Back to the Future: Knowledge Distillation for Human Action Anticipation

Vinh Tran, Yang Wang, Minh Hoai Stony Brook University Stony Brook, NY, 11794

{tquangvinh, wang33, minhhoai}@cs.stonybrook.edu

Abstract

We consider the task of training a neural network to anticipate human actions in video. This task is challenging given the complexity of video data, the stochastic nature of the future, and the limited amount of annotated training data. In this paper, we propose a novel knowledge distillation framework that uses an action recognition network to supervise the training of an action anticipation network, guiding the latter to attend to the relevant information needed for correctly anticipating the future actions. This framework is possible thanks to a novel loss function to account for positional shifts of semantic concepts in a dynamic video. The knowledge distillation framework is a form of self-supervised learning, and it takes advantage of unlabeled data. Experimental results on JHMDB and EPIC-KITCHENS dataset show the effectiveness of our approach.

1. Introduction

Human action anticipation is an important problem, but it is notoriously difficult due to the stochastic nature of the future. Given what is occurring or what can be observed in a video at the current moment, there are multiple possibilities that can happen. Thus, there is a fundamental limit to what we can anticipate, even when we have an infinite amount of training data. In practice, the amount of annotated training data is limited, so anticipation is a much harder problem.

One common approach to address anticipation is to use supervised learning, but learning a direct mapping between distant time steps can be challenging due to the weak correlation between the time steps. Suppose we are interested in anticipation with the lead time τ , we can used supervised learning and train a neural network to map from the video observation *up until* time t (denoted \mathbf{x}_t) to the human action label $y_{t+\tau}$ at time $t+\tau$. That is to use a set of annotated training data pairs $\{\mathbf{x}_t,y_{t+\tau}\}$ to train a network $\mathcal{A}:\mathbf{x}_t\to y_{t+\tau}$. To some extent, the training of the anticipation network \mathcal{A} can be done similarly to the training of a recognition network \mathcal{R} that maps from $\mathbf{x}_{t+\tau}$ to $y_{t+\tau}$, with

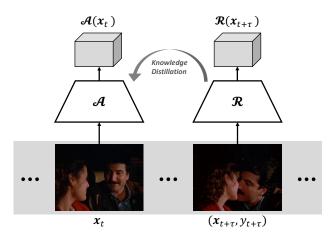


Figure 1: Knowledge Distillation for Action Anticipation. Our goal is to learn an action anticipation network $\mathcal{A}: \mathbf{x}_t \to y_{t+\tau}$ that can map from the current video observation \mathbf{x}_t to the human action category label $y_{t+\tau}$ of the future time step $t+\tau$. We propose to first learn an action recognition model $\mathcal{R}: \mathbf{x}_{t+\tau} \to y_{t+\tau}$, then use \mathcal{R} to supervise the training of \mathcal{A} through a novel knowledge distillation framework.

the only difference being that the input to \mathcal{A} is \mathbf{x}_t while the input to \mathcal{R} is $\mathbf{x}_{t+\tau}$. In general, the correlation between \mathbf{x}_t and $y_{t+\tau}$ is weaker than the correlation between $\mathbf{x}_{t+\tau}$ and $y_{t+\tau}$, so the asymptotic performance of \mathcal{A} is expected to be lower than the asymptotic performance of \mathcal{R} . Furthermore, \mathcal{A} will converge to its asymptotic performance slower than \mathcal{R} . This is due to the complexity of video data, and it will take much training data to separate the relevant features from the irrelevant ones. This separation task is harder for training the anticipation network than for training the recognition network due to the higher ratio of irrelevant features. In general, it will require more training data to get the anticipation network \mathcal{A} to "attend" to the relevant features.

Given that the recognition network is surely better than the anticipation network both asymptotically and nonasymptotically, we propose in this paper a knowledge distillation framework that uses a recognition network \mathcal{R} to guide the training of an anticipation network \mathcal{A} . Our framework trains the anticipation network to attend to the same type of information that is being attended by the recognition network when making classification decisions. Our framework leverages the abundance of (unannotated) data, improving the generalization ability of an anticipation network without requiring additional human annotation.

However, training an anticipation network to attend to the same information as the recognition network is technically challenging. A video is a dynamic environment, so the important concepts/features being attended by a recognition network at time $t+\tau$ might not be at the same location at time t. Thus, we cannot use an L_2 loss to force the one-to-one mapping between elements of two activated feature maps. To address this problem, we propose a novel attention mechanism that does not require pixel-to-pixel correspondence between two input videos or between two activated feature maps.

We perform experimental validation on three datasets: JHMDB, EPIC-Kitchens, and THUMOS. The experimental results show that the proposed knowledge distillation framework can improve the performance of the anticipation network. The level of improvement is consistent with the level of improvement obtained as if the annotated training data is doubled. The proposed approach also achieves the state-of-the-art performance on the JHMDB dataset.

