

FrameExit: Conditional Early Exiting for Efficient Video Recognition

Amir Ghodrati* Babak Ehteshami Bejnordi* Amirhossein Habibian
Qualcomm AI Research[†]

{ghodrati, behtesha, ahabibia}@qti.qualcomm.com

Abstract

In this paper, we propose a conditional early exiting framework for efficient video recognition. While existing works focus on selecting a subset of salient frames to reduce the computation costs, we propose to use a simple sampling strategy combined with conditional early exiting to enable efficient recognition. Our model automatically learns to process fewer frames for simpler videos and more frames for complex ones. To achieve this, we employ a cascade of gating modules to automatically determine the earliest point in processing where an inference is sufficiently reliable. We generate on-the-fly supervision signals to the gates to provide a dynamic trade-off between accuracy and computational cost. Our proposed model outperforms competing methods on three large-scale video benchmarks. In particular, on ActivityNet1.3 and mini-kinetics, we outperform the state-of-the-art efficient video recognition methods with $1.3\times$ and $2.1\times$ less GFLOPs, respectively. Additionally, our method sets a new state of the art for efficient video understanding on the HVU benchmark.

1. Introduction

With the massive growth in the generation of video content comes an increasing demand for efficient and scalable action or event recognition in videos. The astounding performance of deep neural networks for action recognition [5, 43, 58, 12, 42, 56] are obtained by densely applying 2D [56, 31, 58, 13] or 3D [40, 5, 18, 12] models over video frames. Despite demonstrating top-notch performance in recognizing complex and corner case actions, the high data volumes, compute demands, and latency requirements, limit the application of the state-of-the-art video recognition models on resource-constrained devices.

Extensive studies have been conducted to remedy this issue by designing efficient and light-weight architectures [35, 11, 43, 36, 53, 59, 42, 28, 31, 10]. These models have a static computational graph and treat all the videos

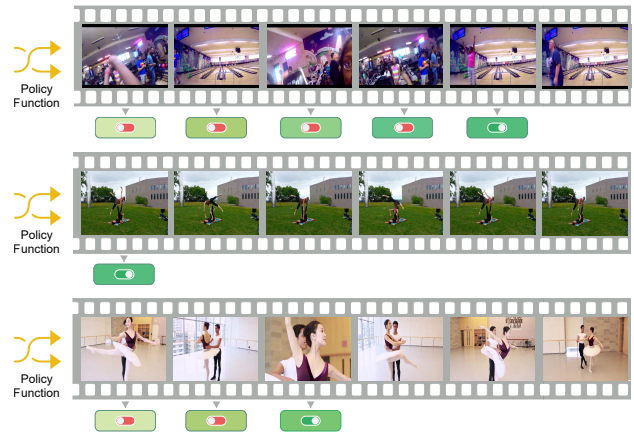


Figure 1: **Efficient video recognition by early exiting.** Our proposed method adjusts the amount of computation to the difficulty of the input, allowing for significant reduction of computational costs. Videos are adopted from [21, 22, 23].

equally regardless of how complex or easy they are for recognition and hence yield sub-optimal results. A large body of research has been focusing on selecting a subset of salient frames to efficiently process the video conditioned on the input [55, 52, 57, 8, 15, 29]. Current methods for frame selection rely on learning a policy function to determine what action should be taken on the selected frame (e.g. process by a heavy recognition model [51], process at a specific spatial resolution [33], etc.). Most of these methods either rely on the assumption that salient frames for the sampler network are also salient for the recognition network [29, 15] or require carefully selected reward functions in case of using policy gradient methods [55, 52, 49]. Moreover, the sampler network may create an additional computational overhead to the model.

An alternative promising direction to reduce the computational complexity of analyzing video content is conditional compute using early exiting. Early exiting has recently been explored for image classification tasks by inserting a cascade of intermediate classifiers throughout the network [30, 39, 54, 24]. In this line of work, the model adjusts the amount of computation to the difficulty of the

*Equal contribution

[†]Qualcomm AI Research is an initiative of Qualcomm Technologies, Inc

input and allows for significant reduction of computational requirements. Inspired by that, we design an efficient video recognition model that performs automatic early exiting by adjusting the computational budget on a per-video basis. Our motivation is that a few frames are sufficient for classifying “easy” samples, and only some “hard” samples need temporally detailed information (see Figure 1).

In this paper, we propose FrameExit, a conditional early exiting framework with learned gating units that decide to stop the computation when an inference is sufficiently reliable. FrameExit has T classifiers accompanied by their associated gates that are attached at different time steps to allow early exiting. The gates are learned in a self-supervised fashion to control the trade-off between model accuracy and total computation costs. We use the recognition loss as a proxy to generate on-the-fly pseudo-labels to train the gates.

Additionally, our early exiting mechanism combined with a simple, deterministic sampling strategy obviates the need for complex sampling policy functions and yet achieves excellent recognition performance. Finally, we propose an accumulated feature pooling module to generate video representations that enable more reliable predictions by the model. Our contributions are as follows:

- We propose a method that employs a simple, deterministic frame sampling strategy, combined with an accumulated feature pooling module to obtain accurate action recognition results.
- We propose a conditional early exiting framework for efficient video recognition. We use a cascade of gating modules to determine when to stop further processing of the video. The gates adapt the amount of computation to the difficulty of an input video, leading to significant reductions in computational costs.
- We show state-of-the-art performance on three large-scale datasets. In all cases, we greatly improve the inference efficiency at a better or comparable recognition accuracy. In particular, on the HVU [7] dataset, we report $5\times$ reduction in computation costs while improving the recognition accuracy upon the state-of-the-art methods.

2. Related work

Efficient Video Recognition: There are two lines of approaches for efficient video recognition. The first focuses on proposing new lightweight video architectures. This can be obtained by decomposing 3D filters into separate 2D spatial and 1D temporal filters [43, 36, 53], extending efficient 2D architectures to their 3D counterparts [42, 28], using shifting operations [31, 9], or exploring neural architecture search for videos [35, 11]. Our method is agnostic to

the network architecture and can be complementary to these types of methods.

