

Suspicious Behavior Detection on Shoplifting Cases for Crime Prevention by Using 3D Convolutional Neural Networks.

Martínez-Mascorro, Guillermo A.

a00824126@itesm.mx

Abreu-Pederzini, José R.

a00793921@itesm.mx

Ortiz-Bayliss, José C.

jcobayliss@tec.mx

Terashima-Marín, Hugo

terashima@tec.mx

Abstract

Crime generates significant losses, both human and economics. Every year, billion of dollars are lost due to attacks, crimes, and scams. Surveillance video camera networks are generating vast amounts of data, and the surveillance staff can not process all the information in real-time. The human sight has its limitations, where the visual focus is among the most critical ones when dealing with surveillance. A crime can occur in a different screen segment or on a distinct monitor, and the staff may not notice it. Our proposal focuses on shoplifting crimes by analyzing special situations that an average person will consider as typical conditions, but may lead to a crime. While other approaches identify the crime itself, we instead model suspicious behavior —the one that may occur before a person commits a crime— by detecting precise segments of a video with a high probability to contain a shoplifting crime. By doing so, we provide the staff with more opportunities to act and prevent crime. We implemented a 3DCNN model as a video feature extractor and tested its performance on a dataset composed of daily-action and shoplifting samples. The results are encouraging since it correctly identifies 75% of the cases where a crime is about to happen.

Keywords— 3D Convolutional Neural Networks, Crime Prevention, Pre-Crime Behavior Analysis, Shoplifting, Suspicious Behavior

1 Introduction

According to the 2018 National Retail Security Survey (NRSS) [1] inventory shrink, a loss of inventory related to theft, shoplifting, error or fraud, had an impact of \$46.8 billion in 2017 on U.S. retail economy. A high number of scams occur every day, from distractions and bar code-switching to booster bags and fake weight strategies, and there is no human power to watch every one of these cases.

The surveillance context is overwhelmed. Vigilance camera networks are generating vast amounts of video screens, and the surveillance staff cannot process all the available information. The more recording devices become available, the more complex the task of monitoring such devices becomes.

Real-time analysis of each camera has become an exhaustive task due to human limitations. The primary human limitation is the Visual Focus of Attention (VFOA) [2]. Human gaze can only concentrate on one specific point at once. Although there are large screens and high-resolution cameras, a person can only regard a small segment of the image at a time. Optical focus is a significant human-related disadvantage in the surveillance context. A crime can occur in a different screen segment or on a different monitor, and the staff may not notice it. Other significant difficulties can be the attention paid, boredom, distractions, lack of experience, among others [3, 4].

Defining what can be considered suspicious behavior is usually tricky, even for psychologists. In this work, the mentioned behavior is related to the commission of a crime, but it does not imply its realization (Figure 1). We define suspicious behavior as a series of actions that happen before a crime occurs. In this context, our proposal focuses on shoplifting crime scenarios, particularly on before-offense situations, that an average person may consider as typical conditions. While existing models identify the crime itself, we model suspicious behavior as a way to anticipate the crime. In other words, we identify behaviors that usually take place before a shoplifting crime. This kind of crime usually occurs in supermarkets, malls, retail stores, and other similar businesses. Many of the models for addressing this problem need the suspect to commit a crime to detect it. Examples of such models include face detection of previous offenders [5, 6] and object analysis in fitting rooms [7]. In this work, we propose an approach that aims at

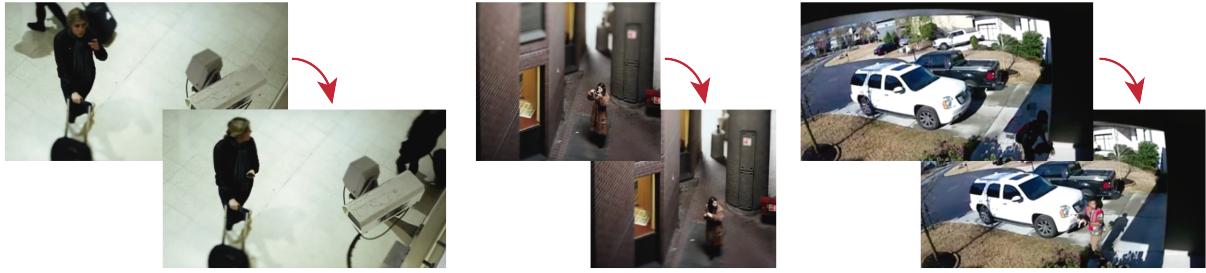


Figure 1: Suspicious behavior is not the crime itself, particular situations will make us distrust of a person.

supporting the monitoring staff to focus their attention on specific screens where crime is more likely to happen. By detecting situations in a video that may indicate that a crime is about to occur, we give the surveillance staff more opportunities to act, prevent, or even respond to such a crime. In the end, it is the security personnel who will decide how to proceed for each situation.

We implement a 3D Convolutional Neural Network (3DCNN) to process criminal videos and extract behavioral features to detect suspicious behavior. We perform the model training by selecting specific videos from the UCF-Crimes dataset [8]. Among the main contributions of this work, we propose a method to extract segments from videos that feed a model based on a 3DCNN and learns to classify suspicious behavior. The model achieves an accuracy of 75% on suspicious behavior detection before committing a crime on a dataset composed of daily-action samples and shoplifting samples. These results suggest that our approach is useful for crime prevention in shoplifting cases.