2. Related Works

Anticipation and Early Recognition. Human action and activity anticipation is an emerging research area. Kong et al. [11] modeled the dynamics of human actions by capturing the temporal evolution of features over time to predict the class label as early as possible. Yuen and Torralba [28] presented a method to identify videos with unusual events in a large collection of short video clips. Soomro et al. [22] divided the video into multiple short segments and used dynamic programming on SVMs scores of each segment to obtain the final action class label. Some recent notable works in this area include: context-fusion [9], within-class Loss [15], ELSTM [20], and Feature Mapping RNN [21].

Note that action anticipation is different from early recognition [8, 13, 14, 16, 19, 20, 26]. Early recognition aims to detect and recognize an ongoing action as soon as possible, while anticipation refers to the ability to forecast the action even before it starts. These two problems are not unrelated, but they also face different technical challenges.

Most of the aforementioned techniques are based on supervised learning, and to some extent are orthogonal to what being presented here. We envision the combination of these techniques with the proposed knowledge distillation framework to obtain complementary benefits.

Distilling knowledge from unlabeled data. tion [7] has been widely used in deep learning as a way to transfer knowledge learned from different models. Le et al. [12] trained a face recognition model using large scale unlabeled images from the Internet. The effectiveness of additional training data has also been observed many times in deep learning for video analysis tasks. For example, models that are pre-trained on large scale video datasets such as Kinetics [2] and Sport1M [24] always have better performance than those that are trained with a limited amount of data. However, collecting such large scale annotated dataset require enormous human effort. The most similar to our work is the method of Vondrick et al. [25] that capitalizes on the temporal structure in unlabeled video to learn to anticipate human actions and objects. However, instead of solely using a single vector as the visual representation, we utilize the output feature map to preserve more local information, and propose a novel approach to handle the positional drift of semantic concepts in video.

3. Future Visual Representation Distillation

In this section, we describe a knowledge distillation framework that uses a recognition network to guide the training of an anticipation network. This framework is designed to leverage the abundance of unlabeled data, and the framework is founded on the assumption that a recognition network is better than an anticipation network both asymptotically and non-asymptotically.

3.1. Framework overview

Suppose the desired anticipation lead time is τ , our goal is to train an anticipation network \mathcal{A} to map from the input video segment at time $t-\tau$ (denoted $\mathbf{x}_{t-\tau}$) to the human action label y_t at time t. That is to train \mathcal{A} so that $\mathcal{A}(\mathbf{x}_{t-\tau}) = y_t$. We assume there is a recognition network \mathcal{R} to recognize the action class of an observed video clip, i.e., predicting y_t from \mathbf{x}_t . In general, the output $\mathcal{R}(\mathbf{x}_t)$ is only an approximate for y_t , but our framework does not require a perfect recognition network. It only assumes that the recognition network is better at predicting y_t than the anticipation network.

We consider here the case where both \mathcal{A} and \mathcal{R} are convolutional neural networks with multiple layers. Without being too specific, let $\bar{\mathcal{A}}(\mathbf{x})$ denote the feature vector/map at a particular layer of the anticipation network for the input video \mathbf{x} ; this layer could be the output or an intermediate layer. Essentially, $\bar{\mathcal{A}}(\mathbf{x})$ can be considered as the feature vector/map that represents the input video \mathbf{x} . Similarly, let $\bar{\mathcal{R}}(\mathbf{x})$ be the feature vector/map of the recognition network for the input video \mathbf{x} .

Let S denote the set of time indexes where the frames are annotated with human action labels; t is in S if y_t is available. One approach for training the anticipation network

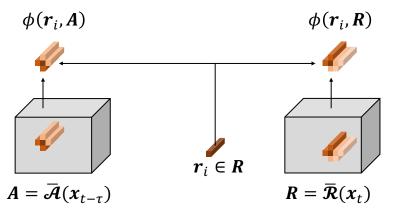


Figure 2: **Knowledge distillation for weakly aligned feature maps.** $\mathbf{R} = \bar{\mathcal{R}}(\mathbf{x}_t)$ encodes the activated features important for the recognition network \mathcal{R} to recognize the human action at time t. For example, in order to recognize a "wash a dish" action, $\mathbf{r}_i \in \mathbf{R}$ might indicate the presence of the dish in the video. Arguably, for the anticipation network \mathcal{A} to successfully anticipate the same action, there must be some activated "dish" features in its feature map $\mathbf{A} = \bar{\mathcal{A}}(\mathbf{x}_{t-\tau})$. Our proposed knowledge distillation framework encourages \mathbf{A} to have the 'same' activated features with \mathbf{R} . However, it would be unreasonable to directly minimize $\|\mathbf{A} - \mathbf{R}\|_2^2$, because the "dish" features in \mathbf{A} and \mathbf{R} might be at different spatiotemporal locations. To perform knowledge distillation for weakly aligned feature maps, we propose an attentional pooling operator ϕ to compute the amount of \mathbf{r}_i within \mathbf{R} and \mathbf{A} respectively, then minimize $\sum_i \|\phi(\mathbf{r}_i, \mathbf{R}) - \phi(\mathbf{r}_i, \mathbf{A})\|_2^2$ instead.

