

DMC-Net: Generating Discriminative Motion Cues for Fast Compressed Video Action Recognition

Zheng Shou^{1,2}Xudong Lin²Yannis Kalantidis¹Laura Sevilla-Lara^{1,3}Marcus Rohrbach¹Shih-Fu Chang²Zhicheng Yan¹¹Facebook AI²Columbia University³University of Edinburgh

Abstract

Motion has shown to be useful for video understanding, where motion is typically represented by optical flow. However, computing flow from video frames is very time-consuming. Recent works directly leverage the motion vectors and residuals readily available in the compressed video to represent motion at no cost. While this avoids flow computation, it also hurts accuracy since the motion vector is noisy and has substantially reduced resolution, which makes it a less discriminative motion representation. To remedy these issues, we propose a lightweight generator network, which reduces noises in motion vectors and captures fine motion details, achieving a more Discriminative Motion Cue (DMC) representation. Since optical flow is a more accurate motion representation, we train the DMC generator to approximate flow using a reconstruction loss and an adversarial loss, jointly with the downstream action classification task. Extensive evaluations on three action recognition benchmarks (HMDB-51, UCF-101, and a subset of Kinetics) confirm the effectiveness of our method. Our full system, consisting of the generator and the classifier, is coined as DMC-Net which obtains high accuracy close to that of using flow and runs two orders of magnitude faster than using optical flow at inference time.

1. Introduction

Video is a rich source of visual content as it not only contains appearance information in individual frames, but also temporal motion information across consecutive frames. Previous work has shown that modeling motion is important to various video analysis tasks, such as action recognition [41, 50, 24], action localization [38, 37, 40, 6, 39, 26, 27] and video summarization [46, 30]. Currently, methods achieving state-of-the-art results usually follow the two-stream network framework [41, 5, 49], which consists of



This work was partially done when Zheng Shou interned at Facebook.

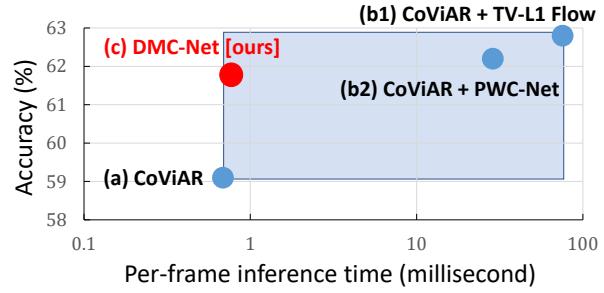


Figure 1: Comparing inference time and accuracy for different methods on HMDB-51. (a) Compressed video based method CoViAR [55] is very fast. (b) But in order to reach high accuracy, CoViAR has to follow two-stream networks to add the costly optical flow computation, either using TV-L1 [58] or PWC-Net [45]. (c) The proposed DMC-Net not only operates exclusively in the compressed domain, but also is able to achieve high accuracy while being two orders of magnitude faster than methods that use optical flow. The blue box denotes the improvement room from CoViAR to CoViAR + TV-L1 Flow; x-axis is in logarithmic scale.

two Convolutional Neural Networks (CNNs), one for the decoded RGB images and one for optical flow, as shown in Figure 2a. These networks can operate on either single frames (2D inputs) or clips (3D inputs) and may utilize 3D spatiotemporal convolutions [47, 49].

Extracting optical flow, however, is very slow and often dominates the overall processing time of video analysis tasks. Recent work [55, 61, 60] avoids optical flow computation by exploiting the motion information from compressed videos encoded by standards like MPEG-4 [25]. Such methods utilize the motion vectors and residuals already present in the compressed video to model motion. The recently proposed CoViAR [55] method, for example, contains three independent CNNs operating over three modalities in the compressed video, i.e. RGB image of I-frame (I), low-resolution Motion Vector (MV) and Residual

I: RGB of I-frame. MV: Motion Vector. R: Residual

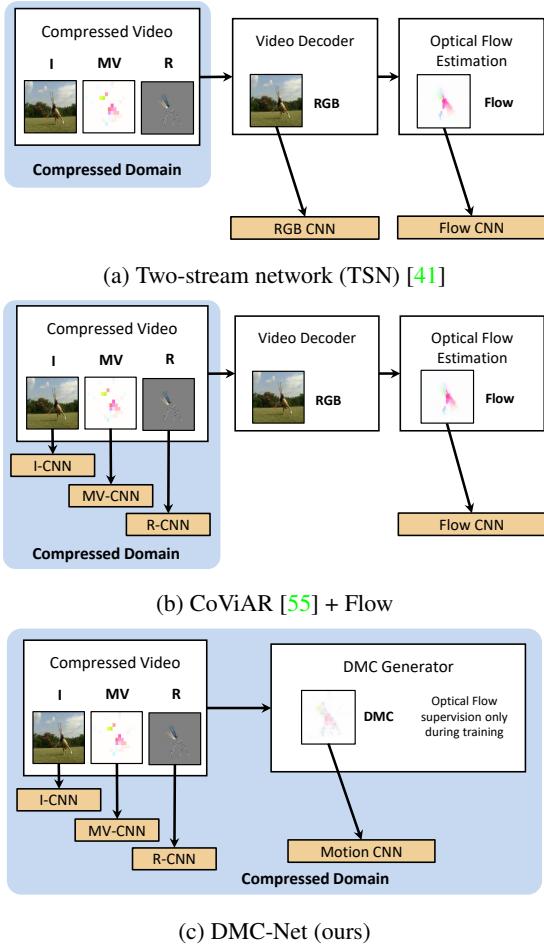


Figure 2: Illustrations of (a) the two-stream network [41], (b) the recent CoViAR [55] method that achieves high accuracy via fusing compressed video data and optical flow, and (c) our proposed DMC-Net. Unlike *CoViAR+Flow* that requires video decoding of RGB images and flow estimation, our DMC-Net operates exclusively in the compressed domain at inference time while using optical flow to learn to capture discriminative motion cues at training time.

(R). The predictions from individual CNNs are combined by late fusion. CoViAR runs extremely fast while modeling motion features (see Figure 2b). However, in order to achieve state-of-the-art accuracy, late fusion with optical flow is further needed (see Figure 1).

This performance gap is due to the motion vector being less informative and discriminative than flow. First, the spatial resolution of the motion vector is substantially reduced (*i.e.* 16x) during video encoding, and fine motion details, which are important to discriminate actions, are permanently lost. Second, employing two CNNs to process motion vectors and residuals separately ignores the strong

interaction between them. Because the residual is computed as the difference between the raw RGB image and its reference frame warped by the motion vector. The residual is often well-aligned with the boundary of moving object, which is more important than the motion at other locations for action recognition according to [35]. Jointly modeling motion vectors and residuals, which can be viewed as coarse-scale and fine-scale motion feature respectively, can exploit the encoded motion information more effectively.

To address those issues, we propose a novel approach to learn to generate a **Discriminative Motion Cue (DMC)** representation by refining the noisy and coarse motion vectors. We develop a lightweight DMC generator network that operates on stacked motion vectors and residuals. This generator requires training signals from different sources to capture discriminative motion cues and incorporate high-level recognition knowledge. In particular, since flow contains high resolution and accurate motion information, we encourage the generated DMC to resemble optical flow by a pixel-level reconstruction loss. We also use an adversarial loss [14] to approximate the distribution of optical flow. Finally, the DMC generator is also supervised by the downstream action recognition classifier in an end-to-end manner, allowing it to learn motion cues that are discriminative for recognition.

During inference, the DMC generator is extremely efficient with merely 0.23 GFLOPs, and takes only 0.106 ms per frame which is negligible compared with the time cost of using flow. In Figure 2c, we call our full model *DMC-Net*. Although optical flow is required during training, our method operates exclusively in the compressed domain at inference time and runs two orders of magnitude faster than methods using optical flow, as shown in Figure 1. Our contributions are summarized as follows:

- We propose DMC-Net, a novel and highly efficient framework that operates exclusively in the compressed video domain and is able to achieve high accuracy without requiring optical flow estimation.
- We design a lightweight generator network that can learn to predict discriminative motion cues by using optical flow as supervision and being trained jointly with action classifier. During inference, it runs two orders of magnitude faster than estimating flow.
- We extensively evaluate DMC-Net on 3 action recognition benchmarks, namely HMDB-51 [23], UCF-101 [42] and a subset of Kinetics [21], and demonstrate that it can significantly shorten the performance gap between state-of-the-art compressed video based methods with and without optical flow.

2. Related Work

Video Action Recognition. Advances in action recognition are largely driven by the success of 2D ConvNets in image recognition. The original Two-Stream Network [41] employs separate 2D ConvNets to process RGB frames and optical flow, and merges their predictions by late fusion. Distinct from image, video possesses temporal structure and motion information which are important for video analysis. This motivates researchers to model them more effectively, such as 3D ConvNets [47, 5], Temporal Segment Network (TSN) [52], dynamic image networks [2], and Non-Local Network [53]. Despite the enormous amount of effort on modeling motion via temporal convolution, 3D ConvNets can still achieve higher accuracy when fused with optical flow [5, 49], which is unfortunately expensive to compute.

