM-LSTM (ED [13])	78.5	78.0	76.3	74.6	73.7	72.7	71.7	71.0	74.5
LSTM-SSP (EFC [13])	78.4	76.9	74.2	72.7	70.7	70.2	69.0	67.9	72.5
LSTM-PPM	78.6	77.5	76.0	75.0	73.5	73.8	72.8	71.8	74.9
TTM-LSTM	79.3	77.5	76.2	75.1	73.3	72.8	71.6	69.9	74.5
TTM-SSP	80.1	77.1	75.3	73.6	72.3	71.7	70.0	68.9	73.6
TTM-PPM (Ours)	81.2	80.3	79.3	77.6	76.9	76.7	76.0	74.9	<b>77.9</b>

(a) TVSeries

Method (A-P)*	Time predicted into the future (seconds)								Avg
	0.25s	0.5s	0.75s	1.0s	1.25s	1.5s	1.75s	2.0s	
LSTM-LSTM (ED [13])	43.8	40.9	38.7	36.8	34.6	33.9	32.5	31.6	36.6
LSTM-SSP (EFC [13])	40.6	39.3	37.2	34.9	33.2	31.5	30.4	28.5	34.4
LSTM-PPM	41.3	40.3	38.9	37.6	35.4	34.6	33.9	33.0	36.8
TTM-LSTM	44.4	43.1	41.3	40.2	38.7	37.9	37.2	36.6	39.9
TTM-SSP	43.9	41.0	38.5	37.1	35.5	32.7	32.2	30.5	36.4
TTM-PPM (Ours)	45.9	43.7	42.4	41.0	39.9	39.4	37.9	37.3	<b>40.9</b>

(b) THUMOS-14

TABLE V: Evaluation of future prediction methods with two-stream features on TVSeries and THUMOS-14. Both TTM and LSTM are evaluated for temporal aggregation. (A-P)\* denotes temporal aggregation and future prediction, respectively.

of the shortcut connection. Specifically, TTM outperforms Conv1D and LSTM by 2.6% (2.7%) and 2.7% (3.9%) with shortcut connection on TVSeries (THUMOS-14), respectively. Second, the shortcut connection to current feature significantly improves all methods on TVSeries. For instance, our TTM degrades from 77.9% to 76.7% after removing the shortcut connection which demonstrates the importance of current feature and the superiority of our design. Last but not the least, the improvements of our TTM over other methods are similar for different time steps, which suggests that TTM provides better aggregated features via attention than the others.

For future prediction, we compare our PPM to Decoder-LSTM and SSP with either Encoder-LSTM or TTM as aggregation method. The results are presented in Table V. Several finds are concluded as follows. First, with both aggregation methods, our PPM consistently outperforms Decoder-LSTM and SSP on both datasets which shows the effectiveness of PPM. Second, our PPM obtains more improvements with our aggregation method TTM than Encoder-LSTM. For instance, TTM-PPM outperforms TTM-LSTM by 3.4% while LSTM-PPM only outperforms LSTM-LSTM by 0.4%. Third, with both aggregation methods, our PPM is significantly superior to SSP on both datasets especially at long-term prediction time steps, which demonstrates the progressive design of our PPM is important.

### G. Importance of Feature Prediction

In order to evaluate the influence of feature prediction for the final action anticipation, we remove the predicted features (w/o FP) by only use the concatenation of action probability and the aggregated historical representation in the PPM. The results are shown in Figure 4. Without considering the predicted features, the performance of the model w/o FP degenerates dramatically. It indicates that only relying on the action probability to predict future actions is not enough and the predicted feature representations are always related to the

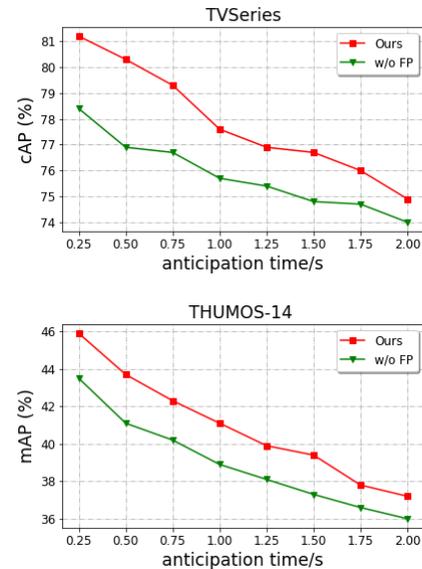


Fig. 4: Evaluation of feature prediction for action anticipation on TVSeries (cAP %) and THUMOS-14 dataset (mAP %) with two-stream features.