is to minimize the following classification loss defined on annotated training data:

$$\sum_{t \in \mathcal{S}} \mathcal{L}_c(\mathcal{A}(\mathbf{x}_{t-\tau}), y_t), \tag{1}$$

Here, \mathcal{L}_c is a loss function that penalizes the difference between the prediction output $\mathcal{A}(\mathbf{x}_{t-\tau})$ and the actual class label y_t , e.g., using the negative log likelihood loss. In general, labeled data is scarce, so it is difficult to obtain good anticipation performance with supervised learning by itself.

Let \mathcal{U} be the set of time indexes t's where y_t is not available (i.e., unlabeled data). Our knowledge distillation framework optimizes the below loss function:

$$\sum_{t \in \mathcal{S}} \mathcal{L}_c(\mathcal{A}(\mathbf{x}_{t-\tau}), y_t) + \sum_{t \in \mathcal{U}} \mathcal{L}_c(\mathcal{A}(\mathbf{x}_{t-\tau}), \mathcal{R}(\mathbf{x}_t))$$
 (2)

$$+ \lambda \sum_{t \in \mathcal{U} \cup \mathcal{S}} \mathcal{L}_d(\bar{\mathcal{A}}(\mathbf{x}_{t-\tau}), \bar{\mathcal{R}}(\mathbf{x}_t)). \tag{3}$$

The above objective function trains the anticipation network \mathcal{A} to output the same output as the recognition network on unlabeled data \mathcal{U} . Furthermore, \mathcal{L}_d is a loss function that measures the discrepancy between two feature maps $\bar{\mathcal{A}}(\mathbf{x}_{t-\tau})$ and $\bar{\mathcal{R}}(\mathbf{x}_t)$. This loss trains the anticipation network to produce the same feature map as the feature map of the recognition network. This formulation uses unlabeled data and the distilled knowledge from the recognition network to guide the anticipation network to attend to the relevant information that is useful for categorizing the future action.

3.2. Distillation loss

In this section, we describe the details of the loss function \mathcal{L}_d for measuring the differences between two activation feature maps. At first glance, a reasonable option for this loss function is to use the sum of squared differences between the elements of the feature maps. However, this loss assumes perfect correspondence between the elements of the feature maps. This turns out to be too restrictive, as will be explained below.

We use convolutional architectures for the anticipation and recognition networks, and the feature maps $\bar{\mathcal{A}}(\mathbf{x})$ and $\bar{\mathcal{R}}(\mathbf{x})$ are typically 4D tensors: $\bar{\mathcal{A}}(\mathbf{x}), \bar{\mathcal{R}}(\mathbf{x}) \in \Re^{l \times h \times w \times d}$. Due to the use of strided convolutional and pooling layers within the network, the dimensions of the feature maps may be different from the sizes of the input video. Usually, l, h, and w can be obtained by dividing the length, the height, and the width of the video \mathbf{x} by their effective convolutional strides respectively. d is the number of channels of the feature map.

Consider a particular video segment \mathbf{x}_t , the feature map $\bar{\mathcal{R}}(\mathbf{x}_t)$ encodes the activated features important for recognizing the human action. For example, in order for the recognition network to recognize a "wash a dish" action, some part of the feature map might indicate the presence of the dish in the video. Arguably, for the anticipation network to successfully anticipate the "wash a dish" action, there must be some activated "dish" features in its feature map. The knowledge distillation framework aims to encourage that, by training the anticipation network to output the 'same' feature map as the recognition network. However,

it would be unreasonable to assume the "dish" feature to stay at the same spatiotemporal location of the feature map. More generally, video is a dynamic environment, where important objects and other semantic entities might not remain at the same locations. Thus, it is unreasonable to use the sum of squared differences to measure the discrepancy between two feature maps of two different time steps.