The second line of works focus on saving compute by selecting salient frames/clips using a policy function parametrized with a neural network (CNN or RNN). This is commonly done by training agents to find which frame to observe next [55, 52, 57, 8], ranking frames based on saliency [29], gating frames [25], skipping RNN states [4], or by using audio to select relevant frames [15, 29]. LiteEval [51] proposes binary gates for selecting coarse or fine features. At each time step, LiteEval computes coarse features with a lightweight CNN to determine whether to examine the current frame more carefully using a heavy CNN. Meng *et al.* [33] propose to adaptively select the input resolution, on a per-frame basis to further balance the accuracy vs compute. However, using policy networks in these methods may come with additional memory and compute overhead. Moreover, optimization is sub-optimal if policy and recognition networks are not trained jointly [29, 15], or requires carefully selected reward functions if policy gradient methods are used [52, 55, 49]. In contrast, our method relies on a single recognition network and does not require complex RL optimization for frame sampling. We formulate the problem in the early-exiting framework and show that a simple sampling strategy, if combined with a proper exiting function, leads to excellent recognition performance.

Conditional compute via early exiting: Conditional computation in neural networks aims to dynamically allocate the components of the model (e.g. layers, sub-networks, etc.) on a per-input basis. ConvNet-AIG [45], SkipNet [48], and BlockDrop [50] exploit the robustness of Residual Networks (ResNets) to layer dropping and activate or deactivate full residual blocks conditioned on the input, to save compute. GaterNet [6], dynamic channel pruned networks [16], batch-shaped channel-gated networks [2], and Conditional channel gated networks for continual learning [1] perform gating in a more fine-grained fashion by turning individual channels of a layer on or off conditioned on the input. Other works focus on adaptive spatial attention for faster inference [14, 46, 47]. Concurrent to our work, AdaFuse [34] dynamically fuses channels from current and history feature maps to reduce the computation. A major research direction in conditional compute is early exiting. Prior works have mainly used early exiting for image classification. To adapt computational costs to the difficulty of the input, Deeply-Supervised Nets [30] and BranchyNet [39] propose architectures that are composed of a cascade of intermediate classifiers. This allows simpler examples to exit early via an intermediate classifier while more difficult samples proceed deeper in the network for more accurate predictions. Multi-scale dense networks [24] and adaptive resolution networks [54] focus on spatial redundancy of input samples and use a multi-scale dense connection archi-

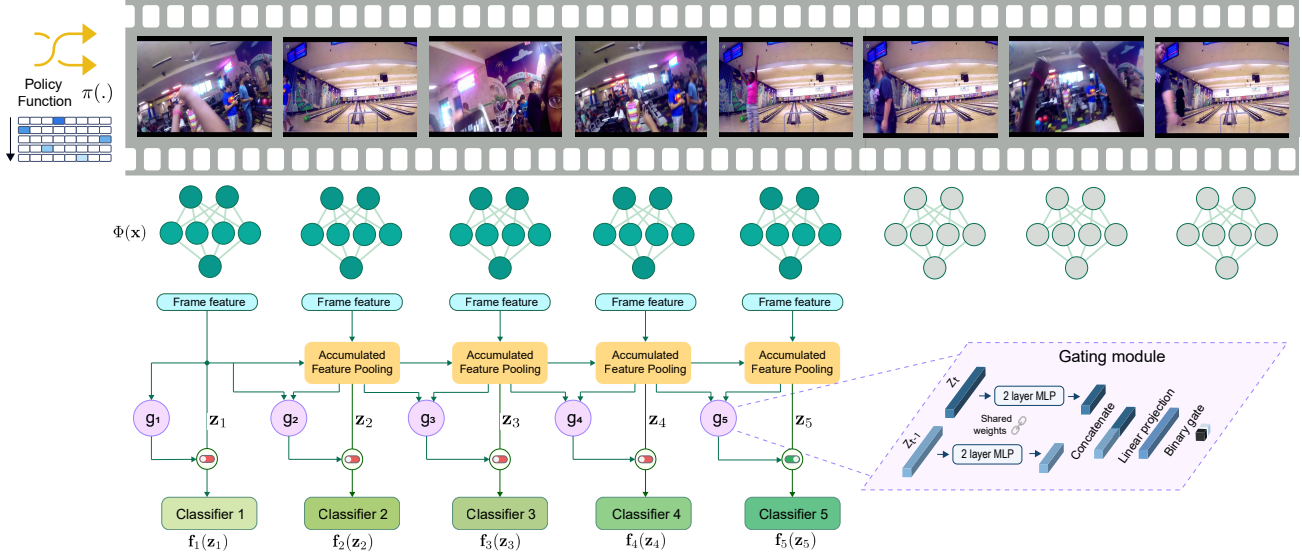


Figure 2: **The overview of FrameExit.** Given a video, at each time step t , we sample a frame from the video using the deterministic policy function π . Each frame is independently represented by the feature extraction network Φ and is aggregated to features of previous time steps using the accumulated feature pooling module (for $t > 1$). The gating modules (g_t) are trained to allow the network to automatically determine the earliest exiting point based on the inferred complexity of the input video. The architecture of the gating module is illustrated in the bottom right corner of the figure. Note that g_1 only receives z_1 as input. The video is adopted from [22].

ture for early exiting.

Our model draws inspiration from these works, but rather than having intermediate classifiers operating on intermediate network features, we focus on early exiting over the temporal dimension for video recognition.

3. FrameExit

Given a set of videos and their labels $\{\mathbf{v}_i, \mathbf{y}_i\}_{i=1}^D$, we aim to classify each video by processing the minimum number of frames. Figure 2 illustrates the overall architecture of our proposed model. Our model consists of *i*) a frame sampling policy π , *ii*) a feature extraction network Φ , *iii*) an accumulated feature pooling module, and *iv*) T classifiers \mathbf{f}_t and their associated exiting gates \mathbf{g}_t , where T is the number of input frames. Given an input video, we extract a partial clip $\mathbf{x}_{1:t}$ by sampling t frames, one at a time, from the video based on the sampling policy π :

$$\mathbf{x}_{1:t} = [\mathbf{x}_{1:t-1}; \mathbf{x}_t], \quad \mathbf{x}_t = \mathbf{v}_{\pi(t)}, \quad (1)$$

where $\mathbf{x}_{1:t-1}$ denotes a partial clip of length $t - 1$ and \mathbf{x}_t is a single video frame. We use the feature extraction network Φ to generate independent representation for each frame \mathbf{x}_t . These representations are then aggregated using the accumulated feature pooling module. The resulting clip level representation, \mathbf{z}_t , is then passed to the classifier \mathbf{f}_t and its associated early exiting gate \mathbf{g}_t .

