

No frame left behind: Full Video Action Recognition

Xin Liu¹ Silvia L. Pinte¹ Fatemeh Karimi Nejadasl² Olaf Booij² Jan C. van Gemert¹
Computer Vision Lab, Delft University of Technology¹ TomTom²

Abstract

Not all video frames are equally informative for recognizing an action. It is computationally infeasible to train deep networks on all video frames when actions develop over hundreds of frames. A common heuristic is uniformly sampling a small number of video frames and using these to recognize the action. Instead, here we propose full video action recognition and consider all video frames. To make this computational tractable, we first cluster all frame activations along the temporal dimension based on their similarity with respect to the classification task, and then temporally aggregate the frames in the clusters into a smaller number of representations. Our method is end-to-end trainable and computationally efficient as it relies on temporally localized clustering in combination with fast Hamming distances in feature space. We evaluate on UCF101, HMDB51, Breakfast, and Something-Something V1 and V2, where we compare favorably to existing heuristic frame sampling methods.

1. Introduction

Videos have arbitrary length with actions occurring at arbitrary moments. Current video recognition methods use CNNs on coarsely sub-sampled frames [2, 24, 28, 31, 38, 41, 45, 47, 48, 50] because using all frames is computationally infeasible. Sub-sampling, however, can miss crucial frames for action recognition. For example, as shown in Fig. 1, sampling the frame with the dish in the pan is crucial for correct recognition. We propose to do away with sub-sampling heuristics and argue for leveraging all video frames: Full video action recognition.

It is worth analyzing why training CNNs on full videos is computationally infeasible in terms of memory and calculations. The calculations in the forward pass yield activations, while the backward pass calculations give gradients which are summed over all frames to update the weights. Many of these calculations can be done in parallel and thus are well-suited for modern GPUs. When treating videos as a large collection of image frames, the amount of calculations are not too different from those on large image datasets [5].

Figure 1. Sub-sampling can miss crucial frames in videos and may cause confusion for action recognition: e.g. compare the two sub-samplings heuristics in row 1 and row 2: Without sampling the dish in the pan it is difficult to classify. Instead, as shown in row 3, we propose to efficiently use all frames during training by clustering frame activations along the temporal dimension and aggregating each cluster to a single representation. The temporal clustering is based on Hamming distances over frame activations, which is computationally fast. With the assumption that similar activations have similar gradients, the aggregated representations approximate the individual frame activations. We efficiently utilize all frames for training without missing important information.

Regarding memory, however, there is a crucial difference between videos and images: The video loss function is not per-frame but on the full video. Hence, to do the backward pass, all activations for each frame, for each filter in each layer need to be stored in memory. This even doubles for storing their gradients. With 10-30 frames per second, this quickly becomes infeasible for even just a few minutes of video. Existing approaches can trade off memory for compute [3, 4, 13] by not storing all intermediate layers, yet they do not scale to video as they would still need to store each frame. The main computational bottleneck for training video CNNs is memory for frame activations.

Here, we propose an efficient method to use all video frames during training. The forward pass computes frame

activations and the backward pass sums the gradients over the frames to update the weights. Now, if only the network was linear, then a huge memory reduction could be gained by first summing all frame activations in the forward pass, which would reduce to just a single update in the backward pass. Yet, deep networks are famously non-linear, and have non-linearities in the activation function and in the loss function. Thus, if all frames were independent, treating the non-linear network as linear would introduce considerable approximation errors. However, subsequent frames in a video are strongly correlated, and it's this correlation that makes it possible for existing approaches to use sub-sampling. Instead of sub-sampling, we propose to process all frames and exploit the frame correlations to create groups of frames where the network is approximately linear. We use the ReLU (Rectified Linear Unit) activation function, which is linear when the signs of two activations agree, to estimate which parts of the video are approximately linear. This allows us to develop an efficient clustering algorithm based on Hamming distances of frame activations as illustrated in Fig. 1. By then aggregating the approximately linear parts in a video in the forward pass, we make large memory savings in the backward pass while still approximating the full video gradient.

We summarize the contributions of our work as follows:

- We propose a method that allows us to use most or even all video frames for training action recognition by approximated individual frame gradients with the gradients of temporally aggregated frame activations;
- We devise an end-to-end trainable approach for efficient grouping of video frames based on temporally localized clustering and Hamming distances;
- Extensive experiments demonstrate that our method compares well to state-of-the-art methods on several benchmark datasets such as UCF101, HMDB51, Breakfast, and Something-Something V1 and V2.

2. Related work

Action recognition architectures. Actions in video involve motion, leading to deep networks which include optical flow [8, 10, 34], 3D convolutions [2, 6, 15, 20] and recurrent connections [10, 35, 40, 39, 46]. Instead of heavy-weight motion representations, a single 2D image can reveal much of an action [19, 22, 34, 41]. 2D CNNs are extremely efficient, and by adding motion information by concatenating a 3D module in ECO [50], modeling temporal relations in TSN [49] or simply shifting filter channels over time in TSM [28] their efficiency is complemented by good accuracy. For this reason, we build on the TSM [28] architecture and modify it for full video action recognition.