Compressed Video Action Recognition. Recently, a number of approaches that utilize the information present in the compressed video domain have been proposed. In the pioneering works [60, 61], Zhang *et al.* replace the optical flow stream in two-stream methods by a motion vector stream, but it still needed to decode RGB image for P-frame and ignored other motion-encoding modalities in compressed videos such as the residual maps. More recently, the CoViAR method [55] proposed to exploit all data modalities in compressed videos, *i.e.* RGB I-frames, motion vectors and residuals to bypass RGB frame decoding. However, CoViAR fails to achieve performance comparable to that of two-stream methods, mainly due to the low-resolution of the motion vectors and the fact that motion vectors and residuals, although highly related, are processed by independent networks. We argue that, when properly exploited, the compressed video modalities have enough signal to allow us to capture more discriminative motion representation. We therefore explicitly learn such representation as opposed to relying on optical flow during inference.

Motion Representation and Optical Flow Estimation. Traditional optical flow estimation methods explicitly model the displacement at each pixel between successive frames [16, 57, 8, 3]. In the last years CNNs have successfully been trained to estimate the optical flow, including FlowNet [9, 18], SpyNet [34] and PWC-Net [45], and achieve low End-Point Error (EPE) on challenging benchmarks, such as MPI Sintel [4] and KITTI 2015 [31]. Im2Flow work [13] also shows optical flow can be hallucinated from still images. Recent work however, shows that accuracy of optical flow does not strongly correlate with accuracy of video recognition [36]. Thus, motion representation learning methods focus more on generating discriminative motion cues. Fan *et al.* [10] proposed to transform TV-L1 optical flow algorithm into a trainable sub-network, which can be jointly trained with downstream recognition network. Ng *et al.* [32] employs fully convolutional ResNet model to generate pixel-wise prediction of optical flow, and

can be jointly trained with recognition network. Unlike optical flow estimation methods, our method does not aim to reduce EPE error. Also different from all above methods of motion representation learning which take decoded RGB frames as input, our method refines motion vectors in the compressed domain, and requires much less model capacity to generate discriminative motion cues.

3. Approach

In this section, we present our approach for generating *Discriminative Motion Cues (DMC)* from compressed video. The overall framework of our proposed **DMC-Net** is illustrated in Figure 3. In Section 3.1, we introduce the basics of compressed video and the notations we use. Then we design the DMC generator network in Section 3.2. Finally we present the training objectives in Section 3.3 and discuss inference in Section 3.4.

3.1. Basics and Notations of Compressed Video

We follow CoViAR [55] and use MPEG-4 Part2 [25] encoded videos where every I-frame is followed by 11 consecutive P-frames. Three data modalities are readily available in MPEG-4 compressed video: (1) RGB image of I-frame (**I**); (2) Motion Vector (**MV**) records the displacement of each macroblock in a P-frame to its reference frame and typically a frame is divided into 16x16 macroblocks during video compression; (3) Residual (**R**) stores the RGB difference between a P-frame and its reference I-frame after motion compensation based on MV. For a frame of height H and width W , I and R have shape $(3, H, W)$ and MV has shape $(2, H, W)$. But note that MV has much lower resolution in effect because its values within the same macroblock are identical.

3.2. The Discriminative Motion Cue Generator

Input of the generator. Existing compressed video based methods directly feed motion vectors into a classifier to model motion information. This strategy is not effective in modeling motion due to the characteristics of MV: (1) MV is computed based on simple block matching, making MV noisy and (2) MV has substantially lower resolution, making MV lacking fine motion details. In order to specifically handle these characteristics of MV, we aim to design a lightweight generation network to reduce noise in MV and capture more fine motion details, outputting DMC as a more discriminative motion representation.

To accomplish this goal, MV alone may not be sufficient. According to [35], the motion nearby object boundary is more important than the motion at other locations for action recognition. We also notice R is often well-aligned with the boundary of moving objects. Moreover, R is strongly correlated with MV as it is computed as the difference between the original frame and its reference I-frame compensated

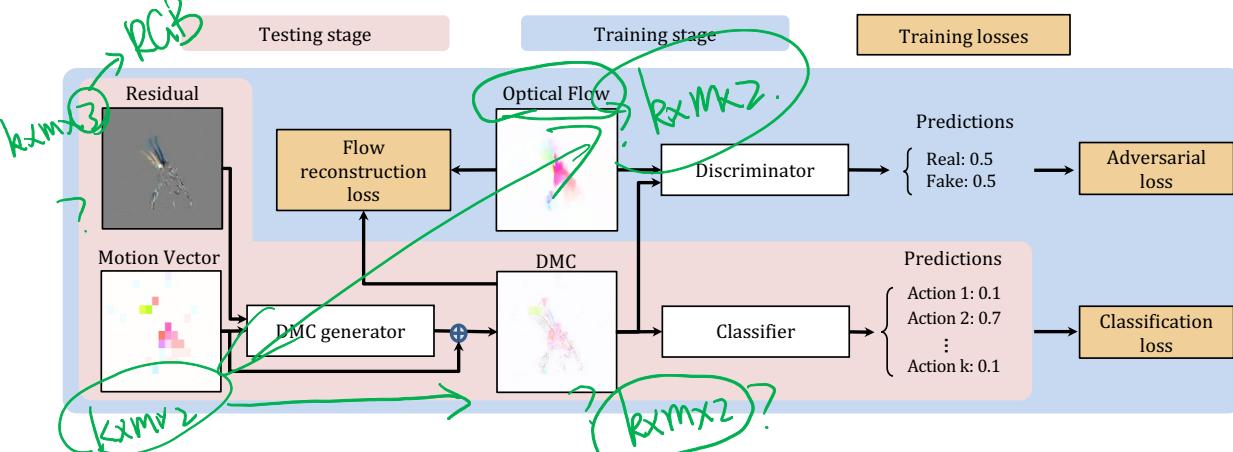


Figure 3: The framework of our Discriminative Motion Cue Network (DMC-Net). Given the stacked residual and motion vector as input, the DMC generator reduces noise in the motion vector and captures more fine motion details, outputting a more discriminative motion cue representation which is used by a small classification network to classify actions. In the training stage, we train the DMC generator and the action classifier jointly using three losses. In the test stage, only the modules highlighted in pink are used.

Network Architecture	GFLOPs
C3D [47]	38.5
Res3D-18 [48]	19.3
ResNet-152 [15]	11.3
ResNet-18 [15]	1.78
DMC generator (PWC-Net [45])	36.15
DMC generator [ours]	0.23

Table 1: Computational complexity of different networks. Input has height 224 and width 224.

Layer	Input size	Output size	Filter config
conv0	5, 224, 224	8, 224, 224	8, 3x3, 1, 1
conv1	13, 224, 224	8, 224, 224	8, 3x3, 1, 1
conv2	21, 224, 224	6, 224, 224	6, 3x3, 1, 1
conv3	27, 224, 224	4, 224, 224	4, 3x3, 1, 1
conv4	31, 224, 224	2, 224, 224	2, 3x3, 1, 1
conv5	33, 224, 224	2, 224, 224	2, 3x3, 1, 1

Table 2: The architecture of our Discriminative Motion Cue (DMC) generator network which takes stacked motion vector and residual as input. Input/output size follows the format of #channels, height, width. Filter configuration follows the format of #filters, kernel size, stride, padding.

using MV. Therefore, we propose to stack MV and R as input into the DMC generator, as shown in Figure 3. This allows utilizing the motion information in MV and R as well as the correlation between them, which cannot be modeled by separate CNNs as in the current compressed video works [55, 61, 60].

Generator network architecture. Quite a few deep generation networks have been proposed for optical flow estimation from RGB images. One of these works is PWC-Net [45], which achieves SoTA performance in terms of both End Point Error (EPE) and inference speed. We there-

fore choose to base our generator design principles on the ones used by PWC-Net. It is worth noting that PWC-Net takes decoded RGB frames as input unlike our proposed method operating only in the compressed domain.

Directly adopting the network architecture of the flow estimator network in PWC-Net for our DMC generator leads to high GFLOPs as indicated in Table 1. To achieve high efficiency, we have conducted detailed architecture search experimentally to reduce the number of filters in each convolutional layer of the flow estimator network in PWC-Net, achieving the balance between accuracy and complexity. Furthermore, since our goal is to refine MV, we propose to add a shortcut connection between the input MV and the output DMC, making the generator to directly predict the refinements which are added on MV to obtain DMC.

Table 2 shows the network architecture of our DMC generator: 6 convolutional layers are stacked sequentially with all convolutional layers densely connected [17]. Every convolutional filter has a 3x3 kernel with stride 1 and padding 1. Each convolutional layer except conv5 is followed by a Leaky ReLU [28] layer, where the negative slope is 0.1.

As shown in Table 1, our DMC generator only requires 0.63% GFLOPs used by the flow estimator in PWC-Net if it were adopted to implement our DMC generator. Also, Table 1 compares our DMC generator with other popular network architectures for video analysis including frame-level models (ResNet-18 and ResNet-152 [15]) and clip-level models (C3D [47] and Res3D [48]). We observe that the complexity of DMC generator is orders of magnitude smaller compared to that of other architectures, which makes it running much faster. In the supplementary material, we explored a strategy of using two consecutive networks to respectively rectify errors in MV and capture fine motion details while this did not achieve better accuracy.