The remainder of this document is organized as follows. In Section 2, we review various approaches that range from psychology to deep learning, to tackle behavior detection. Section 3 presents the methodology followed throughout the experimental process. The results and discussions of the tests are presented in Section 4. Finally, Section 5 presents the conclusions and future works derived from this investigation.

2 Background and Related work

Every surveillance environment must satisfy with a particular set of requirements. Those requirements have promoted the creation of specialized tools, both on equipment and on software, to support the surveillance task. The most common approaches include motion detection [9, 10], face recognition [11, 12, 6, 5], tracking [13, 14, 15], loitering detection [16], abandoned luggage detection [17], crowd behavior [18, 19, 20], and abnormal behavior [21, 22].

Prevention and reaction are two primary aims in the surveillance context. Prevention requires to forestall and deter a crime execution. The monitoring staff must remain alert, watching as much as they can, and alerting the ground personnel. Reaction, on the other hand, involves protocols and measures to respond to a specific event. The security teams take action only after the crime or event has taken place.

Most security-support approaches focus on crime occurrence. [23] present a snatching-detection algorithm, which performs background subtraction and pedestrian tracking, in order to make a decision. That approach divides the frame into eight areas and searches for a speed-shift in one of the tracked persons. The algorithm proposed by [23] can only alert when a person already loses its belongings. [8] present a real-world anomaly detection approach, training thirteen anomalies, such as burglary, fighting, shooting, and vandalism. They label the samples into two categories: normal and anomalous, and use a 3DCNN for feature extraction. Their model includes a ranking loss function and trains a fully-connected neural network for decision making. [24] propose a system to detect loitering people. The system combines several analyzes for decision fusion and final detection, including distance, acceleration, direction-based, and grid-based analysis.

Convolutional Neural Networks (CNN) have shown a remarkable performance in computer vision and different areas in the last recent years. Particularly, 3DCNN —an extension of CNN—, focus on extracting spatial and temporal features from videos. Traditional applications that have been implemented using 3DCNN include object recognition [25], human action recognition [26], gesture recognition [27], and as a specific implementation, Cai et al. [28] used 3DCNN for abnormal behavior detection in examination surveillance within the classroom.

Although all the works mentioned before are based on 3DCNN, each one has a particular architecture, and many parameters —such as depth, number of layers, number of filters on each layer, kernel size, padding, stride— must be adjusted. For example, concerning the number of layers, many approaches rely on simple structures that consist of two or three layers [26, 28, 25], while others require several layers for exhaustive learning [27, 29, 30, 31, 32, 33].

Concerning shoplifting, the current literature is somewhat limited. Surveillance material is, in most cases, a

company’s private property. The latter restricts the amount of data we can get to train and test new surveillance models. For this reason, several approaches focus on training to detect normal behavior. Anything that lies outside the cluster is considered abnormal. In general, surveillance videos contain only a small fraction of crime occurrences. Then, most of the videos in the data are likely to contain normal behavior.

Many approaches have experienced problems regarding the limited availability of samples and their unbalanced category distribution. For this reason, some works have focused on developing models that learn with a minimal amount of data. For example, [26] create a school dataset and test with eight to ten videos, [23] rely on a dataset of nineteen videos (four used for training and fifteen for testing), and [24] work with six videos (one for training and five for testing).

This work aims at developing a support approach for shoplifting crime prevention. Our model detects a person that, according to its behavior, is likely to commit a shoplifting crime. We achieve the latter by analyzing the comportment of the people that appear in the videos before the crime occurs. To the best of our knowledge, this is the first work that analyzes behavior as a means to anticipate a potential shoplifting crime.

3 Methodology

As part of this work, we propose a new methodology to extract segments from videos where people exhibit behaviors that are relevant to the task of preventing shoplifting crime. These behaviors include both normal and suspicious, being the task of the network to classify them. In this section, we will describe the dataset used for experiments and the 3DCNN architecture used for feature extraction and classification.

3.1 Description of the Dataset

We use the UCF-Crime dataset, proposed by [8], to analyze suspicious behavior before a shoplifting crime. The dataset consists of 1900 real-world surveillance videos and provides around 129 hours of videos (with a resolution of 320x240 pixels and not normalized in length). The dataset includes scenarios from several locations and persons that are grouped into thirteen classes such as ‘abuse’, ‘burglary’, and ‘explosion’, among others. From those classes, we extracted the samples used in this work from the ‘shoplifting’ and ‘normal’ classes.

To feed our model, we require videos that show one or more persons and that their activities are visible before the crime is committed. Because of these restrictions, not all the videos in the dataset are useful. Suspicious behavior samples were extracted only from videos that exhibit a shoplifting crime —and these samples omit the crime itself. Normal behavior samples were taken from the ‘normal’ class. Thus, it is important to stress that the model we propose is a suspicious behavior classifier and not a crime classifier.