For brevity, let us reshape the 4D tensors $\mathcal{A}(\mathbf{x}_{t-\tau})$ and $\bar{\mathcal{R}}(\mathbf{x}_t)$ to 2D matrices \mathbf{A} and \mathbf{R} , $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_{lhw}] \in \Re^{d \times lhw}$ and $\mathbf{R} = [\mathbf{r}_1, \dots, \mathbf{r}_{lhw}] \in \Re^{d \times lhw}$. For each vector \mathbf{a}_i , we measure the similarity between \mathbf{a}_i with all vectors in \mathbf{R} and we compute a vector quantity to represent the amount of \mathbf{a}_i in \mathbf{R} :

$$\phi(\mathbf{a}_i, \mathbf{R}) = \sum_{k=1}^{lhw} \omega_k \mathbf{r}_k. \tag{4}$$

with $\omega_k = \frac{1}{Z} \exp(\alpha \mathbf{r}_k^T \mathbf{a}_i)$, where Z is the normalizing constant so that the sum of ω_k 's is 1. Here, α is the hyperparameter that controls the pooling weights. The default value of α is set to $\frac{1}{\sqrt{d}}$, where d is the number of channels of the feature maps. If the value of α is small, the weights associated to each vector are almost equal, and $\phi(\mathbf{a}_i, \mathbf{R})$ is the average pooling of vectors in \mathbf{R} . On the other hand, this operator is similar to max pooling if we use a large value for α , and $\phi(\mathbf{a}_i, \mathbf{R})$ is the vector in \mathbf{R} that is most similar to \mathbf{a}_i . Equivalently, $\phi(\mathbf{a}_i, \mathbf{R})$ can be expressed in the form:

$$\phi(\mathbf{a}_i, \mathbf{R}) = \mathbf{R} \operatorname{softmax}(\alpha \mathbf{R}^T \mathbf{a}_i), \tag{5}$$

Similarly, we can compute a vector to represent the amount of \mathbf{a}_i in \mathbf{A} :

$$\phi(\mathbf{a}_i, \mathbf{A}) = \mathbf{A} \operatorname{softmax}(\alpha \mathbf{A}^T \mathbf{a}_i). \tag{6}$$

We propose to define the loss for the differences between two feature maps ${\bf A}$ and ${\bf R}$ as follows:

$$\mathcal{L}_d(\mathbf{A}, \mathbf{R}) = \tilde{\mathcal{L}}_d(\mathbf{A}, \mathbf{R}) + \tilde{\mathcal{L}}_d(\mathbf{R}, \mathbf{A}), \tag{7}$$

where
$$\tilde{\mathcal{L}}_d(\mathbf{A}, \mathbf{R}) = \sum_{i=1}^{lhw} ||\phi(\mathbf{a}_i, \mathbf{R}) - \phi(\mathbf{a}_i, \mathbf{A})||_2^2$$
, (8)

$$\tilde{\mathcal{L}}_d(\mathbf{R}, \mathbf{A}) = \sum_{i=1}^{lhw} ||\phi(\mathbf{r}_i, \mathbf{A}) - \phi(\mathbf{r}_i, \mathbf{R})||_2^2.$$
 (9)

Note that \mathcal{L}_d is a symmetric loss function, while $\tilde{\mathcal{L}}_d$ is not: $\tilde{\mathcal{L}}_d(\mathbf{A}, \mathbf{R})$ is different from $\tilde{\mathcal{L}}_d(\mathbf{R}, \mathbf{A})$.

4. Experiments

4.1. Datasets

We conducted main experiments on two challenging datasets: JHMDB [10] and EPIC-KITCHENS [4]. We also

performed some controlled experiments on the THUMOS dataset [5] to understand the expected benefits of having extra supervision.

JHMDB. This dataset contains 928 videos from 21 human actions. The videos are collected from movies or from the Internet. The dataset has three splits for training and testing. We follow the standard protocol [21] for evaluation on this dataset, i.e., using only the first 20% of the frames to predict action class labels.

EPIC-KITCHENS. This is the largest Egocentric video dataset for action anticipation. The dataset contains 432 videos recorded by 32 different people in their native kitchen environments. We focus on action anticipation on verb class for this dataset (i.e., 125 action classes). Since the annotations for the test set are not yet available, we follow [1] and split the training set into training and validation subsets. We use the videos recorded by P01 to P25 for training and the remaining videos for validation. In total, there are 23191 clips in the training set and 5281 clips in the validation set. For this dataset, the performance is measured using top-1 and top-5 accuracy values. For brevity, we also refer to top-1 accuracy simply as accuracy.

THUMOS14. This dataset contains 413 videos with annotated action segments. It contains 200 videos for training and 213 videos for testing. This dataset contains untrimmed videos of human actions, so it can be used for anticipation experiments.

4.2. Experiments on JHMDB dataset

We performed several experiments on the JHMDB dataset. We used the I3D network as the backbone for this task. The network was initialized with the weights of a network that was pre-trained on the Kinetics dataset. Training data consisted of pairs of video clips, each with 8 frames. We used spatial crops of 224×224 pixels. The anticipated time between the two clips was set to be $\tau = 15$ frames. We extracted from the Mixed_4f layer of the I3D to produce the $2 \times 14 \times 14 \times 832$ output feature map. During training, we combined both the classification loss \mathcal{L}_c (using the class labels or the predicted class probability) and the attention loss \mathcal{L}_d (using feature maps). We used KL divergence for the classification loss. For the attention loss, we used Huber loss (with $\delta = 1$) instead of mean squared error for better optimization. We set the value α to be $\frac{1}{\sqrt{d}}$. All the models were trained using SGD with momentum 0.9 and weight decay 10^{-5} . The training started with an initial learning rate of 0.01 and stopped after 80 epochs. The learning rate was decreased by a factor of 10 at epoch 20. Both RGB frames and Optical Flow maps were used to train an anticipation network.

