

TinyVIRAT: Low-resolution Video Action Recognition

Ugur Demir*, Yogesh S Rawat[†] and Mubarak Shah[†]

Center for Research in Computer Vision

University of Central Florida, Orlando, Florida, USA

Email: *[ugur]@knights.ucf.edu, [†][yogesh, shah]@crcv.ucf.edu

Abstract—The existing research in action recognition is mostly focused on *high-quality* videos where the action is distinctly visible. In real-world surveillance environments, the actions in videos are captured at a wide range of resolutions. Most activities occur at a distance with a *small resolution* and recognizing such activities is a challenging problem. In this work, we focus on recognizing tiny actions in videos. We introduce a benchmark dataset, *TinyVIRAT*, which contains natural *low-resolution* activities. The actions in *TinyVIRAT* videos have multiple labels and they are extracted from surveillance videos which makes them realistic and more challenging. We propose a novel method for recognizing tiny actions in videos which utilizes a *progressive generative* approach to improve the quality of low-resolution actions. The proposed method also consists of a weakly trained *attention mechanism* which helps in focusing on the activity regions in the video. We perform extensive experiments to benchmark the proposed *TinyVIRAT* dataset and observe that the proposed method significantly improves the action recognition performance over baselines. We also evaluate the proposed approach on synthetically resized action recognition datasets and achieve state-of-the-art results when compared with existing methods. The dataset and code is publicly available at <https://github.com/UgurDemir/Tiny-VIRAT>.

I. INTRODUCTION

Video action recognition has recently seen a good progress, which is mostly due to the availability of large-scale datasets and the success in deep learning. The availability of datasets, such as UCF-101 [1], Kinetics [2], Moments-in-time [3], AVA [4], and Youtube-8M [5], has played an important role in this advancement. Apart from this, there are several novel architectures, such as C3D [6], I3D [7], ResNet-50 [8], and TSN [9], which are proven to be effective for action recognition. However, the performance of these models relies on the quality of the action videos. The videos in all these datasets are of high quality and the action usually covers majority of the video frame. In real-world surveillance environments, the actions in videos are captured at a wide range of resolutions and may appear very tiny, therefore recognizing such actions becomes challenging at a very low-resolution. The existing action recognition models are not designed for low-resolution videos and their performance is still far from satisfactory when the action is not distinctly visible.

In this work, our focus is on action recognition in low-resolution videos. The existing approaches addressing this issue, such as [10], [11], and [12], perform their experiments on *artificially created datasets* where the high-resolution videos

are down-scaled to a smaller resolution to create a low-resolution sample. However, re-scaling a high-resolution video to a lower-resolution does not reflect real world low-resolution video quality. Real world low-resolution videos suffer from grain, camera sensor noise, and other factors, which are not present in the down-scaled videos.

To address this problem, we propose a new benchmark dataset, *TinyVIRAT*, for low-resolution action recognition. The videos in *TinyVIRAT* are realistic and extracted from surveillance videos of VIRAT dataset [13]. This is a multi-label dataset with multiple actions per video clip which makes it even more challenging. The dataset has around 13K video samples from 26 different actions and all the videos are captured at 30fps. The length of the activities vary from sample to sample with an average length of around 3 seconds. It contains arbitrary sized low-resolution videos which ranged from 10x10 pixels to 128x128 pixels with an average of 70x70 pixels. The videos in the proposed dataset are naturally low resolution and they reflect real-life challenges. Some sample video frames from *TinyVIRAT* are shown in Figure 1.

We propose a novel end-to-end deep learning approach to address the problem of tiny action recognition. The proposed approach has three main components; video super-resolution, weakly supervised attention mechanism, and action classification. The video super-resolution network takes a low-resolution video and recovers important appearance and motion details using a Dense Video Super-Resolution network (DVSR), which is trained in a progressive manner. In this set up the foreground and background will have equal importance for the super-resolution task. However, foreground information has more discriminative information for action recognition. Therefore, to make DVSR action aware, we propose a novel attention mechanism which estimates a spatio-temporal map indicating the importance of each pixel for the corresponding action. The attention map is trained in a weakly supervised setting which is guided by the action label of the video without requiring localization bounding boxes. The estimated spatio-temporal map is integrated with the synthesized high quality video to perform action recognition using a classifier.

In summary, this paper makes the following contributions:

- We introduce a tiny action benchmark dataset which is the first dataset for this problem to the best of our knowledge.
- We propose a progressive video super-resolution based approach for tiny action recognition and demonstrate its





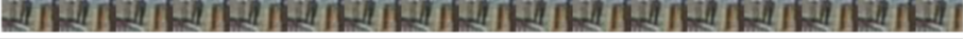


Size	Actions	Frames
20x20	standing	
28x28	carrying walking	
38x38	carrying walking	
40x40	talking standing	
44x44	carrying	
48x48	walking	
52x52	moving riding	

Fig. 1. Some sample video frames for actions from *TinyVIRAT* dataset. The dataset contain low-resolution videos with varying sizes. *TinyVIRAT* is a multi-label dataset and each action video can have multiple action labels

effectiveness on *TinyVIRAT* and artificial low-resolution action recognition benchmarks.

