



## Human Action Recognition in Drone Videos using a Few Aerial Training Examples

Waqas Sultani<sup>a,\*\*</sup>, Mubarak Shah<sup>b</sup>

<sup>a</sup>Intelligent Machine Lab, Information Technology University, Lahore, Pakistan

<sup>b</sup>Center for Research in Computer Vision, University of Central Florida, Orlando, USA

### ABSTRACT

Drones are enabling new forms of human actions surveillance due to their low cost and fast mobility. However, using deep neural networks for automatic aerial action recognition is difficult due to the need for a large number of training aerial human action videos. Collecting a large number of human action aerial videos is costly, time-consuming, and difficult. In this paper, we explore two alternative data sources to improve aerial action classification when only a few training aerial examples are available. As a first data source, we resort to video games. We collect plenty of aerial game action videos using two gaming engines. For the second data source, we leverage conditional Wasserstein Generative Adversarial Networks to generate aerial features from ground videos. Given that both data sources have some limitations, e.g. game videos are biased towards specific actions categories (fighting, shooting, etc.), and it is not easy to generate good discriminative GAN-generated features for all types of actions, we need to efficiently integrate two dataset sources with few available real aerial training videos. To address this challenge of the heterogeneous nature of the data, we propose to use a disjoint multi-task learning framework. We feed the network with real and game, or real and GAN-generated data in an alternating fashion to obtain an improved action classifier. We validate the proposed approach on two aerial action datasets and demonstrate that features from aerial game videos and those generated from GAN can be extremely useful for an improved action recognition in real aerial videos when only a few real aerial training examples are available.

© 2021 Elsevier Ltd. All rights reserved.

### 1. Introduction

Nowadays, drones are ubiquitous and actively being used in several applications such as sports, entertainment, agriculture, forest monitoring, military, and surveillance Dutta and Ekenna (2019); Huang et al. (2018); Zhou et al. (2018). In video surveillance, drones can be much more useful than CCTV cameras due to their freedom of mobility and low cost. One critical task in video surveillance is monitoring human actions using drones.

Automatically recognizing human action in drone videos is

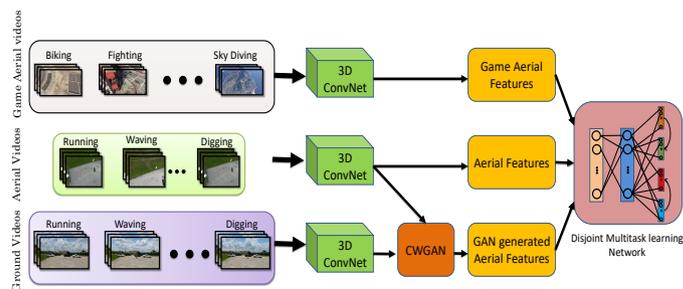
a daunting task. It is challenging due to drone camera motion, small actor size, and most importantly the difficulty of collecting large scale training aerial action videos. Computer vision researchers have tried to detect human action in varieties of videos including sports videos (Soomro et al., 2013), surveillance CCTV videos Sultani et al. (2018), cooking and ego-centric videos Damen et al. (2018). However, despite being very useful and of practical importance, not much research work is done to automatically recognize human action in drone videos.

Deep learning models are data-hungry and need hundreds of training video examples for robust training. However, collecting training dataset is quite challenging in several vision

\*\*Corresponding author: Tel.: +92-3365109108;

e-mail: waqas.sultani@itu.edu.pk (Waqas Sultani)

applications. To address this difficulty of real data collection and its annotations, recently researchers have used games and synthetic images in several computer vision applications such as semantic segmentation Richter et al. (2016), measuring 6D object pose Mercier et al. (2018), and depth image classification Carlucci et al. (2016). Inspired by the use of video games in Richter et al. (2016); Mercier et al. (2018); Carlucci et al. (2016), we propose to collect and use game action videos to improve human action recognition in real-world aerial videos. Recently, computer graphics techniques and gaming technology have improved significantly. For example, GTA (Grand Theft Auto) and FIFA (Federation International Football Association) gaming engines use photo-realistic simulators to render real-world environment, texture, objects (human, bicycle, car, etc) and human actions. Games videos for action recognition are intriguing because 1) without much effort, one can collect a large number of videos containing environment and motion that looks close to real-world, 2) It is easy Richter et al. (2016) to get detailed annotations for action detection and segmentation which are otherwise very expensive to obtain, 3) Most of the gaming engines allow the players to simultaneously capture the same action from the different views (aerial, ground, front, etc.). This means that we can easily collect a large scale multi-view dataset with exact frame-by-frame correspondence. All three advantages make gaming videos quite appealing for aerial action recognition where data collection is difficult and expensive. To the best of our knowledge, we are the first one to use game videos in aerial action recognition research. Another direction to address the scarcity of data is to use GAN-generated video examples generated through generative adversarial networks Goodfellow et al. (2014). Although the quality of images and videos generated by GAN is not yet good enough to train deep networks Xian et al. (2018), GAN generated discriminative features may be still suitable for action classification. Therefore, we propose to employ conditional Wasserstein GAN Arjovsky et al. (2017); Cao et al. (2018) to generate discriminative features. We believe that the GAN-generated aerial examples, when integrated properly with a few real action examples, can



**Figure 1. Summary of the proposed training approach. We propose to utilize game videos and GAN generated aerial features to improve aerial action classification when a few real aerial training examples are available. Our approach does not require the same labels for real and game actions. To tackle different action labels in the game and real dataset, we propose to use disjoint multitask learning framework to efficiently learn robust action classifier.**

help learn a more generalized and robust aerial action classifier.