For processing the videos and extracting the suspicious samples (video segments that exhibit a suspicious behaviour), we propose a new method, the Pre-Crime Behavior (PCB) analysis, which we explain in the next section. Once we obtain the suspicious samples, we applied some transformations to produce several smaller datasets. First, to reduce the computational resources required for training, all the frames in the videos were transformed into grayscale and resized to four testing resolutions: 160×120, 80×60, 40×30, and 32×24 pixels. For organization purposes, all the samples extracted from the videos are indexed. The suspicious samples are indexed as SB_i (where i ranges from 1 to 60) while the samples of normal behavior are indexed as NB_i (where i ranges from 1 to 60). We divided the original sample size by 2, 4, 8, and 10 to explore the performance of each configuration. Table 1 describes how these datasets are conformed. For example, $SBT_balanced_120$ is a set that contains 120 samples with the same number of suspicious and normal samples, 60 of each class (samples SB_1 to SB_{60} and NB_1 to NB_{60}), while $SBT_unbalanced_30s60n$ is a dataset that contains fewer suspicious samples than normal ones (samples SB_1 to SB_{30} and NB_1 to NB_{60}). To increase the number of samples, we applied a flipping procedure that consists of turn over horizontally each frame of the video sample, resulting in a clip where the actions happen in the opposite direction. For example, $SBT_balanced_240$ contains 240 samples (samples SB_1 to SB_{60} , NB_1 to NB_{60} , as well as the flipped versions of SB_1 to SB_{60} and NB_1 to NB_{60}).

Table 1: Datasets description.

Dataset	Suspicious samples	Normal samples
$SBT_balanced_60$	SB_1 to SB_{30}	NB_1 to NB_{30}
$SBT_unbalanced_30s60n$	SB_1 to SB_{30}	NB_1 to NB_{60}
$SBT_balanced_120$	SB_1 to SB_{60}	NB_1 to NB_{60}
$SBT_balanced_120_flip$	Flipped versions of SB_1 to SB_{60}	flipped versions of NB_1 to NB_{60}
$SBT_unbalanced_60s120n$	SB_1 to SB_{60}	NB_1 to NB_{60} + flipped versions of NB_1 to NB_{60}
$SBT_balanced_240$	SB_1 to SB_{60} + flipped versions of SB_1 to SB_{60}	NB_1 to NB_{60} + flipped versions of NB_1 to NB_{60}

3.2 Pre-Crime Behavior

To detect suspicious behavior, the proposed model must focus on what happens before a shoplifting crime is committed. For this purpose, we propose a new method to process surveillance videos. Before we explain our proposal, we introduce some concepts, which are listed below.

- **Strict Crime Moment (SCM).** In a surveillance video, and after being reviewed by a human, the SCM is the segment of video where a person commits shoplifting crime. This moment is the primary evidence to accuse a person of committing a crime.
- **Comprehensive Crime Moment (CCM).** it is the precise moment when an ordinary person can detect the suspect's intentions. He/she started to watch out to go unnoticed and looks for the best moment to commit the crime. Other CCM examples are unsuccessful attempts or reorder things to distract attention. If we isolate this moment, we can doubt the suspect in the video, but there is no clear evidence to know if the suspect steals something.
- **Crime Lapse (CL).** In a video, the CL is the entire segment where a crime takes place. If we remove the CL from the video, it will be impossible to determine that there is a criminal act in the video. The CCM supports the beginning of the CL. It is essential not to leave any trace of the crime to avoid biasing the training.
- **Pre-crime Behavior (PCB).** The PCB contains what happens from the first appearance of the suspect until the CCM begins. These samples have different sizes since each video shows many behaviors. We can find more than one CL per video. The next PCB will start where the previous CL ends and until the next CCM. The result is a video segment in which an ordinary person may not detect that a crime will occur, but we are sure that the sample comes from a video where criminal activity was present.

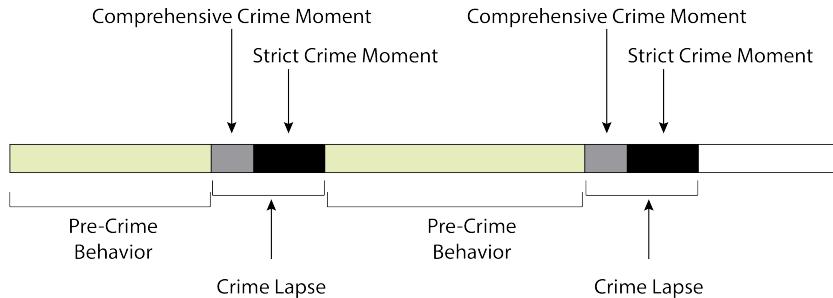


Figure 2: Graphical description of the concepts related to the proposed methodology for suspicious sample extraction.

Figure 2 graphically presents how these concepts interact in one video sample. The sample has two CL, and each one contains its corresponding SCM and CCM. From this video sample, we can extract two PCB training samples: from the beginning of the video to the first CCM, and from the end of the first SCM to the second CCM.

To extract the samples from the videos, we follow the process depicted in Fig. 3. Given a video that contains one or more shoplifting crimes, we identify the precise moment when the offense is committed. After that, we mark the different suspicious moments—moments where a human observer doubts about what the person in the video is doing. Finally, we select the segment where the suspect people are preparing to commit the crime. These segments become the training samples for the Deep Learning (DL) model.

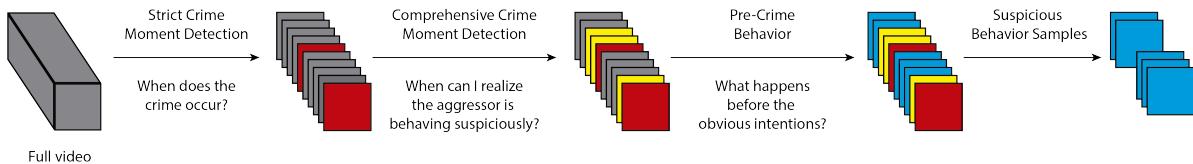


Figure 3: Graphical representation of the process for suspicious sample extraction.

