

Adaptive Focus for Efficient Video Recognition

Yulin Wang*, Zhaoxi Chen*, Haojun Jiang, Shiji Song, Yizeng Han, Gao Huang[†]

Department of Automation, BNRist, Tsinghua University, Beijing, China

{wang-y119, jhj20, hanyz18}@mails.tsinghua.edu.cn, frozen.burning@gmail.com,

{shijis, gaohuang}@tsinghua.edu.cn

Abstract

In this paper, we explore the spatial redundancy in video recognition with the aim to improve the computational efficiency. It is observed that the most informative region in each frame of a video is usually a small image patch, which shifts smoothly across frames. Therefore, we model the patch localization problem as a sequential decision task, and propose a reinforcement learning based approach for efficient spatially adaptive video recognition (AdaFocus). In specific, a light-weighted ConvNet is first adopted to quickly process the full video sequence, whose features are used by a recurrent policy network to localize the most task-relevant regions. Then the selected patches are inferred by a high-capacity network for the final prediction. During offline inference, once the informative patch sequence has been generated, the bulk of computation can be done in parallel, and is efficient on modern GPU devices. In addition, we demonstrate that the proposed method can be easily extended by further considering the temporal redundancy, e.g., dynamically skipping less valuable frames. Extensive experiments on five benchmark datasets, i.e., ActivityNet, FCVID, Mini-Kinetics, Something-Something V1&V2, demonstrate that our method is significantly more efficient than the competitive baselines. Code is available at <https://github.com/blackfeather-wang/AdaFocus>.

1. Introduction

The explosive growth of online videos (e.g., on YouTube or TikTok) has fueled the demands for automatically recognizing human actions, events, or other contents within them, which benefits applications like recommendation [7, 8, 15], surveillance [6, 3] and content-based searching [25]. In the past few years, remarkable success in accurate video recognition has been achieved by leveraging deep networks [11, 59, 12, 2, 40, 19]. However, the impressive perfor-

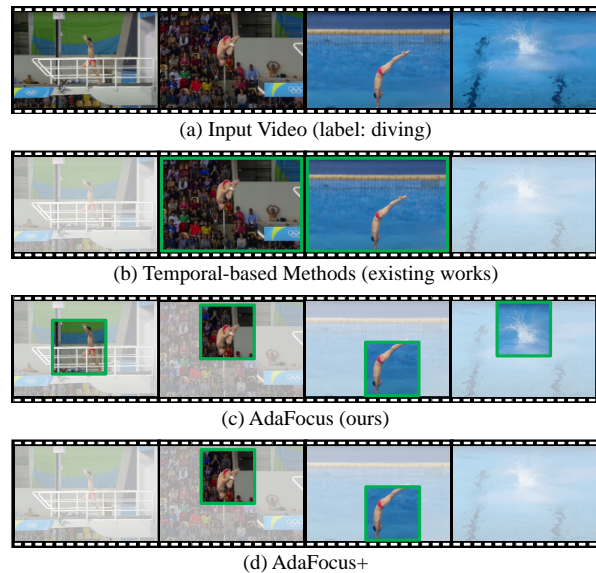


Figure 1. **Comparisons of temporal-based methods and our proposed AdaFocus approach.** Most existing works reduce computational costs by selecting a few informative frames to process, while AdaFocus aims to perform efficient inference by attending to the task-relevant patch of each frame. Importantly, our method is compatible with temporal-based techniques as it can be improved by skipping less important frames (AdaFocus+).

mance of these models usually comes at high computational costs. In real-world scenarios, computation directly translates into power consumption, carbon emission and practical latency, which should be minimized under both economic and safety considerations.

To address this issue, a number of recent works propose to reduce the inherent *temporal redundancy* in video recognition [30, 33, 51, 48, 16, 50, 29]. As shown in Figure 1 (b), it is efficient to focus on the most task-relevant video frames, and allocate the majority of computation to them rather than all frames. However, another important source of redundant computation in image-based data, namely *spatial redundancy*, has rarely been explored in the context of efficient video recognition. In fact, it has been shown in 2D-image classification that convolutional

*Equal contribution.

[†]Corresponding author.

networks (CNNs) are able to produce correct predictions with only a few discriminative regions of the whole image [45, 53, 18, 35, 14, 5]. By performing inference on these relatively small regions, one can dramatically reduce the computational cost of CNNs (e.g., processing a 96x96 patch requires $\sim 18\%$ computation of inferring a 224x224 image).

In this paper, we are interested in whether this *spatial redundancy* can be effectively leveraged to facilitate efficient video recognition. We develop a novel adaptive focus (AdaFocus) approach to dynamically localize and attend to the task-relevant regions of each frame. In specific, our method first takes a quick glance at each frame with a light-weighted CNN to obtain cheap and coarse global information. Then we train a recurrent policy network on its basis to select the most valuable region for recognition. This procedure leverages the reinforcement learning algorithm due to the non-differentiability of localizing task-relevant regions. Finally, we activate a high-capacity deep CNN to process only the selected regions. Since the proposed regions are usually small patches with a reduced size, considerable computational costs can be saved. An illustration of AdaFocus is shown in Figure 1 (c). Our method allocates computation unevenly across the spatial dimension of video frames according to the contributions to the recognition task, leading to a significant improvement in efficiency with preserved accuracy.

The vanilla AdaFocus framework does not model temporal redundancy, i.e., all frames are processed with identical computation, while the only difference lies in the locations of the selected regions. We show that our method is compatible with existing temporal-based techniques, and can be extended via reducing the computation spent on uninformative frames, as presented in Figure 1 (d). This is achieved by introducing an additional policy network that determines whether to skip some less valuable frames. This algorithm is referred to as AdaFocus+.

We evaluate the effectiveness of AdaFocus on five video recognition benchmarks (i.e., ActivityNet, FCVID, MiniKinetics, Something-Something V1&V2). Experimental results show that AdaFocus by itself consistently outperforms all the baselines by large margins, while AdaFocus+ further improves the efficiency. For instance, AdaFocus+ has 2-3x less FLOPs¹ than the recently proposed AR-Net [33] when achieving the same accuracy. We also demonstrate that our method can be deployed on top of the state-of-the-art networks (e.g., TSM [31]) and effectively improve their computational efficiency.

2. Related Works

Video recognition. Significant progress has been made in video recognition with the adoption of convolutional neu-

ral networks (CNNs). One prevalent approach is constructing 3D-CNNs to model the temporal and spatial information jointly, e.g., C3D [40], I3D [2] and ResNet3D [19]. The other line of works first extracts frame-level features, and aggregates the features at different temporal locations via temporal averaging [43], long short-term memory (LSTM) networks [9], channel shifting [31], etc.

Despite the success achieved by the aforementioned works, the expensive computational cost of CNNs, especially 3D-CNNs, usually limits their applicability. Recent research efforts have been made towards improving the efficiency of video recognition, via designing light-weighted architectures [42, 52, 36, 41, 60, 31] or perform dynamic computation on a per-video basis [55, 51, 29, 58, 30, 34]. Our approach shares a similar idea as the latter on reducing the intrinsic redundancy in video data, while with a special focus on spatial redundancy.