In this paper, we propose to utilize game video features and GAN-generated features to improve aerial action classification when a few real aerial training examples are available (see Figure 1). However, one of the key challenges is the disjoint nature of the problem. Video games are designed to address the interest of game playing audience and contain human motions and environments biased towards a few specific human actions. For example, the majority of actions in FIFA games are related to playing a soccer game in a soccer field and the majority of actions in GTA are about fighting. Therefore, it is highly likely that classes of actions in games are different from the types of action classes we are interested to recognize in the real world. Similarly, it is not easy to generate good discriminative GAN-generated features for all types of action. However, our key idea is that despite different classes in games and real videos and the low-quality nature of GAN-generated aerial features, all three data types (games, real and GAN-generated) capture similar local motion patterns, human movements, and human-object interactions, and, if integrated properly, can help learn more accurate aerial action classifiers. To achieve this, we combine games and GAN-generated examples with a few available real training examples using disjoint multitask learning. Specifically, we feed the network with real and game (or GAN-generated) data in an alternating fashion to obtain a more accurate action classi-

---

fier. Note that in this paper, we call the videos as ground action videos if the person recording the videos is on the ground or at side-angle and the aerial videos are the ones that are taken by UAVs. In summary, this paper makes the following contributions:

- We propose to tackle the new problem of drone-based human action recognition when only a few aerial training examples are available.
- To the best of our knowledge, we are the first one to demonstrate the feasibility of game action videos for improving action recognition in real-world aerial videos. Although game imagery has been used before in different computer vision applications, it has not been used for aerial action recognition.
- We show that game and GAN-generated action examples can help to learn a more accurate action classifier through a disjoint multitask learning framework.
- We present two new action datasets: 1) Aerial-Ground game dataset containing seven human actions where for each action we have 100 aerial-ground video pairs, 2) Real aerial dataset containing actions corresponding to eight actions of UCF101.

## 2. Related Work

Human action recognition in videos is one of the most challenging and active vision problems Ali and Shah (2010); Wang and Schmid (2013); Sultani and Saleemi (2014); Carreira and Zisserman (2017); Tran et al. (2015); Wu et al. (2011); Simonyan and Zisserman (2014); Chen et al. (2018); Kataoka and Satoh (2019); Huang et al. (2018); Zhou et al. (2018). Classical approaches used hand-crafted features Ali and Shah (2010); Wang and Schmid (2013) to train generalized human action recognition models that can perform well across different action datasets Cao et al. (2010); Sultani and Saleemi (2014).

With the resurgence of deep learning, several deep learning approaches have been proposed for action recognition. Simonyan et al. Simonyan and Zisserman (2014) proposed RGB and optical flow-based networks for action recognition videos. Both RGB and optical flow networks employ 2D convolution.

Tran et al. Tran et al. (2015) demonstrated the feasibility of 3D convolution for action recognition. In addition to presenting a new large scale action recognition dataset of 400 classes, Carreira et al. Carreira and Zisserman (2017) proposed a two-stream inflated 3D ConvNet (I3D) that is based on 2D convnet inflation and demonstrated state of the art classification accuracy. Recently, an efficient action recognition framework is proposed by Chen et al. Chen et al. (2018). Furthermore, there has been an increased interest to train the generalized action recognition model using multi-task learning. Kataoka et al. Kataoka and Satoh (2019) put forwarded a multi-task approach for the out-of-context action understanding. Similarly, Kim et al. Kim et al. (2018) proposed disjoint multi-task learning to obtain improved video action classification and captioning in a joint framework.

Recently, Zhou et al. Zhou et al. (2018) proposed to analyze human motion using videos that are captured through a drone that orbits around the person. They demonstrated that, as compared to static cameras, videos captured by drones are more suitable for better motion reconstruction. Similarly, Huang et al. Huang et al. (2018) presented a system that can detect cinematic human actions using 3D skeleton points employing a drone.

Although human action recognition is quite an active area of research in computer vision, there does not exist many research works in the literature that deals with aerial action recognition. Wu et al., Wu et al. (2011) proposed to use low-rank optimization to separate objects and moving camera trajectories in aerial videos. UCF-ARG dataset Arjun Nagendran and Shah contains ground, rooftop, and aerial triplets of 10 realistic human actions. This dataset is quite challenging as it contains severe camera motion, non-discriminative backgrounds, and humans in these videos occupy only a few pixels. Perera et al. Perera et al. (2018) proposed to use human pose features to detect gestures in aerial videos. They introduced a dataset that is recorded by a slow and low-altitude (around 10ft) UAV. Although useful, their dataset only contains gestures related to UAV navigation and aircraft handling. Recently, Barekattain et al. Barekattain

et al. (2017) proposed a new video dataset for aerial view concurrent human action detection. It consists of 43 minute-long fully-annotated sequences with 12 action classes. They used a single-shot detection approach Liu et al. (2016) to obtain human bounding boxes and then used features within those bounding boxes for action classification. They have neither addressed the problem for the less number of training videos nor they have used the multiple data sources.

Gathering large-scale datasets and its annotation is expensive and requires hundreds of human hours. To address this challenge, there is an increasing interest in employing synthetic data to train deep neural networks. Josifovski et al. Josifovski et al. (2018) proposed to use annotated synthetic data to train the instance-based object detector and 3D pose estimator. Mercier et al. Mercier et al. (2018) used weakly labeled images and synthetic images to train a deep network for object localization and 6D pose estimation in real-world settings. Carlucci et al. Carlucci et al. (2016) proposed to use synthetic data for depth image classification. Recently, Richter et al., Richter et al. (2016) designed a method to automatically gather ground truth data for semantic segmentation and Hong et al. (2018) presented a GAN based approach to use game annotations for semantic segmentation in real images. Finally, Mueller et al., Mueller et al. (2016) put forwarded photo-realistic simulators to render real-world environment and provide a benchmark for evaluating tracker performance.