Frame sampling for action recognition. Realistic videos contain more frames than can fit in memory. To address this, current methods train by using sub-sampled video frames [2, 28, 41, 50]. Additionally, the SlowFast [7] network also explores the resolution trade-off across temporal, spatial and channel dimension. Rather than using a fixed frame sampling strategy, the sampling can be adaptive [28, 38, 45, 47, 48], or learned to select the best frame [31], or rely on clip sampling [24]. In our work we do not sub-sample frames, but use all frames of the videos, however our clustering is adaptive as it dynamically adapt to the task and the loss function.

Using a subset of frames is computationally more efficient. Using 5-7 frames is sufficient for state-of-the-art action recognition on short videos [32]. Aiming for training efficiency, the work in [42] uses stochastic mini-batch dropping which drops complete batches rather than frames, with a certain probability. Similarly, [44] uses variable mini-batch shapes with different spatio-temporal resolutions varied according to a schedule. Unlike these methods, we do not focus on training efficiency, but propose a method that allows the network to see all video frames during training.

Temporal pooling. To integrate frame-level features, TSN [41] uses average pooling in the late layers of the network. ActionVLAD [11] integrates two-stream networks with a learnable spatio-temporal feature aggregation. Instead of performing temporal pooling or aggregation at a late stage of the network, in [9] RankSVM is used to rank frames temporally and then pool them together. As a follow-up, in [1] a 'dynamic image' is introduced, which is a compact representation of the videos frames using the 'rank pool' operation. In [33, 36] temporal aggregation via pooling and attention is used. Similar to these methods, our proposal performs a temporal pooling of the network activations, however this aggregation is done over clustered activations and it allows us to process all video frames.

Efficient backpropagation. Given that 2/3 of the training computations and memory are spent in the backward pass, existing work focuses on approximations. It is more memory efficient to recompute activations from the previous layer instead of storing them [13], however this comes at the cost of increased training time. In [29] gradient approximations are used where activations are overwritten when new frames are seen without waiting for the backward pass to be performed. Also for efficient backpropagation, randomized automatic differentiation can be used [30], gradients can be reused during training [12], or even quantized during backpropagation [43]. Similar to these works, we use all frames to approximate the full video gradient.

3. Aggregated temporal clusters

3.1. Approximating gradients

We enable the use of all frames of a video during training. To this end, we calculate a single gradient to approximate the gradients of a group of frames. Our hypothesis is that nearby frames in a video are alike, and thus have similar activations, leading to congruent updates. When using the ReLU (Rectified Linear Unit) activation function, we know that for activations with agreeing signs, the activation function is linear. Assuming that similar frames are approximately linear, the standard computation of the sum of gradients over all frames, becomes equivalent to first summing all frame activations and then computing a single gradient. This is computationally and memory efficient. Mathematically, for frames i , this can be formulated as:

$$\sum_i wL(h(x_i w)) = \sum_i wL(h(x_i w)) \quad (1)$$

where x are frame activations, w are the network weights, $h(\cdot)$ is an activation function, and $L(\cdot)$ is the loss function. Note that Eq. (1) only holds in the ideal case when the activation function h is linear for similar frames and the loss function L is also linear. This is not generally the case, and this approximation introduces an error.

With the above ideal scenario in mind, we can use all video frames without calculating the gradient for each frame, by grouping frames that agree in the sign of their activations x . Over these grouped activations we calculate a single gradient $\sum_i wL(h(x_i w))$. However, for similar frames the sign of their activation values may not be in complete agreement. Therefore, we aim to find which frames can be safely grouped together, to minimize the error introduced by our approximation in Eq. (1).

3.2. Error bound for the approximation

For ease of explanation, we consider two input video frames and their activations $x = \{x_1, x_2\}$, and a convolutional operation with parameters w , denoted by xw . The two frames have the same class label, y , since they come from the same video. We consider a multi-class setting using the cross-entropy loss in combination with the softmax function q , which for these two samples is:

$$L(x, y) = -\frac{1}{2} (\log q_y(x_1) + \log q_y(x_2)), \quad (2)$$

where $q_c(x_i) = \frac{\exp(h(x_i w_c))}{\sum_{j=1}^C \exp(h(x_i w_j))}$, $c \in \{1, \dots, C\}$ indexes video classes and $h(\cdot)$ is the ReLU activation function. The gradient of the loss with respect to w is:

$$\sum_i wL(x, y) = \frac{x_1(q_c(x_1) - y_c) + x_2(q_c(x_2) - y_c)}{2}, \quad (3)$$

where y_c is the Dirac function which is 1 when $c=y$. In our method, we first average the two activations after the convolution and before the ReLU. We can do this, because if we assume the activations have agreeing signs $\text{sign}(x_1 w) = \text{sign}(x_2 w)$, then it holds that: $\frac{h(x_1 w) + h(x_2 w)}{2} = h(\frac{x_1 w + x_2 w}{2})$. In this case the cross-entropy loss becomes:

$$\hat{L}(x, y) = -\log q_y \left(\frac{x_1 + x_2}{2} \right). \quad (4)$$

In the backward pass, we calculate a single gradient of the averaged activations as follows:

$$\sum_i w\hat{L}(x, y) = \frac{x_1 + x_2}{2} q_c \left(\frac{x_1 + x_2}{2} \right) - y_c, \quad (5)$$