Reducing temporal redundancy is a popular solution to efficient video recognition, based on the intuition that not all frames contribute equally to the final prediction. In particular, a model could dynamically allocate no/little computation to some less informative or highly correlated frames [18]. This idea has been proven effective by many implementations, including (1) early stopping, i.e., terminating the computation before “watching” the full sequence [10, 49]; (2) conditional computing, e.g., LiteEval [51] adaptively selects an LSTM model with appropriate size at each time step in a recurrent recognition procedure; adaptive resolution network (AR-Net) [33] processes different frames with adaptive resolutions to save unnecessary computation on less important frames; and (3) frame/clip sampling, i.e. dynamically deciding which frames should be skipped without performing any computation [48, 16, 29, 50]. The proposed AdaFocus method differentiates itself from these approaches in that we focus on reducing the *spatial redundancy*, namely allocating the major computation to task-relevant regions of the frames. In addition, our method is compatible with them as it can be improved by further reducing temporal redundancy.

Reducing spatial redundancy. It has been observed that considerable spatial redundancy exists in the process of extracting deep features from image-based data [18, 54]. For example, in 2D-image classification, a number of recent works successfully improve the efficiency of CNNs by attending to some task-relevant or more informative parts of images [53, 45, 13]. In the context of video recognition, existing attention-based methods have demonstrated that different image regions of video frames do not contribute equivalently to the prediction [32]. However, to our best knowledge, the *spatial redundancy* has not been exploited to improve the efficiency of video recognition.

¹In this paper, FLOPs refers to the number of multiply-add operations.

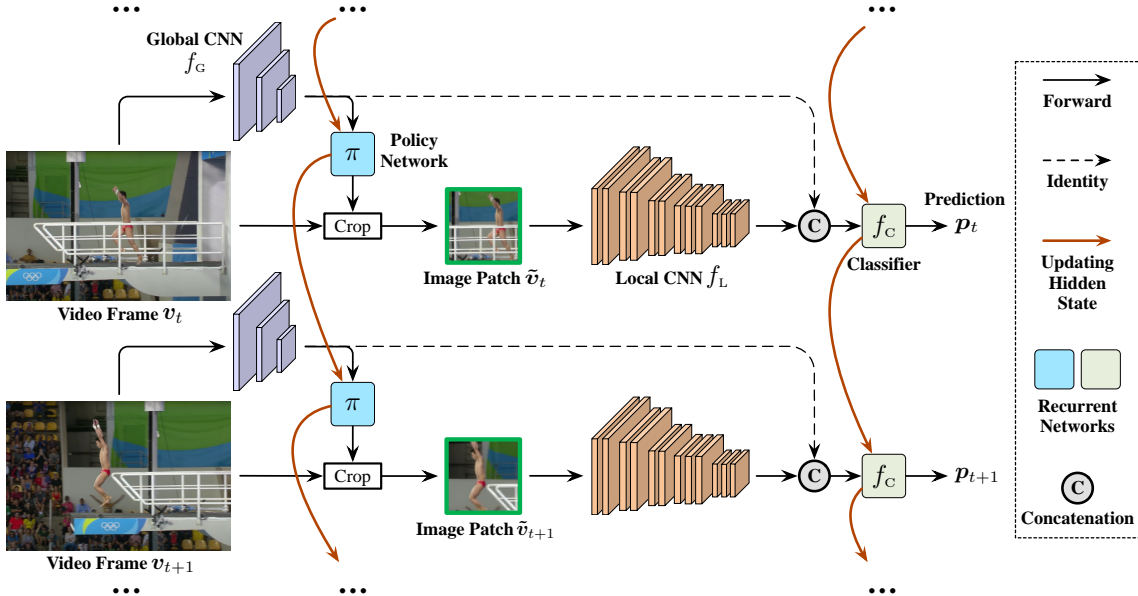


Figure 2. **Overview of AdaFocus.** It first takes a quick glance at each frame v_t using a light-weighted global CNN f_G . Then a recurrent policy network π is built on top of f_G to select the most important image region \tilde{v}_t in terms of recognition. A high-capacity local CNN f_L is adopted to extract features from \tilde{v}_t . Finally, a recurrent classifier aggregates the features across frames to obtain the prediction p_t .

3. Method

Different from most existing works that facilitate efficient video recognition by leveraging the *temporal redundancy*, we seek to save the computation spent on the task-irrelevant regions of video frames, and thus improve the efficiency by reducing the *spatial redundancy*. To this end, we propose an adaptive focus (AdaFocus) framework to adaptively identify and attend to the most informative regions of each frame, such that the computational cost can be significantly reduced without sacrificing accuracy.

In this section, we first describe its components and the correspond training algorithm in Section 3.1 and Section 3.2, respectively. Then we show in Section 3.3 that AdaFocus can be improved by further considering temporal redundancy (e.g., skipping uninformative frames).

3.1. Network Architecture

Overview. We first give an overview of AdaFocus (Figure 2). Consider the online video recognition scenario, where a stream of frames come in sequentially while a prediction may be retrieved after processing any number of frames. At each time step, AdaFocus first takes a quick glance at the full frame with a light-weighted CNN f_G , obtaining cheap and coarse global features. Then the features are fed into a recurrent policy network π to aggregate the information across frames and accordingly determine the location of an image patch to be focused on, under the goal of maximizing its contribution to video recognition. A high-capacity local CNN f_L is then adopted to process the se-

lected patch for more accurate but computationally expensive representations. Finally, a classifier f_C integrates the features of all previous frames to produce a prediction. In the following, we describe these four components in details.

Global CNN f_G and local CNN f_L are backbone networks that both extract deep features from the inputs, but with distinct aims. The former is designed to quickly catch a glimpse of each frame, providing necessary information for determining which region the local CNN f_L should attend to. Therefore, a light-weighted network is adopted for f_G . On the contrary, f_L is leveraged to take full advantage of the selected image regions for learning discriminative representations, and hence we deploy large and accurate models. Since f_L only needs to process a series of relatively small regions instead of the full images, this stage also enjoys high efficiency. We defer the details on the architectures of f_G and f_L to Section 4.