In this paper, in contrast to the above-mentioned methods, we demonstrate the feasibility of game action videos for improving action recognition in real-world *aerial* videos. To evade collecting costly drone training videos, we claim to provide a unified framework employing games and GAN-generated data to achieve improved aerial action recognition. No one has used disjoint multitask learning for aerial action recognition or with three different data types. Note that although we use real ground videos for GAN-generated features, our approach does not require exact ground-aerial pairs. Furthermore, our game dataset collection highlights the built-in multi-view action capturing feature in games that can be used for multi-view action recogni-

tion. Multiple views make the dataset more extensive and open avenues for other researchers to solve novel challenging problems. We believe that our dataset will push the research in joint game-real aerial action recognition.

### 3. Proposed Approach

In this section, we provide the details of our game actions video collection, the method to generated GAN-generated features, and finally disjoint multitask approach where we train the aerial action classifier using game aerial, GAN-generated aerial, and a few real aerial videos in a unified framework.

#### 3.1. Games Action Dataset



**Figure 2.** Two frames of each action for both aerial and ground views from our game action dataset. The first, third, fifth and seventh row represent aerial videos and second, forth, sixth, and eighth row shows the ground videos. For each action, we show the two frames per video.

We employ GTA-5 (Grand Theft Auto) and FIFA (Federation International Football Association) for collecting the game ac-

tion dataset <sup>1</sup>. We ask the players to play the games and record the same action from multiple views. Note that GTA and FIFA allow users to record the actions from multiple angles with real-looking scenes and different realistic camera motions. In total, we collect seven human actions including cycling, fighting, soccer kicking, running, walking, shooting, and skydiving. Due to the availability of plenty of soccer kicking in FIFA games, we collect kicking from FIFA and the rest of the actions are collected from GTA-5. Although in our current approach we are only using aerial game video, for more complete dataset purposes, we capture both ground and aerial video pairs i.e., the same action frames captured from both aerial and ground cameras. Figure 2 shows two frames of each action for both aerial and ground views. These videos will be made publicly available.

For each action, our dataset contains 200 videos (100 ground and 100 aerial) with a total of 1400 videos for seven actions. Note that most of the scenes and interactions in the video games are biased towards actions related to fighting, shooting, walking and running, etc. Therefore, employing game videos to improve action recognition in real-world videos is not trivial. Therefore, in this paper, we propose a unified approach to combine games and real videos employing disjoint multitask learning.

### 3.2. GAN-generated Aerial Examples

We generate GAN-generated aerial video features using Generative Adversarial Networks (GAN) Goodfellow et al. (2014). GAN consists of two networks: Generator and Discriminator. Generator tries to mimic the real data distribution and fools the discriminator by producing realistic looking videos or features while the discriminator job is to robustly classify real and generated video or features. Both generator and discriminator can be simple multi-layer perceptrons. As compared to vanilla-GAN, in conditional GAN Mirza and Osindero (2014), both generator and discriminator are conditioned on auxiliary information. Auxiliary information can be video labels or some other video

features. Our goal is to generate GAN-generated aerial visual features given the real ground features (auxiliary information). Therefore, in our case, the objective function of conditional GAN is given by:

$$\mathcal{L}_{cgan} = \mathbb{E}[\log D(f_{r_a}|f_{r_g})] + \mathbb{E}[\log(1 - D(G(z, f_{r_g})|f_{r_g}))], \quad (1)$$

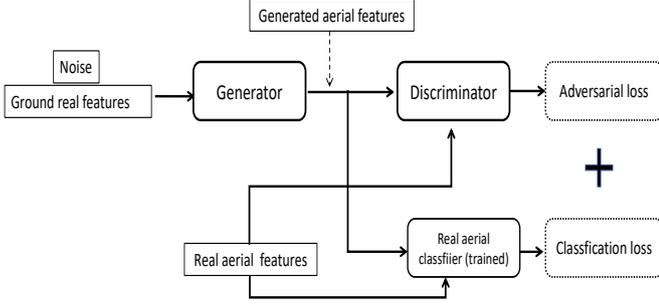
where  $D$  represents discriminator and  $G$  represents generator, in  $D(f_{r_a}|f_{r_g})$ ,  $f_{r_a}$  and  $f_{r_g}$  are real aerial and ground features respectively. These features are randomly sampled from given real aerial and ground features distributions. Note that we do not assume any correspondence between  $f_{r_a}$  and  $f_{r_g}$ . Given the noise vector  $z$  and  $f_{r_g}$ , generator tries to fool discriminator by producing GAN-generated aerial features.

To optimize the above objective function, usually KL or JS divergence is employed to reduce the difference between real and generated data distributions. However, one of the key limitations with KL or JS divergence is that the gradient of divergence decreases with the increase of distance, and the generator learns nothing through gradient descent. To address this limitation, recently Wasserstein GAN is introduced in Arjovsky et al. (2017); Cao et al. (2018), which uses Wasserstein distance. WGAN learns better because it has a smoother gradient everywhere. Finally, to make Wasserstein distance tractable, the 1-Lipschitz constraint is used through gradient penalty loss Gulrajani et al. (2017). The objective function of generating GAN-generated aerial features using conditional Wasserstein GAN (WCGAN-GP) is given by:

$$\mathcal{L}_{cwgan} = \mathbb{E}[D(G(z, f_{r_g})|f_{r_g})] - \mathbb{E}[\log D(f_{r_a}|f_{r_g})] + \mathbb{E}[(\|\nabla_m D(m, (G(z, f_{r_g}))\|_2 - 1)^2], \quad (2)$$

where  $m = tG(z, f_{r_g}) + (1 - t)f_{r_g}$  and  $t$  is uniformly sampled between 0 and 1. Our ultimate goal is to train discriminative action classifiers using GAN-generated features. Although the above objective function generates realistically looking features, it does not guarantee to generate the discriminative features suitable for classification. To accomplish this, we first train soft-max classifiers using a few available real aerial examples. Finally, to enforce WCGAN-GP to produce discriminative features, we use classification loss computed over the

<sup>1</sup>The customized engine such as Lai et al. (2018) can also be used to collect more game videos



**Figure 3. GAN-generated aerial features generation pipeline.** Given ground real features, noise and a few real aerial videos, employing adversarial and classification loss, GAN-generated aerial features are generated.