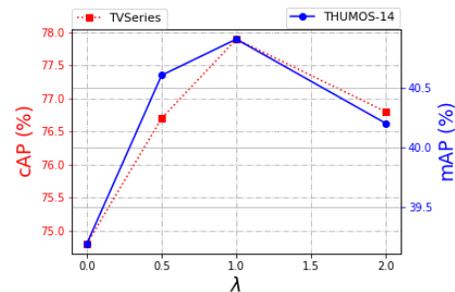


Fig. 5: Evaluation of trade-off weight  $\lambda$  on TVSeries (cAP %) and THUMOS-14 dataset (mAP %) with two-stream features.

action itself and thus could possibly provide some useful information.

### H. Evaluation of Sequence Length and Parameters

In the above experiments, we use a fixed historical length 8 for aggregation, 4 parallel heads and trade-off loss weight  $\lambda = 1.0$  for training by default. To investigate their impacts to the proposed TTPP framework, we evaluate them on both THUMOS-14 and TVSeries.

**The impact of  $\lambda$ .**  $\lambda$  is the weight of feature reconstruction loss in training. Figure 5 shows the results of varied  $\lambda$  on THUMOS-14 and TVSeries. Removing the feature reconstruction loss, *i.e.*  $\lambda = 0$ , degrades performance dramatically on both datasets which suggests the necessary of feature prediction. Increasing the weight from 0 to 1 improves performance, and it gets saturation or slightly hurts performance after 1. This may be explained by that overemphasizing feature reconstruction can hurt the discrimination of predicted features.

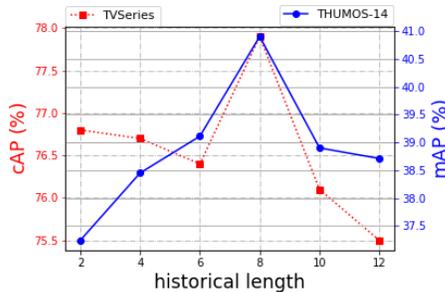


Fig. 6: Evaluation of input sequence length for temporal aggregation on TVSeries (cAP %) and THUMOS-14 dataset (mAP %) with two-stream features.

**Number of heads  $n$ .** We also study performance variations given various number of heads for temporal Transformer. Average prediction performance of our TTPP network with  $n \in \{1, 2, 4, 8, 16\}$  are shown in Table VI. Results in Table VI indicate that our method is not sensitive to parameter  $n$ . The largest performance variation is only within 0.8% on TVSeries and 1.1% on THUMOS-14. On both datasets, we achieve best performance with head number  $n = 4$ .

**Input sequence length.** The length of observed sequence determines how much historical information can be used. Figure 6 illustrates the evaluation results on THUMOS-14 and TVSeries. On both datasets, we achieve the best performance with length 8. Decreasing sequence length leads to insufficient context information and increasing sequence length results to massive background information which are both inferior to the default length.

### I. Efficiency and Visualization

Table VII reports a comparison of parameters, memory footprint, inference time and performance of different models on TVSeries dataset. Compared to the popular Encoder-Decoder LSTM model, our TTPP has 64% fewer parameters, 44% fewer memory footprint and less inference time, while achieves 4.6% higher performance. The efficiency of the proposed TTPP owes to both the Transformer architecture for sequence modeling and the efficient progressive prediction module.

Figure 7 shows some examples of attention weights and action anticipation on TVSeries, THUMOS-14 and TV-Human-Interaction. We find that frames near the current frame usually get higher weights compared to these distant frames since the current frame feature is used as the query. On TVSeries and THUMOS-14, multiple action instances and confusing background frames exist in the videos which lead to incorrect anticipation inevitably.