Formally, given video frames $\{v_1, v_2, \dots\}$ with size $H \times W$, f_G directly takes them as inputs and produces the coarse global feature maps e_t^G :

$$e_t^G = f_G(v_t), \quad t = 1, 2, \dots, \quad (1)$$

where t is the frame index. By contrast, f_L processes $P \times P$ ($P < H, W$) square image patches $\{\tilde{v}_1, \tilde{v}_2, \dots\}$, which are cropped from $\{v_1, v_2, \dots\}$ respectively, and we have

$$e_t^L = f_L(\tilde{v}_t), \quad t = 1, 2, \dots, \quad (2)$$

where e_t^L denotes the fine local feature maps. Importantly, the patch \tilde{v}_t is localized to capture the most informative re-

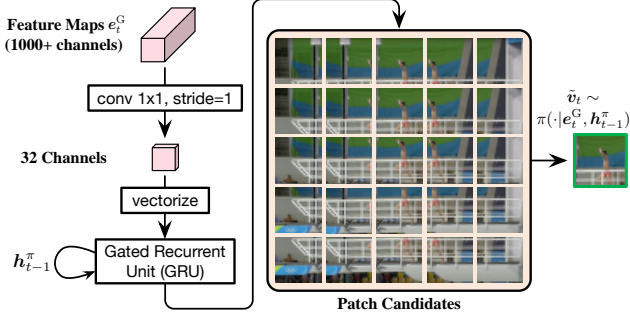


Figure 3. **The architecture of the policy network π .** The global feature maps e_t^G is fed into a 1×1 convolutional layer followed by a gated recurrent unit (GRU) to aggregate temporal information. The outputs of GRU parameterize a categorical distribution $\pi(\cdot|e_t^G, \mathbf{h}_{t-1}^\pi)$ on multiple patch candidates (here we take 25 as an example). During training, we sample \tilde{v}_t from $\pi(\cdot|e_t^G, \mathbf{h}_{t-1}^\pi)$, while at test time, we directly select the patch with the largest softmax probability.

gions for a given task, and this procedure is fulfilled by the policy network π , which is described in the following.

Policy network π is a recurrent network that receives the coarse global features e_t^G from the global CNN f_G , and specifies which region the global CNN f_G should attend to for each frame. Note that both the information of previous and current inputs is used due to the recurrent design. Formally, π determines the locations of images patches $\{\tilde{v}_1, \tilde{v}_2, \dots\}$ to be cropped from the frames. Given that this leads to a non-differentiable operation, we formalize π as an agent and train it with reinforcement learning. In specific, the location of the patch \tilde{v}_t is drawn from the distribution:

$$\tilde{v}_t \sim \pi(\cdot|e_t^G, \mathbf{h}_{t-1}^\pi), \quad (3)$$

where \mathbf{h}_{t-1}^π denotes the hidden states maintained in π that are updated at $(t-1)^{\text{th}}$ frame. In our implementation, we consider multiple candidates (e.g., 36 or 49) uniformly distributed across the images, and establish a categorical distribution on them, which is parameterized by the outputs of π . At test time, we simply adopt the candidate with maximum probability as \tilde{v}_t for a deterministic inference procedure. In addition, note that we do not perform any pooling on the features maps e_t^G since pooling typically corrupts the useful spatial information for localizing \tilde{v}_t . As an alternative, we compress the number of channels with 1×1 convolution to reduce the computational cost of π . An illustration of π is shown in Figure 3.

Classifier f_C is a prediction network aiming to aggregate the information from all the frames that have been processed by the model, and output the current recognition result at each time step. To be specific, we perform global average pooling on the feature maps e_t^G, e_t^L from the two aforementioned CNNs to get feature vectors \bar{e}_t^G, \bar{e}_t^L , and

concatenate them as the inputs of f_C , i.e.,

$$\mathbf{p}_t = f_C([\bar{e}_1^G, \bar{e}_1^L], \dots, [\bar{e}_t^G, \bar{e}_t^L]), \quad (4)$$

where \mathbf{p}_t refers to the softmax prediction at t^{th} step. It is worth noting that we allow e_t^G to be utilized for classification as well, with the aim of facilitating more efficient feature reusing. Such a design leverages previous observations [57, 39] revealing that CNNs are capable of achieving both remarkable localization and recognition performance at the same time. Many of existing methods also adopt similar reusing mechanisms [51, 50, 33, 16]. Besides, there are multiple possible architectures for f_C . In addition to the choice of recurrent networks such as long short-term memory (LSTM) [22] or gated recurrent unit (GRU) [4], f_C can also be set as taking the average of frame-wise predictions, which are typically obtained with a common fully-connected layer, as done in [31, 33, 34].

3.2. Training Algorithm

To ensure the four components function properly, a three-stage training algorithm is introduced.

Stage I: Warming-up. We first initialize f_G, f_L and f_C , but leave the policy network π out at this stage. Then we randomly sample the image patches \tilde{v}_t to minimize the cross-entropy loss $L_{CE}(\cdot)$ over the training set $\mathcal{D}_{\text{train}}$:

$$\underset{f_G, f_L, f_C}{\text{minimize}} \quad \mathbb{E}_{\{v_1, v_2, \dots\} \in \mathcal{D}_{\text{train}}} \left[\frac{1}{T} \sum_{t=1}^T L_{CE}(\mathbf{p}_t, y) \right], \quad (5)$$

$$\tilde{v}_t \sim \text{RandomCrop}(v_t).$$

Here T and y refer to the length and the label corresponding to the video $\{v_1, v_2, \dots\}$, respectively. In this stage, the model learns to extract task-relevant information from an arbitrary sequence of frame patches, laying the basis for training the policy network π .

Stage II: Learning to select informative patches. In this stage, we fix the two CNNs (f_G and f_L) and the classifier f_C obtained in stage I, and evoke a randomly initialized policy network π to be trained with reinforcement learning. Specifically, after sampling a location of \tilde{v}_t from $\pi(\cdot|e_t^G, \mathbf{h}_{t-1}^\pi)$ for the frame v_t (see Eq. (3)), π will receive a reward r_t indicating whether this action is beneficial. We train π to maximize the sum of discounted rewards:

$$\underset{\pi}{\text{maximize}} \quad \mathbb{E}_{\tilde{v}_t \sim \pi(\cdot|e_t^G, \mathbf{h}_{t-1}^\pi)} \left[\sum_{t=1}^T \gamma^{t-1} r_t \right], \quad (6)$$

where $\gamma \in (0, 1)$ is a discount factor for long-term rewards. In our implementation, we fix $\gamma = 0.7$ and solve Eq. (6) using the off-the-shelf proximal policy optimization (PPO) algorithm [38]. Notably, here we directly train π on the basis of the features extracted by f_G , since previous works [57, 39] have demonstrated that CNNs learned for classification generally excel at localizing task-relevant regions with their deep representations.

Ideally, the reward r_t is expected to measure the value of selecting \tilde{v}_t in terms of video recognition. With this aim, we define r_t as:

$$\begin{aligned} r_t(\tilde{v}_t|\tilde{v}_1, \dots, \tilde{v}_{t-1}) \\ = p_{ty}(\tilde{v}_t|\tilde{v}_1, \dots, \tilde{v}_{t-1}) \\ - \mathbb{E}_{\tilde{v}_t \sim \text{RandomCrop}(v_t)} [p_{ty}(\tilde{v}_t|\tilde{v}_1, \dots, \tilde{v}_{t-1})], \end{aligned} \quad (7)$$

where p_{ty} refers to the softmax prediction on y (i.e., confidence on the ground truth label). When computing r_t , we assume all previous patches $\{\tilde{v}_1, \dots, \tilde{v}_{t-1}\}$ have been determined, while only \tilde{v}_t can be changed. The second term in Eq. (7) refers to the expected confidence achieved by the randomly sampled \tilde{v}_t . By introducing it we ensure $\mathbb{E}_{\tilde{v}_t}[r_t] = 0$, which is empirically found to yield a more stable training procedure. In experiments, we estimate this term with a single time of Monte-Carlo sampling. Intuitively, Eq. (7) encourages the model to select the patches that are capable of producing confident predictions on the correct labels with as fewer frames as possible.