Starting from a single-frame clip, we incrementally add more temporal details at each time step until one of the gates

predicts the halt signal. Each gate $\mathbf{g}_t : (\mathbf{z}_{t-1}, \mathbf{z}_t) \rightarrow \{0, 1\}$ is a binary function indicating whether the network has reached a desired confidence level to exit. A gate receives the aggregated features \mathbf{z}_t and \mathbf{z}_{t-1} as input. This allows the gate to make a more informed decisions by considering the agreement between temporal features. For example, a highly confident incorrect exiting prediction solely based on the current frame representation \mathbf{z}_t could potentially be mitigated if there is a significant disagreement with the representation \mathbf{z}_{t-1} . Note that the first gating module in the network only receives \mathbf{z}_1 as input. In the end, the final video label will be predicted as:

$$\mathbf{y} = \mathbf{f}_t(\mathbf{z}_t), \quad \text{where } t = \arg \min_t \{t \mid \mathbf{g}_t = 1\} \quad (2)$$

where t is the earliest frame that meets the gating condition. Note that if none of the gates predict the halt signal, the last classifier will classify the example. In the following sections, we describe our choice for frame sampling policy, accumulated feature pooling module, and our gates for early exiting.

Frame Sampling Policy. We sample T frames, one at a time, using the policy function $\pi(t)$. While in most existing works π is parameterized with a light-weight network and is trained using policy gradient methods [37] or Gumbel reparametrization [32], we use a simple, deterministic, and parameter-free function and show that it performs as well as sophisticated frame selection models. Our sampling

function π follows a coarse-to-fine principle for sampling in temporal dimension. It starts sampling from a coarse temporal scale and gradually samples more frames to add finer details to the temporal structure. Specifically, we sample the first frames from the middle, beginning, and end of the video, respectively, and then repeatedly sample frames from the halves (check Appendix for more details). Compared to sampling sequentially in time, this strategy allows the model to have access to a broader time horizon at each timestamp while mimicking the behaviour of RL-based approaches that jumps forward and backward to seek future informative frames and re-examine past information.

Feature extraction network. Our feature extraction network $\Phi(\mathbf{x}_i; \theta_\Phi)$ is a 2D image representation model, parametrized by θ_Φ , that extracts features for input frame \mathbf{x}_i . We use ResNet-50 [19], EfficientNet-b3 [38], and X3D-S [11] in our experiments.

Accumulated Feature Pooling. This step aims to create a representation for the entire clip $\mathbf{x}_{1:t}$. A major design criteria for our accumulated feature pooling module is efficiency. To limit the computation costs to only the newly sampled frame, the clip representation is incrementally updated. Specifically, given the sampled frame \mathbf{x}_t and features \mathbf{z}_{t-1} , we represent a video clip as:

$$\mathbf{z}_t = \Psi(\mathbf{z}_{t-1}, \Phi(\mathbf{x}_t; \theta_\Phi)), \quad (3)$$

where Ψ is a temporal aggregation function that can be implemented by statistical pooling methods such as average/max pooling, LSTM [20, 17], or self-attention [44].

Early Exiting. While processing the entire frames of a video is computationally expensive, processing a single frame may also restrict the network’s ability to recognize an action. Our conditional early exiting model has T classifiers accompanied by their associated early exiting gates that are attached at different time steps to allow early exiting. Each classifier \mathbf{f}_t receives the clip representation \mathbf{z}_t as input and makes a prediction about the label of the video. During training, we optimize the parameters of the feature extractor network and the classifiers using the following loss function:

$$\mathcal{L}_{cls} = \frac{1}{T} \sum_{t=0}^{t=T} \ell_{cls}(\mathbf{f}_t(\mathbf{z}_t; \theta_f), \mathbf{y}) \quad (4)$$

In our experiments, we use the standard cross-entropy loss for single-label video datasets and binary cross-entropy loss for the multi-label video datasets.

We parameterize each exiting gate \mathbf{g}_t as a multi-layer perceptron, predicting whether the partially observed clip $\mathbf{x}_{1:t}$ is sufficient to accurately classify the entire video. The gates have a very light design to avoid any significant computational overhead. As mentioned earlier, each gate $\mathbf{g}_t(\mathbf{z}_{t-1}, \mathbf{z}_t) \rightarrow \{0, 1\}$ receives as input the aggregated

representations \mathbf{z}_t and \mathbf{z}_{t-1} (See Figure 2 - bottom right). Each of these representations are first passed to two layers of MLP with 64 neurons independently (shared weights between the two streams). The resulting features are then concatenated and linearly projected and fed to a sigmoid function. Note that g_1 only receives \mathbf{z}_1 as input.

During training, gates may naturally learn to postpone exiting so that the last classifier always generates the model output because that may tend to maximize accuracy at the expense of additional processing. However, this sort of training largely defeats the purpose of the early exiting architecture. To overcome this problem we regularize the gates such that they are enforced to early exit. The parameters of the gates θ_g are learned in a self-supervised way by minimizing the binary cross-entropy between the predicted gating output and pseudo labels \mathbf{y}_t^g :

$$\mathcal{L}_{gate} = \frac{1}{T} \sum_{t=0}^{t=T} \text{BCE}(\mathbf{g}_t(\mathbf{z}_{t-1}, \mathbf{z}_t; \theta_g), \mathbf{y}_t^g), \quad (5)$$

We define the pseudo labels for the gate \mathbf{g}_t based on the classification loss:

$$\mathbf{y}_t^g = \begin{cases} 1 & \ell_{cls}(\mathbf{f}_t(\mathbf{z}_t), \mathbf{y}) \leq \epsilon_t \\ 0 & \text{else} \end{cases} \quad (6)$$

where ϵ_t determines the minimum loss required to exit through \mathbf{f}_t . A label 1 indicates the exiting signal while label 0 indicates that the network should proceed with processing further frames. Using this loss, we enforce the gate \mathbf{g}_t to exit when the classifier \mathbf{f}_t generates a sufficiently reliable prediction with a low loss. Provided that the early stage classifiers observe very limited number of frames, we only want to enable exiting when the classifier is highly confident about the prediction, i.e. \mathbf{f}_t has a very low loss. Hence, it is preferred to use smaller ϵ_t for these classifiers. On the other hand, late stage classifiers mostly deal with difficult videos with high loss. Therefore, as we proceed to later stage classifiers, we increase ϵ_t to enable early exiting. In our experiments, we define $\epsilon_t = \beta \exp(\frac{t}{2})$, where β is a hyper-parameter that controls the trade-off between model accuracy and total computation costs. The higher the β , the more computational saving we obtain.