3.3. Flow-guided, Discriminative Motion Cues

Compared to MV, optical flow exhibits more discriminative motion information because: (1) Unlike MV is computed using simple block matching, nowadays dense flow estimation is computed progressively from coarse scales to fine scales [58]. (2) Unlike MV is blocky and thus misses fine details, flow keeps the full resolution of the corresponding frame. Therefore we propose to guide the training of our DMC generator using optical flow. To this end, we have explored different ways and identified three effective training losses as shown in Figure 3 to be presented in the following: a flow reconstruction loss, an adversarial loss, and a downstream classification loss.

3.3.1 Optical Flow Reconstruction Loss

First, we minimize the per-pixel difference between the generated DMC and its corresponding optical flow. Following Im2Flow [13] which approximates flow from a single RGB image, we use the Mean Square Error (MSE) reconstruction loss \mathcal{L}_{mse} defined as:

$$\mathcal{L}_{\text{mse}} = \mathbb{E}_{\mathbf{x} \sim p} \|\mathcal{G}_{\text{DMC}}(\mathbf{x}) - \mathcal{G}_{\text{OF}}(\mathbf{x})\|_2^2, \quad (1)$$

where p denotes the set of P-frames in the training videos, \mathbb{E} stands for computing expectation, $\mathcal{G}_{\text{DMC}}(\mathbf{x})$ and $\mathcal{G}_{\text{OF}}(\mathbf{x})$ respectively denote the DMC and optical flow for the corresponding input frame \mathbf{x} sampled from p . Since only some regions of flow contain discriminative motion cues that are important for action recognition, in the supplementary material we have explored weighting the flow reconstruction loss to encourage attending to the salient regions of flow. But this strategy does not achieve better accuracy.

3.3.2 Adversarial Loss

As pointed out by previous works [29], the MSE loss implicitly assumes that the target data is drawn from a Gaussian distribution and therefore tends to generate smooth and blurry outputs. This in effect results in less sharp motion representations especially around boundaries, making the generated DMC less discriminative. Generative Adversarial Networks (GAN) [14] has been proposed to minimize the Jensen-Shannon divergence between the generative model and the true data distribution, making these two similar. Thus in order to help our DMC generator learn to approximate the distribution of optical flow data, we further introduce an adversarial loss. Note that unlike GAN which samples from random noise, adversarial loss samples from the input dataset, which already has large variability [29].

We relax the notational rigor and use $\mathcal{G}_{\text{OF}}(\mathbf{x})$ to refer to the optical flow corresponding to the frame \mathbf{x} , although for many optical flow algorithms the input would be a pair of frames.

Let our DMC generator \mathcal{G}_{DMC} be the **Generator** in the adversarial learning process. As shown in Figure 3, a **Discriminator** \mathcal{D} is introduced to compete with \mathcal{G}_{DMC} . \mathcal{D} is instantiated by a binary classification network that takes as input either **real** optical flow or **fake** samples generated via our DMC generator. Then \mathcal{D} outputs a two-dimensional vector that is passed through a softmax operation to obtain the probability $P_{\mathcal{D}}$ of the input being *Real*, i.e. flow versus *Fake*, i.e. DMC. \mathcal{G}_{DMC} and \mathcal{D} are trained in an alternating manner: \mathcal{G}_{DMC} is fixed when \mathcal{D} is being optimized, and vice versa.

During training \mathcal{D} , \mathcal{G}_{DMC} is fixed and is only used for inference. \mathcal{D} aims to classify the generated DMC as *Fake* and classify flow as *Real*. Thus the adversarial loss for training \mathcal{D} is:

$$\mathcal{L}_{\text{adv}}^{\mathcal{D}} = \mathbb{E}_{\mathbf{x} \sim p} [-\log P_{\mathcal{D}}(\text{Fake} | \mathcal{G}_{\text{DMC}}(\mathbf{x})) - \log P_{\mathcal{D}}(\text{Real} | \mathcal{G}_{\text{OF}}(\mathbf{x}))], \quad (2)$$

where p denotes the set of P-frames in the training set and $\mathcal{G}_{\text{DMC}}(\mathbf{x})$ and $\mathcal{G}_{\text{OF}}(\mathbf{x})$ respectively represent the DMC and optical flow for each input P-frame \mathbf{x} .

During training \mathcal{G}_{DMC} , \mathcal{D} is fixed. \mathcal{G}_{DMC} is encouraged to generate DMC that is similar and indistinguishable with flow. Thus the adversarial loss for training \mathcal{G}_{DMC} is:

$$\mathcal{L}_{\text{adv}}^{\mathcal{G}} = \mathbb{E}_{\mathbf{x} \sim p} [-\log P_{\mathcal{D}}(\text{Real} | \mathcal{G}_{\text{DMC}}(\mathbf{x}))], \quad (3)$$

which can be trained jointly with the other losses designed for training the DMC generator in an end-to-end fashion, as presented in Section 3.3.3.

Through the adversarial training process, \mathcal{G}_{DMC} learns to approximate the distribution of flow data, generating DMC with more fine details and thus being more similar to flow. Those fine details usually capture discriminative motion cues and are thus important for action recognition. We present details of the discriminator network architecture in the supplementary material.

3.3.3 The Full Training Objective Function

Semantic classification loss. As our final goal is to create motion representation that is discriminative with respect to the downstream action recognition task, it is important to train the generator jointly with the follow-up action classifier. We employ the softmax loss as our action classification loss, denoted as \mathcal{L}_{cls} .

The full training objective. Our whole model is trained with the aforementioned losses putting together in an end-to-end manner. The training process follows the alternating training procedure stated in Section 3.3.2. During training the discriminator, \mathcal{D} is trained while the DMC generator \mathcal{G}_{DMC} and the downstream action classifier are fixed. The full training objective is to minimize the adversarial loss

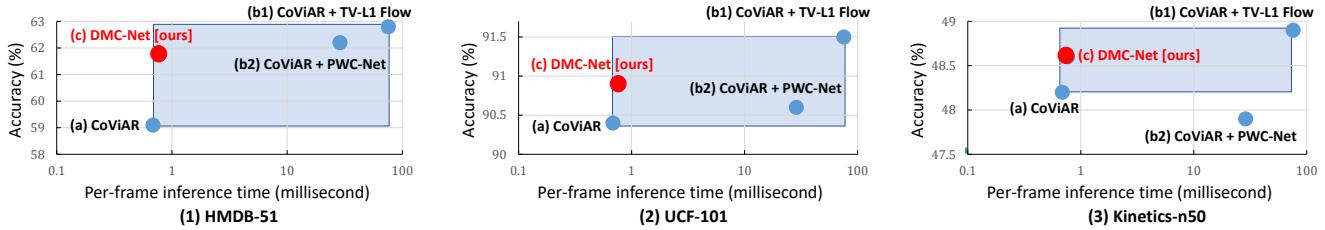


Figure 4: Accuracy vs. speed on 3 benchmarks. Results on UCF-101 and HMDB-51 are averaged over 3 splits. (b1) and (b2) use ResNet-18 to classify flow and (c) also uses ResNet-18 to classify DMC. The proposed *DMC-Net* not only operates exclusively in the compressed domain, but also is able to achieve higher accuracy than (a) while being two orders of magnitude faster than methods that use optical flow. The blue area indicates the improvement room from (a) to (b1).

$\mathcal{L}_{\text{adv}}^D$ in Equation 2. During training the generator \mathcal{G}_{DMC} , \mathcal{D} is fixed while the DMC generator \mathcal{G}_{DMC} and the downstream action classifier are trained jointly with the following full training objective to be minimized:

$$\mathcal{L}_{\text{cls}} + \alpha \cdot \mathcal{L}_{\text{mse}} + \lambda \cdot \mathcal{L}_{\text{adv}}^G, \quad (4)$$

where \mathcal{L}_{mse} is given by Equation 1, $\mathcal{L}_{\text{adv}}^G$ is given by Equation 3, and α, λ are balancing weights.

3.4. Inference

As shown in Figure 3, despite having three losses jointly trained end-to-end, our DMC-Net is actually quite efficient during inference: basically first the generator outputs DMC and then the generated DMC is fed into the classification network to make action class prediction. We compare our inference speed with other methods in Section 4.4.

4. Experiments

In this section, we first detail our experimental setup, present quantitative analysis of our model, and finally compare with state-of-the-art methods.

4.1. Datasets and Evaluation

UCF-101 [43]. This dataset contains 13,320 videos from 101 action categories, along with 3 public train/test splits.

HMDB-51 [23]. This dataset contains 6,766 videos from 51 action categories, along with 3 public train/test splits.

Kinetics-n50. From the original Kinetics-400 dataset [5], we construct a subset referred as **Kinetics-n50** in this paper. We keep all 400 categories. For each class, we randomly sample 30 videos from the original training set as our training videos and randomly sample 20 videos from the original validation set as our testing videos. We evaluate on the full set in the supplementary material.

Evaluation protocol. All videos in the above datasets have single action label out of multiple classes. Thus we evaluate top-1 video-level class prediction accuracy.

4.2. Implementation Details

Training. For I, MV, and R, we follow the exactly same setting as used in CoViAR [55]. Note that I employs ResNet-

152 classifier; MV and R use ResNet-18 classifier. To ensure efficiency, DMC-Net also uses ResNet-18 to classify DMC in the whole paper unless we explicitly point out. To allow apple-to-apple comparisons between DMC and flow, we also choose frame-level ResNet-18 classifier as the flow CNN shown in Figure 2b. TV-L1 [57] is used for extracting optical flow to guide the training of our DMC-Net. All videos are resized to 340×256 . Random cropping of 224×224 and random flipping are used for data augmentation. More details are in the supplementary material.