In a video sample, each moment has its information level importance (see Fig. 4). PCB has less information about the crime itself, but it allows us to analyze the suspect's normal-acting behavior in its first stage, even far from the target. CCM allows us to have a more precise idea about who may commit the crime, but it is not conclusive. Finally, SCM is the doubtless evidence about a person committing a shoplifting crime. If we remove SCM and CCM from the video, the result will be a video containing only people shopping, and there will be no suspicion or evidence if someone commits a crime. That is the importance of the accurate segmentation of the video. From where a Crime Lapse ends until the next SCM, there is new evidence about how a person behaves before committing a new shoplifting crime attempt.

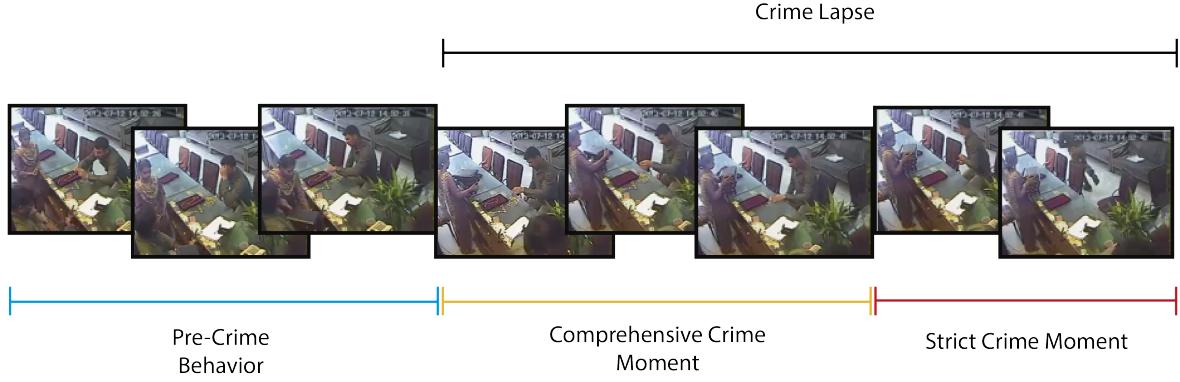


Figure 4: Video segmentation by critical moments.

For experimentation purposes, we only use PCB segments. These segments lack specific criminal behavior and have no information about a transgression. We look to pattern an aggressor's behavior before trying to steal from a store.

3.3 3D Convolutional Neural Networks

We use a 3DCNN for feature extraction and classification. We choose a basic structure to explore the performance of the 3DCNN for suspicious behavior detection task. The architecture of the model consists of four Conv3D layers, two max-pooling layers, and two fully connected layers. As a default configuration, in the first pair of Conv3D layers, we apply 32 filters, and for the second pair, 64 filters. All kernels have a size of $3 \times 3 \times 3$, and the model uses an Adam optimizer and cross-entropy for loss calculation. At the end of the model, it has two dense layers with 512 and 2 neurons, respectively. The output is binary, 1 for Suspicious Behavior and 0 for Normal Behavior. This architecture is selected because it has been used for similar applications [34], and seems suitable as a first approach for behavior detection in surveillance videos.

For handling the model training, we use Google Colaboratory [35]. This free cloud tool allows to write and execute code in cells, runs from a browser, and uses a GPU to train deep learning models. We can upload the datasets to a storage service, link the files, prepare the training environment, and save considerable time during the model training, using a virtual GPU.

4 Experiments and Results

3DCNN is a recent approach for Spatio-temporal analysis, showing a remarkable performance by processing videos in different areas, such as moving objects action recognition [25], gesture recognition [27] and action recognition [26]. We decided to implement 3DCNN in a more challenging context, such as the search for patterns in criminal samples, which lack suspicious and illegal visual behavior. In this section, we present the proposed experiments and their results.

The initial experiment aims at exploring the impact of different values for the parameters of the system. The second experiment focuses on obtaining statistical support that the best configurations obtained from the first experiment are useful for further testing in different situations.

4.1 Exploration of configurations

In this experiment, we explore different values for the parameters of the system. Given different values for the parameters, we estimate the changes in the response due to such configurations. The baseline training uses the most common values for this architecture, such as filters, kernel size, depth, and batch. The rationale behind this first experiment is that by producing small variations on the input parameters, we expect to improve the model performance.

We consider an extensive set of parameters to generate the testing configurations and obtain a total of 22 configurations. For example, we consider a balanced dataset or flipped images. We use unbalanced datasets to simulate real environments where normal behavior is more likely to be present than suspicious ones. For these datasets, we use a sample ratio of 1:2; for each suspicious behavior sample, there are two normal behavior samples. The following is a short description of the nomenclature used to name the datasets so that the reader can understand the differences between each one of the datasets.