We now estimate the relative error introduced by our approximation by comparing equations Eq. (3) and Eq. (5) using Jensen's inequality. We start from the softmax function $q_c(\cdot)$ and we recover back equations Eq. (3) and Eq. (5). The softmax function $q_c(\cdot)$ is convex, therefore we can apply to it Jensen's inequality for the samples x_1 and x_2 : $q_c \left(\frac{x_1 + x_2}{2} \right) \leq \frac{q_c(x_1) + q_c(x_2)}{2}$. We start by considering the case $\frac{(x_1 + x_2)}{2} > 0$. If we multiply both sides of this inequality with $\frac{(x_1 + x_2)}{2}$ we obtain that:

$$\begin{aligned} \frac{(x_1 + x_2)}{2} q_c \left(\frac{x_1 + x_2}{2} \right) &\leq \frac{(x_1 + x_2)}{2} \frac{q_c(x_1) + q_c(x_2)}{2} \\ &= \frac{x_1 q_c(x_1) + x_2 q_c(x_2)}{2} \\ &\quad - \frac{1}{4} (x_1 - x_2) (q_c(x_1) - q_c(x_2)). \end{aligned} \quad (6)$$

In the left hand side of the inequality we recover precisely the $\sum_i w\hat{L}(x, y)$ given by Eq. (5), while in the right hand side we recover Eq. (3) minus the approximation error as $\sum_i wL(x, y) - \frac{1}{4} (x_1 - x_2) (q_c(x_1) - q_c(x_2))$. Note that for the case $y=c$ the additional Dirac terms in x cancel out. We now consider also the case $\frac{(x_1 + x_2)}{2} \leq 0$, which together with the previous case leads to the following bound on the absolute difference between the gradients in Eq. (3) and Eq. (5):

$$|\sum_i wL(x, y) - \sum_i w\hat{L}(x, y)| \leq \frac{1}{4} |(x_1 - x_2) (q_c(x_1) - q_c(x_2))|. \quad (7)$$

Thus, the difference between the two gradient updates is bounded by a function depending on the difference between the activations and their softmax responses. The closer to 0 the difference between the activations the smaller the difference between their gradient updates. We show in the experimental section that, indeed, small differences in the activations entail small differences in the loss.

The inequality in Eq. (6) holds under the condition that the sign of activations agree. Therefore, we want to group frames based on the sign similarity of their activations.

Figure 2. We adopt 2D ResNet-50 with TSM [28] a backbone. The input batch size is n with t frames. We cluster the activations of the first block of size (nt, c, h, w) which groups t frames into g clusters and outputs new activations of size (ng, c, h, w) , as input to the next network blocks. Our full video method efficiently utilizes all frames and is end-to-end trainable.

Figure 3. An illustration of our two clustering algorithms. The numbers on the solid line are pair-wise Hamming distances and the solid line is the cumulative Hamming distance of frame f_1 to f_{10} . For $g=3$ clusters, the *cumulative clustering* groups frames by dividing the total cumulative distance on the y-axis into 3 equally distanced segments, as shown with the dashed lines resulting in the 3 clusters (f_1-f_4) , (f_5-f_7) and (f_8-f_{10}) . The *slope clustering* algorithm is based on the slope of the curve and here selects the top-2 largest slopes, as shown with the solid green lines, which results in the 3 clusters: (f_1-f_6) , (f_7) , (f_8-f_{10}) .

Figure 4. Hamming distances between similar frames and dissimilar frames across 4 blocks of ResNet. The frames are taken from a single Breakfast [25] video. We denote the frames that are similar to their neighbors with circles and the dissimilar ones with squares. Hamming distances are consistent across blocks.

3.3. Temporal clustering and aggregation

Using our proposal in Eq. (5) allows training on all video frames. We group frames based on the sign agreement of their activations. An efficient way to do this, is to binarize the activation values by using the sign function and compute a fast Hamming distance between binarized activations to determine which frames to group.