## V. CONCLUSION

In this paper, we propose a novel deep framework to boost action anticipation by adopting Temporal Transformer with Progressive Prediction, where a TTM is used to aggregate observed information and a PPM to progressively predict

Number of Heads	n=1	n=2	n=4	n=8	n=16
TVSeries	77.1	77.7	<b>77.9</b>	77.5	77.2
THUMOS-14	40.1	40.4	<b>40.9</b>	40.2	39.8

TABLE VI: Comparison between different number of heads on TVSeries and THUMOS-14 with two-stream features.

Model	Parameter (M)	Memory (M)	Inf (s)	cAP (%)
ED [13]	277	6560	212	74.5
TTPP	101	3675	145	77.9

TABLE VII: Comparison of parameter, memory footprint and inference time on TVSeries dataset with two-stream features.

future features and actions. Experimental results on TVSeries, THUMOS-14, and TV-Human-Interaction demonstrate that our framework significantly outperforms the state-of-the-art methods. Extensive ablation studies are conducted to show the effectiveness of each module of our method.

## VI. ACKNOWLEDGEMENTS

This work is partially supported by the National Key Research and Development Program of China (No.2016YFC1400704), and National Natural Science Foundation of China (U1813218, U1713208, 61671125), Shenzhen Basic Research Program (JCYJ20170818164704758, CXB201104220032A), the Joint Lab of CAS-HK, Shenzhen Institute of Artificial Intelligence and Robotics for Society, and the Sichuan Province Key Research and Development Plan (2019YFS0427).

## REFERENCES

- [1] Mohammad Sadegh Aliakbarian, Fatemehsadat Saleh, Mathieu Salzmann, Basura Fernando, and Lars Andersson. Encouraging lstms to anticipate actions very early. In *ICCV*, 2017.
- [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. 2016.
- [3] Fabian Caba, Victor Escorcía, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *CVPR*, 2015.
- [4] João Carreira and Andrew Zisserman. Quo vadis, action recognition? A new model and the kinetics dataset. *CoRR*, abs/1705.07750, 2017.
- [5] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*, 2017.
- [6] Junyoung Chung, Çağlar Gülçehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *CoRR*, abs/1412.3555, 2014.
- [7] Roeland De Geest, Efstratios Gavves, Amir Ghodrati, Zhenyang Li, Cees Snoek, and Tinne Tuytelaars. Online action detection. In *CVPR*, 2016.
- [8] Roeland De Geest and Tinne Tuytelaars. Modeling temporal structure with lstm for online action detection. In *WACV*, pages 1549–1557, 2018.
- [9] Jeff Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. *CoRR*, abs/1411.4389, 2014.
- [10] Tran Du, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *ICCV*, 2015.
- [11] Antonino Furnari and Giovanni Maria Farinella. What would you expect? anticipating egocentric actions with rolling-unrolling lstms and modality attention. In *ICCV*, 2019.
- [12] Jiyang Gao and Ram Nevatia. Revisiting temporal modeling for video-based person reid. In *BMVC*, 2018.