The final objective for training our model is given as:

$$\mathcal{L} = \mathbb{E}_{(\mathbf{v}, \mathbf{y}) \sim D_{train}} [\mathcal{L}_{cls} + \mathcal{L}_{gate}] \quad (7)$$

Note that we use equal weights for the classification and gating loss terms in all of our experiments.

4. Experiments

We conduct extensive experiments to investigate the efficacy of our proposed method using three large-scale

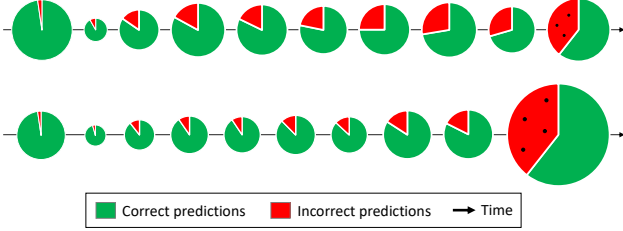


Figure 3: **Watermelon visualization of FrameExit predictions over time.** The exiting statistics of our method over time for $\beta = 1e^{-4}$ (top) and $\beta = 1e^{-6}$ (bottom) on the validation set of ActivityNet. The area of a circle at time step t represents the percentage of samples that exited through classifier t . Easier examples exit earlier from the network with higher accuracy while only hard examples reach to late stage classifiers leading to increased misclassifications. The illustration is inspired from [52].

datasets on two video understanding tasks, namely action recognition and holistic video understanding, as described in Section 4.1. Our experiments demonstrate that FrameExit outperforms the state of the art while significantly reducing computational costs as discussed in Section 4.2. Finally, we present an ablation study of our design choices in Section 4.3

4.1. Experimental setup

Datasets. We conduct experiments on three large-scale datasets: *ActivityNet-v1.3* [3], *Mini-Kinetics* [26], and *HVU* [7]. *ActivityNet-v1.3* [3] is a long-range action recognition dataset consisting of 200 classes with 10,024 videos for training and 4,926 videos for validation. The average length of the videos in this dataset is 167 seconds. *Mini-Kinetics* [26] is a short-range action recognition dataset, provided by [33], that contains 200 classes from Kinetics dataset [26] with 121,215 training and 9,867 testing videos. The average duration of videos is 10 seconds. *HVU* [7] is a large-scale, multi-label dataset for holistic video understanding with 3142 classes of actions, objects, scenes, attributes, events and concepts. It contains 476k training and 31k validation videos. The average duration of videos in this dataset is 8.5 seconds.

Evaluation metrics. Following the literature, we use mean average precision (mAP) and top-1 accuracy to evaluate accuracy for multi-class (*Mini-Kinetics*) and multi-label (*ActivityNet* and *HVU*) classification respectively. We measure the computational cost as giga floating point operations (GFLOPs). As a reference, *ResNet-50* and *EfficientNet-b3* have 4.12 and 1.8 GFLOPs respectively for a single input image of size 224×224 . Moreover, the *X3D-S* architecture has 1.96 GFLOPs for an input size of $13 \times 160 \times 160$. As different baseline methods use different number of frames for recognition, we report per video

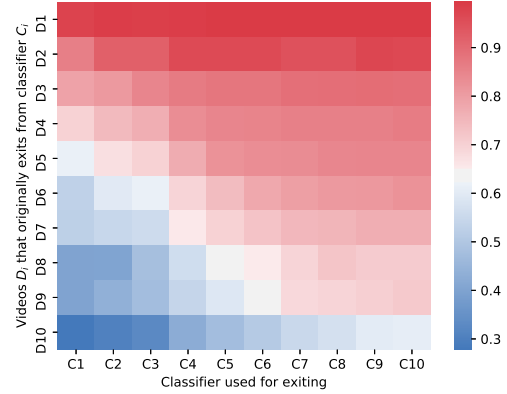


Figure 4: **Illustration of the exiting performance.** The matrix entry (i, j) illustrates the accuracy of the classifier C_j , if it were to classify the video data D_i that originally early exited by the model through C_i . Interestingly, if a video prematurely exits earlier than what was decided, we significantly lose accuracy while if it exits later than the model’s decision, no significant gain is obtained.

GFLOPs in all of our experiments.

Implementation details. We use *ResNet-50* [19] and *EfficientNet-b3* [38] hereafter referred to as *ResNet* and *EfficientNet*, as well as *X3D-S* [11] as our backbone networks. Both backbones are pretrained on *ImageNet*. We remove the last classification layers from the backbone networks, and replace it with a fully-connected layer with 4096 neurons. We set the number of input frames, T , to 10. We use max-pooling as the accumulated feature pooling function, Ψ in Eq. 3, for all the experiments unless otherwise stated. We train the model in two stages: we first train the backbone network and classifiers, then we learn the parameters of the gating modules. We use Adam optimizer [27] with a learning rate of $1e^{-4}$ for both training stages. The first training step runs for 35 epochs whilst dropping the learning rate after 16, and 30 epochs with a factor of 0.1. The second training step runs for 10 epochs and the learning rate is dropped after epochs 5 and 8. We set the hyper-parameter β ranging from $1e^{-6}$ to $1e^{-2}$ to generate varying trade-off points between accuracy and computational costs.

4.2. Results

We first analyze the behaviour of our conditional early exiting method. Then we compare FrameExit with the state of the art in action recognition and holistic video understanding. This section is concluded by reporting several qualitative results.

Conditional early exiting. We analyze the effectiveness of FrameExit in adjusting the amount of computations per video based on its content. Figure 3 illustrates predictions and the number of frames processed by two Frame-

Exit models trained with different values of $\beta = 1e^{-4}$ and $\beta = 1e^{-6}$. As expected, a higher β encourages more video samples to exit early from the network (top row) compared to a lower β (bottom row). A general observation is that the more we proceed to classifiers in later stages, the more inaccurate the predictions become. This may sound counter-intuitive, because if we were to have a model without early exiting, late-stage classifiers produce the most accurate predictions. However, the trend shown in Figure 3 is highly desirable, because it shows that the easier examples have already exited from the network while only hard examples reach to late stage classifiers.