We report the action recognition and anticipation performance of different methods in Table 1. We first trained the

Method	RGB	Flow	Both
Recognition Network	75.3	77.8	83.9
	13.3	17.8	83.9
Anticipation Network			
I3D	69.0	64.4	74.9
I3D + \mathcal{L}_{direct}	69.5	67.1	75.5
I3D + $\tilde{\mathcal{L}}_d(\mathbf{R}, \mathbf{A})$	70.0	67.3	75.7
$I3D + \mathcal{L}_d(\mathbf{A}, \mathbf{R})$	70.2	67.7	75.9

Table 1: Action anticipation results on the JHMDB dataset. The recognition network uses the entire video for classification while the anticipation network only observes the first 20% of the video. All the results are averaged over three splits of the dataset. The improvement on anticipation obtained by our approach, compared to the baseline I3D network, is most evident on the Flow stream (+3.3%). Compared to the direct loss \mathcal{L}_{direct} , the proposed distillation loss $\mathcal{L}_d(\mathbf{A}, \mathbf{R})$ addresses the positional drift problem of the semantic concepts within videos, and achieves better results. We also observe that the combined bidirectional loss $\mathcal{L}_d(\mathbf{A}, \mathbf{R})$ outperforms its unidirectional variant $\tilde{\mathcal{L}}_d(\mathbf{R}, \mathbf{A})$.

action recognition network using all the available frames in the training set. For the action recognition task, we achieved an accuracy of 83.9% (using both RBG and Flow). For reference, this is comparable to the number reported in [3], which achieved 84.1% accuracy using the same model. Directly applying the recognition network on the first 20% frames of the test videos (i.e., using the recognition network for the anticipation task), the accuracy dropped drastically to 74.9%. In the next step, we applied the direct loss (denoted as \mathcal{L}_{direct} as in Table 1) between two feature maps produced by anticipation network and recognition network. With this additional loss, the accuracy is increased to 75.5\%. This is possibly thanks to the small displacement between two feature maps since the dataset contains only action and the 15 frames anticipation is quite short. Hence, the \mathcal{L}_{direct} loss also helps improving the recognition performance on this dataset. Replacing the direct loss function with our distillation loss function that computes the amount of R in A, the performance increase to 75.7%. Finally, we achieved the best performance of 75.9% when using the symmetric bidirectional attention loss $\mathcal{L}_d(A, R)$. As shown in Table 2, we obtained the new state of the art result on the JHMDB dataset.

We illustrate our results on the JHMDB dataset in Figure 3. We visualize the activation at layer Mixed_4f by computing the magnitude of the feature vectors at each spatiotemporal location. The first three columns show the attention maps of different anticipation networks on the same input video (15 frames prior to the start of the action seg-

Method	Accuracy(%)
DP-SVM [23]	5.0
S-SVM [23]	5.0
Where/What [22]	10.0
Context-fusion [9]	28.0
Within-class Loss [15]	33.0
ELSTM [20]	55.0
Future Dynamic Images [17]	61.0
Feature Mapping RNN [21]	73.4
I3D + Knowledge Distillation (Proposed)	75.9

Table 2: Comparison of action anticipation methods on the JHMDB dataset. All methods use the first 20% of the video for prediction. Our approach advances the state-of-the-art performance on this dataset.

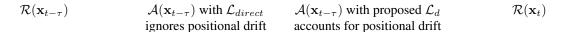
ment). The last column shows the attention map of the recognition network on the action segment. As can be seen from the figure, our method with knowledge distillation is able to 'pay more attention' to locations that are being attended by the recognition network.

4.3. Experiments on Epic-Kitchens dataset

We used I3D network architecture for the experiments described in this section, as in the experiments in the previous section. However, the Epic-Kitchen dataset is significantly larger than the other datasets, while we were constrained by the limited computational resource. For feasible and efficient training on this dataset, we used the feature maps extracted from the MaxPool3d_4a_3x3 layer as input for the network instead of training directly from video frames. To extract the features, we fed the input video clip of 32 frames with a spatial resolution of 256×456 to a pretrained I3D network (that had been trained on the Kinetics dataset). This I3D network generated a $8 \times 16 \times 29$ output feature map with d = 480 channels. During training, we sampled a spatial crop of 14×14 on the feature map and fed it to the remaining layers of the I3D network. For testing, we used a single center crop of the same size. The same type of feature maps were used as inputs for both the anticipation and the recognition networks.

