

SSIVD-Net: A Novel Salient Super Image Classification & Detection Technique for Weaponized Violence

Toluwani Aremu, Li Zhiyuan, Reem Alameeri, Mustaqeem Khan, Abdulmotaleb El Saddik (Fellow, IEEE)
(firstname.lastname, a.elsaddik)@mbzuai.ac.ae
Mohamed Bin Zayed University of Artificial Intelligence (MBZUAI), UAE

Abstract—Detection of violence and weaponized violence in closed-circuit television (CCTV) footage requires a comprehensive approach. In this work, we introduce the *Smart-City CCTV Violence Detection (SCVD)* dataset, specifically designed to facilitate the learning of weapon distribution in surveillance videos. To tackle the complexities of analyzing 3D surveillance video for violence recognition tasks, we propose a novel technique called, *SSIVD-Net* (Salient-Super-Image for Violence Detection). Our method reduces 3D video data complexity, dimensionality, and information loss while improving inference, performance, and explainability through the use of Salient-Super-Image representations. Considering the scalability and sustainability requirements of futuristic smart cities, the authors introduce the *Salient-Classifier*, a novel architecture combining a kernelized approach with a residual learning strategy. We evaluate variations of SSIVD-Net and Salient Classifier on our SCVD dataset and benchmark against state-of-the-art (SOTA) models commonly employed in violence detection. Our approach exhibits significant improvements in detecting both weaponized and non-weaponized violence instances. By advancing the SOTA in violence detection, our work offers a practical and scalable solution suitable for real-world applications. The proposed methodology not only addresses the challenges of violence detection in CCTV footage but also contributes to the understanding of weapon distribution in smart surveillance. Ultimately, our research findings should enable smarter and more secure cities, as well as enhance public safety measures.

Index Terms—Violence Detection, Weaponized Violence Detection, Action Recognition, Signal Processing, Smart Surveillance

I. INTRODUCTION

Violence and gang-related activities can pose a serious threat to a city, particularly when authorities are unable to respond quickly enough to prevent further damage. In some cases, these incidents can result in loss of life and property, especially when weapons are involved. Regrettably, incidents of road rage, gang-related violence, and other spontaneous acts of violent crime frequently occur without prior warning or the ability for authorities to intervene proactively. These events pose a considerable challenge for law enforcement agencies and other relevant authorities. Unfortunately, the reporting of such incidents often occurs after the fact, leaving authorities with limited options for timely intervention and effective prevention.

Although surveillance systems have helped authorities identify instigators and culprits through recordings, it often takes too long to detect, search, and arrest someone after a crime is committed. To reduce turnaround time and increase efficiency, there is a growing need for automated detection and signaling systems. Since the breakthrough of deep learning [1] in the ImageNet 2012 competition, deep neural networks (DNNs) have become the go-to AI technology for automating such tasks. By leveraging DNNs and other AI techniques, smart cities around the world can better detect and respond to instances of violence, safeguarding lives and properties. The benefits of such technologies are clear, for instance, integrated surveillance systems equipped with advanced AI algorithms can analyze real-time video feeds from CCTV cameras to identify and alert authorities to potential violent incidents, enabling swift intervention. Furthermore, AI-powered predictive analytics can analyze various data sources, including social media feeds and sensor data, to identify patterns and trends associated with violence, enabling authorities to allocate resources strategically and prevent outbreaks of violence in specific areas. As these technologies continue to evolve, they will play an increasingly important role in maintaining safety and security in our cities.

Current violence detection methods predominantly rely on spatiotemporal models to identify instances of violent activities within video footage. However, it is essential to recognize that violence encompasses a wide range of behaviors, spanning from physical altercations to gunfights. Treating all violent events equally may not effectively prioritize the severity or potential harm involved. To address this challenge, it becomes crucial to develop methods that specifically focus on detecting weapons in surveillance footage.

Existing research has primarily concentrated on identifying weapons using object detection models, but these approaches often rely on datasets that are limited to specific weapon types, such as knives and guns. While this provides valuable insights into the detection of known weapons, it fails to account for the reality that virtually any object can be utilized as a weapon in an act of violence. Therefore, there

is a pressing need for further research and advancements in open-world weapons detection, which can efficiently identify and classify a broader range of objects that may be employed as weapons. Expanding weapons detection beyond predefined categories equips surveillance systems to recognize potential threats and intervene effectively, improving public safety. Exploring novel approaches and advancements in deep learning, computer vision, and object recognition enables comprehensive open-world weapons detection. These advancements aid in crime prevention and enhance overall community security and well-being.

Furthermore, 3D and spatio-temporal models such as C3D [2], I3D [3] and ConvLSTM [4], as well as object detection models like You-Only-Look-Once (YOLO) [5] and RCNN [6], require a high computational load to achieve state-of-the-art performance in violence and weapons detection. This leads to increased carbon footprints in smart cities. Our objective is to address these challenges by developing a more efficient image classifier capable of accurately detecting violence and weaponized violence while promoting sustainability in smart cities. To achieve this, we contributed the following:

- We address the challenges of detecting violence and weapons in CCTV footage by introducing the *Smart-City CCTV Violence Detection (SCVD)* dataset. Our dataset is designed to facilitate the learning of weapon distribution in surveillance videos, enabling DNNs to effectively detect both weaponized and non-weaponized violence.
- We propose *SSIVD-Net* (Salient-Super-Image for Violence Detection) as a data-centric approach to address the challenges associated with 3D surveillance video in violence recognition tasks. Our approach involves transforming the 3D video data into a Salient-Super-Image representation, resulting in reduced data complexity and dimensionality. This transformation enables faster inference, improved performance, and simplified explainability. In particular, our approach allows for seamless integration with 2D Vision-Classifiers, which are not commonly used in the field.
- The authors introduce a novel architecture called *Salient-Classifier*, that leverages a kernelized approach with residual networks [7]. We evaluate variations of SSIVD-Net and Salient-Classifier on our dataset and benchmark against SOTA models used in violence detection. Additionally, we perform comparative analyses to demonstrate the effectiveness of our model using other violence datasets.

Through our contributions, we aim to advance the SOTA in violence and weaponized violence detection while also providing a practical and scalable solution for real-world applications.

The rest of the paper is structured as follows: Section 2 discusses the relevant literature on violence detection, Sec-

tion 3 introduces the SCVD dataset, Section 4 describes the proposed methodology and its components, Section 5 presents the experimental results and comparative analysis, and Section 6 concludes the paper with potential future directions.