Stage III: Fine-tuning. At the last stage, we fine-tune f_L and f_C (or only f_C) with the learned policy network π from stage II, namely, minimizing Eq. (5) with $\tilde{v}_t \sim \pi(\cdot|e_t^G, h_{t-1}^\pi)$. This stage further improves the performance of our method.

3.3. Reducing Temporal Redundancy

The proposed AdaFocus processes each video frame equivalently with the same amount of computation. In fact, it is compatible with existing methods that focus on reducing the temporal redundancy of videos. To demonstrate this, we propose an extended version of AdaFocus, named as AdaFocus+, which dynamically skips less important frames for the large network f_L .

In specific, we add an additional recurrent policy network π' that has the same inputs and architectures as π , as shown in Figure 4. This new network is trained simultaneously with π in Stage II. For each frame, the outputs of π' parameterize a Bernoulli random variable b_t :

$$b_t \sim \text{Bernoulli}(p_t^B), \quad p_t^B = \pi'(e_t^G, h_{t-1}^{\pi'}) \in (0, 1), \quad (8)$$

which specifies the probability of maintaining the frame v_t , i.e., $\Pr(b_t = 1) = p_t^B$. The hidden state of π' , $h_{t-1}^{\pi'}$, is updated at $(t-1)$ th frame.

During training, we sample b_t according to Eq. (8), and multiply it with the local feature vector \bar{e}_t^L , such that Eq. (4) changes to:

$$p'_t(b_1, \dots, b_t) = f_C([\bar{e}_1^G, b_1 \bar{e}_1^L], \dots, [\bar{e}_t^G, b_t \bar{e}_t^L]). \quad (9)$$

In other words, if $b_t = 1$, the procedure mentioned in Section 3.1 remains unchanged. If $b_t = 0$, we simply do not

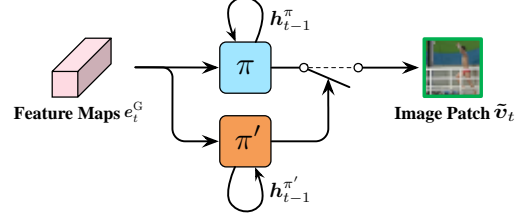


Figure 4. **An illustration of AdaFocus+.** The proposed AdaFocus method is naturally compatible with temporal-based techniques. By involving an additional policy network π' to control whether to attend on each frame (i.e., processing \tilde{v}_t with f_L), we can further reduce the redundant computation spent on less important frames.

feed the image patch \tilde{v}_t into the local CNN f_L , and concatenate \bar{e}_t^G with an all-zero tensor as the inputs of the classifier f_C . Notably, in this case, \bar{e}_t^G will also be fed into π to update its hidden state h_t^π , which introduces negligible computational overhead.

Similarly to π , π' is trained to maximize the sum of discounted rewards (Eq. (6)) as well. Here we define the reward r'_t corresponding to π' as:

$$r'_t(b_t|b_1, \dots, b_{t-1}) = \begin{cases} p'_{ty}(b_1, \dots, b_{t-1}, 1) \\ - p'_{ty}(b_1, \dots, b_{t-1}, 0) - \lambda P^2, & b_t = 1, \\ 0, & b_t = 0, \end{cases} \quad (10)$$

where $p'_{ty}(b_1, \dots, b_{t-1}, 1)$ and $p'_{ty}(b_1, \dots, b_{t-1}, 0)$ refer to the confidence on the ground truth label y with $b_t = 1$ and $b_t = 0$, respectively. The coefficient λ is a pre-defined hyper-parameter, while P is the length (or width) of the patch \tilde{v}_t . We use P^2 to estimate the required computation (FLOPs) of feeding \tilde{v}_t into f_L . When $b_t = 1$, the confidence gain of inferring f_L (i.e., $p'_{ty}(b_1, \dots, b_{t-1}, 1) - p'_{ty}(b_1, \dots, b_{t-1}, 0)$) is compared with the penalty term λP^2 , which reflects its computational costs. Only when this comparison produces positive results will the action of activating f_L be encouraged. Otherwise, π' will be trained to decrease the probability of $b_t = 1$, namely avoiding attending to any local region with the expensive f_L to avoid redundant computation.

During inference, we compare p_t^B of each frame with a fixed threshold $\rho \in (0, 1)$. When $p_t^B \geq \rho$, we process \tilde{v}_t with f_L , and otherwise this patch is skipped. These two cases correspond to $b_t = 1$ and $b_t = 0$ during training, respectively. The value of ρ should be solved on the validation set by only activating f_L for the frames with top $\eta\%$ largest p_t^B ($0 < \eta < 100$). One may vary $\eta\%$ for a flexible trade-off between computational costs and accuracy.

3.4. Offline Video Recognition

All the discussions above are based on the online video recognition setting where the model needs to output a reasonable prediction after seeing each frame. However, we

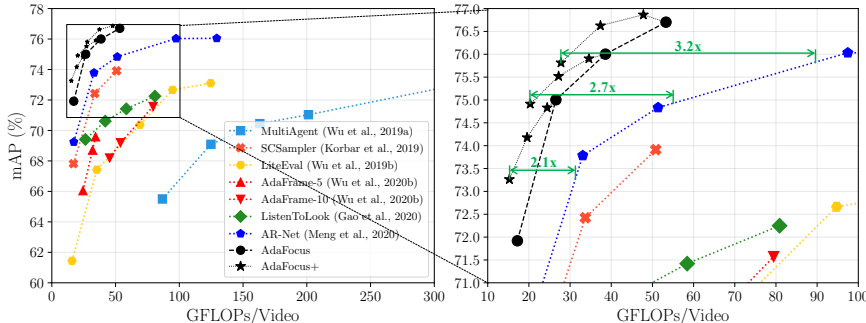


Figure 5. **Offline video recognition results on ActivityNet.** The whole video is provided at a time for a single prediction. Our method is implemented with the patch size $P^2 \in \{96^2, 128^2, 160^2, 192^2\}$. AdaFocus and AdaFocus+ refer to the vanilla spatial redundancy-based AdaFocus and its augmented version by further reducing temporal redundancy.

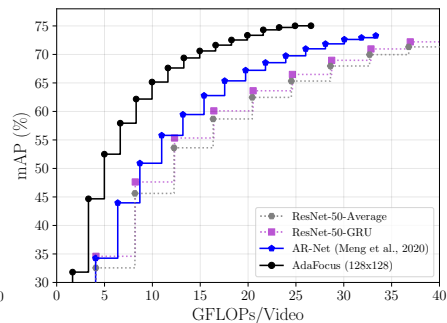


Figure 6. **Online video recognition results on ActivityNet.** Video frames come sequentially, while the models may need to output the prediction after processing any number of frames.