4.3.1 Collecting unlabeled training data

We used the provided action segments to train the action recognition network. The inputs to the network were also the feature maps generated at the MaxPool3d_4a_3x3 layer of the pre-trained I3D network. We trained the recognition network for 30 epochs using SGD optimization with momentum 0.9, weight decay 10^{-5} , and initial learning rate of 0.1. The learning rate was decreased by a factor of 10



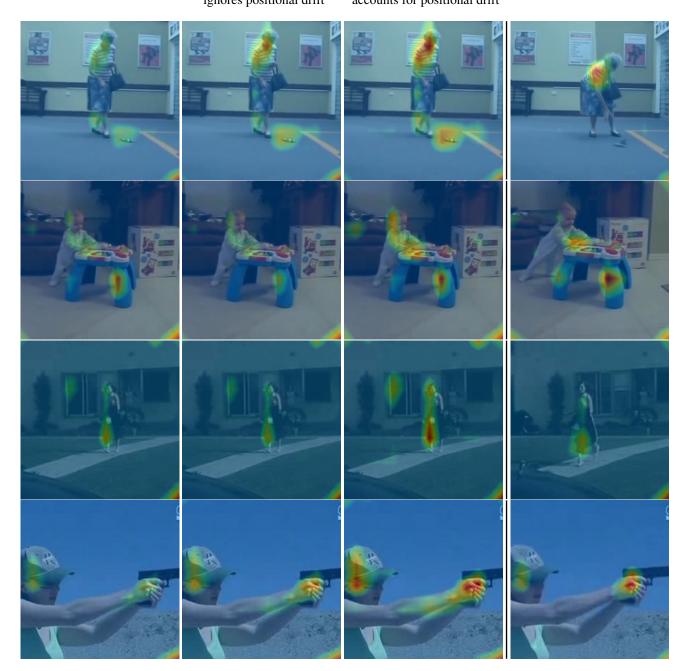


Figure 3: **Visualize informative regions on JHMDB dataset.** From top to bottom, the actions are: *Pick, Push, Run, Shoot Gun.* The last column shows the regions that are being attended by the recognition network \mathcal{R} when it observes (future) action segment. Higher intensity of red color means larger magnitude of activation vector at that location. The first column shows the regions that receive high level of attention when the recognition network \mathcal{R} sees the anticipation segment. The baseline approach visualized in the second column assumes strong alignment between the anticipation segments and the action segments, and uses \mathcal{L}_{direct} to train the anticipation network \mathcal{A} . The third column visualizes the anticipation network \mathcal{A} that is trained with \mathcal{L}_d , the proposed distillation loss for weakly aligned feature maps. As can be observed, the proposed knowledge distillation framework can help the anticipation network \mathcal{A} focus more on the informative regions picked up by the recognition network \mathcal{R} .

after every 5 epoch. For the recognition task on verb action classes, we achieved the recognition accuracy of 46.1% with single center crop evaluation.

We collected video clips from unlabeled segments as follows. First, we randomly took two video segments of 32 frames with anticipation time $\tau=1$ second from unlabeled video segments. To avoid tampering with the test/validation segments, we only collected the data in the 10-second vicinity of the training segments. Second, we used the pretrained I3D (pre-trained on Kinetics) to extract feature maps at the MaxPool3d_4a_3x3 layer for the two video segments. Third, we fed the feature map of the latter segment to the recognition network and computed its visual representation (i.e., both the class probabilities and the feature maps). Finally, the feature map of the first segment and the visual representation of the second segment formed a datapair sample for training the anticipation network.

Using the procedure described above, we collected a total of 26391 unlabeled training samples. Together with the 23191 provided annotations, we have 49582 training samples in total. This is roughly double the amount of the original training data. In theory, the size of unlabeled training data can be made larger, but we were bounded by the computational resource.

During training, if a sample did not have a label, we used the recognition network to provide supervision. Otherwise, we used the provided annotation. We used the KL-Divergence for the classification loss in our knowledge distillation framework to train the anticipation network.

4.3.2 Experimental Results

We trained the I3D network with a mini-batch size of 32 samples. Similar to the optimization for the recognition network, we used SGD with momentum 0.9 and weight decay 10^{-5} for optimization. We trained all the models for 30 epochs. The initial learning rate was set to 0.001, and it was decreased by a factor of 10 after Epoch 20.

We reported all performance values in Table 3. When we trained the I3D network using the annotated training data, we achieved an accuracy of 30.1%. This is better than the performance of the TSN-RGB [27] method, which only achieves 28.5% accuracy on the same dataset. Using the knowledge distillation framework and additional unlabeled data, the obtained anticipation network obtained 1.7% absolute improvement in accuracy ($30.1 \rightarrow 31.8$). The Top-5 accuracy of the TSN-RGB, I3D, and I3D+Knowledge distillation are 73.4, 74.3, and 74.4, respectively.