This observation is more clear in Figure 4. The matrix entry (i, j) illustrates the accuracy of the classifier C_j , if it were to classify the video data D_i that originally early exited by the model through C_i . The diagonal entries $(i, j = i)$ in the matrix represent the actual early-exiting results obtained by the model. In an ideal early exiting scenario, for each row i , it is desired to have a significantly lower performance for C_j compared to C_i when $j < i$ and a similar accuracy to C_i when $j > i$. A similar pattern is observed in Figure 4. As can be seen, if a video prematurely exits earlier than what was decided by the model, we lose accuracy significantly while if it exits later than the model’s decision, no significant gain in accuracy is obtained. FrameExit learns to exit from the optimal classifier to balance computational costs and model accuracy.

Figure 5 illustrates average precision of all categories in the ActivityNet dataset as well as the average number of frames required for each category to confidently exit from the network. As can be seen, the classes with strong visual cues (such as riding bumper car, playing accordion, playing pool, etc.) appear to be simpler and therefore exit in the early stages with a high accuracy. In contrast, actions which involve more complex interactions and require a series of frames to unveil (such as Javelin throw and trimming branches) take more frames to be recognized.

Comparison to state of the art: Action recognition.

We compare our method with an extensive list of recent works on efficient video recognition: AR-Net [33], AdaFrame [52], LiteEval [51], SCSampler [29], MARL [49], and ListenToLook [15]. AR-Net uses MobileNet-V2 as the sampler network and adaptively chooses a ResNet architecture with varying depths as the recognition network. The method is additionally evaluated with variants of EfficientNet [38] as the recognition network. AdaFrame and LiteEval both use MobileNet-V2 as the policy network and ResNet-101 as the recognition network. SCSampler uses MobileNet-V2 as the sampler network and ResNet-50 as the recognition network. MARL uses a ResNet-101 as the recognition network combined with a fully connected layer as a policy network. ListenToLook uses MobileNet-V2 as the sampler network and

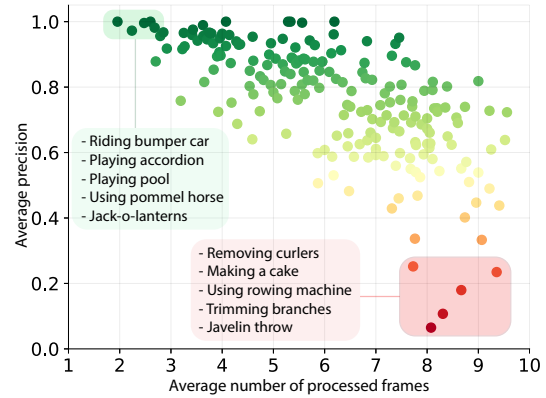


Figure 5: **AP vs number of processed frames per category in the ActivityNet dataset.** Categories with object/scene cues exit in the early stages with high accuracy while complex actions require more frames for recognition.

Table 1: **Comparison to state of the art on action recognition.** FrameExit outperforms competing methods in terms of accuracy and efficiency using ResNet, EfficientNet, and X3D-S backbones. Results of other methods are adopted from [33]. *(IA|IA) denotes the variant of the ListenToLook method that additionally uses audio for sampling and recognition.

	ActivityNet		Mini-kinetics	
	mAP (%)	GFLOPs	Top-1 (%)	GFLOPs
<i>ResNet</i>				
AdaFrame [52]	71.5	79.0	-	-
LiteEval [51]	72.7	95.1	61.0	99.0
ListenToLook [15]	72.3	81.4	-	-
ListenToLook (IA IA)* [15]	75.6	37.5	-	-
SCSampler [29]	72.9	41.9	70.8	41.9
AR-Net [33]	73.8	33.5	71.7	32.0
Ours (w/o exit)	77.3	41.2	73.3	41.2
Ours (FrameExit)	76.1	26.1	72.8	19.7
<i>EfficientNet</i>				
AR-Net [33]	79.7	15.3	74.8	16.3
Ours (w/o exit)	81.1	18.0	75.9	18.0
Ours (FrameExit)	80.0	11.4	75.3	7.8
<i>X3D-S</i>				
Ours (w/o exit)	87.4	19.6	-	-
Ours (FrameExit)	86.0	9.8	-	-

ResNet-101 as the recognition network. Finally, ListenToLook (IA|IA) uses two ResNet-18 for audio and visual modality respectively as the sampler. The same architecture is used for recognition network.

While the competing methods use two networks for sampling and recognizing, FrameExit uses a single network for efficient video recognition. The results for ActivityNet and Mini-Kinetics are shown in Table 1. Our method with a ResNet backbone outperforms all other approaches by obtaining an improved accuracy while using $1.3\times - 5\times$ less GFLOPs. The gain in accuracy is mainly attributed to our accumulated feature pooling module, while the gain in efficiency is attributed to the proposed sampling policy

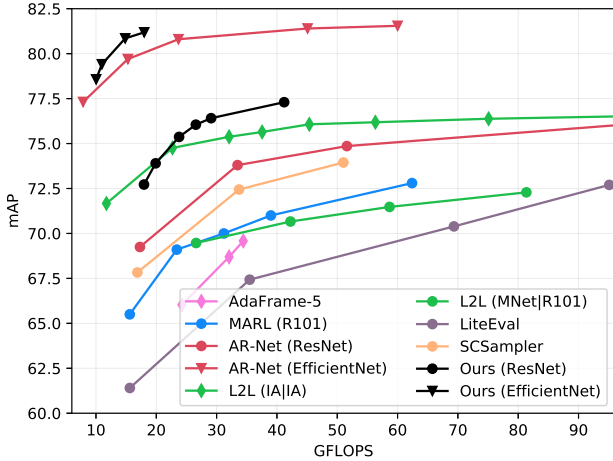


Figure 6: **Accuracy vs. efficiency curves on ActivityNet.** FrameExit performs similar to or better than other methods at a much lower computational cost. Note that L2L (IA|IA) denotes the variant of the ListenToLook method that additionally uses audio for sampling and recognition.

and gating modules for conditional early exiting (see Section 4.3 for detailed analyses). Compared to our model without early exiting, denoted as “Ours (w/o exit)”, FrameExit achieves a comparable accuracy with a significant reduction in GFLOPs (see Appendix for wall-clock timing of FrameExit). In addition, Figure 6 presents accuracy vs. computations trade-off curves on ActivityNet for various methods. Note that except ListenToLook (IA—IA) [15] that uses audio and visual features for sampling and recognition, other methods rely solely on visual features. As shown, FrameExit achieves the same top-performance as other methods at a much lower computational cost.