- We also introduce a weakly supervised foreground attention mechanism that helps a super-resolution network to focus on important regions.
- We perform extensive experiments on the proposed *TinyVIRAT* dataset and quantitatively demonstrate its challenging nature when compared with existing artificially created low-resolution benchmark datasets.

II. RELATED WORK

Action Classification: After deep neural networks became popular for images, they have been successfully applied to the problem of video action recognition [14], [6], [7]. One of the popular deep network architecture C3D [6] showed that using 3D convolution is more suitable to extract spatio-temporal features for video action recognition. Recently, I3D architecture [7] has shown favorable performance on standard benchmarks [1], [15], [2] by employing Inception layers. In [14], deep ResNet [8] architecture variants are investigated for the action recognition task.

For Low Resolution (LR) single image applications, several different approaches have been proposed, where domain adaptation, super-resolution or feature learning are employed to find better representations of LR images [16], [17]. Previous work on this problem is generally motivated by privacy preservation [18], [19], [20]. In [19], a model is proposed which learns a set of different transformation that creates LR videos from the HR training set. It is claimed that action classifiers which are trained on the generated LR dataset learn better decision boundary. In [10], [11], CNN based action classifiers are simultaneously trained with both LR and HR inputs by assuming that models learn similar representations. In [12], the effect of super-resolution on the action recognition task

is analyzed. They compared the existing image and video-based super-resolution networks, and proposed an optical flow guided training approach. However, they only show their performance on artificially created low-resolution videos by downsampling UCF-101 and HMDB-51 dataset to 80x60, which is far from natural tiny actions.

Super-Resolution: One of the seminal CNN based single image SR method is proposed in [21]. After its success, several different CNN architectures have been introduced [22], [23], [24]. Although promising results are shown on single images, these methods cannot capture temporal information in videos if applied frame by frame. In some works, adversarial training has been utilized to obtain more realistic texture in images [25], [26], [27], [28], [29], nonetheless; there is still no consensus about the best training scheme for SR models. Traditional pixel-wise reconstruction losses tend to produce smooth results, however, adversarial training introduces noise and artifacts [30].

One of the most common strategies for Video SR is to incorporate optical flow for capturing motion information in order to synthesize a sequence of video frames [31], [32], [33]. In [12], optical-flow is used to improve super-resolution network performance. The main drawback is that optical flow is computationally expensive, and if motion between frames is high, obtaining reliable estimates becomes difficult. In [34], the 2D CNN network is used for frame SR and then temporal dependency weights are learned, which indicate how to merge input frames to synthesize the final frame. This method synthesizes one frame at a time. In [35] the authors proposed a 2D convolutional progressive network that incorporates short term inter-frame dependencies. There are some 2D convolution-based video SR approaches [36], [37], which only focus on short term temporal relations. In [38], recurrent CNN architecture is used to learn longer temporal

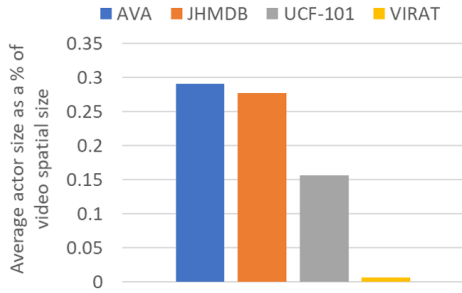


Fig. 2. Average object size ratio comparison.



Fig. 3. Comparison between average actor size and video resolution for different dataset. If the ratio is low enough, actions are naturally tiny.

dependencies for video SR. In some recent works [32], [39], 3D convolution has been explored for effective video super-resolution.

In contrast to these works, we propose a progressively growing architecture using 3D convolutions, where we start from a low-resolution and gradually increase the spatial and temporal resolutions. This progressive approach has been found effective in image generation [40], and we explore it for videos and demonstrate in this paper that it is effective when we have a higher scaling factor.

III. TINYVIRAT DATASET

Most of the existing action recognition datasets contain high resolution, actor centric videos [41], [42], [2], [43], [44], [45], [46], [47], [48], [49]. For example, Kinetics [2], Charades [43], Youtube-8M [44] are collected from Youtube videos where actions cover most of the image regions in every frame of a video. Using these videos to create low-resolution benchmark datasets does not reflect real world situation, and it is not appropriate as they generally contain larger actors.

In the real world, we encounter low-quality actions mostly in surveillance video clips where the camera placed in a distant place. Even though surveillance camera is capable of recording high-quality video, if an action happens far away from the camera, it will suffer from lack of details. Thus, surveillance videos are the perfect candidate for this problem. Figure 2 shows average actor size as a percentage of video spatial size, where most of the action recognition datasets have a significantly larger actor size. If the ratio is large, cropping actions will result similar spatial size with the original video. In comparison, VIRAT dataset has naturally occurring tiny actors which is well suited for low-resolution action recognition task, as can be seen in Figure 3.

We introduced TinyVIRAT dataset which is based on VIRAT [13] dataset for real-life tiny action recognition problems.