Consecutive frames in a video are more likely to be similar in appearance and are thus more likely to have similarly signed activations. Therefore, we explore two variants of a temporal clustering algorithm based on Hamming distances, where we allow a fixed number of clusters g to match the available memory. We employ the temporal order of video frames and calculate Hamming distances only with neighboring frames. Fig. 3 illustrates the two temporal clustering algorithms we consider here: *cumulative clustering* and *slope clustering*. We start by calculating the cumulative Hamming distance $C(x)$ for neighboring frames along the temporal order:

$$C_N(x) = \sum_{i=1}^{N-1} H(x_i, x_{i+1}), \quad (8)$$

where x_i is the binarized activation of frame i , $H(\cdot, \cdot)$ is the Hamming distance, and N is the total number of frames. For *cumulative clustering*, we divide the total cumulative Hamming distance, $C(x)$, into g even segments, where the cluster id for frame i is $\frac{C_i(x)}{C_N(x)}$. For the *slope clustering*, the boundaries of the segments are defined by the frame indexes corresponding to the top- g largest slopes where the cumulative distance increases the most.

For efficiency, we cluster early on in the network, and input to the subsequent layers only aggregated activations. We assume that the sign of the activations corresponding to two similar frames, approximately agree throughout the network. To validate this, we visualize in Fig. 4 the Hamming distance over activations corresponding to similar and dissimilar frames. The distances corresponding to similar frames remain consistent across different layers.

Putting everything together, we input a set of n videos into our TSM-based [28] network architecture. After the first block, we apply temporal clustering and average the activations within each cluster, giving rise to g activations per video. These aggregated activations are input to the subsequent blocks of the network. Our method efficiently utilizes all frames for training and it is end-to-end trainable, as the gradients propagate directly through the aggregated feature-maps. Fig. 2 depicts the outline of our method.

4. Experiments

We evaluate our method on the action recognition datasets Something-Something V1 & V2 [14], UCF-101 [37], HMDB51 [26] and Breakfast [25]. The consistent improvements show the effectiveness and generality of our method. We validate and analyze our method on a fully controlled Move4MNIST dataset we created. We also include ablation studies of the components of our method.

Datasets. Something-Something V1 [14] consists of 86k training videos and 11k validation videos belonging to 174 action categories. The second release V2 of Something-Something increase the number of videos to 220k. The UCF101 [37] dataset contains 101 action classes and 13,320 video clips. The HMDB51 [26] dataset is a varied collection of movies and web videos with 6,766 video clips from 51 action categories. Breakfast [25] has long videos of human cooking activities with 10 categories with 1,712 videos in total, with 1,357 for training and 335 for testing. Our fully controlled dataset Move4MNIST has four action classes $\{move\ up, move\ down, move\ left, move\ right\}$, and 1,800 videos for training and 600 videos for testing. Each video has 32 frames, with a digit from MNIST [27] moving on a UCF-101 video background. To obtain a per-frame ground truth of which frames are relevant we randomly inserted a consecutive chunk of UCF-101 background frames, black frames and frames with MNIST digits that are not part of the target classes. An example is shown in Fig. 7.

Training & Inference. Following the setting in TSM [28], our models are fine-tuned from Kinetics [23] pre-trained weights and we freeze the Batch Normalization [18] layers for HMDB51 [26] and UCF101 [37] datasets. For other datasets, our models are fine-tuned from ImageNet [5] pre-trained weights. To optimize the GPU we train with a fixed number of frames per batch. If the video has less frames, we pad it repeatedly with the last video frame. We compare and cluster the activations of all the frames in each video, and get g average activations for each video, from the first block of our model. We set the number of clusters to $g = \{8, 16\}$ to align with the sub-sampling methods using 8 or 16 frames. During testing, we follow the setting of TSM and sample one clip per video and use the full resolution image with the shorter side 256.

Backbone architecture. For a fair comparison with the state-of-the-art, we evaluate our method on the TSM [28] backbone replying on the ResNet-50 [16] architecture. We use TSM with a ResNet-18 as the backbone for the experiments on our toy dataset Move4MNIST and for model analysis on the Breakfast dataset.

4.1. Are more frames better?

To make it computationally possible to use all individual frames we evaluate on the fully controlled Move4MNIST to test if using more frames during training is better than sub-sampling. We use here the ResNet-18 [16] backbone pretrained on ImageNet [5] and compare with TSM [28]. We evaluate slope clustering and cumulative clustering, and a cluster-free uniform grouping of evenly distributed segments and then aggregating them (*Ours-uniform*).

Model	#Frames	#Clusters	FLOPs /Video	Runtime Mem./Video	Top-1
TSM	8	-	14.56G	1.04GB	90.13 \pm 0.38
TSM	16	-	29.12G	1.72GB	93.78 \pm 0.33
TSM	all	-	58.24G	3.15GB	98.83 \pm 0.16
Ours-uniform	all	8	28.61G	1.56GB	90.25 \pm 0.28
Ours-slope	all	8	28.61G	1.56GB	93.33 \pm 0.19
Ours-cumulative	all	8	28.61G	1.56GB	94.08 \pm 0.25
Ours-uniform	all	16	38.51G	1.79GB	92.73 \pm 0.25
Ours-slope	all	16	38.51G	1.79GB	94.06 \pm 0.18
Ours-cumulative	all	16	38.51G	1.79GB	95.27 \pm 0.21

Table 1. Training with all frames gives best accuracy. Our method with slope or cumulative clustering outperforms the uniform grouping of evenly distributed segments and frame sub-sampling. Our method has less FLOPs and runtime memory usage than TSM training with all frames.