II. RELATED WORK

A. Weapon Detection

There are two main approaches to Object Detection: YOLO [5] and RCNN [6]. YOLO involve taking sliding windows of fixed sizes from the input image at every possible location and feeding them into an image classifier for inference, while RCNN proposes regions to feed into the classifier. Since their inception, research has gone into optimizing these methods to achieve better performance ([8], [9]), faster inference ([10], [11]), or both ([12], [13], [14]).

To detect weapons in surveillance footage, researchers have leveraged the Faster RCNN object detection model for its high accuracy in identifying objects of interest [15], [16]. However, transfer learning has also been applied to pre-trained Faster RCNN models for detecting handheld guns in clustered scenes [15]. In an attempt to classify mostly guns and knives, an ensemble method combining Faster RCNN and Single Shot Detector has been explored [16]. Another study used a pre-trained YOLO-V4 model on similar datasets [17]. However, these methods are limited in their ability to generalize to a wider range of objects that could be used as weapons, and they have not been trained on datasets that include CCTV footage, which limits their ability to detect weapons in various types of video.

B. Violence Detection

While there are multiple attempts on detecting weapons from videos only, there are works that aim to detect violence from videos and CCTV footage. For example, [24] used InceptionNet to detect violence in every frame from sports videos and movies. This results in slower inference and also poor generalization results as their method failed to learn temporal properties connecting frames within similar video embeddings. For their models to not lose temporal information, [20] and [25] used ConvLSTMs to detect violence in CCTVs. In ConvLSTMs, an image classifier is employed for spatial feature extraction while LSTMs learn the temporal information. The authors in [25] used a pre-trained ResNet50 model to extract spatial features from the video frames while [20] leveraged the VGG16 architecture. The extracted features are then concatenated with the latent features from the LSTM [26]. The above techniques and models require a lot of computing resources to get the results.

Our approach differs from previous techniques as we aim to create a novel data-centric approach that enables image classifiers to learn both spatial and temporal features for efficient weaponized violence detection in CCTV videos. We build upon the super image approach, first proposed in [27] and showcase Salient-Super-Image, to minimize the amount of

TABLE I
COMPARISONS BETWEEN THE SCVD AND THE PREVIOUS DATASETS

Dataset	Type	Size	Length/video (sec)	Annotation	Violence	Weapons	Characteristics	Scenario
NTU CCTV-Fights [18]	Video	1000 videos	5-720	Frame	✓	✗	CCTV + Mobile	Natural
Hockey Fight [19]	Video	1000 videos	1-2	Video	✓	✗	Aerial Camera	Hockey Games
RLVS [20]	Video	2000 videos	5-15	Video	✓	✗	CCTV + Mobile	Natural
RWF-2000 [21]	Video	2000 videos	5	Video	✓	✗	CCTV	Surveillance
Sohas [22]	Image	3255 images	N/A	Image	✗	✓	Captured Images	Demonstrations
WVD [23]	Video	168 videos	10-72	Video	✗	✓	Synthetic	Computer Games
SCVD - Ours	Video	500 videos	5-10	Video	✓	✓	CCTV	Surveillance

information lost as a result of rearranging and resizing video frames into 2D images. Our study is the first to investigate open-world weaponized violence detection in surveillance systems.

C. Super-Image

The field of action recognition for video classification has been gaining attention in recent years as it plays a significant role in video understanding. Most approaches in this domain use 3D convolutions to classify videos based on appearance, depth, or body skeletons. However, [27] proposed a different approach using a 2D image classifier (SIFAR). They argued that instead of using deep 3D networks for video action recognition tasks, a simple image classifier could work. To accomplish this, they introduced a technique that involves extracting frames from videos, resizing them, and combining them into a composite image. This super image effectively captures both local spatial information and global temporal dependencies. The experiments of SIFAR showed promising results, demonstrating that 2D image classifiers could achieve comparable results to their 3D counterparts.

In this work, we take this idea to the next level and unleash *Salient-Super-Image*, an innovative and highly effective variant of the Super Image. Our aim is to revolutionize the detection of weaponized violence in surveillance systems by preserving crucial information in a simpler yet more impactful manner. Since there are no datasets currently used for this task, we created **SCVD**, a novel dataset that contains distinctive videos of weaponized, non-weaponized, and normal violence in CCTVs.

III. PROPOSED SMART-CITY CCTV VIOLENCE DETECTION (SCVD) DATASET

The proposed Smart-City CCTV Violence Detection (SCVD) Dataset is a valuable addition to the existing benchmarks for violence detection, as it is specifically focused on weaponized violence scenarios captured strictly by surveillance cameras.

The motivation for creating this dataset is to enable AI models to learn and identify as well as prioritize the degree of danger in chaotic events, particularly those involving weapons, and alert the authorities accordingly. By identifying weapons in violent events, the authorities can respond more quickly and effectively to prevent further damage and harm. This dataset

can also help advance research in the field of AI for violence detection, particularly in the development of models that can effectively identify weapons and distinguish between violent and non-violent events.

Table I provides a summary of the comparison of the proposed Smart-City CCTV Violence Detection (SCVD) dataset with other datasets used to train deep neural networks (DNNs) for detecting violence and weapons in videos. The datasets considered in this comparison include those used for violence detection, such as Nanyang Technological University (NTU) CCTV-Fights [18], Hockey Fight [19], Real-Life Violent Situations (RLVS) [20], and RWF-2000 [21], as well as those for weapons detection, including Sohas [22] and Weapon Violence dataset (WVD) [23].

The datasets used for violence detection contain 1000 to 2000 videos, which are annotated either at the frame [18] or video level [19]–[21]. However, these datasets are not properly collected from surveillance systems. Only RWF-2000 [21] have real-time surveillance videos captured by CCTV surveillance systems, but the dataset only contains generalized violence without a distinctive reference to the type which contains weapons, making it also unsuitable for our tasks, which are addressed in this paper.

To detect weapons in videos, Sohas [22] and Weapon Violence Dataset (WVD) [23] were created. Sohas contains images with bounding box annotations, while WVD contains synthetic videos of length 10 to 72 seconds generated from the GTA-V computer game and annotated at the video level. Both datasets contain guns, knives, and other well-known weapons, but their distribution is not focused on surveillance. Since both datasets focus on only a limited number of weapons as compared to our proposed dataset in this work.

In contrast, our proposed SCVD dataset contains distinctive videos of weaponized and non-weaponized violence scenarios captured from CCTV surveillance systems. In order to enhance the efficiency of weapon recognition, we have opted not to assign specific labels, but instead enable DNN models to statistically comprehend the differentiation between scenes involving weapons and those without. Our objective is to assist AI models in gauging the level of threat in chaotic

events captured by surveillance systems and identifying potential weapons in violent incidents, thereby expediting the intervention of appropriate authorities to prevent further devastation. The SCVD dataset is therefore a suitable benchmark for developing and evaluating AI models for weaponized and non-weaponized violence detection in CCTV footage.