Table 1. **Comparisons of AdaFocus and state-of-the-art efficient video recognition frameworks on ActivityNet-v1.3 and FCVID.** GFLOPs refers to the average computational cost for processing a single video. MN2 and RN denote MobileNet-V2 and ResNet, respectively. The best results are **bold-faced**.

Methods	Backbones	ActivityNet		FCVID	
		mAP	GFLOPs	mAP	GFLOPs
FrameGlimpses [55]	VGG	60.2%	32.9	71.2%	29.9
AdaFrame [50]	MN2+RN	71.5%	79.0	80.2%	75.1
LiteEval [51]	MN2+RN	72.7%	95.1	80.0%	94.3
ListenToLook [16]	MN2+RN	72.3%	81.4	—	—
SCSampler [29]	MN2+RN	72.9%	42.0	81.0%	42.0
AR-Net [33]	MN2+RN	73.8%	33.5	81.3%	35.1
AdaFocus (128x128)	MN2+RN	75.0%	26.6	83.4%	26.6

note that AdaFocus can be straightforwardly adapted to the offline scenario where all frames are given in batch. In specific, one may train a model with the aforementioned approach, but only collect the result correspond to the last frame during inference. Importantly, the feed-forward process of both f_G and f_L , which accounts for the majority of computation, can be executed in parallel, enabling an efficient implementation on GPU devices.

4. Experiment

In this section, we empirically validate our method. We first compare AdaFocus with several recently proposed efficient video recognition frameworks, showing that AdaFocus gives rise to an improved efficiency. Then we implement our method by incorporating state-of-the-art lightweight CNN architectures to demonstrate that AdaFocus complements them and further improves the efficiency. Finally, we provide detailed visualization and ablation results to give additional insights into our method.

Datasets. Our experiments are based on five widely-used video datasets: (1) ActivityNet-v1.3 [1] contains 10,024 training videos and 4,926 validation videos labeled

Table 2. **Performance of AdaFocus and baselines on Mini-Kinetics.** GFLOPs refers to the average computational cost for processing a single video. MN2 and RN denote MobileNet-V2 and ResNet, respectively. The best results are **bold-faced**.

Methods	Backbones	Mini-Kinetics	
		Top-1 Acc.	GFLOPs
LiteEval [51]	MN2+RN	61.0%	99.0
SCSampler [29]	MN2+RN	70.8%	42.0
AR-Net [33]	MN2+RN	71.7%	32.0
AdaFocus (128x128)	MN2+RN	72.2%	26.6
AdaFocus (160x160)	MN2+RN	72.9%	38.6
AdaFocus+ (160x160)	MN2+RN	71.7%	20.3

by 200 action categories. The average duration is 117 seconds; (2) FCVID [26] includes 45,611 training videos and 45,612 validation videos labeled into 239 classes. The average duration is 167 seconds; (3) Mini-Kinetics is a subset of Kinetics [27] introduced by [33, 34]. The dataset consists of 200 classes of videos selected from Kinetics, with 121k videos for training and 10k videos for validation; (4) Something-Something V1&V2 [17] are two large-scale human action datasets, including 98k and 194k videos respectively. We use the official training-validation split.

Data pre-processing. Unless otherwise specified, we uniformly sample 16 frames from each video on ActivityNet, FCVID and Mini-Kinetics, while sample 8 or 12 frames on Something-Something. Following [31, 33], we augment training data by first adopting random scaling followed by 224x224 random cropping, and then performing random flipping on all datasets except for Something-Something V1&V2. During inference, we resize all frames to 256x256 and centre-crop them to 224x224.

4.1. Comparisons with State-of-the-art Efficient Video Recognition Methods

Baselines. In this subsection, AdaFocus is compared with several competitive baselines that focus on facilitat-

Table 3. **Performance of AdaFocus-TSM and other recently proposed efficient CNNs on Something-Something (Sth-Sth).** TSM+ refers to the augmented TSM baseline with the same network architecture as our method except for the policy network π . We uniformly sample 8/12 frames for the MobileNet/ResNet-50 in our models². The latency and throughput are tested on a 2.20GHz Intel Core i7-10870H CPU and a NVIDIA GeForce RTX 2080Ti GPU with the batch size of 1 and 64, respectively. The best results are **bold-faced**.

Method	Backbones	#Frames	Sth-Sth V1		Sth-Sth V2		Latency	Throughput
			Top-1 Acc.	GFLOPs	Top-1 Acc.	GFLOPs	(Intel Core i7, bs=1)	(NVIDIA 2080Ti, bs=64)
I3D [2]	3DResNet50	32×2	41.6%	306	-	-	-	-
I3D+GCN+NL [44]	3DResNet50	32×2	46.1%	606	-	-	-	-
ECO _{En} Lite [60]	BN-Inception + 3DResNet18	92	46.4%	267	-	-	-	-
TSN [43]	ResNet50	8	19.7%	33.2	27.8%	33.2	-	-
TRN _{RGB/Flow} [56]	BN-Inception	8/8	42.0%	32.0	55.5%	32.0	-	-
ECO [60]	BN-Inception+3DResNet18	8	39.6%	32.0	-	-	-	-
AdaFuse [34]	ResNet50	8	46.8%	31.5	59.8%	31.3	-	-
TSM [31]	ResNet50	8	46.1%	32.7	59.1%	32.7	0.32s	128.8 Videos/s
TSM+ [31]	MobileNet-V2+ResNet50	8+8 ²	47.0%	35.1	59.6%	35.1	0.42s	105.0 Videos/s
AdaFocus-TSM (144x144)	MobileNet-V2+ResNet50	8+12	47.0%	23.5 (\downarrow 1.49x)	59.7%	23.5 (\downarrow 1.49x)	0.32s (\downarrow 1.31x)	143.8 Videos/s (\uparrow 1.37x)
AdaFocus-TSM (160x160)	MobileNet-V2+ResNet50	8+12	47.6%	27.5	60.2%	27.5	0.36s	122.1 Videos/s
AdaFocus-TSM (176x176)	MobileNet-V2+ResNet50	8+12	48.1%	33.7	60.7%	33.7	0.42s	104.2 Videos/s

ing efficient video recognition, including MultiAgent [48], SCSampler [29], LiteEval [51], AdaFrame [50], Listen-to-look [16] and AR-Net [33]. Due to spatial limitations, we briefly introduce them in Appendix A.

Implementation details. We deploy MobileNet-V2 [37] and ResNet-50 [20] as the global CNN f_G and local CNN f_L in AdaFocus. A one-layer gated recurrent unit (GRU) [4] with a hidden size of 1024 is used in both the policy network π and the classifier f_C . The number of patch candidates is set to 49 (uniformly distributed in 7x7). Due to the limited space, training details are deferred to Appendix B.