Table 1 shows that TSM trained on all the 32 frames of a video in Move4MNIST significantly outperforms TSM trained on 8 and 16 sub-sampled frames. The uniform grouping of evenly distributed segments does not much better than random sub-sampling, and uniform grouping performs worse than random sub-sampling when the frame and cluster numbers increased from 8 to 16. This can be explained since the videos in the Move4MNIST contain black frames, UCF-101 background frames, and frames containing another digits at random positions, which can make sub-sampling miss frames related to the task and evenly distributed segments group frames wrongly. Both our clustering approaches with 8 and 16 clusters do better than evenly distributed segments or sub-sampling with 8 or 16 frames as they can adapt to the content and dynamically choose which frames to group. In addition, our method has significantly reduced FLOPs and runtime memory when compared to the baseline on all frames.

4.2. Do similar frames have similar gradients?

In this experiment, we evaluate our assumption that similar frame activations have similar gradients. The activa-

Figure 5. An illustration of activation distance versus gradient distance for frames from three videos in Move4MNIST dataset. For frames that are similar with respect to recognizing the action, the activation distance and the gradients distance between them have a nearly linear relation for both the Euclidean distance and the Hamming distance. Our assumption that frames having similar activations with respect to the task have similar gradients is validated.

tions and gradients are taken from the 1st block of ResNet-18. We show the Euclidean and the Hamming activation distance versus the gradient Euclidean distance between all $32 \times 31/2 = 496$ frame pairs for three videos in Move4MNIST in Fig. 5. For both the Euclidean distance and the Hamming distance the relation between activations and gradients is close to linear. It validates our assumption that frames having similar activations with respect to the task have similar gradients.

We compare the ground truth gradients when training truly on all frames to our efficient approximation. We use 16 clusters and compare our approximate gradients to the real gradients which are from 3rd block of ResNet-18 for a video in Move4MNIST. We compare the results of our method with cumulative clustering, slope clustering and uniform grouping. Results in Fig. 6 show that compared to uniform grouping, cumulative clustering and slope clustering give smaller Euclidean distance between the single gradient from each cluster and the sum of gradients of frames in the corresponding cluster. And cumulative clustering gives even smaller gradients differences than slope clustering. In other words, it means that our method with cumulative clustering (the right hand side of Eq. (1)) approximates the standard gradients calculation (the left hand side of Eq. (1)) in the network with a small difference.

4.3. Analyzing model properties

We evaluate the clustering methods, the number of clusters, and the training time efficiency on Breakfast and Move4MNIST with a ResNet-18 backbone.

Different temporal clustering methods. We compare slope clustering, cumulative clustering, and uniform grouping where the videos are split into equal segments. From Table 2, cumulative clustering outperforms slope clustering, while uniform grouping has the lowest top-1 accuracy. This is because equal temporal grouping merges non similar frames together leading to linear approximations of non-

Figure 6. Comparing the Euclidean distance between gradients of the ground truth of truly using all frames to our efficient approximation per cluster for cumulative clustering, slope clustering and uniform grouping on Move4MNIST. Compared to uniform grouping and slope clustering, cumulative clustering results in smaller gradients difference and thus a better approximation.

Model	#Frames	#Clusters	Tr. sec/epoch	Top-1
TSM	8	-	97.6	59.1
TSM	16	-	113.7	61.4
Ours-uniform	all	8	100.1	58.3
Ours-slope	all	8	99.6	60.7
Ours-cumulative	all	8	101.3	62.0
Ours-uniform	all	16	114.0	60.2
Ours-slope	all	16	114.5	63.7
Ours-cumulative	all	16	115.2	64.4

Table 2. With 8 and 16 clusters we consistently outperform TSM with 8 and 16 frames for comparable training time on the Breakfast dataset.

linear information and incorrect network updates, resulting in a low action recognition accuracy. A similar trend is also visible on the Move4MNIST dataset in Table 1.

Figure 7. Temporal clustering results for a video in Move4MNIST. Cumulative temporal clustering groups frames more accurately than slope temporal clustering.

Figure 8. Cumulative temporal clustering results over epochs for six videos in the Breakfast dataset. Each cluster is shown in a different color. Clusters contains segments with different lengths. Our cumulative temporal clustering groups frames with similar activations together. The cluster lengths change according to the changes in the frame activations during training.

In Fig. 7, we show the temporal clustering results for a small number of frames of a Move4MNIST video. Cumulative clustering correctly groups similar frames together, while slope clustering groups moving zero frames and black frames together.