To verify that the performance of our model is not limited to one architecture, we conduct similar experiments with EfficientNet and X3D backbones. Using EfficientNet-b3, FrameExit obtains 3.9% further absolute gain in mAP on ActivityNet and 2.5% on Mini-Kinetics, while consuming $2.3\times$ and $2.5\times$ less compute. In particular, on ActivityNet and mini-kinetics, we outperform AR-Net [33], the leading method among competitors, with $1.3\times$ and $2.1\times$ less GFLOPs, respectively. When using the highly efficient X3D-S as our backbone, at a similar computational cost, FrameExit achieves 86.0% which is 6.0% higher than the best 2D model. This demonstrates the superiority of our method for efficient video understanding using both 2D and 3D backbones.

Comparison to state of the art: Holistic understanding.

We evaluate our method for the task of holistic video understanding, where the goal is to recognize various semantic aspects including objects, scenes, and actions within a video. To this end, we conduct experiments on the

Table 2: **Comparison to state of the art on holistic video understand.** FrameExit significantly improves mAP on HVU dataset while saving compute. *indicates GFLOPs for one input clip.

	mAP (%)	GFLOPs
<i>2D/3D ResNet</i>		
Uniform-10	44.7	41.2
Random-10	43.6	41.2
3D-ResNet18 [41]	35.4	38.6*
HATNet [7]	39.6	41.8*
FrameExit ($\beta = 1e^{-3}$)	45.7	8.6
FrameExit ($\beta = 1e^{-5}$)	49.2	18.7
<i>EfficientNet</i>		
FrameExit ($\beta = 1e^{-3}$)	47.7	11.7
FrameExit ($\beta = 1e^{-2}$)	46.1	5.7

large-scale, multi-label HVU dataset [7]. we use average pooling in our accumulated feature pooling module, Ψ in Eq. 3, as it outperforms max pooling. We compare our method with the commonly-used uniform and random sampling baselines. These baselines sample K frames (uniformly/randomly) and average frame-level predictions to generate a final video-level prediction. We additionally, compare our method with 3D-ResNet18 [41] and mixed 2D/3D ResNet [7] models. As shown in Table 2, our method consistently outperforms other methods in terms of accuracy while saving significant computation. In particular, we report two variants of FrameExit with the ResNet backbone by changing accuracy-efficiency trade-off parameter β . FrameExit trained with $\beta = 1e^{-3}$ only uses 8.6 GFLOPs to obtain 45.7% mAP, which is highly suitable for low-budget requirements. FrameExit trained with $\beta = 1e^{-5}$ uses 18.7 GFLOPs on average to obtain 49.2% mAP, which is suitable for high-accuracy requirements. Interestingly, FrameExit outperforms 3D models [41, 7] in this task. Given that HVU requires recognizing static 2D concepts such as scenes (8% of the class labels) and objects (55% of the class labels), 3D-CNN methods are not necessarily the optimal choice for this task. The most budget-friendly trained model is FrameExit with the EfficientNet backbone that uses only 5.7 GFLOPs and achieves a mAP of 46.1.

This experiment implies that our method can be used in a wider range of efficient video recognition tasks. Our results set a new state of the art on the HVU dataset for future research.

4.3. Ablation study

In this section we inspect different aspects of FrameExit. For all ablations, we follow the same training procedure explained in Section 4.1 and use the Activitynet-v1.3 dataset.

Impact of policy function. To validate the impact of our proposed deterministic policy function for frame sampling,

Table 3: Impact of sampling policy. we report the results over 5 runs.

	Train-random	Test-random	mAP (%)	GFLOPs
(1)	✗	✗	73.7	27.6
(2)	✓	✗	74.3 ± 0.9	26.7 ± 0.33
(3)	✗	✓	74.4 ± 0.02	27.0 ± 0.19
(4)	✓	✓	74.4 ± 0.34	25.0 ± 0.37
FrameExit	-	-	76.1	26.1

we make comparisons with two commonly used sampling approaches namely “Sequential” and “Random”, both during training and evaluation. For “Sequential”, we keep original frame order while in “Random”, we randomly sample frames. Results for all combinations are shown in Table 3. Using the original frame ordering during training and testing results in a lower mAP and a higher GFLOPs. We speculate that during training, earlier classifiers receive less informative gradients due to not observing a holistic view of the entire video (row 1 vs. row 2). Similarly, the model consumes higher GFLOPs because it needs to observe more frames to infer a confident prediction (row 1 vs. row 3). Random sampling both during training and testing improves the results because early stage classifiers have a higher chance of obtaining more useful information reflecting the entire video time-span. The best result, however, is obtained by our simple deterministic sampling policy. This is because jumping forward and backward enabled by our frame sampling policy effectively improves the chance of picking up informative frames for recognition.

Impact of accumulated feature pooling. We evaluate the impact of accumulated feature pooling module by comparing the performance of FrameExit with and without feature pooling as reported in Table 4. For the FrameExit without feature pooling, we train each classifier and its associated gating module only on the currently sampled frame. We then average frame-level predictions made by all classifiers up until the exiting classifier to obtain the video-level prediction. Therefore, the major difference in the design of these two settings relates to the use of the pooling operation in feature space or in the prediction space. The results in Table 4 demonstrates that pooling over features is much more effective than over output predictions.

Number of input frames. We train FrameExit with various number of input frames $T = \{4, 8, 10, 16\}$. During inference, the gates decide when to stop processing. As shown in Table 5, the performance of the model increases

Table 4: Impact of accumulated feature pooling.

Feature pooling	mAP (%)	GFLOPs
✓	76.1	26.1
✗	67.8	27.6

Table 5: Impact of number of frames on FrameExit.

	$T = 4$	$T = 8$	$T = 10$	$T = 16$
mAP (%)	66.2	74.3	76.1	76.1
GFLOPs	13.0	22.4	26.1	35.1

Table 6: Impact of adaptive exiting. *for each column, results are reported over an average no. of processed frames.

No. processed frames	3	4	6	7
Fixed budget	60.8	67.5	74.0	75.2
FrameExit*	67.5	70.1	75.8	76.4

as the number of input frames increases, but up to a certain limit. This is a reasonable observation, as certain actions/videos may require more frames to be recognized. However, the reason why increasing the number of frames after a limit does not further improve the performance could be mainly attributed to the limited capacity of 2D convolutional networks in leveraging temporal information. As a result, to keep the balance between accuracy and efficiency, we set $T = 10$ in our experiments.