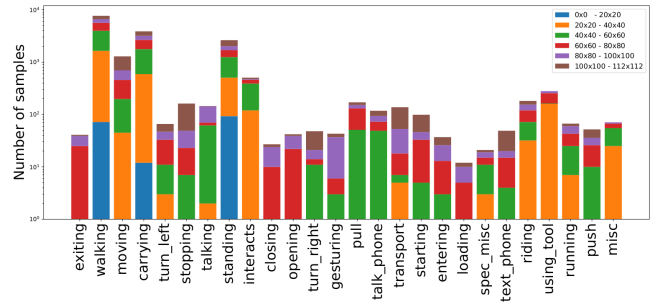


Fig. 4. Number of samples per action labels and resolution. Numbers on y-axis are shown in log scale.

TABLE I
DATASET STATISTICS. ANF: AVERAGE NUMBER OF FRAMES, ML: MULTI-LABEL, NC: NUMBER OF CLASSES, AND NV: NUMBER OF VIDEOS.

Dataset	Resolution	ANF	ML	NC	NV
UCF-101	320x240	186.50	No	101	13320
JHMDB-51	320x240	94.49	No	51	7000
AVA	264x440 - 360x640	127081.66	Yes	80	272
TinyVIRAT	10x10 - 128x128	93.93	Yes	26	12829

VIRAT dataset is a natural candidate for low-resolution actions but it contains a large variety of different actor sizes and it is a very complex since actions can happen any time in any spatial position. To focus only on low-resolution action recognition problem, we crop small action clips from VIRAT videos.

In VIRAT dataset actors can perform multiple actions and temporally actions can start and end at different times. Before deciding which actions are tiny, we merged spatio-temporally overlapping actions and created multi-label action clips. We split these clips if the labels are changing temporally. This steps makes sure that created clips are trimmed. We selected clips that are spatially smaller than 128x128. Finally, long videos are split into smaller chunks and actions which do not have enough samples are removed from the dataset. We use the same train and validation split from the VIRAT dataset.

TODO TinyVIRAT has 7,663 training and 5,166 testing videos with 26 action labels. Table I shows statistics from TinyVIRAT and several other datasets. Figure 4 shows the number of samples per action class and the distribution of the videos by spatial size.

IV. METHOD

The proposed method focuses on learning to enhance the quality of low-resolution videos to improve action classification performance. The action classifier network is trained with super-resolved videos instead of raw low-resolution video clips. Our approach consists of two main parts: (i) super-resolution and (ii) action classifier networks which can be seen in Figure 5 The first module, **Super-resolution Network** (SR network), is a novel deep convolutional neural network, which takes a low-resolution video clip and introduces sharp appearance and motion details to synthesize a high-resolution

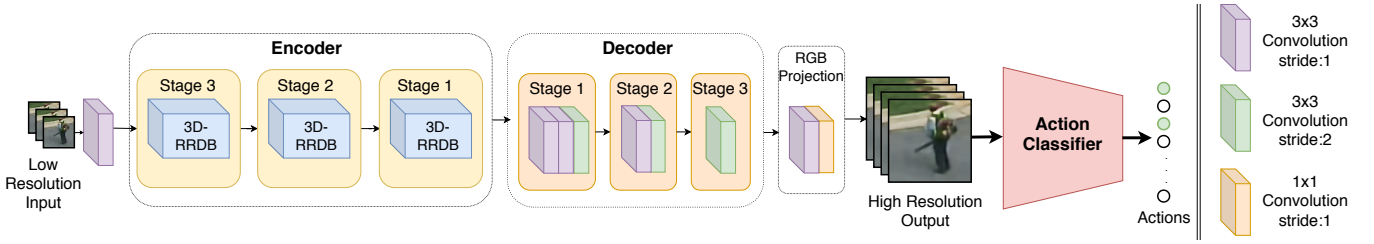


Fig. 5. Overview of the progressive video generation and action classification approach. During the training process, we are introducing new blocks to Progressive DVSR network architecture at each stage. After video synthesis is completed, action classifier process the video to predict actions.

counterpart. The second module, **Action Classification Network**, takes the generated SR video clip and recognizes the action in the video.

We propose progressive training strategy for the SR network which improves the reconstruction that helps action classification task. Improving texture quality leads us to get better performance but super-resolution network still does not know any task specific information. Focusing on the regions where the actions are happening is more beneficial for our main goal. Since, background of the videos does not have the importance as much as foreground for the action recognition task, guiding SR network is more crucial. We introduce a novel weakly-supervised foreground attention branch which guides our network to focus on important regions and learn features that are important for action recognition. Figure 6 shows weakly this supervised attention branch.

A. Video Super-Resolution

Video SR can be defined as finding sharp appearance and motion details from a low-resolution (LR) video to generate high-resolution (HR) video. We introduce a 3D convolution-based dense video SR (DVSR) network to solve this problem. The problem can be formulated as video-to-video translation, $\hat{V}_{HR} = G(V_{LR})$, where V_{LR} is the low-resolution video clip, \hat{V}_{HR} is the generated high-resolution video output and G is the generator network, termed as DVSR network.