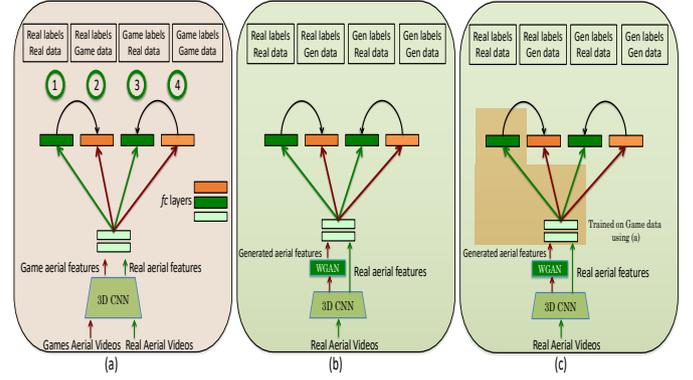
GAN-generated aerial examples given as:

$$\mathcal{L}_{cl} = -E[\log P(y_{r_g}|G(z, f_g); \theta)], \quad (3)$$

where  $P(y_{r_g}|G(z, f_g))$  denotes the probability of correct label prediction of generated examples. Since labels for real ground and GAN-generated aerial examples are the same, we use the labels of real ground videos ( $y_{r_g}$ ) as ground truth. Finally, the overall objective function for GAN-generated aerial examples generation is given by combination of Eq. 2 and Eq. 3. Figure 3 summarize the complete GAN-generated aerial feature generation scheme.

### 3.3. Aerial videos classification using Disjoint Multi-Task learning

Multitask learning improves the generalization capabilities of the model by effectively learning multiple related tasks. It has been used in several computer vision problems to learn the joint model such as; simultaneous object detection and segmentation Hariharan et al. (2014), surface normal estimation, and pixel labeling Misra et al. (2016) and joint pose estimation and action recognition Gkioxari et al. (2014). One of the limitations of multitask learning is the requirement of the availability of multiple labels for each task for the *same* data. However, most of existing action datasets do not have such labels and hence restricting multitask learning on these datasets. To address this, recently disjoint multitask learning Kim et al. (2018) is introduced. Since, in our approach, the two data sources (games and real) are different, and secondly, we do not assume any common



**Figure 4. Disjoint multitasking framework for aerial, games and GAN-generated examples.** (a) Disjoint multitask learning (DML) using games and a few real aerial videos. (b) DML using GAN-generated aerial and real aerial video features (c) DML using real, GAN-generated, and game data. GAN-generated aerial features are abbreviated as Gen data.

action classes, this fits well in the context of disjoint multitask learning.

Figure 4 demonstrates the model overview. We first compute deep features of a few available real aerial videos and game videos using a 3D convolutional neural network Chen et al. (2018); Carreira and Zisserman (2017); Hara et al. (2018). Secondly, we obtain GAN-generated aerial features using the method described in Section 3.2. We use two fully connected layers shared between all tasks and one dedicated fully connected layer for each task. Figure 4.a shows training using real and game visual features, Figure 4.b shows training using real and GAN-generated visual features, and Figure 4.c demonstrates training using the real, game and GAN-generated data.

**Joint learning using real and game videos:** We denote the real and game aerial visual features as  $r_a \in \mathcal{R}_a$ ,  $g_a \in \mathcal{G}_a$  respectively. Note that we do not assume the type or the number of actions in both datasets to be the same. We train all four branches of the network for the classification using softmax as a final activation function along with cross-entropy loss.

As shown in Figure 4 (a), we have real and game labels available for real and game data (branch ① and ④). However, we do not have real labels for the game data and the game labels for the real data (branch ② and ③) due to the disjoint nature of two datasets. We train different branches of the multi-task

framework using the aerial real and game iteratively. First, we train real and game classification branches (① and ④) for which we have corresponding ground truth labels available. Next, we train branch ② and ③ which predicts the real labels for the game data and game labels for real data respectively. However, due to unavailability of labels for ② and ③, we use the prediction from ① and ④ as a ground truth labels for ② and ③ respectively. The overall objective function of the framework is given by:

$$\begin{aligned} & \min_{\Theta} \sum_{r_a \in \mathcal{R}_a} \overbrace{\mathcal{L}(y_{r_a}, P(y_{r_a}|r_a))}^{\textcircled{1}} + \overbrace{\mathcal{L}(y_{g_a}^{\hat{}}, P(y_{g_a}|r_a))}^{\textcircled{3}} \\ & + \min_{\Theta} \sum_{g_a \in \mathcal{G}_a} \overbrace{\mathcal{L}(y_{r_a}^{\hat{}}, P(y_{r_a}|g_a))}^{\textcircled{2}} + \overbrace{\mathcal{L}(y_{g_a}, P(y_{g_a}|g_a))}^{\textcircled{4}} \end{aligned} \quad (4)$$

where  $y_{r_a}$  and  $P(y_{r_a}|r_a)$  represents ground truth labels and predicted labels of real aerial videos,  $r_a, y_{g_a}^{\hat{}}$  are the labels obtained from ④ (the layer trained with game ground truth labels) and  $P(y_{g_a}|r_a)$  are predicted game action labels for real videos. Similarly,  $y_{g_a}$  and  $P(y_{g_a}|g_a)$  are ground truth and predicted action labels of game aerial videos,  $y_{r_a}^{\hat{}}$  are obtained from ① (the layer trained with real aerial ground truth labels) and  $P(y_{r_a}|g_a)$  is predicted real action labels for game videos. Finally  $\Theta$  represents network parameters. We repeat the above procedure for several epochs and fine-tune the parameters on the validation data.