Testing. For I, MV, and R, we follow the exactly same setting as in CoViAR [55]: 25 frames are uniformly sampled for each video; each sampled frame has 5 crops augmented with flipping; all 250 ($25 \times 2 \times 5$) score predictions are averaged to obtain one video-level prediction. For DMC, we following the same setting except that we do not use cropping and flipping, which shows comparable accuracy but requires less computations. Finally, we follow CoViAR [55] to obtain the final prediction via fusing prediction scores from all modalities (*i.e.* I, MV, R, and DMC).

4.3. Model Analysis

How much gain DMC-Net can improve over CoViAR?

Figure 4 reports accuracy on all three datasets. **CoViAR + TV-L1** and **CoViAR + PWC-Net** follow two-stream methods to include an optical flow stream computed by TV-L1 [58] and PWC-Net [45] respectively. **CoViAR + TV-L1** can be regard as our upper bound for improving accuracy because TV-L1 flow is used to guide the training of **DMC-Net**. By only introducing a lightweight DMC generator, our **DMC-Net** significantly improves the accuracy of **CoViAR** to approach **CoViAR + Flow**. Figure 5 shows that the generated DMC has less noisy signals such as those in the background area and DMC captures fine and sharp details of motion boundary, leading to the accuracy gain over **CoViAR**.

How effectiveness is each proposed loss? On HMDB-51, when only using the classification loss, the accuracy of DMC-Net is 60.5%; when using the classification loss and the flow reconstruction loss, the accuracy is improved to 61.5%; when further including the adversarial training loss, DMC-Net eventually achieves 61.8% accuracy. As in-

	Two-Stream Method (RGB+Flow)		Compressed Video Based Methods		Generator Time (ms) / FPS	Generator + Cls. Time (ms) / FPS
	BN-Inception	ResNet152	CoViAR	DMC-Net [ours]		
Time (ms)	Preprocess	75.0	75.0	0.46	0.46	1449.2 / 0.7
	CNN (S)	1.6	7.5	0.59	0.89	220.8 / 4.5
	Total (S)	76.6	82.5	1.05	1.35	83.3 / 12.0
	CNN (C)	0.9	4.0	0.22	0.30	28.6 / 35.0
FPS	Total (C)	75.9	79.0	0.68	0.76	28.8 / 34.8
	CNN (C)	1111.1	250.0	4545.4	3333.3	0.1 / 9433.9
Total (C)						
(b) DMC-Net vs. flow estimation methods						

(a) DMC-Net vs. Two-stream methods and CoViAR

Table 3: Comparisons of per-frame inference speed. **(a)** Comparing our DMC-Net to the two-stream methods [19, 15] and the CoViAR method [55]. We consider two scenarios of forwarding multiple CNNs sequentially and concurrently, denoted by **S** and **C** respectively. We measure CoViAR’s CNN forwarding time using our own implementation as mentioned in Section 4.4 and numbers are comparable to those reported in [55]. **(b)** Comparing our DMC-Net to deep network based optical flow estimation and motion representation learning methods, whose numbers are quoted from [10]. CNNs in DMC-Net are forwarded concurrently. All networks have batch size set to 1. For the classifier (denoted as **Cl.**), all methods use ResNet-18.

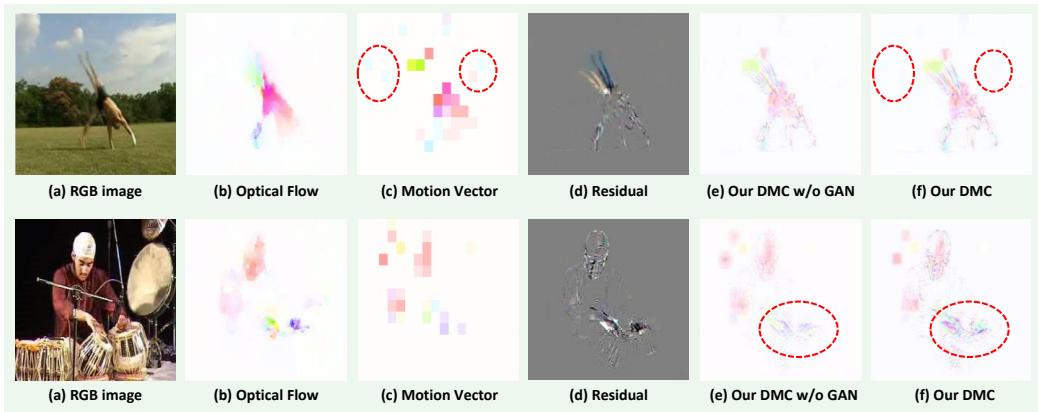


Figure 5: A Cartwheel example (top) and a PlayingTabla (bottom) example. All images in one row correspond to the same frame. For the Cartwheel example, these noisy blocks in the background (highlighted by two red circles) are reduced in our DMC. For the PlayingTabla example, our DMC exhibits sharper and more discriminative motion cues around hands (highlighted by the red circle) than our DMC w/o adversarial loss during training. Better viewed in color.

dicated by previous literature [20], using an adversarial loss without a reconstruction loss often introduces artifacts.

4.4. Inference Speed

Following [55], we measure the average per-frame running time, which consists of the time for data pre-processing and the time for CNN forward pass. For the CNN forward pass, both the scenarios of forwarding multiple CNNs sequentially and concurrently are considered. Detailed results can be found in Table 3 (a). Results of **two-stream methods** are quoted from [55]. Due to the need of decoding compressed video into RGB frames and then computing optical flow, its pre-process takes much longer time than compressed video based methods. **DMC-Net** accepts the same inputs as **CoViAR** and thus **CoViAR** and **DMC-Net** have the same pre-processing time. As for the CNN forward-

ing time of compressed video based methods, we measure **CoViAR** and **DMC-Net** using the exactly same implementation as stated in Section 4.2 and the same experimental setup: we use one NVIDIA GeForce GTX 1080 Ti and set the batch size of each CNN to 1 while in practice the speed can be further improved to utilize larger batch size. Despite adding little computational overhead on CoViAR, DMC-Net is still significantly faster than the conventional **two-stream methods**.

Deepflow [54], **Flownet** [18] and **PWC-Net** [45] have been proposed to accelerate optical flow estimation by using deep networks. **TVNet** [10] was proposed to generate even better motion representation than flow with fast speed. Those estimated flow or generated motion representation can replace optical flow used in two-stream methods to go through a CNN for classification. We combine these meth-

	HMDB-51	UCF-101
Compressed video based methods		
EMV-CNN [60]	51.2 (split1)	86.4
DTMV-CNN [61]	55.3	87.5
CoViAR [55]	59.1	90.4
DMC-Net (ResNet-18) [ours]	62.8	90.9
DMC-Net (I3D) [ours]	71.8	92.3
Decoded video based methods (RGB only)		
<i>Frame-level classification</i>		
ResNet-50 [15]	48.9	82.3
ResNet-152 [15]	46.7	83.4
<i>Motion representation learning</i>		
ActionFlowNet (2-frames) [32]	42.6	71.0
ActionFlowNet [32]	56.4	83.9
PWC-Net (ResNet-18) + CoViAR [45]	62.2	90.6
TVNet [10]	71.0	94.5
<i>Spatio-temporal modeling</i>		
C3D [47]	51.6	82.3
Res3D [48]	54.9	85.8
ARTNet [51]	70.9	94.3
MF-Net [7]	74.6	96.0
S3D [56]	75.9	96.8
I3D RGB [5]	74.8	95.6
I3D RGB + DMC-Net (I3D) [ours]	77.8	96.5
Decoded video based methods (RGB + Flow)		
Two-stream [41]	59.4	88.0
Two-Stream fusion [12]	65.4	92.5
I3D [5]	80.7	98.0
R(2+1)D [49]	78.7	97.3

Table 4: Accuracy averaged over all three splits on HMDB-51 and UCF-101 for both state-of-the-art compressed video based methods and decoded video based methods.

ods with a ResNet-18 classifier in Table 3 (b). We can see that our DMC generator runs much faster than these state-of-the-art motion representation learning methods.

4.5. Comparisons with Compressed Video Methods

As shown in the top section of Table 4, **DMC-Net** outperforms all other methods that operate in the compressed video domain, *i.e.* **CoViAR** [55], **EMV-CNN** [60] and **DTMV-CNN** [61]. Our method outperforms methods like [60, 61] that the output of the MV classifier is trained to approximate the output of the optical flow classifier. We believe this is because of the fact that approximating the classification output directly is not ideal, as it does not explicitly address the issues that MV is noisy and low-resolutional. By generating a more discriminative motion representation DMC, we are able to get features that are highly discriminative for the downstream recognition task. Furthermore, our DMC-Net can be combined with these classification networks of high capacity and trained in an end-to-end manner. **DMC-Net (I3D)** replaces the classifier from ResNet-18 to I3D, achieving significantly higher accuracy and outperforming a number of methods that require video decoding.

Our supplementary material discusses the speed of I3D.