Our dataset was obtained by running an automated script that crawled YouTube and downloaded related videos using geographical prompts to ensure that the dataset is not limited to a single location but instead generalized to most parts of the world. The videos were recorded both indoors and outdoors and from various CCTV sensors. The downloaded videos were in 720p resolution with frame dimensions of 1280×720 . After cleaning the dataset by avoiding noisy and staged videos, the remaining videos' quantity is around 500. These videos were trimmed to a length of 5-10 seconds and annotated at the video level. The final dataset contains three classes: Violence (V), Normal (N), and Weaponized-Violence (WV), with the latter containing videos with an arbitrary number of possible weapons to improve DNNs' performance on weaponized violence detection.

IV. METHODOLOGY

Traditional Convolutional Neural Networks (CNNs) architectures are designed to process individual images, with the object of interest usually located at the center of the image [6], [27]. However, these approaches are not suitable for video input as each frame in a video is related to the preceding and succeeding frames. To account for the temporal information in videos, researchers have proposed different methods. One of the approaches is to process each frame individually, neglecting the interdependence between frames, which leads to temporal information loss. Another approach is to extract features from each frame and then use an LSTM network [4], [28] to learn the temporal dependencies between these features. However, this method can be computationally expensive and slow.

A more efficient and effective approach is to use 3D Convolutional Neural Networks (3D CNNs), which can directly process spatio-temporal information from videos [2], [3], [21]. In 3D CNNs, the filters used for convolution are applied both spatially and temporally, allowing the network to learn features from the entire video sequence. This approach has been successfully applied in various video understanding tasks such as action recognition and video captioning.

A. *Salient-Super-Image*

Frame Arrangements: With video understanding tasks in mind, [27] explored the possibility of merging multiple input frames from a video to a super image that contains both spatial and temporal information before processing by CNN. They found that a square formation, i.e. a 4×4 arrangement, had the best performance compared to linear or

rectangular formations, i.e. a 16×1 arrangement. To achieve this, the frames were resized to 224×224 , combined to form a super image, and then finally resized again to obtain a 224×224 image for input into their Swin-Transformer-based SIFAR classifier. However, the super image approach presents certain limitations, particularly in dealing with the original aspect ratio of surveillance videos. Since most videos have a wide/non-square aspect ratio, resizing each frame with dimensions such as 1280×720 to a 224×224 can introduce noise, result in loss of aspect information by converting the wide frame into a square shape, and potentially hinder performance.

To address this issue, we propose a novel approach called *Salient-Super-Image*. The main idea of our approach is to maintain as much aspect information as possible and reduce information loss in super images by selecting an optimal sample size, sampler, aspect ratio, and spatial arrangement (grid shape) that rearranges the frames while maintaining the aspect ratio of the original frames i.e., a 1280×720 , which have rectangular dimensions ($h \times w$, where $w > h$). Our goal is to achieve a final salient image with h as close as possible to w . In mathematical terms, we aim to minimize the difference between the width of n frames and the height of m frames:

$$\min(w_n - h_m) \quad (1)$$

so that:

$$w_n \approx h_m \quad (2)$$

By combining frames using this optimal grid selection, we reduced the information loss and obtain a 224×224 salient image for input into the classifier, while still preserving the temporal context of the video (in Figure 1, the six frames are uniformly sampled, and the $144p_B$ aspect ratio (192×144) was used to resize the images before arranging them in a 3×2 grid shape).

Optimal Frame Selection and Salient Arrangement: In order to process the videos from our SCVD dataset, which have frame dimensions of 1280×720 , we needed to find an optimal way to merge multiple frames into a single image for input into our *Salient-Classifer*. To achieve this, we considered several factors:

- 1) **Sample Size:** The sample size is an important factor that determines the final grid size for arranging the processed frames. It refers to the number of frames that would like to extract per second, dependent on the selected sampler, and denoted by k . It is worth noting that the value of k should be within the range of 3 and $len(fpgs)$, where $fpgs$ is the total number of frames that can be extracted from a given second(s).
- 2) **Sampler:** The type of sampler utilized is also a crucial factor in determining the performance of the classifier. To this end, we developed seven types of samplers

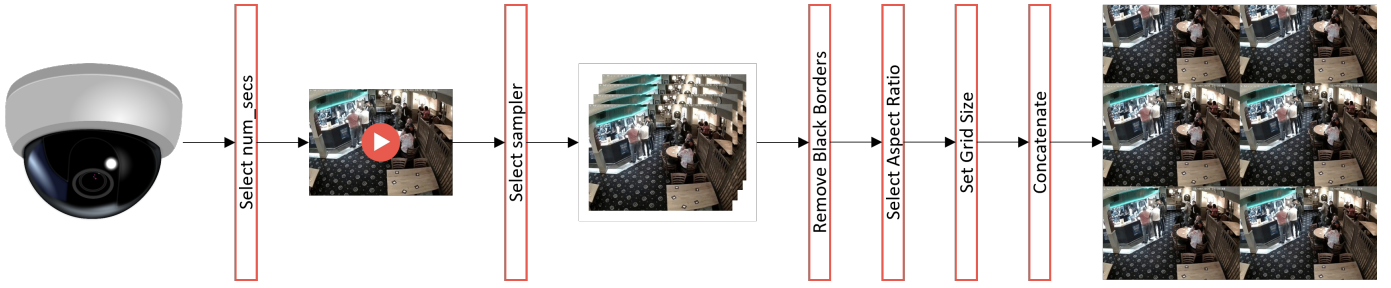


Fig. 1. Saliency-Super-Image: A sequence of video frames gotten from a CCTV surveillance system are rearranged into a saliency-super-image based on given factors such as sample size, sampler, aspect ratio, and spatial arrangement (grid shape).

inspired by existing literature, each designed to select a fixed number of k sample sizes from a given array of frames.

- **Uniform:** This sampler selects frames uniformly from the video, where the number of frames to select is specified by the parameter k . It first calculates the stride between the selected frames and then chooses k equally spaced indices according to (Algorithm 1).

Algorithm 1: Uniform Sampler

```

1  $stride \leftarrow num\_frames // k$ ;
2  $indices \leftarrow$  evenly spaced integers from
   $[0, num\_frames - 1, stride]$ ;
3  $sampled\_frames \leftarrow []$ ;
4 for  $index$  in  $indices$  do
5    $sampled\_frames.append(frames[index])$ ;
6 end

```

- **Random:** This sampler randomly selects k frames from the video without replacement. It uses the function named; `np.random.choice()` to generate a list of k with unique random indices, which correspond to the selected frames according to (Algorithm 2).