**Joint learning using GAN-generated and real videos:** For joint learning using real and GAN-generated video features, we repeat the same approach as discussed above for real and game videos. Specifically, to obtain an improved action classifier, we feed the network with real and GAN-generated aerial features in an alternating fashion to a disjoint multitask learning framework. Figure 4(b) illustrates the joint learning using real and GAN-generated video features. Note that since we use real ground features to generate GAN-generated aerial examples (see Figure 3), we use the labels of real ground videos ( $y_{r_g}$ ) as GAN-generated aerial examples labels, abbreviated as ‘Gen labels’ in the third and fourth branch in Figure 4(b). As shown in Figure 4(a) and Figure 4(b), we use the same network architecture for joint learning using games or GAN-generated data.

#### Joint learning using real, GAN-generated and game videos:

Finally, we combine all three data types i.e., real, GAN-generated, and game videos in a single framework (see Figure 4(c)). In this case, instead of training the network from scratch, we initialize the network with the weights obtained through training using game data. Specifically, networks weights of Figure 4(a) are used to initialize weights of two backbone ( $f_c$ ) layers and ( $f_c$ ) layer being trained using real data and real labels (the layers are shown in the brown block in Figure 4(c)). In our case, since the numbers of classes in real and GAN-generated features are the same but that of games are different, we initialize the rest of ( $f_c$ ) layers from scratch. Finally, the network is fine-tuned iteratively by feeding real and GAN-generated data.

## 4. Experiments

The main goal of our experiments is to quantitatively evaluate the proposed approach and analyze the different components. For evaluation, we use two aerial action datasets: UCF-ARG-Aerial Arjun Nagendran and Shah (publicly available) and YouTube-Aerial (will be publicly released). We perform experiments with and without games and GAN-generated data (Table 1 and Table 2), with and without disjoint multi-task learning (Fig 3). We also performed K-shot learning experiments (Fig 7), and analyzed robustness of proposed approach across three different visual features (Table 1 and Table 2). Finally, action-wise performance on two datasets are given in Table 4 and Table 5 and confusion matrices for the different components of our approach are shown in Fig 8.

### 4.1. Datasets

**UCF-ARG** Arjun Nagendran and Shah: UCF-ARG dataset contain 10 human actions. This dataset includes: boxing, carrying, clapping, digging, jogging, open-close trunk, running, throwing, walking, and waving. This is a multi-view dataset where videos are collected from an aerial camera mounted on a Helium balloon, ground camera, and rooftop camera. All videos are of high resolution  $1920 \times 1080$  and recorded at 60fps. The aerial videos contain severe camera shake and large



Figure 5. Examples of videos from the YouTube-Aerial dataset. In each video, different human action is being performed. We aim to automatically recognize human action in these videos when only a few training aerial examples are available.

camera motion. On average, each action contains 48 videos. The dataset partition includes 60% of videos of each action for training, 10% for validation, and 30% for testing. Figure 6 shows some of the videos from the UCF-ARG dataset. Note that the testing experiments are done on the aerial part of the UCF-ARG dataset (named as UCF-ARG-Aerial).

**YouTube-Aerial Dataset:** We collect this new dataset ourselves from the drones videos available on YouTube. This dataset contains actions corresponding to eight actions of UCF101 Soomro et al. (2013). The actions include band marching, biking, cliff-diving, golf-swing, horse-riding, kayaking, skateboarding, and surfing. The videos in this dataset contain large and fast camera motion and aerial videos are captured at variable heights. A few examples of videos in this dataset are shown in Figure 5. Each action contains 50 videos. Similar to the UCF-ARG dataset, the dataset partition includes 60%, 10%, and 30% of videos for training, validation, and testing respectively.

#### 4.2. Implementation details

We use five aerial videos of each action (named as a few available training examples in the above sections). For visual features computations, we use three recently proposed video features; namely 3D multi-fiber network (MFN-3D) Chen et al. (2018), 3D Inception network (I3D) Carreira and Zisserman (2017), and 3D residual network (Resnet-3D) Hara

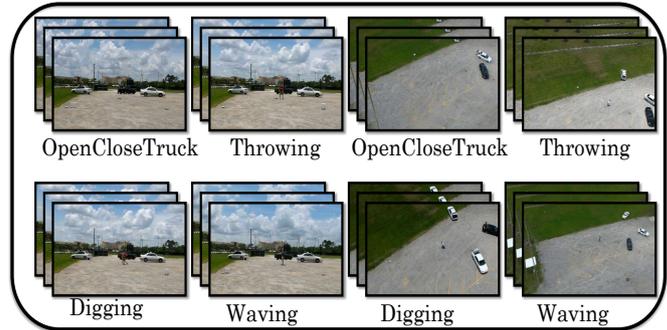


Figure 6. Examples of videos from the UCF-ARG dataset. The first two columns show the videos captured by the ground camera while the last two columns show the same actions captured by a UAV.

et al. (2018). Authors in Chen et al. (2018) showed that a multi-fiber network provides state-of-the-art results on several competitive datasets and is the order of magnitude faster than several other video features networks. It achieves high computational efficiency by dividing the complex neural network into small lightweight networks. We extract the features (768D) for all videos from the second last layer of the network. I3D features were proposed in Carreira and Zisserman (2017), where authors suggested a novel technique to inflate 2D ConvNets into 3D and bootstrap 3D filters from 2D filters. We extract the features (1024D) from the global average pooling layer using 128-frame snippets. Experiments are performed using RGB stream only. Similarly, we extract 512 dimension features from the last layer of 3D-Resnet34 Hara et al. (2018).