4.6. Comparisons with Decoded Video Methods

In this section we compare DMC-Net to approaches that require decoding all RGB images from compressed video. Some only use the RGB images, while others adopt the two-stream method [41] and further require computing flow.

RGB only. As shown in Table 4, decoded video methods only based on RGB images can be further divided into three categories. **(1) Frame-level classification:** 2D CNNs like ResNet-50 and ResNet-152 [15] have been experimented in [11] to classify each frame individually and then employ simple averaging to obtain the video-level prediction. Due to lacking motion information, frame-level classification underperforms **DMC-Net**. **(2) Motion representation learning:** In Table 4, we evaluate **PWC-Net (ResNet-18) + CoViAR** which feeds estimated optical flow into a ResNet-18 classifier and then fuses the prediction with **CoViAR**. The accuracy of **PWC-Net (ResNet-18) + CoViAR** is not as good as **DMC-Net (ResNet-18)** because our generated DMC contains more discriminative motion cues that are complementary to MV. For TVNet [10], the authors used BN-Inception [19] to classify the generated motion representation and then fuse the prediction with a RGB CNN. The accuracy of TVNet is better **DMC-Net (ResNet-18)** thanks to using a strong classifier but is worse than our **DMC-Net (I3D)**. **(3) Spatio-temporal modeling:** There are also a lot of works using CNN to model the spatio-temporal patterns across multiple RGB frames to implicitly capture motion patterns. It turns out that our **DMC-Net** discovers motion cues that are complementary to such spatio-temporal patterns: **I3D RGB + DMC-Net (I3D)** improves **I3D RGB** via incorporating predictions from our **DMC-Net (I3D)**.

RGB + Flow. As shown in Table 4, the state-of-the-art accuracy is belonging to the two-stream methods [21, 49], which combine predictions made from a RGB CNN and an optical flow CNN. But as discussed in Section 4.4, extracting optical flow is quite time-consuming and thus these two-stream methods are much slower than our **DMC-Net**.

5. Conclusion

In this paper, we introduce **DMC-Net**, a highly efficient deep model for video action recognition in the compressed video domain. Evaluations on 3 action recognition benchmarks lead to substantial gains in accuracy over prior work, without the assistance of computationally expensive flow. The supplementary materials can be found in the following appendix.

6. Acknowledgment

Zheng Shou would like to thank the support from Wei Family Private Foundation when Zheng was at Columbia.

7. Appendix

7.1. Data Modalities in the Compressed Domain

Prevailing video compression standards employs the Group Of Pictures (GOP) structure to encode the raw video into successive, non-overlapping GOPs. Frames or pictures within one GOP are compressed together. Each GOP begins with an I-frame (intra coded frame) whose RGB pixel values are stored. I-frame can be decoded independently with other frames.

The rest of frames within a GOP are P-frame (predictive coded frame) and/or B-frame (bi-predictive coded frame), containing motion-compensated difference information relative to the previously decoded frames. Each P-frame can only reference one frame which could be either I-frame or P-frame while each B frame can only reference two frames. In this thesis, we follow [55] to focus on the low-latency scenario which only involves P-frame without B-frame. Each P-frame stores motion vectors and residual errors: during encoding, the video codec divides a P-frame into macroblocks of size such as 16x16 and find the most similar image patch in the reference frame for each macroblock; the displacement between a macroblock in P-frame and its most similar image patch in the reference frame is regarded as the corresponding motion vector, which will be used in motion compensation during decoding; the pixel differences between a macroblock in P-frame and its most similar image patch in the reference frame are denoted as residual errors. During the decoding of a P-frame, the video codec performs motion compensation which effectively warps the reference frame using the motion vectors and then adds the residual errors to the motion-compensated reference frame to reconstruct the P-frame.

Consequently, three data modalities in the compression domain are available: (1) RGB values of I-frame; (2) motion vectors and (3) residual errors of P-frame. We refer readers to [25] for more details.

7.2. More implementation details

In addition to Section 4.2 in the main paper, here we present more implementation details. Our model is implemented using PyTorch [33]. We first train our DMC-Net with the adversarial loss and then train it with all losses together. We elaborate these two steps separately in the following. On all three datasets (*i.e.* HMDB-51, UCF-101 and Kinetics-n50), we found a generic settings can work well. We use Adam optimizer [22].

Configuration of the training with the flow reconstruction loss and the classification loss. We first train the DMC generator for 1 epoch using the flow reconstruction loss only with the classification network fixed. Then we include the classification loss to train both the generator and

classifier end-to-end for 49 epochs. In the total loss (*i.e.* the Equation 4 in the main paper), we set α to 10 to balance weights. The overall learning rate is set to 0.01 and it is divided by 10 whenever the total training loss plateaus. All layers in the classification network except its last layer have the learning rate set to be 100x smaller.

Configuration of the training with all losses including the adversarial loss. Then we use the above trained model as the initialization for training our whole model with all three losses including the adversarial loss. Our whole model consists of the generator, the classifier and the discriminator now. In the total loss (*i.e.* the Equation 4 in the main paper), we set α to 10 and set λ to 1. The overall learning rate is set to 0.01 and it is divided by 10 whenever the total training loss plateaus. All layers in the classification network except its last layer have the learning rate set to be 100x smaller. Based on the network architectures for the discriminator used in a popular GAN implementation repository , we experimented with various number of filters in each layer and various number of layers. Finally we identified a network architecture for implementing our discriminator which achieves accuracy comparable to more complicated architectures. This discriminator’s architecture consists of a stack of 2D convolutional layers with a two-way Fully Connected layer at the end, as shown in the following Figure 6.

7.3. Other early fusion possibilities

As shown in Figure 7, we explore other early fusion possibilities: we duplicate the first convolution layer (*i.e.* conv0) of our DMC generator as conv0_mv and conv0_r to respectively process MV and R independently. Their outputs are fused before feeding into conv1 and two fusion methods are studied: element-wise addition (denoted as **Add**) and channel-wise concatenation (denoted as **Concat**). On HMDB-51, our method (*i.e.* directly stacking MV and R) achieves accuracy **61.80%**, which is better than **Add** (61.32%) and **Concat** (61.36%). We believe this is because MV and R are strongly correlated in the original pixel space before convolution.

7.4. Attention-weighted flow reconstruction loss

In this section we describe a way to attend to the discriminative regions of optical flow during generating DMC. However, in our experiments we found that this idea does not offer quantitative benefit beyond the GAN method on the datasets we experimented with. Thus this idea was not included in the main paper.

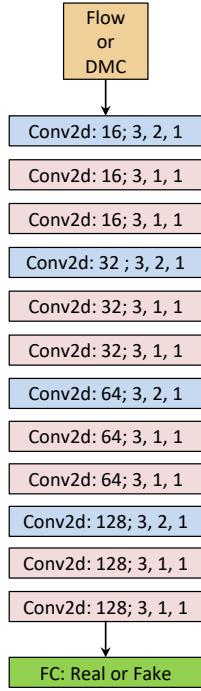


Figure 6: The network architecture of our discriminator. We denote each 2D convolutional layer in the format of #filters; kernel size, stride, padding.

7.4.1 Approach

The Mean Square Error (MSE) loss penalizes errors evenly over the whole image. In many cases, parts of optical flow contain noises, e.g. motions corresponding to background or camera motion. When reconstructing the optical flow, our DMC generator would ideally focus only on parts of the flow that contain motion cues discriminative with respect to the downstream action recognition task. Because these would be the regions of optical flow that are important for action recognition, and the regions where we would want a better reconstruction. Conversely, the reconstruction error in other regions of the optical flow, such as background, may not be important or even could be misleading.

This motivates us to try to create an adaptive MSE loss, where a weight is assigned for each location of the optical flow, based on the discriminative ability of that location. To get such a set of weights for each optical flow, we utilize recent related works on network interpretation [59], including the Class Activation Map [62] method and the Guided Back-Propagation [44] method. Such methods were proposed with a view to highlighting discriminative regions of the input data with respect to the classification outputs and are able to calculate *attention*-like weights for every location of the input data.