Algorithm 2: Random Sampler

```

1  $num\_frames \leftarrow len(frames)$ ;
2  $indices \leftarrow$  randomly select  $k$  indices from  $[0, num\_frames - 1]$ 
  without replacement;
3  $sampled\_frames \leftarrow [frames[i] \text{ for } i \text{ in } indices]$ ;

```

- **Continuous:** This sampler selects k frames that are evenly spaced across the entire video. Initially calculates the stride between the selected frames and then chooses k indices such that they are evenly spaced across the entire video according to (Algorithm 3).
- **Mean absolute difference:** This sampler selects k frames that have the smallest average absolute difference between adjacent frames. It calculates the absolute differences between adjacent frames and then selects the k frames with the smallest

Algorithm 3: Continuous Sampler

```

1  $stride \leftarrow \frac{num\_frames-1}{k-1}$ ;
2  $indices = []$ ;
3 for  $i \leftarrow 0$  to  $k-1$  do
4    $index_i \leftarrow [i \times stride]$ ;
5    $indices.append(index_i)$ ;
6 end
7  $sampled\_frames \leftarrow [frames[i] \text{ for } i \text{ in } indices]$ ;

```

average absolute difference and returns them in a list according to (Algorithm 4).

Algorithm 4: Mean Absolute Difference Sampler

```

1  $diffs \leftarrow abs(diff(frames))$ ;
2  $avg\_difs \leftarrow mean(diffs)$ ;
3  $indices \leftarrow argsort(avg\_difs)[:k]$ ;
4  $sampled\_frames \leftarrow [frames[i] \text{ for } i \text{ in } indices]$ ;

```

- **Lucas-Kanade:** This sampler uses the Lucas-Kanade algorithm to compute optical flow between adjacent frames. It then selects k frames such that they have the largest amount of motion and computes optical flow between adjacent frames and then selects the k frames with the largest amount of motion according to (Algorithm 5).

Algorithm 5: Lucas-Kanade Sampler

```

// Compute optical flow using Lucas-Kanade
algorithm
1  $gray\_frames \leftarrow [gray(frame) \text{ for } frame \text{ in } frames]$ ;
2  $sampled\_frames \leftarrow [frames[0]]$ ;
3  $prev\_frame \leftarrow gray\_frames[0]$ ;
4 for  $i \leftarrow 1$  to  $k-1$  do
5    $next\_frame \leftarrow$ 
      $gray\_frames[int((i/k) \times (num\_frames - 1))]$ ;
6    $flow \leftarrow$ 
      $opt\_flow\_farneback(prev\_frame, next\_frame)$ ;
7    $mag, ang \leftarrow cartToPolar(flow[...], 0], flow[...], 1]$ ;
8    $mag \leftarrow normalize(mag)$ ;
9    $sampled\_frames.append(frames[int((i/k) *$ 
      $(num\_frames - 1))])$ ;
10   $prev\_frame \leftarrow next\_frame$ 
11 end

```

- **Centered:** This sampler selects k frames that are centered around the middle of the video. It first

selects the middle frame of the video and then selects $k/2$ frames from the first half of the video and $k/2$ frames from the second half of the video. The selected frames are then returned in a list according to (Algorithm 6).

Algorithm 6: Centered Sampler

```

1 mid ← num_frames // 2
2 half_k ← k // 2
3 stride ← mid // half_k
4 indices ← [i * stride | i ∈ [0, half_k)]
5 sampled_frames += [frames[i] | i ∈ indices]
6 indices ← [i * stride + mid | i ∈ [0, half_k)]
7 sampled_frames += [frames[i] | i ∈ indices]

```

- **Consecutive:** The consecutive sampler selects a fixed number of consecutive frames from the video frames. For example, if we want to select k consecutive frames from a video with num_frames total frames, we can start at frame index i and select k frames from that index according to (Algorithm 7).

Algorithm 7: Consecutive Sampler

```

1 num_frames ← len(frames);
2 sampled_frames ← [];
3 start ← random.randint(0, num_frames - k);
4 for i ← start to start+k-1 do
5 | sampled_frames.append(frames[i]);
6 end

```

- 3) **Aspect Ratio:** The *aspect_ratio* parameter is a tuple of integers that specifies the desired size (width, height) of the cropped and resized frames. These aspect ratios correspond to commonly used video resolutions and dimensions and can be used to specify the desired output size for video frames in a standardized way. The available options for *aspect_ratio* include:

- 144p_A: 192×144 pixels
- 240p_A: 320×240 pixels
- 360p_A: 480×360 pixels
- 480p_A: 640×480 pixels
- 144p_B: 256×144 pixels
- 240p_B: 426×240 pixels
- 360p_B: 640×360 pixels
- 480p_B: 852×480 pixels
- square: 360×360 pixels
- vertical: 270×450 pixels

- 4) **Grid Shape:** Optimizing salient performance also involves integrating an essential factor, namely a tuple of integers whose multiplication should be equivalent to the sample size k . When selecting an optimal grid shape, it is crucial to determine suitable row and column coefficients (r, c) based on the width w and height h of the salient-super-image. Ideally, $w \times c$ should be approximately equal to $h \times r$.

In the case of our dataset comprising surveillance videos with a 1280×720 aspect ratio, the choice of r and c

depends on the frame’s height h and width w according to the following guideline:

$$\mathbb{O}_{mn} = \begin{cases} r > c, & \text{if } w > h \\ r < c, & \text{otherwise} \end{cases} \quad (3)$$

This ensures that the optimal selection of r and c aligns with the aspect ratio of the frames.

B. Salient-Classifer

We developed a novel deep convolutional neural network, called Salient-Classifer, to learn the embeddings from the proposed Salient-Super-Images. Salient-Classifer draws inspiration from the ResNet architecture [7], and Kervolution [29]. Specifically, the Salient-Classifer takes an input image and passes the signal tensors through a residual-based network architecture similar to ResNet, but with several key differences. The input layer, KConv2D, is based on Kervolution’s kernel-based convolutional layer. Additionally, we proposed SaliNet, which adds a minimal block and could serve as a sustainable and energy-saving alternative to ResNet’s basic and bottleneck blocks.