- **Balance.** The dataset has the same number of samples of each class. The values this parameter can take are balanced and unbalanced.
- **Ratio.** The proportion of samples of each class in the dataset. The different configurations for balanced sets include 60, 120 and 240 samples. For unbalanced datasets the ratio is 1:2 for suspicious (s) and normal (n) class, respectively, with a total of 90 and 180 samples.
- **Test size.** The percentage of the dataset intended to the testset. The possible percentage are 20, 30 or 40 percent.
- **Depth.** The number of consecutive frames used for a 3D convolution. The values this parameter can take are 10, 30 and 90.
- **Resolution.** The size of the input images. 160×120 , 80×60 , 40×30 or 32×24 .
- **Flip.** If the word ‘flip’ appear in the dataset name, the frames in the videos have turned over horizontally.

By following the previous description, a dataset named *SBT_unbalanced_60s120n_30t_30f_40x30_flip* refers to a dataset which has fewer samples of suspicious behavior than normal behavior, sixty and one hundred and twenty respectively. It destines thirty percent of the dataset for the test, and it uses thirty frames to perform a 3D convolution. Finally, the input images have a resolution of 40x30 pixels, and they were flipped horizontally. The tests focus on comparing different depths, test set sizes, the balance in the number of samples, which image resolution is optimal between time processing and image detail, and the data-augmentation technique of flip the images. The objective of the exploratory experiment is to find a suitable configuration to model suspicious behavior. As previously mentioned, we analyzed 22 configurations, which are tested in four different resolutions that run three times each.

The depth sizes (number of consecutive frames) considered for the test are 10, 30 and 90. Table 2 shows the results of these runs with the four resolutions. Based on the results, using 10 and 30 frames achieves the best classification results, 69.4% to 83.3% and 69.4% to 75%, respectively. Table 3 presents the results of the test set size comparison. We select values of 20%, 30%, and 40% of the complete dataset for testing purposes. Although the first case (20% test set size) uses more information to train, this proportion did not get the best results. It produced outcomes between 47.2% and 72.2%. The second case obtained results between 68.5% and 75.9%, and the third one between 61.1% and 70.3%. 30% of the total dataset has the best results to define the test set size.

Table 2: Results of depth comparison.

Dataset	Resolution			
	32x24	40x30	80x60	160x120
SBT_balanced_60_20t_10f	83.3%	72.2%	77.7%	69.4%
SBT_balanced_60_20t_30f	69.4%	75.0%	69.4%	69.4%
SBT_balanced_60_20t_90f	69.4%	63.9%	61.1%	58.3%

To deal with unbalanced training, we create three datasets with sixty normal samples, thirty suspicious ones and three different depths (Table 4). We are aware that our model requires more samples to provide a better performance. However, in this test, the results reveal a similar performance, around 80%, between 30 frames and 90 frames depth. 3DCNN can handle unbalanced datasets. The difference may relay in the training time of each depth.

Data augmentation techniques are an option to take advantage of small datasets. For this reason, we test the model performance using original and flipped images in different runs. The used test set has a size of 30% and 40%. The tests throw accuracy results between 70% and 80% (Table 5). Therefore, we consider that both orientations can effectively be used as samples to train the model.

Table 3: Results of test set size comparison

Dataset	Resolution			
	32x24	40x30	80x60	160x120
SBT_balanced_60_20t_10f	72.2%	72.2%	66.6%	47.2%
SBT_balanced_60_30t_10f	68.5%	74.0%	68.5%	75.9%
SBT_balanced_60_40t_10f	63.9%	68.0%	61.1%	70.3%

Table 4: Results of unbalanced dataset test

Dataset	Resolution			
	32x24	40x30	80x60	160x120
SBT_unbalanced_30s60n_30t_10f	66.6%	67.8%	68.8%	79.0%
SBT_unbalanced_30s60n_30t_30f	69.1%	65.4%	80.2%	80.2%
SBT_unbalanced_30s60n_30t_90f	69.1%	65.4%	81.4%	66.6%

Table 5: Results of flipped images comparison.

Dataset	Resolution			
	32x24	40x30	80x60	160x120
SBT_balanced_120_40t_10f	71.5%	71.5%	77.0%	77.0%
SBT_balanced_120_40t_10f_flip	73.6%	77.0%	83.3%	70.8%
SBT_balanced_120_30t_10f	75.9%	71.3%	71.3%	79.6%
SBT_balanced_120_30t_10f_flip	76.8%	72.2%	81.5%	78.6%

Finally, we create datasets with balanced 240 samples and unbalanced 180 samples. Each type, balanced and unbalanced, combines three different depths and four resolutions, for a total of 24 datasets. Table 6 shows both the best result and the average result from three runs, for each resolution. It is essential to clarify better results by resolution may be achieved by using different depths. In this table, the value inside the parenthesis indicates the depth value.

The presented results demonstrate that the best results are obtained through higher resolutions and using unbalanced datasets. Most of the results were achieved using depths of ten or thirty frames. The next experiments explore a more in-depth analysis of the best configurations, their performance, and statistical validation.

4.2 Statistical Validation

Once the exploration tests end, we analyze the results to decide which parameters improve the classification and select configurations with the best performance. As a second experiment, the prominent configurations were run thirty times, using cross-validation, to give statistical support to the results that previously presented. For this experiment, the configurations train with the largest datasets we already create (SBT_balanced_240_30t and SBT_unbalanced_60120_30t). Complementing with cross-validation, we use ten folds of the dataset for train and test.

From previous results, four configurations were trained with 240 balanced-samples and 180 unbalanced-samples datasets. Fixed parameters were 100 epochs for training, 70% samples for training and 30% for testing, both datasets use the original and the flipped images. The ratio and number of samples per class can be inferred from the balance parameter (see section 4.1, balance and ratio). For this test, we perform thirty runs per configuration and use ten dataset folds for cross-validation.