Adaptive vs fixed exiting. To show the merits of adaptive early exiting versus a fixed budget exiting, we conduct an ablation in two settings. In the first setting, we use our conditional early exiting model and in the second setting, we assume a fixed number of frames is processed for each test video. Table 6 shows that our conditional early exiting method consistently outperforms fixed exiting.

5. Conclusions

In this paper, we presented FrameExit, a conditional early exiting method for efficient video recognition. Our proposed method uses gating modules, trained to allow the network to automatically determine the earliest exiting point based on the inferred complexity of the input video. To enable gates to make reliable decisions we use an effective video representation obtained using accumulated feature pooling. We showed that our early exiting mechanism combined with a simple, deterministic sampling strategy obviates the need for complex sampling policy techniques. Our proposed method is model-agnostic and can be used with various network architectures. Comprehensive experiments show the power of our method in balancing accuracy versus efficiency in action recognition as well as holistic video understanding tasks.

Acknowledgements We thank Fatih Porikli, Michael Hofmann, Haitam Ben Yahia, Arash Behboodi and Iliia Karmenov for their feedback and discussions.

References

- [1] Davide Abati, Jakub Tomczak, Tijmen Blankevoort, Simone Calderara, Rita Cucchiara, and Babak Ehteshami Bejnordi. Conditional channel gated networks for task-aware continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3931–3940, 2020. [2](#)
- [2] Babak Ehteshami Bejnordi, Tijmen Blankevoort, and Max Welling. Batch-shaping for learning conditional channel gated networks. In *International Conference on Learning Representations*, 2019. [2](#)
- [3] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 961–970, 2015. [5](#)
- [4] Víctor Campos, Brendan Jou, Xavier Giró-i Nieto, Jordi Torres, and Shih-Fu Chang. Skip rnn: Learning to skip state updates in recurrent neural networks. *arXiv preprint arXiv:1708.06834*, 2017. [2](#)
- [5] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*, 2017. [1](#)
- [6] Zhouong Chen, Yang Li, Samy Bengio, and Si Si. Gaternet: Dynamic filter selection in convolutional neural network via a dedicated global gating network. *arXiv preprint arXiv:1811.11205*, 2018. [2](#)
- [7] Ali Diba, Mohsen Fayyaz, Vivek Sharma, Manohar Paluri, Jürgen Gall, Rainer Stiefelhagen, and Luc Van Gool. Large scale holistic video understanding. In *European Conference on Computer Vision*, pages 593–610. Springer, 2020. [2](#), [5](#), [7](#)
- [8] Hehe Fan, Zhongwen Xu, Linchao Zhu, Chenggang Yan, Jianjun Ge, and Yi Yang. Watching a small portion could be as good as watching all: Towards efficient video classification. In *IJCAI International Joint Conference on Artificial Intelligence*, 2018. [1](#), [2](#)
- [9] Linxi Fan, Shyamal Buch, Guanzhi Wang, Ryan Cao, Yuke Zhu, Juan Carlos Niebles, and Li Fei-Fei. Rubiksnet: Learnable 3d-shift for efficient video action recognition. In *European Conference on Computer Vision*, pages 505–521. Springer, 2020. [2](#)
- [10] Quanfu Fan, Chun-Fu Chen, Hilde Kuehne, Marco Pistoia, and David Cox. More is less: Learning efficient video representations by big-little network and depthwise temporal aggregation. *arXiv preprint arXiv:1912.00869*, 2019. [1](#)
- [11] Christoph Feichtenhofer. X3d: Expanding architectures for efficient video recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 203–213, 2020. [1](#), [2](#), [4](#), [5](#)
- [12] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 6202–6211, 2019. [1](#)
- [13] Basura Fernando, Efstratios Gavves, José Oramas, Amir Ghodrati, and Tinne Tuytelaars. Rank pooling for action recognition. *IEEE transactions on pattern analysis and machine intelligence*, 39(4):773–787, 2016. [1](#)
- [14] Michael Figurnov, Maxwell D Collins, Yukun Zhu, Li Zhang, Jonathan Huang, Dmitry Vetrov, and Ruslan Salakhutdinov. Spatially adaptive computation time for residual networks. In *CVPR*, 2017. [2](#)
- [15] Ruohan Gao, Tae-Hyun Oh, Kristen Grauman, and Lorenzo Torresani. Listen to look: Action recognition by previewing audio. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10457–10467, 2020. [1](#), [2](#), [6](#), [7](#)
- [16] Xitong Gao, Yiren Zhao, Lukasz Dudziak, Robert Mullins, and Cheng-Zhong Xu. Dynamic channel pruning: Feature boosting and suppression. *arXiv preprint arXiv:1810.05331*, 2018. [2](#)
- [17] Amir Ghodrati, Efstratios Gavves, and Cees GM Snoek. Video time: Properties, encoders and evaluation. *arXiv preprint arXiv:1807.06980*, 2018. [4](#)
- [18] Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. Learning spatio-temporal features with 3d residual networks for action recognition. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 3154–3160, 2017. [1](#)
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016. [4](#), [5](#)
- [20] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997. [4](#)
- [21] [Sophie and Greg playing acro yoga](#) by [Gregology](#) is licensed under [CC BY](#). [1](#)
- [22] [SPORK! Exclusive: Pathways Waveland Bowling](#) by [SPORK! NFP](#) is licensed under [CC BY](#). [1](#), [3](#)
- [23] [Canada's National Ballet School: Training in the Professional Ballet Program](#) by [Canada's National Ballet School](#) is licensed under [CC BY](#). [1](#)
- [24] Gao Huang, Danlu Chen, Tianhong Li, Felix Wu, Laurens van der Maaten, and Kilian Weinberger. Multi-scale dense networks for resource efficient image classification. In *International Conference on Learning Representations*, 2018. [1](#), [2](#)
- [25] Noureldien Hussein, Mihir Jain, and Babak Ehteshami Bejnordi. Timegate: Conditional gating of segments in long-range activities. *arXiv preprint arXiv:2004.01808*, 2020. [2](#)
- [26] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. [5](#)
- [27] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2014. [5](#)
- [28] Okan Köpüklü, Neslihan Kose, Ahmet Gunduz, and Gerhard Rigoll. Resource efficient 3d convolutional neural networks. In *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*, pages 1910–1919. IEEE, 2019. [1](#), [2](#)
- [29] Bruno Korbar, Du Tran, and Lorenzo Torresani. Scsampler: Sampling salient clips from video for efficient action recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6232–6242, 2019. [1](#), [2](#), [6](#)