Offline video recognition. We first implement AdaFocus under the offline recognition setting, where our method produces a single prediction after processing the whole video. This setting is adopted by the papers of most baselines as well. The results on ActivityNet and FCVID are presented in Table 1. We use a patch size of $P^2 = 128^2$ in AdaFocus, and evaluate the performance of different methods via mean average precision (mAP) following the common practice [50, 51, 16, 33] on these two datasets. It can be observed that our method outperforms the alternative baselines by large margins in terms of efficiency. For example, on FCVID, AdaFocus achieves 2.1% higher mAP (83.4% v.s. 81.3%) than the strongest baseline, AR-Net, with 1.3x less computation (26.6 GFLOPs v.s. 33.5 GFLOPs).

Results of varying patch sizes are presented in Figure 5. We change the patch size within $P^2 \in \{96 \times 96, 128 \times 128, 160 \times 160, 192 \times 192\}$, and plot the corresponding mAP v.s. FLOPs relationships in black dots. We also present the variants of baselines with various computational costs. One can observe that AdaFocus leads to a considerably better trade-off between efficiency and accuracy.

Improvements from further reducing temporal re-

²In fact, we can also sample 8/12 frames for TSM+, but this increases computational costs dramatically (\sim 1.5x). Hence, we do not consider it.

dundancy. Then we test extending AdaFocus by skipping less informative frames, as stated in Section 3.3. The results are presented as AdaFocus+ (black stars) in Figure 5. The coefficient λ is set to $1e - 6$, while the skipping proportion $\eta\%$ is varied within $\{0.9, 0.7, 0.5\}$. For the ease of implementation, we solve the threshold ρ on the training set. We find this achieves almost the same performance as using a validation set. It is clear that further reducing temporal redundancy leads to a significantly better efficiency. With a given mAP, the number of required GFLOPs per video for AdaFocus+ is approximately 2.1-3.2x less than AR-Net.

Results on Mini-Kinetics are presented in Table 2. The observations here are similar to ActivityNet/FCVID. AdaFocus+ reduces the required computation to reach 71.7% accuracy by 1.6x (20.3 GFLOPs v.s. 32.0 GFLOPs).

Online video recognition results are shown in Figure 6. Note that here we assume a stream of video frames come in sequentially for processing and the model may need to output a prediction at any time. In specific, we take a fixed number of frames from the beginning of videos, feed them into networks to evaluate the results, and change the number of frames to obtain the mAP-FLOPs trade-off. We consider two additional baselines: (1) ResNet-50-Average averages the frame-level predictions of a ResNet-50 with the full inputs; (2) ResNet-50-GRU augments (1) by aggregating the features across frames using a GRU classifier. Figure 6 shows that our method is able to obtain much better performance given the same number of FLOPs, which enables accurate and fast recognition in real-time applications.

4.2. Building on Top of Efficient CNNs

Setup. In this subsection, we implement AdaFocus on top of the recently proposed efficient network architecture, CNNs with temporal shift module (TSM) [31], to demonstrate that our method can effectively improve the efficiency



Figure 7. **Visualization results (zoom in for details).** The green boxes indicate the locations of the image patches selected by AdaFocus.

Table 4. **Effects of Reusing e_t^G for Recognition.**

Reusing e_t^G for Recognition	mAP			
	96x96	128x128	160x160	192x192
✗	70.2%	73.4%	75.0%	75.9%
✓	71.9%	75.0%	76.0%	76.7%

of such state-of-the-art light-weighted models. Specifically, we still use MobileNet-V2 and ResNet-50 as f_G and f_L , but add TSM to them. A fully-connected layer is deployed as the classifier f_C , and we average the frame-wise predictions as the output, following the design of TSM [31]. For a fair comparison, we augment the vanilla TSM by introducing the same two backbone networks as ours (named as TSM+), where their output features are also concatenated to be fed into a linear classifier. In other words, TSM+ differentiates itself from AdaFocus only in that it feeds the whole frames into ResNet-50, while we feed the selected image patches.

The setting of offline video recognition on Something-Something V1&V2 is used here. Notably, as the videos in the two datasets are very short (average duration \approx 4s), we find the networks require the visual angle of adjacent inputs (frames/patches) to be similar for high generalization performance. We also observe that the locations of task-relevant regions do not significantly change across the frames in the same video. Therefore, here we let AdaFocus generate a single patch location for the whole video after aggregates the information of all frames. Importantly, such a simplification does not affect the main idea of our method since different videos have varying patch locations. The architecture of the policy network π and the training algorithm remain unchanged. Details of training hyper-parameters are deferred to Appendix B.

Results on Something-Something are reported in Table 3. One can observe that by reducing the input size of the relatively expensive ResNet-50 network, AdaFocus en-

ables TSM to process more frames in the task-relevant region of each video using the same computation, leading to a significantly improved efficiency. For example, AdaFocus achieves the same performance as TSM+ with 1.5x less GFLOPs on Something-Something V1.

Practical efficiency. In Table 3, we also test the actual inference speed of AdaFocus-TSM on both an Intel i7 CPU and a NVIDIA 2080Ti GPU, with the batch size of 1 and 64, respectively, which are sufficient to saturate the two devices. It can be observed that our practical speedup is significant as well, with a slight drop compared with theoretical results. We tentatively attribute this to the inadequate hardware-oriented optimization in our implementation.

4.3. Analytical Results

Visualization. In Figure 7, we visualize the regions selected by our proposed AdaFocus. Here we uniformly sample 8 video frames from ActivityNet. One can observe that our method effectively guides the expensive local CNN f_L to attend to the task-relevant regions of each frame, such as the sailboard, bicycle and flute.

Importance of reusing e_t^G for recognition. As aforementioned, our method effectively leverages the coarse global feature e_t^G for both localizing the task-relevant patch \tilde{v}_t and recognition. As shown in Table 4, only using e_t^G for localization degrades the mAP by 1 – 1.5%, which demonstrates the effects of this reusing mechanism.

Effectiveness of the learned patch selection policy is validated in Table 5. We consider its three alternatives: (1) *randomly* sampling patches, (2) cropping patches from the *centres* of the frames, and (3) sampling patches from a standard *gaussian* distribution centred the frame. In addition, we test altering the reward function for reinforcement learning to: (1) *confidence reward* directly uses the confidence on ground truth labels as rewards, and (2) *increments reward* utilizes the increments of confidence as rewards. For a

Table 5. **Comparisons of various patch selection policies.** Fixed policies are pre-defined without leveraging reinforcement learning (RL). For RL-based policy, we change the design of rewards.

Ablation		mAP			
		96x96	128x128	160x160	192x192
Fixed Policy	Random Policy	65.8%	70.7%	73.1%	74.8%
	Central Policy	61.9%	68.7%	72.4%	74.8%
	Gaussian Policy	64.7%	70.6%	73.5%	74.9%
Learned Policy	Confidence Reward	68.5%	72.3%	74.1%	75.5%
	Increments Reward	69.4%	72.7%	74.4%	75.6%
by RL	AdaFocus (ours)	70.2%	73.4%	75.0%	75.9%

clear comparison, here we do not reuse e_t^G for recognition. An interesting phenomenon is that random policy appears strong and outperforms the central policy, which may be attributed to the spatial similarity between frames. That is to say, adjacent central patches might have repetitive contents, while randomly sampling is likely to collect more comprehensive information. Besides, it is shown that the learned policies have considerably better performance, and our proposed reward function significantly outperforms others.