Table 7 presents average accuracy and the standard deviation of each configuration tested. Most of the results are around 70% accuracy. There is not significative deviation on each training group. The results seem very similar between them. We analyze the confusion matrices to search for biased results. Although we find cases where the classification results are biased to a particular class, we discover good results.

In this investigation, the 80x60 resolution has the best results in suspicious behavior detection task. It achieves accuracy rates above 85% both balanced and unbalanced datasets, preferably with ten frames depth. The best result in a single run is 92.50% of accuracy. This performance was obtained in the thirtieth run, using the unbalanced

Table 6: Best results comparison

Dataset	Resolution							
	32x24		40x30		80x60		160x120	
	Individual	Average	Individual	Average	Individual	Average	Individual	Average
Balanced Dataset	83.3% (90f)	75.9% (30f)	74.0% (90f)	73.4% (90f)	87.0% (30f)	79.6% (10f)	87.0% (30f)	79.6% (10f)
Unbalanced Dataset	84.7% (10f)	77.7% (30f)	86.1% (10f)	76.3% (30f)	91.6% (10f)	87.0% (10f)	90.2% (30f)	86.1% (30f)

Table 7: Thirty runs training results.

Resolution	Dataset	Avg Accuracy	Std Deviation
160x120	unb_60s120n_30t_10f	75.7%	0.0638
	unb_60s120n_30t_30f	73.9%	0.0543
	bal_240_30t_10f	73.1%	0.0661
	bal_240_30t_30f	71.6%	0.0999
80x60	unb_60s120n_30t_10f	75.0%	0.0689
	unb_60s120n_30t_30f	74.8%	0.0500
	bal_240_30t_10f	73.0%	0.0717
	bal_240_30t_30f	73.6%	0.0821
40x30	unb_60s120n_30t_10f	68.7%	0.0569
	unb_60s120n_30t_30f	69.1%	0.0576
	bal_240_30t_10f	71.8%	0.0468
	bal_240_30t_30f	71.9%	0.0555
32x24	unb_60s120n_30t_10f	69.4%	0.0686
	unb_60s120n_30t_30f	71.6%	0.0533
	bal_240_30t_10f	70.3%	0.0476
	bal_240_30t_30f	70.1%	0.0574

dataset and ten frames depth. After 30 runs, the model obtain an average accuracy of 75%.

Table 8 exhibits the best results and their confusion matrices. Even in the confusion matrices, the accuracy per class is above 90% for suspicious-behavior class and around 80% for normal-behavior class.

4.3 Discussion

As the first experiment in this work, we select a 3D Convolutional Neural Network with a basic configuration as a base model, and then we perform a parameter tuning, searching for network model improvement.

From the parameter exploration, we found that 80x60 and 160x120 resolutions deliver better results than a commonly used low resolution or. This experiment was limited to a maximum resolution of 160x120 due to processing resources.

Another significant aspect is the "depth" parameter. This parameter describes the number of consecutive frames used to perform the 3D convolution. After testing different values, we observed that low values, between 10 and 30 frames, have a good relationship between image detail and processing time. These two factors impact the network model training and the correct classification of the samples.

Also, the proposed model can correctly handle flipped images and unbalanced datasets. We performed a more realistic simulation where the dataset has more normal-behavior samples than suspicious-behavior examples. The unbalanced datasets were also correctly classified.

For the second experiment, we use the configurations with the best performance and test them with bigger datasets. We performed 30 runs for each configuration, applying cross-validation, with 10 and 30 frames values for depth and using a 240-samples balanced dataset and a 180-samples unbalanced dataset.

From this experimentation, we found that 80x60 resolution reports better accuracy for the four scenarios we test. Table 7 presents the average accuracy for each configuration. The average performance for the four scenarios in 80x60 resolution is 74.1%, while the 160x120 resolution obtains 73.5%. Also, in a single training, 80x60 resolution performance achieves over 90%, 92.5% for a balanced dataset and 91.6% for the unbalanced dataset.

Finally, when comparing the base model against the proposed one (Table 9), we obtained that our model is capable of improving the classification results by 1.3% and 4.5% on average for balanced and unbalanced datasets,

Table 8: Confusion Matrix of best results
Dataset: **unb_60s120n_30t_10f_80x60**

Accuracy: **92.5%**

	Suspicious	Normal	Accuracy
Suspicious	18	0	100%
Normal	4	32	88.9%

Dataset: **bal_240_30t_10f_80x60**

Accuracy: **91.6%**

	Suspicious	Normal	Accuracy
Suspicious	36	0	100%
Normal	6	30	83.3%

Dataset: **unb_60s120n_30t_10f_80x60**

Accuracy: **90.7%**

	Suspicious	Normal	Accuracy
Suspicious	18	0	100%
Normal	5	31	86.0%

Dataset: **bal_240_30t_10f_80x60**

Accuracy: **90.2%**

	Suspicious	Normal	Accuracy
Suspicious	36	0	100%
Normal	7	29	80.6%

respectively. In the best-single-training comparison, the proposed architecture exceeds 90% accuracy in both cases. The confusion matrices show that for both balanced and unbalanced datasets, the proposed architecture successfully classifies 100% of suspicious samples and obtains a low number of false positives from normal-behavior samples.