All methods mentioned above require a trained classifier

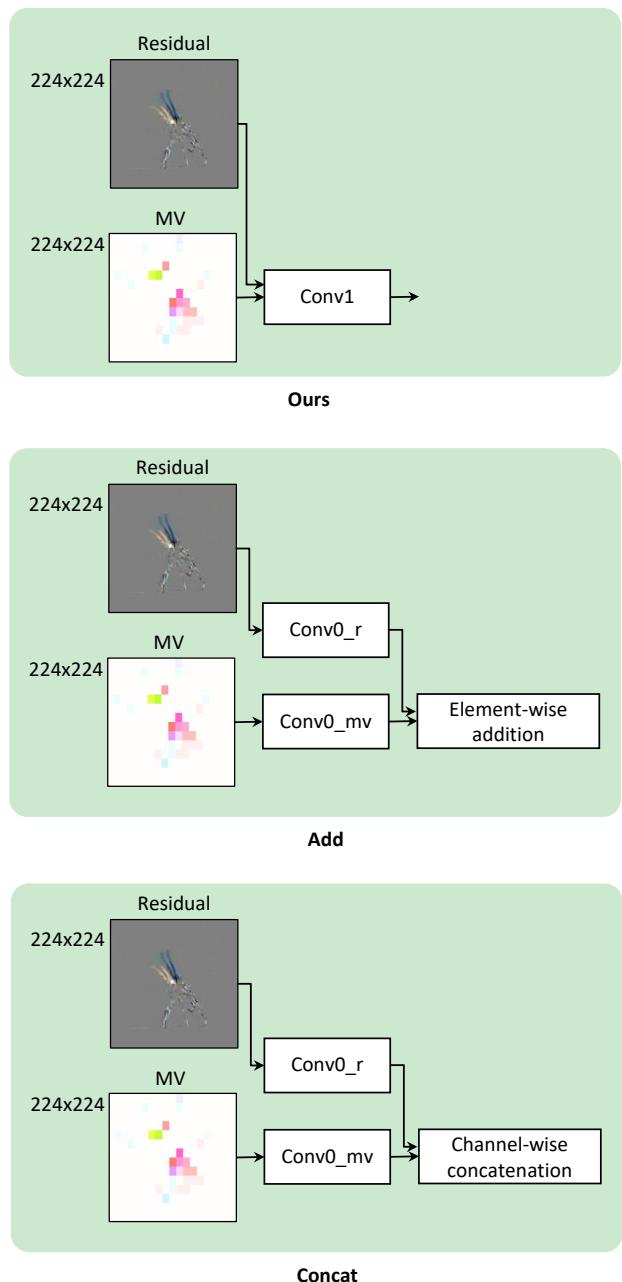


Figure 7: Different early fusion possibilities.

to inspect. We therefore first train a ResNet-18 classifier network for action recognition using optical flow as input. We then use network interpretation methods to output a set of attention-like weights $A \in \mathbb{R}^{H \times W}$ for each input optical flow of height H and width W . These attention-like weights can be computed before training our DMC generator and then be utilized during the training of our DMC-Net. Specifically, we can extend the optical flow reconstruction loss of Equation 1 in the main paper to take into

account the location-specific weights and derive the adaptively weighted flow reconstruction loss $\mathcal{L}_{\text{MSE}-\alpha}$:

$$\begin{aligned} \mathcal{L}_{\text{MSE}-\alpha} &= \mathbb{E}_{\mathbf{x} \sim p} \frac{1}{H \times W} \\ &\sum_{h=1}^H \sum_{w=1}^W A_{h,w} \cdot \|\phi_{h,w} - \sigma_{h,w}\|_2^2, \end{aligned} \quad (5)$$

where p denotes the set of P-frames in the training videos, ϕ is set to the generated DMC denoted as $\mathcal{G}_{\text{DMC}}(\mathbf{x})$, σ is set to the corresponding optical flow denoted as $\mathcal{G}_{\text{OF}}(\mathbf{x})$, \mathbb{E} stands for computing expectation, $A_{h,w}$ denotes the learned weight for location h, w . In order to obtain A , we have explored two widely used network interpretation techniques as presented in the following.

Class Activation Map (CAM) [62]. The CAM method practically switches the order of the last two layers of the trained flow classifier, *i.e.* the fully connected classification layer and the Global Average Pooling layer pool15. This way, the fully connected classifier can be re-purposed as a convolutional layer f_{cls} to slide over every location of the conv5 (*i.e.* the layer right before the pool15)'s output, effectively producing a classification score at each location. As the output of f_{cls} is of low spatial dimension and each location has a wide receptive field with respect to the input flow, the high activations effectively focus on the most discriminative salient regions of the input flow. We choose the activation map corresponding to the ground truth action class as the attention map A . Finally we deal with the negative values in A via passing A through a ReLU operation, which leads to the best accuracy compared to other common normalization methods according to our experimental explorations. Note that in our experiments discussed in the following Section 7.4.2, we resize the input flow from 224x224 to 448x448 before feeding it into the classifier so that we can obtain the attention map A of higher spatial resolution (*i.e.* 14x14), covering more details. Further, we upsample A back to 224x224 via bilinear interpolation so that A has the same size as the generated DMC. As shown in Figure 8 (c), the attention map generated by the CAM method can indeed highlight the salient regions of flow such as the player's hands and head. The flow values along the x direction and the y direction at the same spatial location share the same attention weight.

Guided Back-Propagation (GBP) [44]. Rather than finding the salient regions, some methods [59, 44] have been proposed to determine the contribution from the input's each value to the final classification output. Since the input in our case is optical flow, the higher the contribution of a value, the more discriminative motion information this value contains. Therefore, we can obtain an attention map

A of the data shape as the same as the flow (*i.e.* 2x224x224 in the following Section 7.4.2). Each value in A stands for the contribution of the corresponding flow's value at the same location. Specifically, we utilize the GBP [44] method, which improves the De-conv method [59] by combining it with the regular back-propagation pass. Concretely, we set the classification output as a one-hot vector with the ground truth class indicated and then we back-propagate the one-hot vector back to the input flow. Note that following the conventional back-propagation can only generate a generic attention map independent to the input rather than a map that is related with a specific input flow. To address this issue, GBP further integrates the De-conv method into the conventional back-propagation pass: basically whenever back-propagating gradients through a ReLU layer, GBP sets the negative gradients to 0. Finally, we pass the obtained A through a ReLU operation to set its negative values to 0. Figure 8 (d) and (e) show the attention maps generated by the GBP method, highlighting the pixels whose values are sensitive for classifying the input optical flow as PlayingTable.

	Accuracy
DMC-Net	61.5
DMC-Net w/ Att (CAM)	61.4
DMC-Net w/ Att (GBP)	61.5
DMC-Net w/ GAN	61.8
DMC-Net w/ GAN w/ Att (CAM)	61.5
DMC-Net w/ GAN w/ Att (GBP)	61.0

Table 5: Accuracy on HMDB-51 averaged over 3 splits for the study of the effectiveness of attending to the discriminative regions of optical flow during training our DMC-Net.

7.4.2 Experimental results

Although it is reasonable and intuitive to attend to the discriminative regions of optical flow during generating DMC, this idea does not offer benefit beyond the GAN method proposed in the main paper. In Table 5, **DMC-Net** is only trained with the flow reconstruction loss and the classification loss; **DMC-Net w/ Att (CAM)** is replacing the flow reconstruction loss in **DMC-Net** by the above attention-weighted flow reconstruction loss based on the attention map generated by the CAM method; **DMC-Net w/ Att (GBP)** is replacing the flow reconstruction loss in **DMC-Net** by the above attention-weighted flow reconstruction loss based on the attention map generated by the GBP method. We can see that **DMC-Net** achieves accuracy comparable with **DMC-Net w/ Att (GBP)** and **DMC-Net w/ Att (CAM)**. But if we equip **DMC-Net** with the generative adversarial loss, denoted as **DMC-Net w/ GAN**, the highest

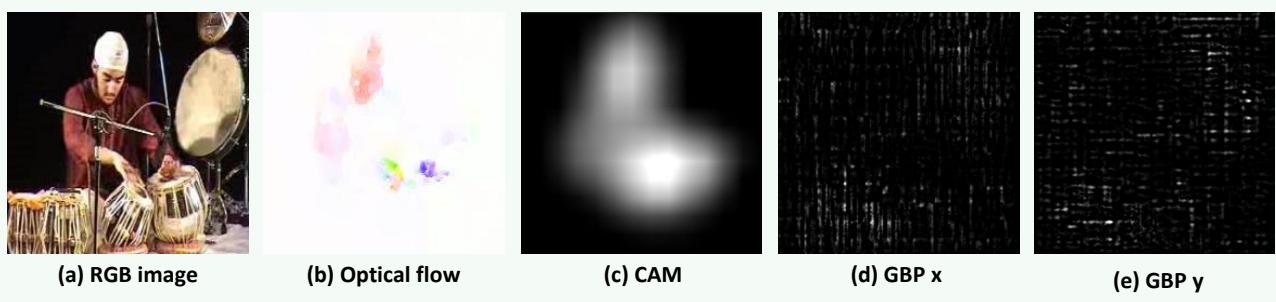


Figure 8: Illustrations of the attention maps generated by CAM [62] and GBP [44] for a PlayingTabla example. (a) shows the RGB image. (b) shows the corresponding optical flow. (c) is the attention map generated using the CAM method. (d) and (e) show the attention maps generated using the GBP method for the input flow respectively along the x direction and the y direction. Better viewed in color.

accuracy can be achieved.

Furthermore, we explore whether this strategy of making flow reconstruction loss attending to the discriminative regions of flow is complementary to the proposed adversarial loss. **DMC-Net w/ GAN w/ Att (CAM)** and **DMC-Net w/ GAN w/ Att (GBP)** respectively use the attention map generated by the CAM method and the GBP method to weight the flow reconstruction loss in **DMC-Net w/ GAN**. But this strategy of attention-weighted flow reconstruction loss hurts the accuracy of **DMC-Net w/ GAN** and thus is not complementary to the idea of using the adversarial loss. We believe this is because the DMC generator trained with the original flow reconstruction loss, the classification loss and the GAN loss can already capture sufficient motion information that can be learned from approximating flow and thus explicitly focusing on the discriminative regions of flow does not offer additional benefits. Consequently, we opt to do not use the attention-weighted flow reconstruction loss in the main paper.