1) *Dense Video Super-Resolution (DVSR) Architecture*: The proposed DVSR network consists of three main components, encoder, decoder, and a projection module. The encoder is responsible for feature extraction, the decoder part focuses on increasing the resolution of the features and the projection module generates the HR videos using those features. The use of residual blocks [8] for image super-resolution tasks has been found effective in improving the image quality due to the similarities between the input and output [22], [50], [26]. Motivated by this, we introduce 3D convolution-based residual-in-residual dense block (3D-RRDB). The 3D-RRDB module consists of a sequence of 3D-RDB modules integrated along with a residual skip connection (Figure 5). Each 3D-RDB module has densely connected five convolution layers [51]. The input of 3D-RDB and the output is merged by a residual connection. 3D-RRDB contains three 3D-RDB modules and a skip connection from the first block to the last block. A detailed architecture of DVSR is shown in Figure 5.

The encoder takes low-resolution video frames and passes them through a 3D convolution layer and several 3D-RDB modules to extract important video features. The decoder part takes LR spatio-temporal features and projects them to HR space. Depending on the scale factor, the feature maps are up-scaled by linear interpolation followed by 3D convolutions. Instead of transposed convolution (fractionally strided convolution), using this strategy prevents checkerboard artifacts [52]. Each up-scaling layer increases the spatial size by a factor of two. After completing spatial enhancement, obtained HR features are given to the projection module which consists of a sequence of 3x3x3 and 1x1x1 convolutions.

2) *Progressive DVSR*: The proposed Progressive DVSR approach learns to increase the resolution in steps, with one scale at a time (Figure 5). We start with a shallow variant of DVSR architecture at the beginning of the training, which only increases the resolution by a factor of two. After it converges, we increase the depth of the encoder and decoder part of DVSR by adding new blocks, so that it can learn to generate features at a scale factor of four. This process is repeated until the desired resolution is obtained. This approach simplifies the problem by dividing it into multiple smaller steps.

The encoder starts with a 3D convolution followed by a 3D-RRDB module. The decoder has one 3D upsampling along with 3D convolutions. Each step uses its projection layer. After progressing to the next step, previous projection layers are omitted. The network takes 16 RGB frames with 14x14 spatial resolution. In the first step, the network produces 4 frames of size 28x28. In the next step, a new 3D-RRDB module is added to the encoder and a new decoder module is added on top of the previous decoder. The network tries to synthesize 8 frames with a resolution of 56x56. We increase the temporal extent along with the spatial resolution. Therefore, the network learns to focus on necessary parts out of 16 frames. This progressive process continues until we have the desired resolution.

Each newly added layer causes a huge degradation in generated video quality. In order to make a smooth transition between progressive steps, we apply a fade-in operation to encoder and decoder separately. After each step, we keep using previous layer outputs along with new block outputs, since the new blocks produce noisy features. We decrease the effect of previous block outputs very smoothly until the new block is trained enough. The fade-in parameter α is set to 0

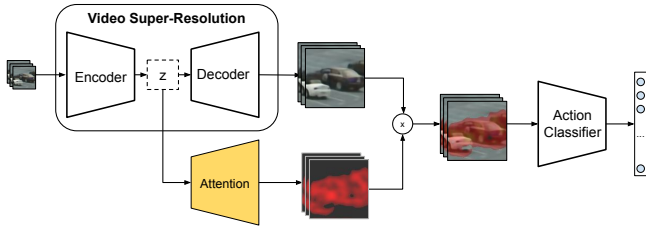


Fig. 6. Overview of weakly-supervised foreground attention branch. It takes encoded features from the super-resolution network and predicts spatio-temporal importance weights for each pixel. During the training super-resolved video clip and the attention map is multiplied and given to the action classifier. The classifier guides the attention branch to highlight foreground regions.

and gradually increased to 1 as the training progresses. In the encoder part, features from the previous layer are added to the new block output. In the decoder part, the output video clips are faded-in. Since each up-sampling block works on a different resolution, the previous projection layer’s output is increased using linear interpolation. The final network architecture after completing the transition is similar to the end-to-end DVSR network.

3) *Foreground Attention*: Super-resolution networks generally focus on texture quality and reconstruction of the whole scene. They do not have knowledge about the foreground or background without guidance. Intuitively, we know that the foreground has much more importance for the action recognition task. To force our DVSR network to attach importance to prominent regions, we add a foreground attention branch to our DVSR network. It takes intermediate features from the encoder and predicts a spatio-temporal importance map for each frame. The predicted weights are used as a weight in DVSR training.

The foreground attention branch is weakly supervised by using an action classifier network. During the training, generated HR video is masked with a predicted attention map and given to the action classifier. If the attention weights cluster around the foreground, action classifier should be able to classify the video. Otherwise, the attention branch will be penalized during the back-propagation. Figure 6 shows the training setup for the attention branch.

Training foreground attention branch by using only gradients from an action classification network can lead to a trivial solution that gives equally high importance to all of the pixels. If weights are equally important and high, all the pixels will be sent to the action classifier without filtering and there will be no feedback to the attention branch. To prevent that we use the L1 norm of the predicted attention map as a penalty so that large region predictions will be discouraged. Figure 7 demonstrates some of our weakly-supervised foreground predictions.