KConv2D: In Figure 2, the Salient-Super-Image output from SSIVD-Net is fed into a Kernelized Convolution layer (*KConv2D*). As demonstrated in [29], the *KConv2D* can effectively learn discriminative features from a larger latent space obtained by expanding the linear outputs $\mathbf{x}^T \mathbf{w}$ of the vanilla convolution using a given kernel type. However, the proposed SCVD dataset faces a challenge due to the high similarity between the distributions of the violence and weaponized violence classes. This similarity makes it difficult for any Deep Neural Network (DNN) to learn from the dataset. To address this issue, we take inspiration from Support Vector Machines (SVMs) which use kernels such as Sigmoid, RBF, Polynomial, etc., to find a hyperplane in a high-dimensional space that separates data points of different classes with a maximum margin. Thus, this study utilizes two kernels to replace the original convolution kernel in order to find the maximum distance between these two classes in our dataset. These kernels are:

- **Polynomial Kernel:** In Support Vector Machines (SVMs), the polynomial kernel is a widely-used kernel for non-linear classification. It maps data points to a higher-dimensional feature space using a polynomial degree that controls the complexity of the decision boundary. In this work, the authors employ a polynomial kernel in place of a linear kernel in convolution layers. Specifically, we compute the polynomial kernel as:

$$K(x, w) = (x^T w + c)^d \quad (4)$$

where c is a learnable balance factor, and d is the degree of the polynomial. Notably, we make c learnable instead of a constant in the original polynomial kernel to control the influence of individual data points and help with an expressive translation to improve and shift the decision

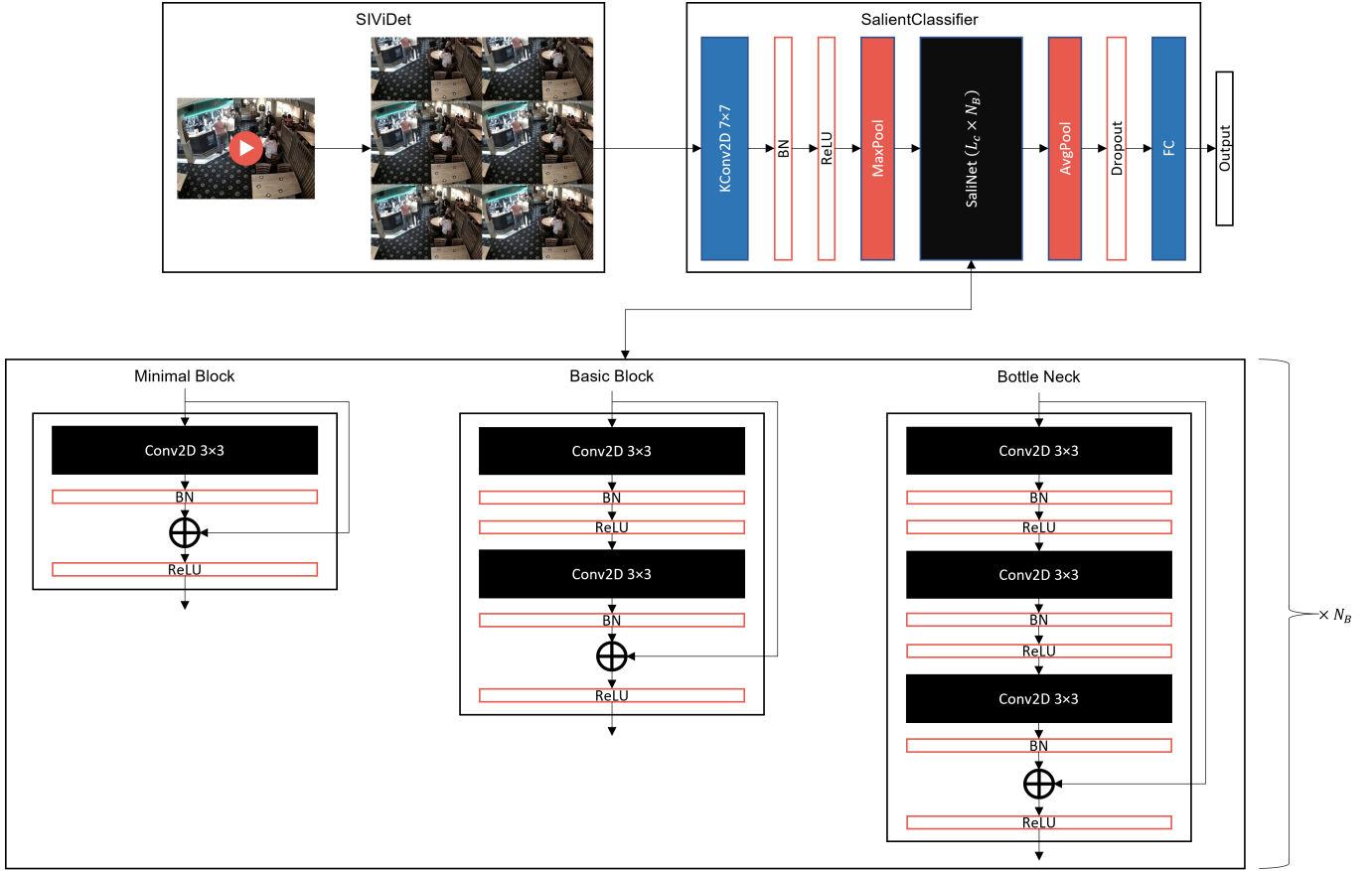


Fig. 2. Salient Classifier: The Salient-Super-Images produced in Figure 1 are fed into our sustainable Salient-Classifier which follows the residual learning strategy.

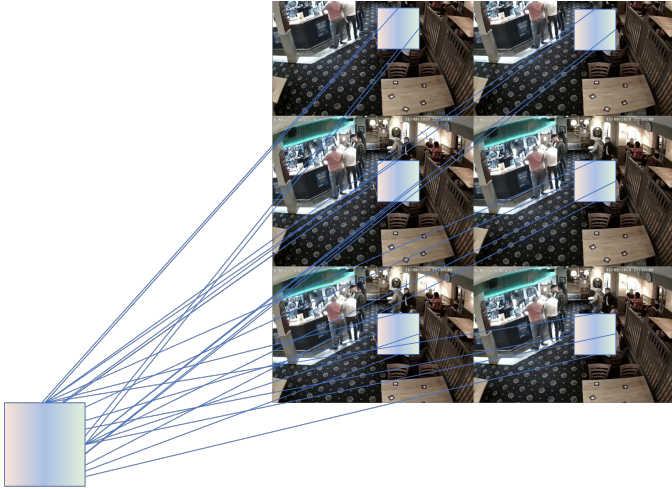


Fig. 3. Our model tries to learn similar spatial patterns in the concatenated salient-super-image to avoid the temporal information lost.

boundary between close classes. The kernel function takes two inputs x and w , representing data points in the input space, and maps them to a higher-dimensional feature space. In this feature space, the polynomial kernel

computes the dot product of the mapped features. The degree d determines the polynomial and controls the complexity of the decision boundary. Nevertheless, if the degree is too high, it can lead to overfitting.