7.5. More discussions about the speed

7.5.1 Speed of DMC-Net (I3D)

Table 4 in the main paper shows that our DMC-Net implemented using a I3D classifier, denoted as **DMC-Net (I3D)**, achieves much better accuracy than using a ResNet-18 classifier, denoted as **DMC-Net (ResNet-18)**. Note that the speed of **DMC-Net (I3D)** and the speed of **DMC-Net (ResNet-18)** are not directly comparable. ResNet-18 is a frame-level classifier: given an input frame, **DMC-Net (ResNet-18)** can classify it with the speed at 0.76ms as reported in the Table 3 in the main paper.

However, I3D is a clip-level classifier: during testing, we follow [5] to feed 250 frames concurrently into a I3D classifier to obtain one action class prediction. The per-frame inference time of **DMC-Net (I3D)** is 0.79ms which is slightly slower yet very close to **DMC-Net (ResNet-18)**.

(i.e. 0.76ms). But in order to make one action prediction, **DMC-Net (I3D)** needs to take $0.79 \times 250 = 197.5$ ms while **DMC-Net (ResNet-18)** only takes 0.76ms with the need of only one input frame.

7.6. Generalize DMC-Net to different compressed video standards

It is worthwhile pointing out that although we follow CoViAR [55] to specifically use MPEG-4 video [25], in real applications it would be interesting to develop methods that can handle different video encoding formats. In the worst case, we can always convert the input video of arbitrary format into MPEG-4 first. On HMDB-51, FFmpeg [1] takes 1.13ms in average to convert one frame when processing each video sequentially, still being much faster than extracting flow for the two-stream method. Table 3 in the main paper shows that the per-frame inference speed of DMC-Net is 0.76ms and that of the two-stream method is more than 75ms.

7.7. More ablation studies

In addition to the model analysis in the main paper’s Section 4.3, to further validate our design choices, here we present more ablation studies and some strategies that are alternative to the current settings used in the main paper.

7.7.1 End-to-end learning

In the main paper, we train the generator and the classifier in an end-to-end manner with the gradients from the classification loss propagated to not only the classifier but also the generator. An alternative training strategy is to separate the training of the generator and the training of the classifier. Concretely, we can first train the generator without the classifier and the classification loss. Then we fix the generator only for doing inference and then feed the generated

DMC into the classification network to train the classifier using the classification loss only.

7.7.2 Decomposed two-stage DMC generation

In the main paper, we design a lightweight network to refine Motion Vector (denoted as MV) to generate DMC. Note that MV has 224x224 spatial size but MV is composed of 16x16 macroblocks in which every pixel has the identical value. If we downsample MV by a factor of 16 (denoted as MV_d), the same amounts of motion information are still preserved. Thus the effects of our DMC generator can be considered to be two-fold: (1) correcting errors and reducing noises in MV and (2) generating fine details of discriminative motion cue during the process of upsampling MV_d .

Consequently, an alternative way of designing the DMC generator is to first have an error correction network to rectify noises in MV_d and then have another network to conduct upsampling from 14x14 to 224x224. As shown in Figure 9 (b), given the stacked residual and MV_d both of size 14x14, we have an error correction network to generate MV'_d of size 14x14. Then the generated MV'_d is resized from 14x14 to 224x224 via bilinear interpolation to obtain the MV' . Finally, we feed the stacked residual and MV' into an up-sampling network to generate DMC of more fine motion details. Note that in Figure 9 (b) we not only measure the flow reconstruction loss between the generated DMC and the corresponding flow but also measure the flow reconstruction loss between the MV'_d and the downsampled flow of size 14x14.

7.8. Smoothing Motion Vector via bilinear interpolation before fed into the DMC generator

As shown in Figure 9 (a), the DMC generator used in the main paper accepts the blocky MV of size 224x224 as input. Since the optical flow extract by TV-L1 is smooth rather than blocky, smoothing MV before feeding it into the generator can generate DMC of much less blocky artifacts which do not exhibit useful motion information. Therefore, instead of directly feeding the blocky MV into the DMC generator, we can make the input MV more smooth by first downampling MV of size 224x224 to MV_d of size 14x14 and then resizing MV_d back to 224x224 via bilinear interpolation. The rest process follows the main paper.

7.8.1 Experimental results

To investigate the effectiveness of the above strategies, we explore different scenarios during the training of DMC-Net using the flow reconstruction loss and the classification loss. We denote the scenario used in the main paper as **Ours**, which trains the generator and the classifier in an end-to-end manner, generates the DMC in one single stage, and takes the blocky MV as input for the DMC generator.

First, we compare **Ours** to **Ours w/o end-to-end** which follows the above Section 7.7.1 to separate the training of the DMC generator and the training of the classifier. Table 6 confirms the effectiveness of the end-to-end learning strategy and therefore we use it in the main paper.

Second, we compare **Ours** to **Ours w/ two-stage** which follows the above Section 7.7.2 to decompose the DMC generation into a two-stage process. Table 6 shows that decomposing the DMC generation into two-stage does not offer benefit in terms of accuracy. Thus we opt to use the single network in the main paper to generate DMC via jointly correcting errors and generating fine motion details in one single step.

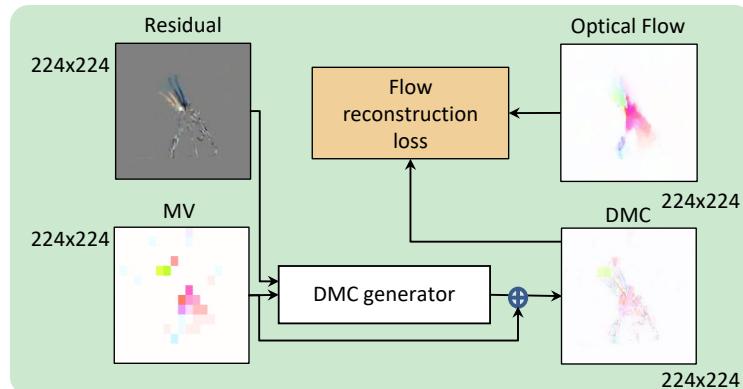
Third, we compare **Ours** to **Ours w/ bilinear interp** which follows the above Section 7.8 to first smooth MV via bilinear interpolation before feeding it into the DMC generator. It turns out that **Ours** and **Ours w/ bilinear interp** can generate DMC of comparably good motion cues that lead to similar accuracy. Therefore in the main paper we directly feed the blocky MV into the DMC generator.

	Accuracy
Ours w/o end-to-end	59.3
Ours w/ two-stage	60.6
Ours w/ bilinear interp	61.4
Ours	61.5

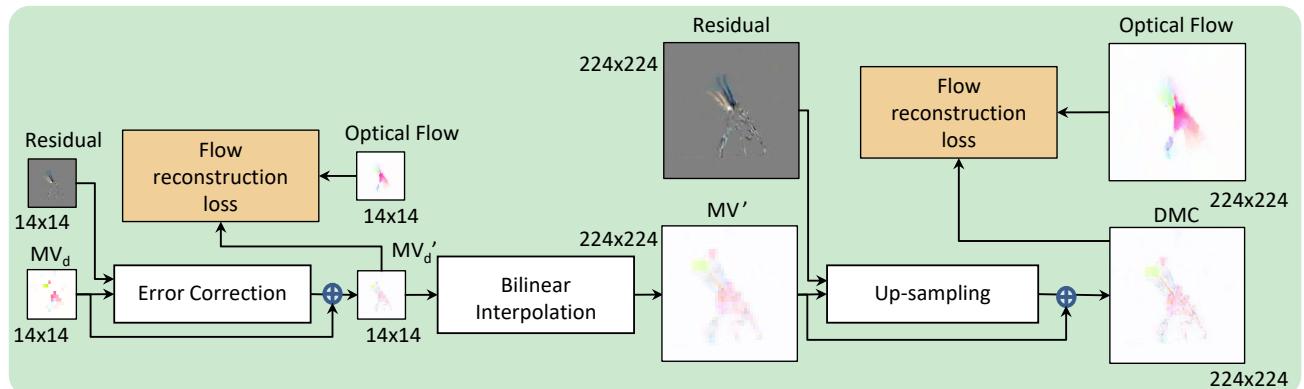
Table 6: Accuracy on HMDB-51 averaged over 3 splits when our DMC-Net is trained with the flow reconstruction loss and the classification loss.

7.9. Results on the full Kinetics dataset

Due to the extremely long training time on the full Kinetics dataset using one single GPU, we directly adopt the training hyper-parameters used for the Kinetics-n50 subset. The accuracy of **CoViAR** is 65.37%; the accuracy of **CoViAR + TV-L1 Flow** is 65.43%; the accuracy of **DMC-Net (ours)** is 65.42%. We can observe that **DMC-Net (ours)** still improves **CoViAR** to match the performance of **CoViAR + TV-L1 Flow** but the performances are very close. We conjecture this is because when training on such a large-scale dataset, the models for I-frame and Residual have already seen training data of large variance and thus motion information cannot offer significantly complementary cues for distinguishing different action categories.



(a) Single one-stage DMC generation



(b) Decomposed two-stage DMC generation

Figure 9: Illustrations for (a) the strategy of single one-stage DMC generation used in the main paper and (b) the strategy of decomposed two-stage DMC generation.