4.4 Processing Time

As mentioned before, we use Google Colaboratory to perform the experiments. This tool is based on Jupyter Notebooks and allows the free use of a GPU. The speed of each training depends on the tool demand. Most of the network trainings end in less than an hour, but a higher GPU demand may impact the training time. We are not able to establish a relationship between resolutions and training time, but we have an approximate correlation between different depths.

Table 10 shows the average training time of the final tests, running one hundred epochs, described in section 4.1 (accuracy results of these experiments are shown in Table 6). Comparing the training time from using 10-frames against 30-frames, it increases approximately three times in the second training. We find the same increase's relation when comparing trainings with 30-frames and 90-frames. Some ninety-frames trainings, with a hundred epochs and high resolution, have reached a duration of up to four hours. Training duration is an essential factor due to the size of the used dataset is considerably small.

Another point to consider is the system's accuracy against the training time. Although the training time increases approximately three times, using the same number of epochs, the accuracy is usually lower when using 90-frames depth, in most of the cases. In some instances, we get a higher precision using 90-frames, but accuracy reached by smaller depths was not far from the best one, and the training time was considerably lower.

5 Conclusion

For this work, we focus on the behavior performed by a person before committing a shoplifting crime. The neural network model identifies the previous conduct, looking for suspicious behavior and not to recognize the crime itself. This behavior analysis is the principal reason why we remove the committed crime segment from the video samples, to allow the artificial model to focus on decisive conduct and not in the offense. We implement a 3D Convolutional Neural Network due to its capability to obtain abstract features from signals and images, based on previous approaches for action recognition and movement detection.

Table 9: Comparison between base model and the proposed one.

	Balanced		Unbalanced	
	Base model	Proposed	Base model	Proposed
Avg Acc	71.7%	73.0%	70.5%	75.0%
Std Dev	0.06	0.07	0.05	0.07
Best Result	81.9%	91.6%	88.8%	92.5%

	Base Model		Proposed	
			Susp	Norm
	Susp	Norm	Susp	Norm
Balanced	34	2	36	0
	11	25	6	30
Unbalanced	16	2	18	0
	4	32	4	32

Table 10: Average training times in seconds comparison between different depths and resolutions.

Dataset		Resolution			
		32x24	40x30	80x60	160x120
SBT_unbalanced_60120...	10f	96	126	369	1,356
	30f	196	279	1,027	3,918
	90f	518	758	2,929	11,655
SBT_balanced_240...	10f	118	157	475	1,714
	30f	257	364	1,304	4,952
	90f	688	1,011	3,879	15,415

Based on the results obtained from the conducted experimentation, a 75% accuracy in suspicious behavior detection, we can state that it is possible to model the suspicious behavior of a person in the shoplifting context. Through the presented experimentation, we found which parameters fit better for behavior analysis, particularly for the shoplifting context. We explore different parameters and configurations, and in the end, we compare our results against a reference 3D Convolutional architecture. The proposed model demonstrates a better performance with balanced and unbalanced datasets and using the particular configuration obtained from previous experiments.

The final intention of this experimentation is the development of a tool capable of supporting the surveillance staff, presenting visual behavioral cues, and this work is a first step to achieve the mentioned goal. From this point, we will explore different aspects that will contribute to the project development, such as bigger datasets, adding more criminal contexts that present suspicious behavior and real-time tests.

5.1 Future Work

For these experiments, we use a selected number of videos from the UCF-Crimes dataset. To test in a more realistic simulation, we have to increase the number of samples, preferably the normal-behavior ones, to create a bigger sample-unbalance between classes.

Another interesting aspect of the developing of this project is to expand our behavior detection model to other contexts. It exists many situations where we can find suspicious behavior, such as stealing, arson intents, burglary, among others. We will gather videos of different contexts to strengthen the capability to detect suspicious behavior.

References

- [1] National Retail Federation. *2018 National Retail Security Survey*. National Retail Federation, June 2018.
- [2] S. O. Ba and J. Odobez. Recognizing visual focus of attention from head pose in natural meetings. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 39(1):16–33, Feb 2009.