4) *Super-resolution training*: The DVSR network takes LR input and synthesizes HR output. The difference between ground truth and generated video is used as a loss value to update network parameters. We employ a two-stage training

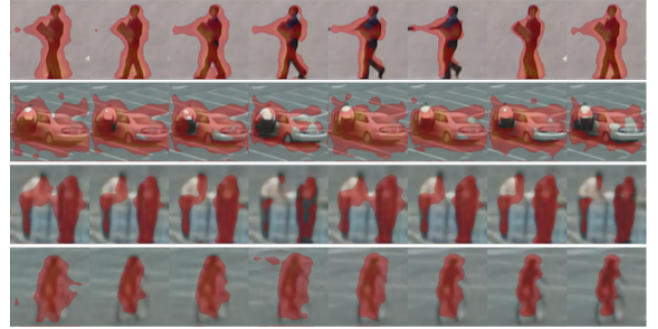


Fig. 7. Examples of attention maps predictions by our weakly-supervised foreground attention branch. Attention maps are concentrated around the foreground actors. The predicted attention maps are used in super-resolution training to weight reconstruction loss. They successfully highlight the important regions for different cases; single actor, multiple actor, person object interaction at very low resolution settings.

strategy for the DVSR network. In the first stage, we pre-train our network progressive or end-to-end approach by using standard reconstruction loss. Afterward, we add the foreground attention branch and use the attention map to weight reconstruction loss.

Reconstruction loss is pixel-wise L1 distance between a ground truth video, V_{HR}^i , and the generated video, \hat{V}_{HR}^i . It is defined as:

$$\mathcal{L}_{rec} = \frac{1}{N} \sum_{i=1}^N \frac{1}{CTHW} \|\hat{V}_{HR}^i - V_{HR}^i\|_1 \odot F_{att}, \quad (1)$$

where N is number of samples in a batch, C, T, H, W are channel size, number of input frame, height and width respectively. F_{att} is the foreground prediction from the attention branch. For the first phase training F_{att} is set to 1.

Foreground attention training The foreground attention branch is trained by an action classifier network. The predicted foreground attention map is applied to the reconstructed \hat{V}_{HR}^i and it is used as an input to the action classifier. We use binary cross-entropy and cross-entropy for multi-label and single-label action classification tasks respectively. Also, we add L1 sparsity constraint to the predicted attention map and calculate the loss for attention branch by

$$\mathcal{L}_{att} = -\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^K y_c^i \log(A(\hat{V}_{HR}^i \odot F_{att})) + \lambda_{att} \|F_{att}\|_1, \quad (2)$$

where A is the action classifier, λ_{att} is the coefficient for the regularization term and y_c^i is the ground truth action labels.

The overall loss \mathcal{L} is computed by combining both reconstruction and attention loss,

$$\mathcal{L} = \lambda_{rec} \mathcal{L}_{rec} + \lambda_{att} \mathcal{L}_{att}, \quad (3)$$

such that each individual loss function is governed by a λ coefficient.

B. Action classification

We use existing action classifier networks I3D [7] and ResNet variants [14] as our backbone structures. We combine

TABLE II
VIDEO ACTION CLASSIFICATION RESULTS FOR TINY VIRAT DATASET.

Method	F1-Score
I3D	28.73
I3D + Prog. DVSR	32.55
I3D + Prog. DVSR + Att.	34.49
ResNet-50	29.08
ResNet-50 + Prog. DVSR	29.81
ResNet-50 + Prog. DVSR + Att.	30.80
WideResNet	32.66
WideResNet + Prog. DVSR	34.05
WideResNet + Prog. DVSR + Attn.	35.07

our DVSR networks and classification network and use them as an end-to-end prediction model. Figure 5 shows the final architecture. The super-resolution part takes the low-resolution videos and increase the spatial size while introducing new details. The synthesized high-resolution video is passed through the action classifier backbone to get final action prediction.

The main idea is that instead of using primitive interpolation methods on tiny action videos, our DVSR network is applied to improve LR video quality. Moreover, we utilize a weakly-supervised foreground attention prediction approach to highlight important features for the action classifiers. To achieve this, we use action classification as an auxiliary task for the super-resolution network during the joint training. After that phase classifier is trained by using final DVSR network outputs without foreground attention branch. We use cross-entropy loss to train the action classifier network,

$$\mathcal{L}_{act} = -\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^K y_c^i \log(A(G(V_{LR}^i))), \quad (4)$$

where V_{LR} is Low-Resolution video, G is the super-resolution network, A is the action classifier, K is the number of classes and y_c^i is the indicator function which is 1 if c is equal to the given video label, otherwise, it will be zero.

V. EXPERIMENTS

A. Datasets and Metrics

We evaluate our approach on the proposed TinyVIRAT; TinyVIRAT has multi-label videos so performance is evaluated by F1 score. For super-resolution training, we crop the video clips from VIRAT following a similar strategy with TinyVIRAT but only allowing larger than 112x112 clips. To create low and high-resolution video pairs for super-resolution training, we down-scale and then up-scale the videos by using bicubic interpolation.

In addition, we evaluate the performance of our system on two publicly available action dataset UCF-101 [1] and HMDB-51 [15] in order to compare with existing methods.