Number of clusters. We conduct experiments using 8 and 16 clusters for our method, which follows the protocol of TSM with 8 and 16 frames for training. Table 2 shows that using 16 clusters consistently outperforms using 8 clusters for all clustering methods. A larger number of clusters improves accuracy. In the extreme case, the cluster numbers equal the number of frames in a video, which is equivalent with using all frames for training. From the table we can also see that our cumulative temporal clustering implementation improves the top-1 accuracy by 2.9% and 3.0%, separately for 8 clusters and 16 clusters comparing to TSM with 8 and 16 frames.

To show that our cumulative temporal clustering algorithm is different from the naive uniform grouping, we vi-

Model	Backbone	#Frames	#Clusters	Top-1
ResNet-152[17]	ResNet152	64	-	41.1%
ActionVLAD [17]	ResNet152	64	-	55.5%
VideoGraph [17]	ResNet152	64	-	59.1%
TSM [28] (our impl.)	ResNet50	16	-	72.1%
Ours-slope	ResNet50	all	16	74.9%
Ours-cumulative	ResNet50	all	16	76.6%

Table 3. Our method using either slope temporal clustering or cumulative temporal clustering compared to existing works on the Breakfast dataset. Our proposal outperforms TSM, and significantly exceeds in top-1 accuracy methods using the deeper backbone architecture, ResNet-152. By using all frames our method has an advantage on long-term video action recognition.

sualize the 8 clusters obtained from cumulative temporal clustering for six videos over different epochs in the Breakfast dataset in Fig. 8. Different videos have different segment lengths in the cumulative temporal clustering, which takes the similarity of frame activations into consideration. In Fig. 8, we also show that the cluster length changes over epochs during training, since the activations change during training.

Efficiency of training time. Table 2 gives the training time per epoch for all the models. Our method with 8 clusters and 16 clusters only has an increase of 3.7 seconds and 1.5 seconds in training time per epoch, when compared to TSM with 8 frames and 16 frames. The results show that our method is efficient during training time, while using all video frames.

4.4. Comparison with the state-of-the-art

We compare our method with the state-of-the-art on Something-Something V1&V2, Breakfast, UCF-101 and HMDB51. All methods use ResNet-50 pre-trained on ImageNet as a backbone, unless specified otherwise.

Comparison on the Breakfast dataset. We compare our method with existing work on the Breakfast dataset, which contains long action videos. Our method using either slope temporal clustering or cumulative temporal clustering largely outperforms the three methods using ResNet-152 as a backbone, in Table 3. Compared to TSM using 16 sub-sampled frames, our method improves the top-1 accuracy by 2.8% and 4.5% with slope temporal clustering and cumulative temporal clustering, respectively. Methods using sub-sampling can easily miss important frames for the recognition task on long action videos. Our method has an advantage on the long videos for action recognition, by efficiently utilizing all the frames.

Comparison on the Something-Something dataset. In Table 5, we list the results of our method compared to other methods on the Something-Something V1&V2 datasets. We achieve state-of-the-art performance on both V1 and

Model	Backbone	Pre-train	#Frames	#Clusters	Top-1 UCF-101	Top-1 HMDB51
TSM [28] (our impl.)	ResNet50	Kinetics	1	-	91.2%	65.1%
TSN [28]	ResNet50	Kinetics	8	-	91.7%	64.7%
SI+DI+OF+DOF [1]	ResNeXt50	Imagenet	dynamic images	-	95.0%	71.5%
TSM [28]	ResNet50	Kinetics	8	-	95.9%	73.5%
STM [21]	ResNet50	ImageNet+Kinetics	16	-	96.2%	72.2%
Ours-slope	TSM-ResNet50	Kinetics	all	8	96.2%	73.3%
Ours-cumulative	TSM-ResNet50	Kinetics	all	8	96.4%	73.4%

Table 4. Top-1 accuracy on UCF-101 and HMDB51. Our method performs only slightly better than the state-of-the-art on the scene-related datasets UCF-101 and HMDB51. These datasets do not have much frame diversity per video, thus, the improvement of our method over sampling methods is limited.

Model	#Frames	#Clusters	Top-1 V1	Top-1 V2
TSN [28]	8	-	19.7%	30.0%
TRN-Multiscale [28]	8	-	38.9%	48.8%
TSM [28]	8	-	45.6%	59.1%
TSM [21]	16	-	47.2%	63.4%
STM [21]	8	-	49.2%	62.3%
STM [21]	16	-	50.7%	64.2%
Ours-slope	all	8	46.7%	60.2%
Ours-cumulative	all	8	49.5%	62.7%
Ours-cumulative	all	16	51.4%	65.1%

Table 5. Top-1 accuracy on Something-Something V1 and V2 datasets. Our method using cumulative temporal clustering outperforms the state-of-the-art methods on both Something-Something V1 and V2. Our method achieves limited accuracy improvement for shorter videos.