- [30] Chen-Yu Lee, Saining Xie, Patrick Gallagher, Zhengyou Zhang, and Zhuowen Tu. Deeply-supervised nets. In *Artificial intelligence and statistics*, pages 562–570, 2015. 1, 2
- [31] Ji Lin, Chuang Gan, and Song Han. Tsm: Temporal shift module for efficient video understanding. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7083–7093, 2019. 1, 2
- [32] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. In *ICLR*, 2017. 3
- [33] Yue Meng, Chung-Ching Lin, Rameswar Panda, Prasanna Sattigeri, Leonid Karlinsky, Aude Oliva, Kate Saenko, and Rogerio Feris. Ar-net: Adaptive frame resolution for efficient action recognition. *arXiv preprint arXiv:2007.15796*, 2020. 1, 2, 5, 6, 7
- [34] Yue Meng, Rameswar Panda, Chung-Ching Lin, Prasanna Sattigeri, Leonid Karlinsky, Kate Saenko, Aude Oliva, and Rogerio Feris. Adafuse: Adaptive temporal fusion network for efficient action recognition. In *International Conference on Learning Representations*, 2021. 2
- [35] AJ Piergiovanni, Anelia Angelova, and Michael S Ryoo. Tiny video networks. *arXiv preprint arXiv:1910.06961*, 2019. 1, 2
- [36] Zhaofan Qiu, Ting Yao, and Tao Mei. Learning spatio-temporal representation with pseudo-3d residual networks. In *proceedings of the IEEE International Conference on Computer Vision*, pages 5533–5541, 2017. 1, 2
- [37] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018. 3
- [38] Mingxing Tan and Quoc V Le. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv preprint arXiv:1905.11946*, 2019. 4, 5, 6
- [39] S. Teerapittayanon, B. McDanel, and H. T. Kung. Branchynet: Fast inference via early exiting from deep neural networks. In *2016 23rd International Conference on Pattern Recognition (ICPR)*, pages 2464–2469, Dec 2016. 1, 2
- [40] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 4489–4497, 2015. 1
- [41] Du Tran, Jamie Ray, Zheng Shou, Shih-Fu Chang, and Manohar Paluri. Convnet architecture search for spatiotemporal feature learning. *arXiv preprint arXiv:1708.05038*, 2017. 7
- [42] Du Tran, Heng Wang, Lorenzo Torresani, and Matt Feiszli. Video classification with channel-separated convolutional networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5552–5561, 2019. 1, 2
- [43] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 6450–6459, 2018. 1, 2
- [44] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 4
- [45] Andreas Veit and Serge Belongie. Convolutional networks with adaptive inference graphs. In *European Conference on Computer Vision (ECCV)*, 2018. 2
- [46] Thomas Verelst and Tinne Tuytelaars. Dynamic convolutions: Exploiting spatial sparsity for faster inference. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2320–2329, 2020. 2
- [47] Longguang Wang, Xiaoyu Dong, Yingqian Wang, Xinyi Ying, Zaiping Lin, Wei An, and Yulan Guo. Learning sparse masks for efficient image super-resolution. *arXiv preprint arXiv:2006.09603*, 2020. 2
- [48] Xin Wang, Fisher Yu, Zi-Yi Dou, Trevor Darrell, and Joseph E. Gonzalez. Skipnet: Learning dynamic routing in convolutional networks. In *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part XIII*, pages 420–436, 2018. 2
- [49] Wenhao Wu, Dongliang He, Xiao Tan, Shifeng Chen, and Shilei Wen. Multi-agent reinforcement learning based frame sampling for effective untrimmed video recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6222–6231, 2019. 1, 2, 6
- [50] Zuxuan Wu, Tushar Nagarajan, Abhishek Kumar, Steven Rennie, Larry S Davis, Kristen Grauman, and Rogerio Feris. Blockdrop: Dynamic inference paths in residual networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8817–8826, 2018. 2
- [51] Zuxuan Wu, Caiming Xiong, Yu-Gang Jiang, and Larry S Davis. Liteeval: A coarse-to-fine framework for resource efficient video recognition. In *Advances in Neural Information Processing Systems*, pages 7780–7789, 2019. 1, 2, 6
- [52] Zuxuan Wu, Caiming Xiong, Chih-Yao Ma, Richard Socher, and Larry S Davis. Adaframe: Adaptive frame selection for fast video recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1278–1287, 2019. 1, 2, 5, 6
- [53] Saining Xie, Chen Sun, Jonathan Huang, Zhuowen Tu, and Kevin Murphy. Rethinking spatiotemporal feature learning for video understanding. *arXiv preprint arXiv:1712.04851*, 1(2):5, 2017. 1, 2
- [54] Le Yang, Yizeng Han, Xi Chen, Shiji Song, Jifeng Dai, and Gao Huang. Resolution adaptive networks for efficient inference. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2369–2378, 2020. 1, 2
- [55] Serena Yeung, Olga Russakovsky, Greg Mori, and Li Fei-Fei. End-to-end learning of action detection from frame glimpses in videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2678–2687, 2016. 1, 2
- [56] Joe Yue-Hei Ng, Matthew Hausknecht, Sudheendra Vijayanarasimhan, Oriol Vinyals, Rajat Monga, and George Toderici. Beyond short snippets: Deep networks for video classification. In *Proceedings of the IEEE conference on*

computer vision and pattern recognition, pages 4694–4702, 2015. [1](#)

- [57] Yin-Dong Zheng, Zhaoyang Liu, Tong Lu, and Limin Wang. Dynamic sampling networks for efficient action recognition in videos. *IEEE Transactions on Image Processing*, 29:7970–7983, 2020. [1](#), [2](#)
- [58] Bolei Zhou, Alex Andonian, Aude Oliva, and Antonio Torralba. Temporal relational reasoning in videos. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 803–818, 2018. [1](#)
- [59] Mohammadreza Zolfaghari, Kamaljeet Singh, and Thomas Brox. Eco: Efficient convolutional network for online video understanding. In *Proceedings of the European conference on computer vision (ECCV)*, pages 695–712, 2018. [1](#)