4.4. Experiments on THUMOS14 dataset

We propose in this paper a method that distills knowledge from a recognition network to supervise the training of an anticipation network. As shown in the previous two

Method	Accuracy(%)
TSN-RGB [27]	28.5
I3D [2]	30.1
I3D + Knowledge Distillation (Proposed)	31.8

Table 3: The accuracy of anticipation methods on the EPIC-KITCHENS dataset (for anticipating the verb actions). All methods reported here are implemented by us, trained with the same amount of labeled data.

subsections, this method yields an absolute improvement in anticipation accuracy between 1%-1.7%. The improvement on the Flow stream is even more evident (+3.3% on JHMDB dataset). This level of improvement is significant, and our method outperforms the previous state-of-the-art on the JHMDB dataset. At the same time, this level of improvement is not as high as one would expect. In this section, we perform some controlled experiments to rectify the expected level of improvement.

In the proposed framework, the recognition network acts as a teacher, supervising the training of the anticipation network. To some extent, the recognition network plays the role of an annotator that provides some form of annotation for the unlabeled data. This form of annotation arguably has lower quality than the actual labels. Given this link between knowledge distillation and extra annotation, we performed some controlled experiments to analyze how the size of the annotated training set affects the performance of an anticipation network trained with supervised learning.

We used videos from the THUMOS14 action detection challenge [5] to create a dataset for action anticipation. We first identified the temporal location of an action segment. Interested in the anticipation lead time of one second, we moved back one second and extracted a one-second clip ending at that location. This one-second clip will be used as the input to the anticipation network. Using this strategy, we can compile a training dataset of multiple one-second clips. From the full training set, we randomly dropped 50% of the data to create another smaller training set with only half of the full training data. We then trained two anticipation networks, one using the full training set and the other using the smaller training set with 50% of the data.

We experimented with both 2D and 3D ConvNet architectures to see how the size of the training set affects the anticipation performance. To train a 2D network, we used Resnet152 [6] (pre-trained on ImageNet/ILSVRC [18]) to extract a feature embedding vector for each frame. Then we used a multilayer perceptron to predict the action class on top of the extracted features. For evaluation, we computed the classification score for each frame and averaged the prediction. To train a 3D network, we used the I3D network to train on the input clips of 16 frames with spatial crops of

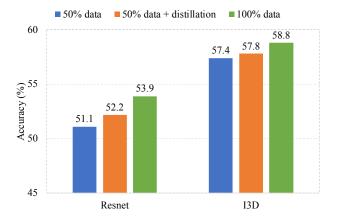


Figure 4: Expected performance gain when doubling the amount of training data. This shows the performance of anticipation networks trained with supervised learning using either all training data or 50% of the training data. This experiment is done on the THUMOS dataset using either ResNet or I3D features. In general, doubling the amount of training data yields a moderate improvement from 1% to 3%, and the gain is larger for weaker features. Using half of the labeled data and the other half as unlabeled data for knowledge distillation leads to some improvement, but the level of improvement is not as good as having actual ground truth labels. After all, the recognition network acts in the role of an imperfect teacher rather than an oracle.

 224×224 pixels. For evaluation using the I3D network, we use the center crops of the test clip.

The anticipation performance of these networks is plotted in Figure 4. As can be seen, using more data clearly improves the performance of an anticipation network. When doubling the amount of annotated training data, the gain in accuracy of the two networks (for two types of features) are 1.4% and 2.8%. This experiment leads us to believe that doubling the amount of annotated training data will only yield moderate improvement in anticipation accuracy, perhaps somewhere from 1% to 3%. In the experiments described in the previous subsections regarding the performance of the proposed knowledge distillation framework, the amount of unlabeled data used is roughly the same with the amount of labeled data, and the total amount of training data is roughly doubled. To some extent, the level of improvement obtained by the knowledge distillation framework is consistent with the level of improvement observed in this experiment when we double the amount of annotated training data.

Of course, the quantity of unlabeled data is not the only factor that determines the amount of performance gain for using the proposed knowledge distillation framework. Other important factors are the accuracy of the recognition

network and also the quality of the unlabeled data. Figure 4 also shows the performance of the anticipation networks trained with knowledge distillation. In this experiment, we only use half of the labeled training data, and we use the other half as unlabeled data for knowledge distillation. As can be seen, some performance gain is obtained by using the knowledge distillation framework. The level of improvement is not as good as having actual ground truth annotations, and this can probably be attributed to the imperfection of the recognition network. After all, the recognition network acts in the role of an imperfect teacher rather than an oracle.