V2, with outperforming STM of 8 frames by 0.3% and 0.4% for V1 and V2, and STM of 16 frames by 0.7% and 0.9% for V1 and V2 respectively. Comparing to TSM, we significantly improve the top-1 accuracy of 8 frames by 3.9% and 3.5%, and the top-1 accuracy of 16 frames by 4.2% and 1.7% for the V1 and V2 datasets. Although the Something-Something dataset is characterized by temporal variations, the video clips are short compared to the Breakfast dataset. The methods using frame sampling heuristics can capture the main movement in videos. Therefore, our accuracy improvement is not as pronounced as for the Breakfast dataset.

Comparison on the UCF-101 and HMDB51 datasets. We train with 8 clusters and evaluate over three splits and report averaged results in Table 4. Our performance is on par with state-of-the-art methods on both datasets. The UCF-101 and HMDB51 have a scene-bias, where motion plays a limited role and just a few number of frames –or even a single frame– is sufficient. Thus, methods relying on sampling heuristics can correctly classify the actions and our method using all frames is not expected to improve results. To test this, we show results with a single frame in Table 4 which shows that TSM with 1 frame achieves comparable accuracy to TSN with 8 frames on UCF-101 and outperforms TSN with 8 frames on HMDB51. For scene-biased datasets, using all frames does not bring accuracy benefits.

5. Conclusion

We propose an efficient method for training action recognition deep networks without relying on sampling heuristics. Our work offers a solution to using all video frames during training based on the assumption that similar frames have similar gradients, leading to similar parameter updates. To this end, we efficiently find frames that are similar with respect to the classification task, by using a cumulative temporal clustering algorithm based on Hamming distances. The clustering based on Hamming distances enforces that activations in a cluster agree in signs, which is a requirement entailed by our assumption that we can approximate the gradients of multiple frames with a single gradient of an aggregated frame. We accumulate the activations within each cluster to create new representations used to classify the actions. Our proposed method shows competitive results on large datasets when compared to existing work.

Despite our state of the art results, we identify several limitations. One limitation is that the number of clusters is fixed and thus not well-suited for inhomogeneous videos with more semantic (shot) changes than clusters. This could create a dependency for action proposals or other approaches to pre-segment a video in homogeneous segments which somewhat counters the philosophy of using full video action recognition. Another limitation is that for grouping frames the only non-linearity we consider is the activation function and do not use the non-linearity in the loss. This limitation seems insurmountable, as memory constraints prevent us to store all frame activations for when the loss is computed. Nevertheless, with our current results and analysis, we make a first move for action recognition to go full video.

Acknowledgments This work is part of the research program Efficient Deep Learning (EDL), which is (partly) financed by the Dutch Research Council (NWO).

References

- [1] H. Bilen, B. Fernando, E. Gavves, and A. Vedaldi. Action recognition with dynamic image networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(12):2799–2813, 2018. [2](#), [8](#)
- [2] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6299–6308, 2017. [1](#), [2](#)
- [3] Bo Chang, Lili Meng, Eldad Haber, Lars Ruthotto, David Begert, and Elliot Holtham. Reversible architectures for arbitrarily deep residual neural networks. In *AAAI*, 2018. [1](#)
- [4] Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear memory cost. *arXiv preprint arXiv:1604.06174*, 2016. [1](#)
- [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255, 2009. [1](#), [5](#)
- [6] Tran Du, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. C3d: Generic features for video analysis. *Corr*, 2(8), 2014. [2](#)
- [7] 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. [2](#)
- [8] Christoph Feichtenhofer, Axel Pinz, and Andrew Zisserman. Convolutional two-stream network fusion for video action recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1933–1941, 2016. [2](#)
- [9] 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. [2](#)
- [10] Harshala Gammulle, Simon Denman, Sridha Sridharan, and Clinton Fookes. Two stream lstm: A deep fusion framework for human action recognition. In *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 177–186. IEEE, 2017. [2](#)
- [11] Rohit Girdhar, Deva Ramanan, Abhinav Gupta, Josef Sivic, and Bryan Russell. Actionvlad: Learning spatio-temporal aggregation for action classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 971–980, 2017. [2](#)
- [12] Negar Goli and Tor M Aamodt. Resprop: Reuse sparsified backpropagation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1548–1558, 2020. [2](#)
- [13] Aidan N Gomez, Mengye Ren, Raquel Urtasun, and Roger B Grosse. The reversible residual network: Backpropagation without storing activations. In *Advances in neural information processing systems*, pages 2214–2224, 2017. [1](#), [2](#)
- [14] 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*, volume 1, page 5, 2017. [5](#)
- [15] Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet? In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 6546–6555, 2018. [2](#)
- [16] 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*, pages 770–778, 2016. [5](#)
- [17] Nouredien Hussein, Efstratios Gavves, and Arnold WM Smeulders. Videograph: Recognizing minutes-long human activities in videos. *ICCV 2019, Workshop on Scene Graph Representation and Learning*, 2019. [7](#)
- [18] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *Proceedings of the 32nd International Conference on Machine Learning*, 2015. [5](#)
- [19] Mihir Jain, Jan C Van Gemert, and Cees GM Snoek. What do 15,000 object categories tell us about classifying and localizing actions? In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 46–55, 2015. [2](#)
- [20] Shuiwang Ji, Wei Xu, Ming Yang, and Kai Yu. 3d convolutional neural networks for human action recognition. *IEEE transactions on pattern analysis and machine intelligence*, 35(1):221–231, 2012. [2](#)
- [21] Boyuan Jiang, MengMeng Wang, Weihao Gan, Wei Wu, and Junjie Yan. Stm: Spatiotemporal and motion encoding for action recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2000–2009, 2019. [8](#)
- [22] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *CVPR*, 2014. [2](#)
- [23] 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. *CoRR*, 2017. [5](#)
- [24] 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](#)
- [25] H. Kuehne, A. B. Arslan, and T. Serre. The language of actions: Recovering the syntax and semantics of goal-directed human activities. In *Proceedings of Computer Vision and Pattern Recognition Conference (CVPR)*, 2014. [4](#), [5](#)
- [26] Hilde Kuehne, Hueihan Jhuang, E. Garrote, T. Poggio, and Thomas Serre. Hmdb: A large video database for human motion recognition. *2011 International Conference on Computer Vision*, pages 2556–2563, 2011. [5](#)
- [27] Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. 2010. [5](#)
- [28] 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](#), [4](#), [5](#), [7](#), [8](#)