5. Conclusion

This paper has proposed a *spatial redundancy* based approach for efficient video recognition, AdaFocus. Inspired by the fact that not all image regions in video frames are task-relevant, AdaFocus reduces computational costs by inferring the high-capacity network only on a small but informative patch of each frame, which is adaptively localized with reinforcement learning. We further show that our method can be extended by dynamically skipping less valuable frames. Extensive experiments demonstrate that our method outperforms existing works in terms of both theoretical computational efficiency and actual inference speed.

Acknowledgements

This work is supported in part by the National Science and Technology Major Project of the Ministry of Science and Technology of China under Grants 2018AAA0100701, the National Natural Science Foundation of China under Grants 61906106 and 62022048, the Institute for Guo Qiang of Tsinghua University and Beijing Academy of Artificial Intelligence.

References

[1] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *CVPR*, pages 961–970, 2015. [6](#)

[2] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*, pages 6299–6308, 2017. [1](#), [2](#), [7](#)

[3] Jianguo Chen, Kenli Li, Qingying Deng, Keqin Li, and S Yu Philip. Distributed deep learning model for intelligent video surveillance systems with edge computing. *IEEE Transactions on Industrial Informatics*, 2019. [1](#)

[4] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *EMNLP*, pages 1724–1734, Doha, Qatar, Oct. 2014. Association for Computational Linguistics. [4](#), [7](#)

[5] Wen-Hsuan Chu, Yu-Jhe Li, Jing-Cheng Chang, and Yu-Chiang Frank Wang. Spot and learn: A maximum-entropy patch sampler for few-shot image classification. In *CVPR*, pages 6251–6260, 2019. [2](#)

[6] Robert T Collins, Alan J Lipton, Takeo Kanade, Hironobu Fujiyoshi, David Duggins, Yanghai Tsin, David Tolliver, Nobuyoshi Enomoto, Osamu Hasegawa, Peter Burt, et al. A system for video surveillance and monitoring. *VSAM final report*, 2000(1-68):1, 2000. [1](#)

[7] James Davidson, Benjamin Liebald, Junning Liu, Palash Nandy, Taylor Van Vleet, Ullas Gargi, Sujoy Gupta, Yu He, Mike Lambert, Blake Livingston, et al. The youtube video recommendation system. In *Proceedings of the fourth ACM conference on Recommender systems*, pages 293–296, 2010. [1](#)

[8] Yashar Deldjoo, Mehdi Elahi, Paolo Cremonesi, Franca Garzotto, Pietro Piazzolla, and Massimo Quadrana. Content-based video recommendation system based on stylistic visual features. *Journal on Data Semantics*, 5(2):99–113, 2016. [1](#)

[9] Jeffrey 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. In *CVPR*, pages 2625–2634, 2015. [2](#)

[10] 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*, 2018. [2](#)

[11] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *ICCV*, pages 6202–6211, 2019. [1](#)

[12] Christoph Feichtenhofer, Axel Pinz, and Andrew Zisserman. Convolutional two-stream network fusion for video action recognition. In *CVPR*, pages 1933–1941, 2016. [1](#)

[13] 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*, pages 1039–1048, 2017. [2](#)

[14] Jianlong Fu, Heliang Zheng, and Tao Mei. Look closer to see better: Recurrent attention convolutional neural network for fine-grained image recognition. In *CVPR*, pages 4438–4446, 2017. [2](#)

[15] Junyu Gao, Tianzhu Zhang, and Changsheng Xu. A unified personalized video recommendation via dynamic recurrent neural networks. In *ACM MM*, pages 127–135, 2017. [1](#)

[16] Ruohan Gao, Tae-Hyun Oh, Kristen Grauman, and Lorenzo Torresani. Listen to look: Action recognition by previewing audio. In *CVPR*, pages 10457–10467, 2020. [1](#), [2](#), [4](#), [6](#), [7](#), [11](#)

- [17] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The "something something" video database for learning and evaluating visual common sense. In *ICCV*, pages 5842–5850, 2017. 6
- [18] Yizeng Han, Gao Huang, Shiji Song, Le Yang, Honghui Wang, and Yulin Wang. Dynamic neural networks: A survey. *arXiv preprint arXiv:2102.04906*, 2021. 2
- [19] Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet? In *CVPR*, pages 6546–6555, 2018. 1, 2
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 7, 11
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 11
- [22] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997. 4
- [23] Gao Huang, Yixuan Li, Geoff Pleiss, Zhuang Liu, John E Hopcroft, and Kilian Q Weinberger. Snapshot ensembles: Train 1, get m for free. In *ICLR*, 2017. 11
- [24] Gao Huang, Zhuang Liu, Geoff Pleiss, Laurens Van Der Maaten, and Kilian Weinberger. Convolutional networks with dense connectivity. *IEEE transactions on pattern analysis and machine intelligence*, 2019. 11
- [25] Nazli Ikizler and David Forsyth. Searching video for complex activities with finite state models. In *CVPR*, pages 1–8. IEEE, 2007. 1
- [26] Y.-G. Jiang, Z. Wu, J. Wang, X. Xue, and S.-F. Chang. Exploiting feature and class relationships in video categorization with regularized deep neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(2):352–364, 2018. 6
- [27] 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. 6
- [28] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 11
- [29] Bruno Korbar, Du Tran, and Lorenzo Torresani. Scsampler: Sampling salient clips from video for efficient action recognition. In *ICCV*, pages 6232–6242, 2019. 1, 2, 6, 7, 11
- [30] Hengduo Li, Zuxuan Wu, Abhinav Shrivastava, and Larry S Davis. 2d or not 2d? adaptive 3d convolution selection for efficient video recognition. *arXiv preprint arXiv:2012.14950*, 2020. 1, 2
- [31] Ji Lin, Chuang Gan, and Song Han. Tsm: Temporal shift module for efficient video understanding. In *ICCV*, pages 7083–7093, 2019. 2, 4, 6, 7, 8, 11
- [32] Lili Meng, Bo Zhao, Bo Chang, Gao Huang, Wei Sun, Frederick Tung, and Leonid Sigal. Interpretable spatio-temporal attention for video action recognition. In *ICCV Workshops*, pages 0–0, 2019. 2
- [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. In *ECCV*, pages 86–104. Springer, 2020. 1, 2, 4, 6, 7, 11
- [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 *ICLR*, 2021. 2, 4, 6, 7
- [35] Volodymyr Mnih, Nicolas Heess, Alex Graves, et al. Recurrent models of visual attention. In *NeurIPS*, pages 2204–2212, 2014. 2
- [36] Bowen Pan, Rameswar Panda, Camilo Fosco, Chung-Ching Lin, Alex Andonian, Yue Meng, Kate Saenko, Aude Oliva, and Rogerio Feris. Va-red²: Video adaptive redundancy reduction. *arXiv preprint arXiv:2102.07887*, 2021. 2
- [37] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *CVPR*, pages 4510–4520, 2018. 7, 11
- [38] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. 4
- [39] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *ICCV*, pages 618–626, 2017. 4
- [40] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *ICCV*, pages 4489–4497, 2015. 1, 2
- [41] Du Tran, Heng Wang, Lorenzo Torresani, and Matt Feiszli. Video classification with channel-separated convolutional networks. In *ICCV*, pages 5552–5561, 2019. 2
- [42] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *CVPR*, pages 6450–6459, 2018. 2
- [43] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *ECCV*, pages 20–36. Springer, 2016. 2, 7
- [44] Xiaolong Wang and Abhinav Gupta. Videos as space-time region graphs. In *ECCV*, pages 399–417, 2018. 7
- [45] Yulin Wang, Kangchen Lv, Rui Huang, Shiji Song, Le Yang, and Gao Huang. Glance and focus: a dynamic approach to reducing spatial redundancy in image classification. In *NeurIPS*, 2020. 2, 11
- [46] Yulin Wang, Zanlin Ni, Shiji Song, Le Yang, and Gao Huang. Revisiting locally supervised learning: an alternative to end-to-end training. In *ICLR*, 2021. 11
- [47] Yulin Wang, Xuran Pan, Shiji Song, Hong Zhang, Gao Huang, and Cheng Wu. Implicit semantic data augmentation for deep networks. In *NeurIPS*, pages 12635–12644, 2019. 11