- **Gaussian Kernel:** In SVMs, the RBF (Radial Basis Function) kernel maps the input space to a higher dimensional feature space, where it is more likely to be linearly separable. The feature space is defined by the distance between the input data points and a set of reference points, also known as support vectors. The RBF kernel function measures the similarity between the input data points and these support vectors based on their distance in this higher dimensional space. The RBF kernel function is defined as:

$$K(x, x') = \exp(-\gamma \|x - x'\|^2) \quad (5)$$

where x and x' are two input data points, γ is a parameter that controls the width of the RBF kernel, and $\|x - x'\|^2$ is the squared Euclidean distance between x and x' . The kernel function outputs a similarity score between the two input data points.

For the proposed network, the Gaussian kernel which is

based on the RBF definition is written as:

$$K(x, w) = \exp(-\gamma \|x - w\|^2) \quad (6)$$

Here, we replace the reference points with the weights.

TABLE II

THE SALIENT-CLASSIFIER’S SALINET HAS THREE BLOCK STYLES: THE BASIC BLOCK(B) AND BOTTLENECK(N) WHICH WERE ADAPTED FROM THE RESNET ARCHITECTURE, AS WELL AS OUR NOVEL MINIMAL BLOCK(M). THIS TABLE SHOWS THE LAYER ARRANGEMENT FOR EACH SALINET BLOCK STYLE, AND THE NUMBER OF PARAMETERS BASED ON THE CHOSEN ARRANGEMENT AND STYLE.

Classifier	Layer Arrangement	Number of Params (M)		
		Minimal Block(m)	Basic Block(b)	Bottle Neck(n)
SaliNet-2	1,1,0,0	1.8	4.9	8.0
SaliNet-4	1,1,1,1	1.8	4.9	8.0
SaliNet-8	2, 2, 2, 2	4.9	11.2	14.0
SaliNet-16	3, 4, 6, 3	10.0	21.3	23.5

SaliNet: We propose the *Minimal-Block* (m), while retaining the *Basic-Block* (b) and *Bottle-Neck* (n) of ResNets to make up our Salient Block (*SaliNet*). While the BottleNeck and Basic Block have three and two convolution layers respectively, the minimal only requires one which uses a 3×3 filter followed by batch normalization before the residual comes in. The output of all of these computations is finally activated using the ReLU function. Also, we proposed new layer arrangements for each adopted block type to satisfy our requirement of creating networks requiring less energy and computation. Table II shows the information about these blocks.

With our design approach, the Salient-Classifier harnesses the power of KConv2D to extract similar spatial patterns within each Salient-Super-Image (see Figure 3). These learned dynamics are subsequently fed into the SaliNet blocks, enabling the model to acquire more valuable features. This design enhances the accuracy of the classifier and facilitates efficient inference performance.

V. EXPERIMENTS AND RESULTS

A. Environment

Setup: The following details show the environment setup for the experiments conducted for this research:

- System used: NVIDIA Quadro RTX6000
- GPU Memory: 24GB GDDR6
- Operating System: Ubuntu 21.04
- Libraries: Tensorflow¹, Keras², Pytorch³

Training Details: To train all networks including baselines on the SCVD dataset, we used the Stochastic Gradient Descent (SGD optimizer) with a learning rate of $1e-3$ and a momentum of 0.9. The choice of the optimizer is based on the general knowledge that SGDM ensures full convergence, compared to Adam. The SGD optimizer for parameter update is:

$$w_{t+1} = w_t - \alpha \frac{\partial L}{\partial w_t} \quad (7)$$

where w_t denotes the weight, w_{t+1} denotes the updated parameter, α denotes the learning rate, and $\frac{\partial L}{\partial w_t}$ denotes the partial derivative of the gradient. To ensure fairness in comparison, all Salient-Classifier variants were initialized with the same weights and trained from scratch for 30 epochs. We recorded the average test accuracy and average precision for each network.

B. Results

Baseline Results: We conducted a comprehensive hyperparameters search for the Salient-Classifier module using the smallest and fastest block variant, SaliNet-2m. Our iterative approach involved exhaustively experimenting with all possible hyperparameters to identify the best-performing option. This allowed us to eliminate inferior options and ultimately obtain the optimal hyperparameters for subsequent experiments. Table III shows the results.

We began by utilizing various samplers to select 4 frames from the input videos using a square aspect ratio. Among these samplers, the uniform sampler (1) outperformed the others, achieving an accuracy of 78.4% and an average precision of 80.5%, and was thus used in all subsequent experiments. We then proceeded to use 6 frames with different aspect ratios, with 480p_A and 480p_B achieving the best results, achieving an accuracy of 83.0% and an average precision of 83.4%. Of these, 480p_A achieved a slightly faster average inference time per salient-super-image than 480p_B.

Finally, we used the uniform sampler and 480p_A aspect ratio to select 12 and 15 frames from the input videos, for which we achieved the best accuracy of 86.6% and average precision of 89.6% on 12 frames. Furthermore, our results demonstrate the superiority of our *salient-super-image* approach over the original *superimage* method [27]. The authors of the *superimage* approach used a consecutive frame sampler with a square aspect ratio. This supports our hypothesis that the use of *superimage* is limiting when analyzing wide-aspect videos extracted from surveillance systems.

Comparison with SOTA approaches on the SCVD dataset:

Table IV compares the performances of various models on the SCVD dataset. We can observe that the proposed *Salient-Classifier* variants achieve highly competitive results compared to existing methods.

Initially, we observe that the Flow-Gated-Network (FGN) model [21], which incorporates 3D convolutions, achieved a maximum accuracy of 74.4%. The LSTM-based networks, including ConvLSTM [4] and SepConvLSTM [28], attained an accuracy of 71.6% and 78.4%, respectively. These LSTM

¹<https://github.com/tensorflow/tensorflow>

²<https://github.com/keras-team/keras>

³<https://github.com/pytorch/pytorch>

TABLE III

TABLE SHOWING RESULTS OF DIFFERENT HYPERPARAMETERS FOR THE SALINET-2M. WE EMPLOYED AN ELIMINATION ROUTINE TO SELECT THE BEST HYPERPARAMETERS FOR THE PROPOSED SCVD DATASET.

k - grid_shape	Sampler	Aspect Ratio	Accuracy(%)	AP(%)	Inference time (s)
4 - 2x2	uniform	square	78.4	80.5	0.04
4 - 2x2	random	square	75.5	79.9	0.05
4 - 2x2	continuous	square	74.4	76.2	0.04
4 - 2x2	mean_abs	square	71.1	77.7	0.15
4 - 2x2	LK	square	69.6	78.2	0.21
4 - 2x2	centered	square	73.2	78.7	0.04
4 - 2x2	consecutive	square	70.4	79.4	0.04
6 - 3x2	uniform	144p_A	78.9	81.2	0.05
6 - 3x2	uniform	144p_B	79.7	81.9	0.05
6 - 3x2	uniform	240p_A	80.9	84.0	0.05
6 - 3x2	uniform	240p_B	81.3	84.2	0.05
6 - 3x2	uniform	360p_A	78.4	81.9	0.05
6 - 3x2	uniform	360p_B	82.4	83.8	0.05
6 - 3x2	uniform	480p_A	83.0	83.4	0.05
6 - 3x2	uniform	480p_B	83.0	83.4	0.06
9 - 3x3	uniform	square	84.7	85.0	0.06
12 - 4x3	uniform	480p_A	86.6	89.6	0.07
15 - 5x3	uniform	480p_A	84.3	86.8	0.08

TABLE IV

COMPARISONS BETWEEN CURRENT SOTA MODELS IN VIOLENCE DETECTION AND VARIANTS OF SALIENT-CLASSIFIERS ON THE SCVD DATASET.