For disjoint multitask learning, we have two shared fully connected ( $f_c$ ) layers (512 and 256 units respectively). We have four task-specific layers: two  $f_c$  layers with the number of units equal to the number of actions in the real dataset and two  $f_c$  layers with the number of units equal to the number of actions in the game dataset. Similarly, for training without DML, we use only five aerial videos for each action. We employ three fully connected ( $f_c$ ) layers. The first two ( $f_c$ ) layers have 512 and 256 units respectively and the third one has the number of units equal to the number of actions in the dataset. The network is trained using the negative log-likelihood loss. We use the Adam optimizer with learning rate of 0.001, beta1=0.5, and beta2=0.999.

To generate GAN-generated examples, both our generator and discriminator contain four fully connected ( $f_c$ ) layers where the first three  $f_c$  layers have Leaky ReLU activation. In the case of the generator, the last  $f_c$  has ReLU activation. The noise vector  $z$  (312D) is drawn from unit Gaussian. For all networks, we use Adam optimizer. Since the network is already trained on game data, for joint learning from GAN-generated and game data (Figure 4(c)), we reduce the weight of loss for the branches being trained on the GAN-generated data. We ran the experiments several times with random initialization of the network ( $f_c$  layers) and report the average results.

### 4.3. Experiments Results

Table 1 and Table 2 demonstrates the experimental results of proposed approach on YouTube-Aerial and UCFARG-Aerial datasets using three visual features Chen et al. (2018); Carreira and Zisserman (2017); Hara et al. (2018). All the classification results are on real aerial videos. As compared to the YouTube-Aerial dataset, all visual features have lower classification accuracy on UCF-ARG-Aerial. This is mainly due to two reasons; firstly visual features networks (MFN-3D, I3D, Resnet-3D) were initially pre-trained on YouTube videos (Kinetics Carreira and Zisserman (2017), Sports-1M Karpathy et al. (2014) datasets) which are similar to YouTube-Aerial videos, secondly, UCF-ARG dataset is more challenging due to non-discriminative backgrounds and very small actors size.

**Component-wise accuracy:** In Table 1 and Table 2, trained using ‘Ground videos’ demonstrates classification results when training is done on ground camera videos only. Note that the UCF-ARG dataset contains ground cameras videos for the corresponding aerial action videos. For the YouTube-Aerial dataset, we use the videos of eight actions from UCF101 ground camera videos. We use ground videos from UCF101 (instead of collecting new ground videos ourselves) because our approach does not require pair-wise correspondence between aerial and ground videos. Note that in our approach ground videos features are only used to generate GAN-generated aerial videos features (see Figure 3). Training using ‘GAN-generated with DML’ demonstrates the experimental results when the network

is trained using disjoint multi-task learning employing GAN-generated visual features. Training using ‘Games with DML’ shows the results using game data and finally, training using ‘Games + GAN-generated with DML’ depicts classification results using both game and GAN-generated data. The experimental results in Table 1 and Table 2 show that the proposed approach results in improved aerial action classification. The results emphasize the strength of the proposed approach and suggest that given a few aerial videos (five in our case), games and GAN-generated aerial features can improve the classification accuracy when integrated properly using disjoint multitask learning.

**Table 1. Results of the proposed approach on YouTube-Aerial dataset**

Trained using	MFN-3D	I3D	Resnet-3D
Ground videos	49.7	50.7	53.5
GAN-generated with DML (Fig 4.b)	64.2	65.6	58.3
Games with DML (Fig 4.a)	62.9	64.8	56.7
Games + GAN-generated with DML (Fig 4.c)	<b>68.2</b>	<b>67.0</b>	<b>58.6</b>

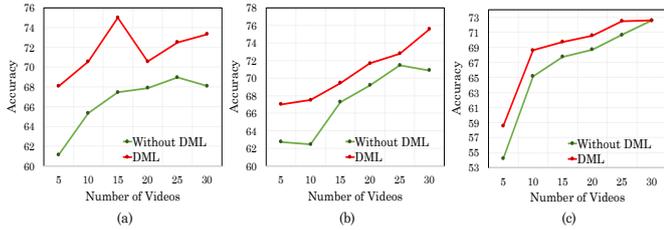
**Table 2. Results of the proposed approach on UCF-ARG-Aerial dataset for three visual features. The results in Table 1 and Table 2 demonstrate that each component of our approach is important. Both games and GAN-generated visual features are useful for improved aerial classification and combining them further improve the classification results.**

Trained using	MFN-3D	I3D	Resnet-3D
Ground videos	21.3	11.3	9.7
GAN-generated with DML (Fig 4.b)	32.1	15.6	12.4
Games with DML (Fig 4.a)	34.4	<b>16.8</b>	13.7
Games + GAN-generated with DML (Fig 4.c)	<b>35.9</b>	16.3	<b>15.1</b>

**Impact of disjoint multitask learning:** To verify the usefulness of disjoint multitask learning in our approach, in Table 3, we show the classification accuracy with and without training using disjoint multitask learning. We use the same experimental settings in both experiments and use the same number (five) of aerial videos. In experiments without DML, we use five aerial videos and in experiments with DML, we use games and GAN-generated data along with five aerial videos. The results demonstrate that integrating games and GAN-generated data through disjoint multitask learning significantly outperforms the training without disjoint multitask learning specifically when the

**Table 3. Classification results on YouTube-Aerial dataset when training is done without and with employing disjoint multitask learning. In training Without/with disjoint multitask learning, we have used the same five real aerial videos.**

Training	MFN-3D	I3D	Resnet-3D
Without Disjoint Multitask Learning	61.1	62.7	54.2
With Disjoint Multitask Learning	<b>68.2</b>	<b>67.0</b>	<b>58.6</b>

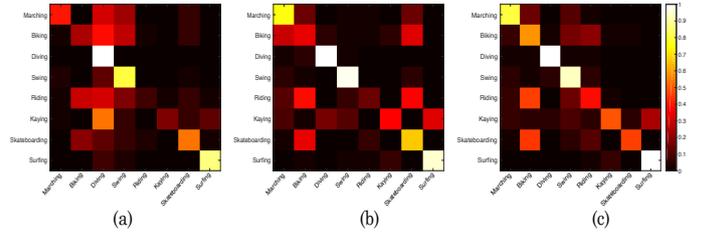


**Figure 7. Accuracy of the proposed approach for the different number of videos. (a), (b), and (c) show the results on YouTube-Aerial dataset with MFN-3D Chen et al. (2018), I3D Carreira and Zisserman (2017) and Resnet-3D Hara et al. (2018) features respectively. The red curves show the proposed approach while the green curves show the results when training is done without disjoint multitask learning.**

number of available training videos is small.