- [48] Wenhao Wu, Dongliang He, Xiao Tan, Shifeng Chen, and Shilei Wen. Multi-agent reinforcement learning based frame sampling for effective untrimmed video recognition. In *ICCV*, pages 6222–6231, 2019a. [1](#), [2](#), [7](#), [11](#)
- [49] Wenhao Wu, Dongliang He, Xiao Tan, Shifeng Chen, Yi Yang, and Shilei Wen. Dynamic inference: A new approach toward efficient video action recognition. In *CVPR Workshops*, pages 676–677, 2020a. [2](#)
- [50] Zuxuan Wu, Hengduo Li, Caiming Xiong, Yu-Gang Jiang, and Larry Steven Davis. A dynamic frame selection framework for fast video recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020b. [1](#), [2](#), [4](#), [6](#), [7](#), [11](#)
- [51] Zuxuan Wu, Caiming Xiong, Yu-Gang Jiang, and Larry S Davis. Liteeval: A coarse-to-fine framework for resource efficient video recognition. In *NeurIPS*, 2019b. [1](#), [2](#), [4](#), [6](#), [7](#), [11](#)
- [52] Saining Xie, Chen Sun, Jonathan Huang, Zhuowen Tu, and Kevin Murphy. Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification. In *ECCV*, pages 305–321, 2018. [2](#)
- [53] Zhenda Xie, Zheng Zhang, Xizhou Zhu, Gao Huang, and Stephen Lin. Spatially adaptive inference with stochastic feature sampling and interpolation. In *ECCV*, pages 531–548. Springer, 2020. [2](#)
- [54] Le Yang, Yizeng Han, Xi Chen, Shiji Song, Jifeng Dai, and Gao Huang. Resolution adaptive networks for efficient inference. In *CVPR*, pages 2369–2378, 2020. [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 *CVPR*, pages 2678–2687, 2016. [2](#), [6](#)
- [56] Bolei Zhou, Alex Andonian, Aude Oliva, and Antonio Torralba. Temporal relational reasoning in videos. In *ECCV*, pages 803–818, 2018. [7](#)
- [57] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *CVPR*, pages 2921–2929, 2016. [4](#)
- [58] Sijie Zhu, Taojiannan Yang, Matias Mendieta, and Chen Chen. A3d: Adaptive 3d networks for video action recognition. *arXiv preprint arXiv:2011.12384*, 2020. [2](#)
- [59] Xizhou Zhu, Yuwen Xiong, Jifeng Dai, Lu Yuan, and Yichen Wei. Deep feature flow for video recognition. In *CVPR*, pages 2349–2358, 2017. [1](#)
- [60] Mohammadreza Zolfaghari, Kamaljeet Singh, and Thomas Brox. Eco: Efficient convolutional network for online video understanding. In *ECCV*, pages 695–712, 2018. [2](#), [7](#)

Appendix for “Adaptive Focus for Efficient Video Recognition”

A. Introduction of Baselines

AdaFocus is compared with several competitive baselines that focus on facilitating efficient video recognition, including MultiAgent [48], SCSampler [29], LiteEval [51], AdaFrame [50], Listen-to-look [16] and AR-Net [33].

- MultiAgent [48] proposes to learn to select important frames with multi-agent reinforcement learning.
- SCSampler [29] introduces a light-weighted framework to efficiently identify the most salient temporal clips within a long video. We follow the implementation of [33].
- LiteEval [51] combines a coarse LSTM and a fine LSTM to adaptively allocate computation based on the importance of frames.
- AdaFrame [50] learns to dynamically select informative frames with reinforcement learning and performs adaptive inference.
- Listen-to-look [16] fuses image and audio information to select the key clips within a video. As we do not leverage the audio of videos, for a fair comparison, we adopt its image-based version introduced in their paper.
- AR-Net [33] dynamically identifies the importance of video frames, and processes them with different resolutions accordingly.

B. Implementation Details

B.1. Training Hyper-parameters for Section 4.1

In our implementation, we always train f_G , f_L and f_C using a SGD optimizer with cosine learning rate annealing and a Nesterov momentum of 0.9 [21, 24, 23, 47, 46, 31, 33]. The size of the mini-batch is set to 64, while the L2 regularization coefficient is set to $1e-4$. We initialize f_G and f_L by fine-tuning the ImageNet pre-trained MobileNet-V2 [37] and ResNet-50 [20] (we use the official models provided by PyTorch) using full inputs for 15 epochs with an initial learning rate of 0.01. In stage I, we train f_L and f_C using randomly sampled patches for 50 epochs with an initial learning rate of $5e-4$ and 0.05, respectively. Here we do not train f_G as we find this does not significantly improve the performance, but increases the training time. In stage II, we train π/π' with an Adam optimizer [28] for 50/10 epochs. The same training hyper-parameters as [45] are adopted. In stage III, we only fine-tune f_C with the learned policy for 10 epochs, since we find further fine-tuning f_L leads to trivial improvements but prolongs the training time. The initial learning rates are set to $5e-4$ and $5e-3$ for Mini-Kinetics and ActivityNet/FCVID, respectively.

B.2. Training Hyper-parameters for Section 4.2

Here we initialize f_G and f_L by training them using the same configuration as [31]. The training procedure of AdaFocus is the same as Section 4.1 except for the following changes. In stage I, we use the initial learning rate of

1e-5 and 0.01 for f_L and f_C , respectively, and train them for 10 epochs. In stage III, we use an initial learning rate of 5e-4 for f_C . Note that TSM+ follows exactly the same training procedure as our method. The only difference is that TSM+ does not train the policy network π , since it adopts full frames as inputs.