Model	Num_Params (M)	Accuracy (%)
FGN [21]	0.3	74.4
Conv-LSTM [4]	47.4	71.6
Sep-Conv-LSTM [28]	0.4	78.4
SaliNet-2m	1.8	86.6
SaliNet-4m	1.8	83.1
SaliNet-8m	4.9	77.8
SaliNet-2b	4.9	75.9
SaliNet-2n	8.0	78.8

models utilize dynamic 2D convolution filters from pre-trained models to extract spatial features from each frame, which are then combined to obtain spatio-temporal features for inference. While these models exhibit a relatively lower parameter count, they encounter challenges in capturing the nuanced temporal dynamics that differentiate between weaponized and non-weaponized violence classes.

In contrast to the aforementioned models, our proposed Salient-Classifier variants, which employ the SaliNet module with different block styles, demonstrated superior performance. Notably, the minimal block variants exhibited a greater capacity to capture critical information, enabling effective differentiation between occluding classes. Among our Salient-Classifier models, those utilizing the SaliNet-2m and SaliNet-4m achieved remarkable accuracy of 86.6% and 83.1% respectively, while maintaining compact sizes of only 1.8 million parameters each. It is worth mentioning that the remaining variants displayed performances comparable to the current SOTA models i.e., Flow-Gated-Net [21] and Separable-ConvLSTMs [28].

Moreover, we note that increasing the model complexity through the utilization of larger block styles, such as SaliNet-8m, SaliNet-2b, and SaliNet-2n, does not necessarily improve the Salient-Classifier’s performance. This observation indicates that the optimal trade-off between model capacity and complexity is achieved with SaliNet-2m. However, further investigation is required in future work to understand the underlying reasons for this phenomenon. Nevertheless, these results highlight the effectiveness of our proposed network in capturing pertinent spatial features within intricate classes in our dataset. Additionally, the compact size of the model contributes to its computational efficiency, rendering it suitable for real-world applications.

Comparison with SOTA approaches on benchmark datasets: We conducted a comprehensive performance comparison among various variants of our classifiers, including 2m, 2b, and 2n, alongside well-established 3D CNN architectures such as C3D, I3D, and FGN. Additionally, we evaluated architectures based on CNN-LSTM hybrids, namely Conv-LSTM, Bi-Conv-LSTM, and Sep-Conv-LSTM. The evaluation encompassed three distinct datasets: MovieFight [19], HockeyFight [19], and RWF-2000 [21]. The detailed comparative analysis of these techniques with the proposed is mentioned in Table V

We trained our proposed classifiers on the MovieFight and HockeyFight datasets [19] using the 480p_A aspect ratio. Since these datasets do not have a consistent number of frames per second (fps) across all videos, we utilized the uniform sampler to extract 12 fps with a grid shape of (4, 3). On the other hand, for the SCVD dataset, all videos have a fixed frame rate of 30 fps. Therefore, we utilized all the frames to create a Salient-Super-Image with a 144p_A aspect

TABLE V
COMPARISON BETWEEN CURRENT SOTA MODELS IN VIOLENCE
DETECTION AND VARIANTS OF SALIENT-CLASSIFIERS ON BENCHMARK
DATASETS

Method	Model	MovieFight	HockeyFight	SCVD
3D-CNNs	C3D	100.0	96.5	82.8
	I3D	100.0	98.5	85.8
	FGN	100.0	98.0	87.3
Conv-LSTM	Conv-LSTM	100.0	97.1	77.0
	Bi-Conv-LSTM	100.0	98.1	-
	Sep-Conv-LSTM	100.0	99.5	89.3
Salient-Classifiers	SaliNet-2m	100.0	100.0	88.5
	SaliNet-2b	100.0	100.0	89.7
	SaliNet-2n	100.0	100.0	90.3

ratio and a grid shape of (6, 5). We conducted training for 30 epochs on the MovieFight and HockeyFight datasets and 50 epochs on the SCVD dataset to ensure sufficient model convergence.

Based on the results presented in Table V, our classifiers utilizing the Salient-Super-Image approach outperformed the other approaches across multiple datasets. Notably, our classifiers achieved a perfect accuracy of 100% on the MovieFight dataset, demonstrating their robustness in capturing the distinguishing characteristics of fight scenes in movies. Additionally, our classifiers achieved a perfect accuracy of 100% on the HockeyFight dataset, indicating their effectiveness in accurately detecting fights in hockey videos.

On the more challenging SCVD dataset, our best-performing model achieved an impressive accuracy of 90.3%. This showcases the superior performance of our approach in identifying violence in surveillance videos, where the presence of occlusions, varying lighting conditions, and other factors make the task more complex.

These results highlight the efficacy of our Salient-Super-Image approach in enhancing the performance of violence detection classifiers across diverse datasets. The high accuracies achieved by the Salient-Classifier’s variations demonstrate the potential of our approach for real-world applications in smart surveillance systems.

VI. CONCLUSION

The purpose of this study is to identify both weaponized and non-weaponized violence in surveillance systems. Recognizing the scarcity of video datasets specifically focusing on weaponized violence classes, we have introduced a novel dataset called the *Smart City Violence Detection (SCVD)* Dataset. To enhance the efficiency of violence detection, we have presented a technique called *Salient-Super-Image*, which transforms the 3D video understanding task into a 2D perspective, enabling the utilization of selected image classifiers.

Our experiments demonstrate that the Salient-Super-Image technique effectively mitigates information loss compared to the Super-Image approach [27]. Moreover, the proposed *Salient-Classifiers* exhibit superior performance on the SCVD dataset when compared to the current SOTA methods. Extensive evaluations on benchmark datasets further validate the effectiveness of our approaches in surpassing the current SOTA techniques in violence detection.