5. Summary

In this paper, we have presented a framework for knowledge distillation. This framework uses the action recognition network to supervise the training of an action anticipation network. With a novel knowledge distillation technique to account for the positional drift of semantic concepts in video, the action recognition network acts as a teacher guiding the anticipation network to attend to the relevant information needed for predicting the future action. Using this framework, we are able to leverage unlabeled data to train the anticipation network in a self-supervised manner. The experimental results on the JHMDB and EPIC-KITCHENS datasets show the benefits of our proposed method.

References

- [1] F. Baradel, N. Neverova, C. Wolf, J. Mille, and G. Mori. Object level visual reasoning in videos. In *Proc. ECCV*, 2018.
- [2] J. Carreira and A. Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *Proc. CVPR*. IEEE, 2017.
- [3] V. Choutas, P. Weinzaepfel, J. Revaud, and C. Schmid. Potion: Pose motion representation for action recognition. In Proc. CVPR, 2018.
- [4] D. Damen, H. Doughty, G. M. Farinella, S. Fidler, A. Furnari, E. Kazakos, D. Moltisanti, J. Munro, T. Perrett, W. Price, and M. Wray. Scaling egocentric vision: The epickitchens dataset. In *Proc. ECCV*, 2018.
- [5] A. Gorban, H. Idrees, Y.-G. Jiang, A. Roshan Zamir, I. Laptev, M. Shah, and R. Sukthankar. THUMOS challenge: Action recognition with a large number of classes. http://www.thumos.info, 2015.
- [6] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proc. CVPR*, 2016.
- [7] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [8] M. Hoai and F. De la Torre. Max-margin early event detectors. *IJCV*, 107(2):191–202, 2014.
- [9] A. Jain, A. Singh, H. S. Koppula, S. Soh, and A. Saxena. Recurrent neural networks for driver activity anticipation via sensory-fusion architecture. In *Proc. Intl. Conf. on Robotics* and Automation. IEEE, 2016.

- [10] H. Jhuang, J. Gall, S. Zuffi, C. Schmid, and M. J. Black. Towards understanding action recognition. In *Proc. ICCV*, 2013
- [11] Y. Kong, D. Kit, and Y. Fu. A discriminative model with multiple temporal scales for action prediction. In *Proc. ECCV*, 2014.
- [12] Q. V. Le, M. Ranzato, R. Monga, M. Devin, K. Chen, G. S. Corrado, J. Dean, and A. Y. Ng. Building high-level features using large scale unsupervised learning. In *Proc. ICML*, 2012.
- [13] K. Li and Y. Fu. Prediction of human activity by discovering temporal sequence patterns. *IEEE PAMI*, 36(8):1644–1657, 2014
- [14] S. Ma, L. Sigal, and S. Sclaroff. Learning activity progression in lstms for activity detection and early detection. In Proc. CVPR, 2016.
- [15] S. Ma, L. Sigal, and S. Sclaroff. Learning activity progression in lstms for activity detection and early detection. In Proc. CVPR, 2016.
- [16] M. Pei, Y. Jia, and S.-C. Zhu. Parsing video events with goal inference and intent prediction. In *Proc. ICCV*, 2011.
- [17] C. Rodriguez, B. Fernando, and H. Li. Action anticipation by predicting future dynamic images. In ECCV'18 workshop on Anticipating Human Behavior, 2018.
- [18] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. Imagenet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015.
- [19] M. Ryoo. Human activity prediction: Early recognition of ongoing activities from streaming videos. In *Proc. ICCV*, 2011.
- [20] M. Sadegh Aliakbarian, F. Sadat Saleh, M. Salzmann, B. Fernando, L. Petersson, and L. Andersson. Encouraging lstms to anticipate actions very early. In *Proc. ICCV*, 2017.
- [21] Y. Shi, B. Fernando, and R. Hartley. Action anticipation with rbf kernelized feature mapping rnn. In *Proc. ECCV*, 2018.
- [22] K. Soomro, H. Idrees, and M. Shah. Predicting the where and what of actors and actions through online action localization. In *Proc. CVPR*, 2016.
- [23] K. Soomro, H. Idrees, and M. Shah. Online localization and prediction of actions and interactions. *IEEE PAMI*, 41(2), 2019
- [24] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features with 3D convolutional networks. In *Proc. ICCV*, 2015.
- [25] C. Vondrick, H. Pirsiavash, and A. Torralba. Anticipating visual representations from unlabeled video. In *Proc. CVPR*, 2016
- [26] B. Wang and M. Hoai. Back to the beginning: Starting point detection for early recognition of ongoing human actions. In CVIU, 2018.
- [27] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Val Gool. Temporal segment networks: Towards good practices for deep action recognition. In *Proc. ECCV*, 2016.
- [28] J. Yuen and A. Torralba. A data-driven approach for event prediction. In *Proc. ECCV*, 2010.