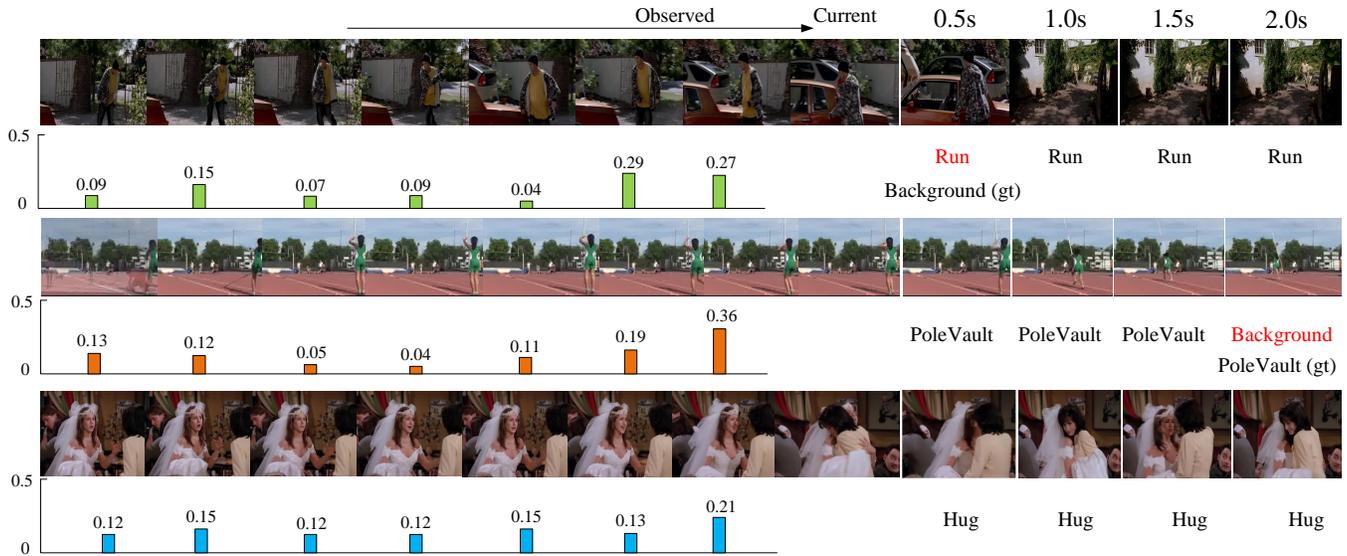


Fig. 7: Visualization of attention weights and action anticipation on TVSeries (1st row), THUMOS-14 (2nd row), and TV-Human-Interaction (3rd row). Incorrect anticipation results are shown in red.

[13] Jiyang Gao, Zhenheng Yang, and Ram Nevatia. RED: reinforced encoder-decoder networks for action anticipation. In *BMVC*, 2017.

[14] Rohit Girdhar, Joao Carreira, Carl Doersch, and Andrew Zisserman. Video action transformer network. In *CVPR*, 2019.

[15] Alex Graves. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.

[16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015.

[17] Minh Hoai and Fernando De La Torre. Max-margin early event detectors. In *CVPR*, 2012.

[18] Sepp Hochreiter and Jrgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.

[19] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. 2015.

[20] Y-G. Jiang, J. Liu, A. Roshan Zamir, G. Toderici, I. Laptev, M. Shah, and R. Sukthankar. THUMOS challenge: Action recognition with a large number of classes. 2014.

[21] Xu Jing, Zhao Rui, Zhu Feng, Huaming Wang, and Wanli Ouyang. Attention-aware compositional network for person re-identification. In *CVPR*, 2018.

[22] Rafal Józefowicz, Oriol Vinyals, Mike Schuster, Noam Shazeer, and Yonghui Wu. Exploring the limits of language modeling. *CoRR*, abs/1602.02410, 2016.

[23] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, and Fei Fei Li. Large-scale video classification with convolutional neural networks. In *2014 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.

[24] Will Kay, João Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The kinetics human action video dataset. *CoRR*, abs/1705.06950, 2017.

[25] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, and Andrew Zisserman. The kinetics human action video dataset. 2017.

[26] Qihong Ke, Mohammed Bannamoun, Senjian An, Farid Boussaid, and Ferdous Sohel. Human interaction prediction using deep temporal features. In *ECCV*, 2016.

[27] Qihong Ke, Mario Fritz, and Bernt Schiele. Time-conditioned action anticipation in one shot. In *CVPR*, 2019.

[28] Y. Kong, Z. Tao, and Y. Fu. Deep sequential context networks for action prediction. In *CVPR*, 2017.

[29] Yu Kong, Shangqian Gao, Bin Sun, and Yun Fu. Action prediction from videos via memorizing hard-to-predict samples. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

[30] Yu Kong, Dmitry Kit, and Yun Fu. A discriminative model with multiple temporal scales for action prediction. In *ECCV*, 2014.

[31] Tian Lan, Tsung-Chuan Chen, and Silvio Savarese. A hierarchical representation for future action prediction. In David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 689–704, Cham, 2014. Springer International Publishing.

[32] Ivan Laptev, Marcin Marszalek, Cordelia Schmid, and Benjamin Rozenfeld. Learning realistic human actions from movies. In *CVPR*, 2008.

[33] Kang Li and Yun Fu. Prediction of human activity by discovering temporal sequence patterns. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(8):1644–1657, 2014.

[34] Jun Liu, Gang Wang, Ling Yu Duan, Kamila Abdiyeva, and Alex C. Kot. Skeleton based human action recognition with global context-aware attention lstm networks. *IEEE Transactions on Image Processing*, PP(99):1–1, 2018.