- [29] Mateusz Malinowski, Grzegorz Swirszcz, Joao Carreira, and Viorica Patraucean. Sideways: Depth-parallel training of video models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11834–11843, 2020. 2
- [30] Deniz Oktay, Nick McGreivy, Joshua Aduol, Alex Beatson, and Ryan P Adams. Randomized automatic differentiation. *CoRR*, 2020. 2
- [31] Jian Ren, Xiaohui Shen, Zhe Lin, and Radomir Mech. Best frame selection in a short video. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 3212–3221, 2020. 1, 2
- [32] Konrad Schindler and Luc Van Gool. Action snippets: How many frames does human action recognition require? In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, 2008. 2
- [33] Fadime Sener, Dipika Singhania, and Angela Yao. Temporal aggregate representations for long-range video understanding. In *European Conference on Computer Vision*, pages 154–171, 2020. 2
- [34] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In *Advances in neural information processing systems*, pages 568–576, 2014. 2
- [35] Bharat Singh, Tim K Marks, Michael Jones, Oncel Tuzel, and Ming Shao. A multi-stream bi-directional recurrent neural network for fine-grained action detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1961–1970, 2016. 2
- [36] Sibo Song, Ngai-Man Cheung, Vijay Chandrasekhar, and Bappaditya Mandal. Deep adaptive temporal pooling for activity recognition. In *Proceedings of the 26th ACM international conference on Multimedia*, pages 1829–1837, 2018. 2
- [37] Khurram Soomro, Amir Roshan Zamir, and M Shah. A dataset of 101 human action classes from videos in the wild. *Center for Research in Computer Vision*, 2(11), 2012. 5
- [38] Swathikiran Sudhakaran, Sergio Escalera, and Oswald Lanz. Gate-shift networks for video action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1102–1111, 2020. 1, 2
- [39] Amin Ullah, Jamil Ahmad, Khan Muhammad, Muhammad Sajjad, and Sung Wook Baik. Action recognition in video sequences using deep bi-directional lstm with cnn features. *IEEE Access*, 6:1155–1166, 2017. 2
- [40] Vivek Veeriah, Naifan Zhuang, and Guo-Jun Qi. Differential recurrent neural networks for action recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 4041–4049, 2015. 2
- [41] 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 *European conference on computer vision*, pages 20–36, 2016. 1, 2
- [42] Yue Wang, Ziyu Jiang, Xiaohan Chen, Pengfei Xu, Yang Zhao, Yingyan Lin, and Zhangyang Wang. E2-train: Training state-of-the-art cnns with over 80% energy savings. In *Advances in Neural Information Processing Systems*, pages 5138–5150, 2019. 2
- [43] Simon Wiedemann, Temesgen Mehari, Kevin Kepp, and Wojciech Samek. Dithered backprop: A sparse and quantized backpropagation algorithm for more efficient deep neural network training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 720–721, 2020. 2
- [44] Chao-Yuan Wu, Ross Girshick, Kaiming He, Christoph Feichtenhofer, and Philipp Krahenbuhl. A multigrid method for efficiently training video models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 153–162, 2020. 2
- [45] 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 International Conference on Computer Vision*, pages 6222–6231, 2019. 1, 2
- [46] 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. 2
- [47] 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/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 1, 2
- [48] 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
- [49] 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. 2
- [50] 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, 2