B. Implementation Details

DVSR is trained by setting λ_{rec} to 1 and λ_{att} to 0.5 in Equation 3 during the joint training. Both DVSR and action classification networks are trained with Adam optimizer [53] and we use β_1 as 0.5 and β_2 as 0.9. Both the DVSR network

TABLE III
VIDEO ACTION CLASSIFICATION COMPARISON FOR HMDB-51 DATASET.

Method	Input	Accuracy %
I3D	112x112	52.61
SoSR [12]	80x60	54.77
Bicubic - I3D	14x14	10.59
Privacy-Preserv [19]	12x16	28.68
F. Coupled [32]	12x16	39.15
DVSR	14x14	41.24
Prog. DVSR	14x14	41.63
Bicubic - I3D	28x28	46.97
Privacy-Preserv [19]	24x32	32.15
DVSR	28x28	53.66
Prog. DVSR	28x28	55.95

and the action classifier are trained with a learning rate of 0.0002. The super-resolution network is trained without using any pre-trained model weights. The value of α for layer transition in a progressive approach is set to 0 initially and increased by $5e-3$ after each iteration. This step size is set empirically and can be estimated based on the batch size and the number of epochs required for convergence.

For the action classification task, the I3D action classifier network is used with pre-trained weights which were obtained by training I3D on Charades dataset [43]. ResNet architectures are pre-trained on Kinetics [2] dataset, we obtain the model parameters from [14].

C. Quantitative Results

We first train the action classifier networks standalone in order to set the baselines for TinyVIRAT. Spatial size of the low-resolution videos is resized to 112x112 by using bicubic interpolation. Then we apply our progressively trained DVSR network and foreground attention approaches to demonstrate improvement. For each experiment, the backbone action classifiers are initialized with the same pre-trained weights.

Table II shows action classification performance baselines and our approach on TinyVIRAT dataset. Our DVSR network shows favorable improvement comparing to baseline network results. We experimented with different action classifier architectures to show that our DVSR is not biased to a certain network but improves the results consistently for all the models. After introducing our weakly-supervised foreground attention approach we improve the baseline scores by a large margin.

We also compare our method with previous work [19] [32] and [12] on public datasets. For a fair comparison, we compare our method with other non-optical flow based methods. Table III and Table IV show the results on HMDB-51 and UCF-101 respectively.

D. Ablation study

We have already shown that using super-resolution improves the action classification results. We experiment with different strategies to study the variation in the performance of the action classification task. Table V shows the performance for different training strategies. For each experiment, the action classifier part uses the same action classification loss and for

TABLE IV
VIDEO ACTION CLASSIFICATION COMPARISON FOR UCF-101 DATASET.

Method	Input	Accuracy %
I3D	112x112	84.72
SoSR [12]	80x60	83.92
Bicubic - I3D	14x14	14.14
DVSR	14x14	68.17
Prog. DVSR	14x14	70.55
Bicubic - I3D	28x28	66.72
DVSR	28x28	82.37%
Prog. DVSR	28x28	82.87

SR experiments we use DVSR network. We only change the generator training strategy while training.

1) *Non-progressive Super-Resolution Training*: We show that using the standard end-to-end super-resolution approach, action classifier performance can be boosted. In Table V, the first two rows show the effect of using a standard super-resolution approach.

2) *Progressive Training*: Progressive training strategy simplifies the optimization process for the super-resolution problem. Instead of learning the mapping between low-resolution to high-resolution, it breaks down the process into smaller tasks. As we expected the proposed progressive DVSR model provides a significant improvement over the baseline approach. It also improves the standard super-resolution performance as we can see in the third row of Table V.

3) *Foreground Attention Branch*: Using weakly-supervised foreground attention branch while super-resolution network training gives us the best result. It guides DVSR network to focus on foreground regions that have more meaningful information for the action recognition task. Figure 7 shows qualitative results for super-resolution and attention maps. From the visual results, we can see that when the quality of the video is low, distinguishing the background and foreground becomes more difficult but our attention prediction successfully concentrates around the actors. The improvement we get from attention guidance is orthogonal to progressive training strategy, and using both of them together lead us the best performance as we can see in the last row of Table V.

VI. CONCLUSION

We introduce a new tiny action recognition benchmark dataset *TinyVIRAT* which consists of natural low-resolution videos. We propose a novel tiny video action classification framework which incorporates progressively growing video super-resolution network to improve tiny action recognition performance. We utilize a 3D convolution-based dense residual network and weakly-supervised foreground attention branch which helps in learning effective appearance and motion features from low-resolution videos. We also demonstrate that the enhanced videos using the progressive DVSR network learn important appearance and motion features which is beneficial for action recognition. We perform experiments on two artificial benchmark datasets and demonstrate that the proposed approach leads to a better action recognition performance on both artificially created and natural low-resolution videos. Our

TABLE V
ABLATION RESULTS. OUR NOVEL PROGRESSIVE TRAINING STRATEGY AND WEAKLY-SUPERVISED ATTENTION APPROACH IMPROVES THE PERFORMANCE OF STANDARD ACTION CLASSIFIERS. WE USE I3D ARCHITECTURE FOR THESE EXPERIMENTS.