[35] Jun Liu, Gang Wang, Ping Hu, Ling Yu Duan, and Alex C. Kot. Global context-aware attention lstm networks for 3d action recognition. In *CVPR*, 2017.

[36] Minh Thang Luong, Hieu Pham, and Christopher D. Manning. Effective approaches to attention-based neural machine translation. *Computer Science*, 2015.

[37] Tahmida Mahmud, Mahmudul Hasan, and Amit K. Roy-Chowdhury. Joint prediction of activity labels and starting times in untrimmed videos. In *CVPR*, 2017.

[38] Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010.

[39] Alonso Patron, Marcin Marszalek, Andrew Zisserman, and Ian Reid. High five: Recognising human interactions in tv shows. In *BMVC*, 2010.

[40] S. Ren, K. He, R Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(6):1137–1149, 2017.

[41] Cristian Rodriguez, Basura Fernando, and Hongdong Li. Action anticipation by predicting future dynamic images.

[42] M. S. Ryoo. Human activity prediction: Early recognition of ongoing activities from streaming videos. In *CVPR*, 2012.

[43] Shikhar Sharma, Ryan Kiros, and Ruslan Salakhutdinov. Action recognition using visual attention. In *ICLR*, 2015.

[44] Shikhar Sharma, Ryan Kiros, and Ruslan Salakhutdinov. Action recognition using visual attention. *Computer Science*, 2017.

[45] Zheng Shou, Junting Pan, Jonathan Chan, Kazuyuki Miyazawa, Hassan Mansour, Anthony Vetro, Xavier Giro-I-Nieto, and Shih Fu Chang. Online detection of action start in untrimmed, streaming videos. In *CVPR*, 2018.

[46] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In *NIPS*, 2014.

[47] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *Computer Science*, 2014.

- [48] Sijie Song, Cuiling Lan, Junliang Xing, Wenjun Zeng, and Jiaying Liu. An end-to-end spatio-temporal attention model for human action recognition from skeleton data. In *AAAI*, 2016.
- [49] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *Computer Science*, 2012.
- [50] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [51] Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew W. Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *CoRR*, abs/1609.03499, 2016.
- [52] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017.
- [53] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Anticipating the future by watching unlabeled video. *CoRR*, abs/1504.08023, 2015.
- [54] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Anticipating visual representations from unlabeled video. In *CVPR*, 2016.
- [55] Heng Wang, Alexander Klser, Cordelia Schmid, and Cheng Lin Liu. Dense trajectories and motion boundary descriptors for action recognition. *International Journal of Computer Vision*, 103(1):60–79, 2013.
- [56] Hongsong Wang and Jiashi Feng. Delving into 3d action anticipation from streaming videos. 2019.
- [57] Limin Wang, Yuanjun Xiong, Wang Zhe, and Qiao Yu. Towards good practices for very deep two-stream convnets. *Computer Science*, 2015.
- [58] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*, 2018.
- [59] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, and Klaus Macherey. Google’s neural machine translation system: Bridging the gap between human and machine translation. 2016.
- [60] Zuxuan Wu, Xi Wang, Yu-Gang Jiang, Hao Ye, and Xiangyang Xue. Modeling spatial-temporal clues in a hybrid deep learning framework for video classification. *CoRR*, abs/1504.01561, 2015.
- [61] Yuanjun Xiong, Limin Wang, Zhe Wang, Bowen Zhang, Hang Song, Wei Li, Dahua Lin, Yu Qiao, Luc Van Gool, and Xiaoou Tang. Cuhk and ethz and siat submission to activitynet challenge 2016.
- [62] Mingze Xu, Mingfei Gao, Yi-Ting Chen, Larry S. Davis, and David J. Crandall. Temporal recurrent networks for online action detection. In *ICCV*, 2019.
- [63] Fan Yang, Ke Yan, Shijian Lu, Huizhu Jia, Xiaodong Xie, and Wen Gao. Attention driven person re-identification. *Pattern Recognition*, pages S0031320318303133–, 2018.
- [64] Cao Yu, Daniel Barrett, Andrei Barbu, Siddharth Narayanaswamy, and Wang Song. Recognize human activities from partially observed videos. In *CVPR*, 2013.