By creating the SCVD dataset and developing innovative techniques like Salient-Super-Image and Salient-Classifiers, we have made significant contributions to the field of violence detection in surveillance systems. Our findings not only highlight the advantages of our approaches but also emphasize the importance of addressing the specific challenges associated with weaponized violence detection. We believe that these advancements pave the way for more accurate and efficient violence detection systems, which have the potential to improve public safety and contribute to the development of smarter and safer cities.

REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” *Communications of the ACM*, June 2017.
- [2] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, “Learning spatiotemporal features with 3d convolutional networks,” *arXiv.org*, 2014. [Online]. Available: <https://arxiv.org/abs/1412.0767>
- [3] J. Carreira and A. Zisserman, “Quo vadis, action recognition? a new model and the kinetics dataset,” 2018.
- [4] S. Sudhakaran and O. Lanz, “Learning to detect violent videos using convolutional long short-term memory,” 2017.
- [5] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, “You only look once: Unified, real-time object detection,” *arXiv.org*, 2015, [Online]. Accessed on 17-Sep-2022. [Online]. Available: <https://arxiv.org/abs/1506.02640>
- [6] R. Girshick, J. Donahue, T. Darrell, and J. Malik, “Rich feature hierarchies for accurate object detection and semantic segmentation,” *arXiv.org*, 2013, [Online]. Accessed on 17-Sep-2022. [Online]. Available: <https://arxiv.org/abs/1311.2524>
- [7] 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, pp. 770–778.
- [8] K. He, G. Gkioxari, P. Dollár, and R. Girshick, “Mask R-CNN,” *arXiv preprint arXiv:1703.06870*, 2017. [Online]. Available: <https://arxiv.org/abs/1703.06870>
- [9] J. Redmon and A. Farhadi, “YOLOv3: An incremental improvement,” *arXiv preprint arXiv:1804.02767*, 2018. [Online]. Available: <https://arxiv.org/abs/1804.02767>
- [10] S. Ren, K. He, R. Girshick, and J. Sun, “Faster R-CNN: Towards real-time object detection with region proposal networks,” *arXiv.org*, 2015. [Online]. Available: <https://arxiv.org/abs/1506.01497>
- [11] J. Redmon and A. Farhadi, “YOLO9000: Better, faster, stronger,” *arXiv.org*, 2016. [Online]. Available: <https://arxiv.org/abs/1612.08242>
- [12] R. Girshick, “Fast R-CNN,” 2015. [Online]. Available: <https://arxiv.org/abs/1504.08083>
- [13] A. Bochkovskiy, C.-Y. Wang, and H.-Y. Liao, “YOLOv4: Optimal speed and accuracy of object detection,” 2020. [Online]. Available: <https://arxiv.org/abs/2004.10934>
- [14] C.-Y. Wang, A. Bochkovskiy, and H.-Y. Liao, “YOLOv7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors,” 2022. [Online]. Available: <https://arxiv.org/abs/2207.02696>
- [15] G. K. Verma and A. Dhillon, “A handheld gun detection using faster r-cnn deep learning,” in *Proceedings of the 7th International Conference on Computer and Communication Technology - ICCCT-2017*, 2017.

- [16] H. Jain, A. Vikram, Mohana, A. Kashyap, and A. Jain, "Weapon detection using artificial intelligence and deep learning for security applications," in *2020 International Conference on Electronics and Sustainable Communication Systems (ICESC)*, 2020.
- [17] M. T. Bhatti, M. G. Khan, M. Aslam, and M. J. Fiaz, "Weapon detection in real-time CCTV videos using deep learning," *IEEE Access*, vol. 9, pp. 34 366–34 382, 2021.
- [18] M. Perez, A. C. Kot, and A. Rocha, "Detection of real-world fights in surveillance videos," in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 2662–2666.
- [19] E. R. Nievas, D. S. Oscar, B. G. Gloria, and S. Rahul, "Hockey fight detection dataset," in *Computer Analysis of Images and Patterns*. Springer, 2011, pp. 332–339.
- [20] M. M. Soliman, M. H. Kamal, M. A. El-Massih Nashed, Y. M. Mostafa, B. S. Chawky, and D. Khattab, "Violence recognition from videos using deep learning techniques," in *2019 Ninth International Conference on Intelligent Computing and Information Systems (ICICIS)*, 2019, pp. 80–85.
- [21] M. Cheng, K. Cai, and M. Li, "Rwf-2000: An open large scale video database for violence detection," in *2020 25th International Conference on Pattern Recognition (ICPR)*, 2021, pp. 4183–4190.
- [22] F. Pérez-Hernández, S. Tabik, A. Lamas, R. Olmos, H. Fujita, and F. Herrera, "Object detection binary classifiers methodology based on deep learning to identify small objects handled similarly: Application in video surveillance," *Knowledge-Based Systems*, vol. 194, p. 105590, 2020.
- [23] M. S. Nadeem, V. N. Franqueira, F. Kurugollu, and X. Zhai, "Wvd: A new synthetic dataset for video-based violence detection," in *Lecture Notes in Computer Science*, 2019, pp. 158–164.
- [24] A. Mumtaz, A. B. Sargano, and Z. Habib, "Violence detection in surveillance videos with deep network using transfer learning," in *2018 2nd European Conference on Electrical Engineering and Computer Science (EECS)*, 2018, pp. 558–563.
- [25] M. Sharma and R. Baghel, "Video surveillance for violence detection using deep learning," *Advances in Data Science and Management*, pp. 411–420, 2020. [Online]. Available: https://link.springer.com/chapter/10.1007/978-981-15-0058-9_35
- [26] H. Sak, A. W. Senior, and F. Beaufays, "Long short-term memory based recurrent neural network architectures for large vocabulary speech recognition," *CoRR*, vol. abs/1402.1128, 2014. [Online]. Available: <http://arxiv.org/abs/1402.1128>
- [27] Q. Fan, C.-F. Chen, and R. Panda, "An image classifier can suffice for video understanding," *arXiv.org*, Jun 2021. [Online]. Available: <https://arxiv.org/abs/2106.14104>
- [28] Z. Islam, M. Rukonuzzaman, R. Ahmed, M. H. Kabir, and M. Farazi, "Efficient two-stream network for violence detection using separable convolutional LSTM," in *2021 International Joint Conference on Neural Networks (IJCNN)*. IEEE, jul 2021. [Online]. Available: <https://doi.org/10.1109%2Fijcnn52387.2021.9534280>
- [29] C. Wang, J. Yang, L. Xie, and J. Yuan, "Kervolutional neural networks," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, jun 2019. [Online]. Available: <https://doi.org/10.1109%2Fcvpr.2019.00012>