Method	F1-Score
w/o DVSR	28.73
DVSR	30.45
Progressive DVSR	32.55
Progressive DVSR + Attention	34.49

super-resolution based tiny action classification framework can be integrated into any video analysis pipeline without much effort for other problems such as semantic segmentation, object localization, and object tracking.

VII. ACKNOWLEDGEMENT

This research is based upon work supported by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), via IARPA R&D Contract No. D17PC00345. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the ODNI, IARPA, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon.

REFERENCES

- [1] K. Soomro, A. R. Zamir, and M. Shah, "UCF101: A dataset of 101 human actions classes from videos in the wild," *CoRR*, vol. abs/1212.0402, 2012.
- [2] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, M. Suleyman, and A. Zisserman, "The kinetics human action video dataset," *CoRR*, vol. abs/1705.06950, 2017.
- [3] M. Monfort, A. Andonian, B. Zhou, K. Ramakrishnan, S. A. Bargal, T. Yan, L. Brown, Q. Fan, D. Gutfreund, C. Vondrick *et al.*, "Moments in time dataset: one million videos for event understanding," *arXiv preprint arXiv:1801.03150*, 2018.
- [4] C. Gu, C. Sun, D. A. Ross, C. Vondrick, C. Pantofaru, Y. Li, S. Vijayanarasimhan, G. Toderici, S. Ricco, R. Sukthankar *et al.*, "Ava: A video dataset of spatio-temporally localized atomic visual actions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6047–6056.
- [5] S. Abu-El-Haija, N. Kothari, J. Lee, P. Natsev, G. Toderici, B. Varadarajan, and S. Vijayanarasimhan, "Youtube-8m: A large-scale video classification benchmark," *arXiv preprint arXiv:1609.08675*, 2016.
- [6] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3d convolutional networks," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 4489–4497.
- [7] J. Carreira and A. Zisserman, "Quo vadis, action recognition? A new model and the kinetics dataset," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, 2017, pp. 4724–4733.
- [8] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, 2016, pp. 770–778.
- [9] 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 *European conference on computer vision*. Springer, 2016, pp. 20–36.