**K-shot learning:** In Figure 7, we demonstrate the accuracy of the proposed approach when training is done using different numbers of aerial videos on the YouTube-Aerial dataset. For the better analysis, we quantitatively compare the proposed approach against training without the disjoint multitask learning framework. The experiments are done for all three visual features. The classification results for the different number of training videos suggest that the proposed approach is not only useful when training data is less but is also beneficial with the increased training data.

**Class-wise accuracy:** Table 4 and Table 5 show the class-wise accuracy of proposed approach using MFN-3D features Chen et al. (2018) for UCF-ARG and YouTube-Aerial datasets. For YouTube-Aerial datasets, in five out of eight classes, the proposed approach significantly outperforms the classifier trained only on ground videos. A similar trend can be seen in six out of ten classes of the UCF-ARG dataset. Furthermore, in both datasets, for the majority of classes, combining GAN-



**Figure 8. This figure shows confusion matrixes averaged over all three visual features (MFN-3D, I3D, and Resnet-3D) on the YouTube-Aerial dataset. Confusion matrixes are (a) for the network trained using ground videos, (b) the network trained using aerial videos without DML, (c), and the network trained using aerial videos with DML employing both GAN-generated and game data. It can be seen that confusion between actions reduces through DML training.**

generated and game data either improve the accuracy or keep the best of both. The proposed approach works better for the actions which have discriminative motion patterns such as Biking, Swing, Kayaking, Carry, Clap, and Running, etc. However, our approach has limitations for the actions which have strong background scene biases water or mountains in Diving class) or contain less human body part motion (Skateboarding).

**Confusion matrices:** Figure 8 shows the confusion matrix averaged over all three visual features on the YouTube-Aerial dataset. The proposed disjoint multi-task learning framework significantly reduced the confusion between different actions as the classification accuracy increases from 51.3 (Figure8.a) to 59.4 (Figure8.b) to 64.5 (Figure8.c).

## 5. Conclusion

Recently, low cost and lightweight hardware make drones a good candidate for monitoring human actions. However, training the deep neural network for action recognition needs lots of training examples that are difficult to collect. In this paper, we explore two alternative data sources to obtain more accurate neural network classifiers. To tackle the different types of actions in the game and real action datasets, we propose to use disjoint multitask learning. Our experimental results and thorough analysis demonstrated that game action and GAN-generated examples, when integrated properly, can help to get improved aerial classification accuracy. The future works will aim at spatio-temporal localization of actors in drone videos,

**Table 4. Quantitative results on YouTube-Aerial dataset. The top row shows class-wise action recognition accuracy on aerial testing videos when trained on done on ground videos. The second, third and fourth rows demonstrate accuracy when training is done using DML employing GAN-generated features, game features and both respectively.**

Trained Using	March- ing	Biking	Diving	Swing	Riding	Kayak	Skate- board	Surfing	Avg
Ground Videos	26.7	0	100	80	26.7	40	53.3	73.3	49.7
GAN-generated with DML	53.3	46.7	73.3	100	40	86.7	20	86.7	64.3
Games with DML	60	73.3	86.7	100	26.7	86.7	13.3	53.3	62.9
Games + GAN-generated with DML	60	73.3	86.7	100	13.3	80	46.70	86.7	68.2

**Table 5. Quantitative results for UCF-ARG dataset. Similar to the Table 4, the top row shows class-wise action recognition accuracy on aerial testing videos when trained on done on ground videos. The second, third and fourth rows demonstrate accuracy when training is done using DML employing GAN-generated features, game features and both respectively.**

Trained Using	Boying	Carry	Clap	Dig	Jog	Trunk	Run	Throw	Walk	Wave	Avg
Ground Videos	0	0	0	0	0	0	40	33.3	80	60	21.33
GAN-generated with DML	60	40	20	13.3	13.3	58.3	26.7	6.70	13.3	66.7	31.83
Games with DML	53.3	60	40	00.0	6.70	83.3	40	6.70	13.3	40.0	34.33
Games + GAN-generated with DML	73.3	60.0	26.7	6.7	13.3	58.3	26.7	13.3	13.3	66.7	35.92

which will need attention based deep features. In our current work, we only use aerial game videos. However, it would be useful to use ground and aerial game videos jointly to learn the transformations between two ground and aerial views. Finally, one of the limitations of drones is their limited battery life. Future work could include designing the algorithms which work on low power devices.