References

- [1] Ffmpeg: A complete, cross-platform solution to record, convert and stream audio and video. <https://www.ffmpeg.org/>. 12
- [2] Hakan Bilen, Basura Fernando, Efstratios Gavves, and Andrea Vedaldi. Action recognition with dynamic image networks. *IEEE transactions on pattern analysis and machine intelligence*, 2018. 3
- [3] Andrés Bruhn, Joachim Weickert, and Christoph Schnörr. Lucas/kanade meets horn/schunck: Combining local and global optic flow methods. *International journal of computer vision*, 61(3):211–231, 2005. 3
- [4] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black. A naturalistic open source movie for optical flow evaluation. In A. Fitzgibbon et al. (Eds.), editor, *European Conf. on Computer Vision (ECCV)*, Part IV, LNCS 7577, pages 611–625. Springer-Verlag, Oct. 2012. 3
- [5] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*, 2017. 1, 3, 6, 8, 12
- [6] Yu-Wei Chao, Sudheendra Vijayanarasimhan, Bryan Seybold, David A Ross, Jia Deng, and Rahul Sukthankar. Rethinking the faster r-cnn architecture for temporal action localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1130–1139, 2018. 1
- [7] Yunpeng Chen, Yannis Kalantidis, Jianshu Li, Shuicheng Yan, and Jiashi Feng. Multi-fiber networks for video recognition. In *ECCV*, 2018. 8
- [8] J Lewis M Black D Sun, S Roth. Learning optical flow. In *ECCV*, 2008. 3
- [9] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. Flownet: Learning optical flow with convolutional networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2758–2766, 2015. 3
- [10] Lijie Fan, Wenbing Huang, Stefano Ermon Chuang Gan, Boqing Gong, and Junzhou Huang. End-to-end learning of motion representation for video understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6016–6025, 2018. 3, 7, 8
- [11] Christoph Feichtenhofer, Axel Pinz, and Richard P Wildes. Spatiotemporal multiplier networks for video action recognition. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 8
- [12] 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*, 2016. 8
- [13] Ruohan Gao, Bo Xiong, and Kristen Grauman. Im2flow: Motion hallucination from static images for action recognition. In *CVPR*, 2018. 3, 5
- [14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, 2014. 2, 5
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 4, 7, 8
- [16] Berthold KP Horn and Brian G Schunck. Determining optical flow. *Artificial intelligence*, 17(1-3):185–203, 1981. 3
- [17] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *CVPR*, 2017. 4
- [18] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. In *IEEE conference on computer vision and pattern recognition (CVPR)*, volume 2, page 6, 2017. 3, 7
- [19] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 7, 8
- [20] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. *arXiv preprint*, 2017. 7
- [21] 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. 2, 8

- [22] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 9
- [23] Hildegarde Kuehne, Hueihan Jhuang, Estibaliz Garrote, Tomaso Poggio, and Thomas Serre. Hmdb: a large video database for human motion recognition. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, pages 2556–2563. IEEE, 2011. 2, 6
- [24] Ivan Laptev, Marcin Marszałek, Cordelia Schmid, and Benjamin Rozenfeld. Learning realistic human actions from movies. In *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, pages 1–8. IEEE, 2008. 1
- [25] Didier Le Gall. Mpeg: A video compression standard for multimedia applications. *Communications of the ACM*, 1991. 1, 3, 9, 12
- [26] Tianwei Lin, Xu Zhao, and Zheng Shou. Single shot temporal action detection. In *Proceedings of the 2017 ACM on Multimedia Conference*, pages 988–996. ACM, 2017. 1
- [27] Tianwei Lin, Xu Zhao, and Zheng Shou. Temporal convolution based action proposal: Submission to activitynet 2017. *arXiv preprint arXiv:1707.06750*, 2017. 1
- [28] Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml*, volume 30, page 3, 2013. 4
- [29] Michael Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond mean square error. *ICLR*, 2016. 5
- [30] Engin Mendi, Hélio B Clemente, and Coskun Bayrak. Sports video summarization based on motion analysis. *Computers & Electrical Engineering*, 39(3):790–796, 2013. 1
- [31] Moritz Menze, Christian Heipke, and Andreas Geiger. Joint 3d estimation of vehicles and scene flow. In *ISPRS Workshop on Image Sequence Analysis (ISA)*, 2015. 3
- [32] Joe Yue-Hei Ng, Jonghyun Choi, Jan Neumann, and Larry S Davis. Actionflownet: Learning motion representation for action recognition. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1616–1624. IEEE, 2018. 3, 8
- [33] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 9
- [34] Anurag Ranjan and Michael J Black. Optical flow estimation using a spatial pyramid network. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 2, page 2. IEEE, 2017. 3
- [35] Laura Sevilla-Lara, Yiyi Liao, Fatma Guney, Varun Jampani, Andreas Geiger, and Michael J Black. On the integration of optical flow and action recognition. *arXiv preprint arXiv:1712.08416*, 2017. 2, 3
- [36] Laura Sevilla-Lara, Yiyi Liao, Fatma Guney, Varun Jampani, Andreas Geiger, and Michael J. Black. On the integration of optical flow and action recognition. In *German Conference on Pattern Recognition (GCPR)*, Oct. 2018. 3
- [37] Zheng Shou, Jonathan Chan, Alireza Zareian, Kazuyuki Miyazawa, and Shih-Fu Chang. Cdc: Convolutional-deconvolutional networks for precise temporal action localization in untrimmed videos. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 1417–1426. IEEE, 2017. 1
- [38] Zheng Shou, Hang Gao, Lei Zhang, Kazuyuki Miyazawa, and Shih-Fu Chang. Autoloc: Weaklysupervised temporal action localization in untrimmed videos. In *ECCV*, 2018. 1
- [39] Zheng Shou, Junting Pan, Jonathan Chan, Kazuyuki Miyazawa, Hassan Mansour, Anthony Vetro, Xavier Giro-i Nieto, and Shih-Fu Chang. Online action detection in untrimmed, streaming videos-modeling and evaluation. In *ECCV*, 2018. 1
- [40] Zheng Shou, Dongang Wang, and Shih-Fu Chang. Temporal action localization in untrimmed videos via multi-stage cnns. In *CVPR*, 2016. 1
- [41] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. In *NIPS*, 2014. 1, 2, 3, 8
- [42] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012. 2
- [43] K. Soomro, A. R. Zamir, and M. Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. In *CRCV-TR-12-01*, 2012. 6
- [44] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014. 10, 11, 12
- [45] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. PWC-Net: CNNs for optical flow using pyramid, warping, and cost volume. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1, 3, 4, 6, 7, 8
- [46] Antonio Tejero-de Pablos, Yuta Nakashima, Tomokazu Sato, Naokazu Yokoya, Marko Linna, and Esa Rahtu. Summarization of user-generated sports video by using deep action recognition features. *IEEE Transactions on Multimedia*, 20(8):2000–2011, 2018. 1
- [47] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 4489–4497, 2015. 1, 3, 4, 8
- [48] Du Tran, Jamie Ray, Zheng Shou, Shih-Fu Chang, and Manohar Paluri. Convnet architecture search for spatiotemporal feature learning. *arXiv preprint arXiv:1708.05038*, 2017. 4, 8
- [49] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6450–6459, 2018. 1, 3, 8
- [50] Heng Wang and Cordelia Schmid. Action recognition with improved trajectories. In *Proceedings of the IEEE international conference on computer vision*, pages 3551–3558, 2013. 1

- [51] Limin Wang, Wei Li, Wen Li, and Luc Van Gool. Appearance-and-relation networks for video classification. In *CVPR*, 2018. 8
- [52] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *ECCV*, 2016. 3
- [53] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*, 2018. 3
- [54] Philippe Weinzaepfel, Jerome Revaud, Zaid Harchaoui, and Cordelia Schmid. Deepflow: Large displacement optical flow with deep matching. In *Proceedings of the IEEE International Conference on Computer Vision*, 2013. 7
- [55] Chao-Yuan Wu, Manzil Zaheer, Hexiang Hu, R Manmatha, Alexander J Smola, and Philipp Krähenbühl. Compressed video action recognition. In *CVPR*, 2018. 1, 2, 3, 4, 6, 7, 8, 9, 12
- [56] Saining Xie, Chen Sun, Jonathan Huang, Zhuowen Tu, and Kevin Murphy. Rethinking spatiotemporal feature learning for video understanding. In *ECCV*, 2018. 8
- [57] Christopher Zach, Thomas Pock, and Horst Bischof. A duality based approach for realtime tv-l 1 optical flow. In *Joint Pattern Recognition Symposium*, pages 214–223. Springer, 2007. 3, 6
- [58] Christopher Zach, Thomas Pock, and Horst Bischof. A duality based approach for realtime tv-l1 optical flow. In *Joint Pattern Recognition Symposium*, 2007. 1, 5, 6
- [59] M.D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *ECCV*, 2014. 10, 11
- [60] Bowen Zhang, Limin Wang, Zhe Wang, Yu Qiao, and Hanli Wang. Real-time action recognition with enhanced motion vector cnns. In *CVPR*, 2016. 1, 3, 4, 8
- [61] Bowen Zhang, Limin Wang, Zhe Wang, Yu Qiao, and Hanli Wang. Real-time action recognition with deeply transferred motion vector cnns. *IEEE Transactions on Image Processing*, 2018. 1, 3, 4, 8
- [62] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *CVPR*, 2016. 10, 11, 12