- [10] J. Chen, J. Wu, J. Konrad, and P. Ishwar, "Semi-coupled two-stream fusion convnets for action recognition at extremely low resolutions," in *2017 IEEE Winter Conference on Applications of Computer Vision, WACV 2017, CA, USA, March 24-31, 2017*, 2017, pp. 139–147.
- [11] M. Xu, A. Sharghi, X. Chen, and D. J. Crandall, "Fully-coupled two-stream spatiotemporal networks for extremely low resolution action recognition," in *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, vol. 00, Mar 2018, pp. 1607–1615.
- [12] H. Zhang, D. Liu, and Z. Xiong, "Two-stream action recognition-oriented video super-resolution," in *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
- [13] S. Oh, A. Hoogs, A. Perera, N. Cuntoor, C.-C. Chen, J. T. Lee, S. Mukherjee, J. Aggarwal, H. Lee, L. Davis *et al.*, "A large-scale benchmark dataset for event recognition in surveillance video," in *CVPR 2011*. IEEE, 2011, pp. 3153–3160.
- [14] K. Hara, H. Kataoka, and Y. 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 (CVPR)*, 2018, pp. 6546–6555.
- [15] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre, "HMDB: a large video database for human motion recognition," in *Proceedings of the International Conference on Computer Vision (ICCV)*, 2011.
- [16] Z. Wang, S. Chang, Y. Yang, D. Liu, and T. S. Huang, "Studying very low resolution recognition using deep networks," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [17] Y. Bai, Y. Zhang, M. Ding, and B. Ghanem, "Finding tiny faces in the wild with generative adversarial network," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [18] J. Dai, B. Saghaei, J. Wu, J. Konrad, and P. Ishwar, "Towards privacy-preserving recognition of human activities," *2015 IEEE International Conference on Image Processing (ICIP)*, pp. 4238–4242, 2015.
- [19] M. S. Ryoo, B. Rothrock, C. Fleming, and H. J. Yang, "Privacy-preserving human activity recognition from extreme low resolution," in *AAAI*, 2017.
- [20] M. S. Ryoo, K. Kim, and H. J. Yang, "Extreme low resolution activity recognition with multi-siamese embedding learning," in *AAAI*, 2018.
- [21] C. Dong, C. C. Loy, K. He, and X. Tang, "Image super-resolution using deep convolutional networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 2, pp. 295–307, Feb. 2016.
- [22] J. Kim, J. K. Lee, and K. M. Lee, "Accurate image super-resolution using very deep convolutional networks," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR Oral)*, June 2016.
- [23] W.-S. Lai, J.-B. Huang, N. Ahuja, and M.-H. Yang, "Deep laplacian pyramid networks for fast and accurate super-resolution," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [24] B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee, "Enhanced deep residual networks for single image super-resolution," in *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2017, July 21-26, 2017*, 2017, pp. 1132–1140.
- [25] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi, "Photo-realistic single image super-resolution using a generative adversarial network," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [26] M. S. M. Sajjadi, B. Schölkopf, and M. Hirsch, "EnhanceNet: Single image super-resolution through automated texture synthesis," in *Proceedings IEEE International Conference on Computer Vision (ICCV)*. Piscataway, NJ, USA: IEEE, Oct. 2017, pp. 4501–4510.
- [27] T. Tong, G. Li, X. Liu, and Q. Gao, "Image super-resolution using dense skip connections," in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [28] J. Li, F. Fang, K. Mei, and G. Zhang, "Multi-scale residual network for image super-resolution," in *The European Conference on Computer Vision (ECCV)*, September 2018.
- [29] P. Abolghasemi, A. Mazaheri, M. Shah, and L. Boloni, "Pay attention!-robustifying a deep visuomotor policy through task-focused visual attention," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4254–4262.
- [30] X. Deng, R. Yang, M. Xu, and P. L. Dragotti, "Wavelet domain style transfer for an effective perception-distortion tradeoff in single image super-resolution," in *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
- [31] M. S. M. Sajjadi, R. Vemulapalli, and M. Brown, "Frame-Recurrent Video Super-Resolution," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [32] J. Caballero, C. Ledig, A. P. Aitken, A. Acosta, J. Totz, Z. Wang, and W. Shi, "Real-time video super-resolution with spatio-temporal networks and motion compensation," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2848–2857, 2017.
- [33] X. Tao, H. Gao, R. Liao, J. Wang, and J. Jia, "Detail-revealing deep video super-resolution," in *IEEE International Conference on Computer Vision, ICCV 2017, October 22-29, 2017*, 2017, pp. 4482–4490.
- [34] D. Liu, Z. Wang, Y. Fan, X. Liu, Z. Wang, S. Chang, and T. Huang, "Robust video super-resolution with learned temporal dynamics," in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [35] P. Yi, Z. Wang, K. Jiang, J. Jiang, and J. Ma, "Progressive fusion video super-resolution network via exploiting non-local spatio-temporal correlations," in *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
- [36] A. Kappeler, S. Yoo, Q. Dai, and A. Katsaggelos, "Video super-resolution with convolutional neural networks," *IEEE Transactions on Computational Imaging*, vol. 2, pp. 1–1, 06 2016.
- [37] W. Shi, J. Caballero, F. Huszar, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang, "Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [38] Y. Huang, W. Wang, and L. Wang, "Bidirectional recurrent convolutional networks for multi-frame super-resolution," in *Advances in Neural Information Processing Systems 28*, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, Eds., 2015, pp. 235–243.
- [39] Y. Jo, S. Wug Oh, J. Kang, and S. Joo Kim, "Deep video super-resolution network using dynamic upsampling filters without explicit motion compensation," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [40] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of gans for improved quality, stability, and variation," in *ICLR*. OpenReview.net, 2018.
- [41] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei, "Large-scale video classification with convolutional neural networks," in *CVPR*, 2014.
- [42] C. Xu, S.-H. Hsieh, C. Xiong, and J. J. Corso, "Can humans fly? Action understanding with multiple classes of actors," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2015.
- [43] G. A. Sigurdsson, G. Varol, X. Wang, A. Farhadi, I. Laptev, and A. Gupta, "Hollywood in homes: Crowdsourcing data collection for activity understanding," in *European Conference on Computer Vision*, 2016.
- [44] S. Abu-El-Haija, N. Kothari, J. Lee, P. Natsev, G. Toderici, B. Varadarajan, and S. Vijayanarasimhan, "Youtube-8m: A large-scale video classification benchmark," *arXiv preprint arXiv:1609.08675*, 2016.
- [45] B. Schiele, "A database for fine grained activity detection of cooking activities," in *Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, ser. CVPR '12, 2012, pp. 1194–1201.
- [46] C. Gu, C. Sun, D. A. Ross, C. Vondrick, C. Pantofaru, Y. Li, S. Vijayanarasimhan, G. Toderici, S. Ricco, R. Sukthankar, C. Schmid, and J. Malik, "Ava: A video dataset of spatio-temporally localized atomic visual actions," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [47] B. G. Fabian Caba Heilbron, Victor Escorcia and J. C. Niebles, "Activitynet: A large-scale video benchmark for human activity understanding," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 961–970.
- [48] Y.-G. Jiang, J. Liu, A. Roshan Zamir, G. Toderici, I. Laptev, M. Shah, and R. Sukthankar, "THUMOS challenge: Action recognition with a large number of classes," <http://csrc.ucf.edu/THUMOS14/>, 2014.
- [49] D. Damen, H. Doughty, G. M. Farinella, S. Fidler, A. Furnari, E. Kazakos, D. Moltisanti, J. Munro, T. Perrett, W. Price, and M. Wray, "Scaling egocentric vision: The epic-kitchens dataset," in *European Conference on Computer Vision (ECCV)*, 2018.
- [50] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. C. Loy, "Esrgan: Enhanced super-resolution generative adversarial networks," in *The European Conference on Computer Vision Workshops (ECCVW)*, September 2018.

- [51] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, 2017, pp. 2261–2269.
- [52] A. Odena, V. Dumoulin, and C. Olah, "Deconvolution and checkerboard artifacts," *Distill*, 2016. [Online]. Available: <http://distill.pub/2016/deconv-checkerboard>
- [53] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *ICLR*, 2015.