## References

- Ali, S., Shah, M., 2010. Human action recognition in videos using kinematic features and multiple instance learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Arjovsky, M., Chintala, S., Bottou, L., 2017. Wasserstein generative adversarial networks, in: *Proceedings of the 34th International Conference on Machine Learning*.
- Arjun Nagendran, D.H., Shah, M., . Ucf-arg data set. URL: <https://www.crcv.ucf.edu/research/data-sets/ucf-arg/>.
- Barekatin, M., Martí, M., Shih, H., Murray, S., Nakayama, K., Matsuo, Y., Prendinger, H., 2017. Okutama-action: An aerial view video dataset for concurrent human action detection, in: *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2017, Honolulu, HI, USA, July 21-26, 2017*.
- Cao, L., Liu, Z., Huang, T.S., 2010. Cross-dataset action detection, in: *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*.
- Cao, Y., Liu, B., Long, M., Wang, J., 2018. Hashgan: Deep learning to hash with pair conditional wasserstein gan, in: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Carlucci, F.M., Russo, P., Caputo, B., 2016. A deep representation for depth images from synthetic data. *2017 IEEE International Conference on Robotics and Automation (ICRA)*, 1362–1369.
- Carreira, J., Zisserman, A., 2017. Quo vadis, action recognition? a new model and the kinetics dataset. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 4724–4733.
- Chen, Y., Kalantidis, Y., Li, J., Yan, S., Feng, J., 2018. Multi-fiber networks for video recognition, in: *European Conference on Computer Vision (ECCV)*.
- Damen, D., Doughty, H., Farinella, G.M., Fidler, S., Furnari, A., Kazakos, E., Moltisanti, D., Munro, J., Perrett, T., Price, W., Wray, M., 2018. Scaling egocentric vision: The epic-kitchens dataset, in: *European Conference on Computer Vision (ECCV)*.
- Dutta, S., Ekenna, C., 2019. Air-to-ground surveillance using predictive pursuit. *2019 International Conference on Robotics and Automation (ICRA)*, 8234–8240.
- Gkioxari, G., Hariharan, B., Girshick, R., Malik, J., 2014. R-cnns for pose estimation and action detection. *arXiv:1406.5212*.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y., 2014. Generative adversarial nets, in: *Advances in Neural Information Processing Systems 27*.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., Courville, A.C., 2017.

- Improved training of wasserstein gans, in: *Advances in Neural Information Processing Systems* 30.
- Hara, K., Kataoka, H., Satoh, Y., 2018. 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)*, pp. 6546–6555.
- Hariharan, B., Arbeláez, P., Girshick, R., Malik, J., 2014. Simultaneous detection and segmentation, in: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (Eds.), *Computer Vision – ECCV 2014*.
- Hong, W., Wang, Z., Yang, M., Yuan, J., 2018. Conditional generative adversarial network for structured domain adaptation, in: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Huang, C., Gao, F., Pan, J., Yang, Z., Qiu, W., Chen, P., Yang, X., Shen, S., Cheng, K.T., 2018. Act: An autonomous drone cinematography system for action scenes. 2018 IEEE International Conference on Robotics and Automation (ICRA) .
- Josifovski, J., Kerzel, M., Pregizer, C., Posniak, L., Wermter, S., 2018. Object detection and pose estimation based on convolutional neural networks trained with synthetic data. 2018 IEEE International Conference on Intelligent Robots and Systems (IROS) .
- Karpathy, A., Toderici, G., Shetty, S., Leung, T., Sukthankar, R., Fei-Fei, L., 2014. Large-scale video classification with convolutional neural networks, in: *CVPR*.
- Kataoka, H., Satoh, Y., 2019. Unsupervised out-of-context action understanding, in: 2019 International Conference on Robotics and Automation (ICRA).
- Kim, D.J., Choi, J., Oh, T.H., Yoon, Y., Kweon, I., 2018. Disjoint multi-task learning between heterogeneous human-centric tasks, in: *Winter Conference on Application of Computer Vision*.
- Lai, K.T., Lin, C.C., Kang, C.Y., Liao, M.E., Chen, M.S., 2018. Vivid: Virtual environment for visual deep learning.
- Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C., 2016. Ssd: Single shot multibox detector. URL: <http://arxiv.org/abs/1512.02325>.
- Mercier, J.P., Mitash, C., Giguère, P., Boularias, A., 2018. Learning object localization and 6d pose estimation from simulation and weakly labeled real images. 2019 International Conference on Robotics and Automation (ICRA) , 3500–3506.
- Mirza, M., Osindero, S., 2014. Conditional generative adversarial nets. CoRR URL: <http://arxiv.org/abs/1411.1784>.
- Misra, I., Shrivastava, A., Gupta, A., Hebert, M., 2016. Cross-stitch networks for multi-task learning, in: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- Mueller, M., Smith, N., Ghanem, B., 2016. A benchmark and simulator for uav tracking, in: *Proc. of the European Conference on Computer Vision (ECCV)*.
- Perera, A., Law, Y.W., Chahl, J., 2018. Uav-gesture: A dataset for uav control and gesture recognition, in: *UAVision workshop, ECCV*.
- Richter, S.R., Vineet, V., Roth, S., Koltun, V., 2016. Playing for data: Ground truth from computer games. ArXiv abs/1608.02192.
- Simonyan, K., Zisserman, A., 2014. Two-stream convolutional networks for action recognition in videos, in: *Advances in Neural Information Processing Systems* 27.
- Soomro, K., Zamir, R., Shah, M., 2013. Ucf101: A dataset of 101 human actions classes from videos in the wild, in: *ICCV*.
- Sultani, W., Chen, C., Shah, M., 2018. Real-world anomaly detection in surveillance videos, in: 2018 IEEE Conference on Computer Vision and Pattern Recognition, pp. 6479–6488.
- Sultani, W., Saleemi, I., 2014. Human action recognition across datasets by foreground-weighted histogram decomposition, in: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Tran, D., Bourdev, L., Fergus, R., Torresani, L., Paluri, M., 2015. Learning spatiotemporal features with 3d convolutional networks, in: *Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV)*.
- Wang, H., Schmid, C., 2013. Action recognition by dense trajectories, in: *ICCV*.
- Wu, S., Oreifej, O., Shah, M., 2011. Action recognition in videos acquired by a moving camera using motion decomposition of lagrangian particle trajectories, in: *Proceedings of the 2011 International Conference on Computer Vision*.
- Xian, Y., Lorenz, T., Schiele, B., Akata, Z., 2018. Feature generating networks for zero-shot learning, in: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zhou, X., Liu, S., Pavlakos, G., Kumar, V.S.A., Daniilidis, K., 2018. Human motion capture using a drone. 2018 IEEE International Conference on Robotics and Automation (ICRA) , 2027–2033.