- [3] Nandita M. Nayak, Ricky J. Sethi, Bi Song, and Amit K. Roy-Chowdhury. *Modeling and Recognition of Complex Human Activities*, pages 289–309. Springer London, London, 2011.
- [4] S. Rankin, N. Cohen, K. MacLennan-Brown, and K. Sage. Cctv operator performance benchmarking. In *2012 IEEE International Carnahan Conference on Security Technology (ICCST)*, pages 325–330, 2012.
- [5] FaceFirst. Official website, 2019. (accessed 6 April 2019).
- [6] DeepCam. Official website, 2018. (accessed 6 April 2019).
- [7] Xin Geng, Gang Li, Yangdong Ye, Yiqing Tu, and Honghua Dai. Abnormal behavior detection for early warning of terrorist attack. In Abdul Sattar and Byeong-ho Kang, editors, *AI 2006: Advances in Artificial Intelligence*, pages 1002–1009, Berlin, Heidelberg, 2006. Springer Berlin Heidelberg.
- [8] Waqas Sultani, Chen Chen, and Mubarak Shah. Real-world anomaly detection in surveillance videos. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6479–6488, 06 2018.
- [9] K. K. Hati, P. K. Sa, and B. Majhi. Lobs: Local background subtracter for video surveillance. In *2012 Asia Pacific Conference on Postgraduate Research in Microelectronics and Electronics*, pages 29–34, 2012.
- [10] D. Berjon, C. Cuevas, F. Moran, and N. Garcia. GPU-based implementation of an optimized nonparametric background modeling for real-time moving object detection. *IEEE Transactions on Consumer Electronics*, 59(2):361–369, 2013.
- [11] Joshila Grace.L.K and K. Reshma. Face recognition in surveillance system. In *2015 International Conference on Innovations in Information, Embedded and Communication Systems (ICIIECS)*, pages 1–5, March 2015.
- [12] Ade Nurhopipah and Agus Harjoko. Motion detection and face recognition for cctv surveillance system. *IJCCS (Indonesian Journal of Computing and Cybernetics Systems)*, 12:107, 07 2018.
- [13] Tang Sze Ling, Liang Kim Meng, Mei Lim, Kadim Zulaikha, Ahmed A. Bahaa'a, and Ahmed Al-Obaidi. Colour-based object tracking in surveillance application. In *Proceedings of the International MultiConference of Engineers and Computer Scientists 2009, IMECS 2009*, volume 1, page No pages, 03 2009.
- [14] J. S. Kim, D. H. Yeom, and Y. H. Joo. Fast and robust algorithm of tracking multiple moving objects for intelligent video surveillance systems. *IEEE Transactions on Consumer Electronics*, 57(3):1165–1170, 2011.
- [15] L. Hou, W. Wan, K. Han, R. Muhammad, and M. Yang. Human detection and tracking over camera networks: A review. In *2016 International Conference on Audio, Language and Image Processing (ICALIP)*, pages 574–580, 2016.
- [16] Joohyung Kang and Sooyeong Kwak. Loitering detection solution for cctv security system. *Journal of Korea Multimedia Society*, 17, 01 2014.
- [17] Jing-Ying Chang, Huei-Hung Liao, and Liang-Gee Chen. Localized detection of abandoned luggage. *EURASIP Journal on Advances in Signal Processing*, 2010(1):675784, 2010.
- [18] S. Wu, H. Wong, and Z. Yu. A bayesian model for crowd escape behavior detection. *IEEE Transactions on Circuits and Systems for Video Technology*, 24(1):85–98, Jan 2014.
- [19] Manuel Alvar, Andrea Torsello, Alvaro Sanchez-Miralles, and José María Armingol. Abnormal behavior detection using dominant sets. *Machine Vision and Applications*, 25(5):1351–1368, 2014.
- [20] Tian Wang, Meina Qiao, Yingjun Deng, Yi Zhou, Huan Wang, Qi Lyu, and Hichem Snoussi. Abnormal event detection based on analysis of movement information of video sequence. *Optik-International Journal for Light and Electron Optics*, 152:50–60, 2018.
- [21] Kan Ouvirach, Shashi Gharti, and Matthew N. Dailey. Automatic suspicious behavior detection from a small bootstrap set. In *Proceedings of the International Conference on Computer Vision Theory and Applications(VISAPP-2012)*, pages 655–658, 2012.
- [22] Mohammad Sabokrou, Mahmood Fathy, Zahra Moayed, and Reinhard Klette. Fast and accurate detection and localization of abnormal behavior in crowded scenes. *Machine Vision and Applications*, 28(8):965–985, 2017.
- [23] Hiroaki Tsushita and Thi Thi Zin. A study on detection of abnormal behavior by a surveillance camera image. In Thi Thi Zin and Jerry Chun-Wei Lin, editors, *Big Data Analysis and Deep Learning Applications*, pages 284–291, Singapore, 2019. Springer Singapore.
- [24] Tatsuya Ishikawa and Thi Thi Zin. A study on detection of suspicious persons for intelligent monitoring system. In Thi Thi Zin and Jerry Chun-Wei Lin, editors, *Big Data Analysis and Deep Learning Applications*, pages 292–301, Singapore, 2019. Springer Singapore.
- [25] Tao He, Hua Mao, and Zhang Yi. Moving object recognition using multi-view three-dimensional convolutional neural networks. *Neural Computing and Applications*, 28(12):3827–3835, Dec 2017.

- [26] S. Ji, W. Xu, M. Yang, and K. Yu. 3D convolutional neural networks for human action recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(1):221–231, Jan 2013.
- [27] L. Zhang, G. Zhu, P. Shen, and J. Song. Learning spatiotemporal features using 3DCNN and convolutional lstm for gesture recognition. In *2017 IEEE International Conference on Computer Vision Workshops (ICCVW)*, pages 3120–3128, 2017.
- [28] X. Cai, F. Hu, and L. Ding. Detecting abnormal behavior in examination surveillance video with 3D convolutional neural networks. In *2016 6th International Conference on Digital Home (ICDH)*, pages 20–24, 2016.
- [29] G. Varol, I. Laptev, and C. Schmid. Long-term temporal convolutions for action recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(6):1510–1517, June 2018.
- [30] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1*, NIPS’12, pages 1097–1105, USA, 2012. Curran Associates Inc.
- [31] C. Szegedy, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–9, June 2015.
- [32] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.
- [33] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. *arXiv e-prints*, page arXiv:1512.03385, Dec 2015.
- [34] Fujimoto Lab. 3D convolutional neural network for video classification, code repository, 2017. accessed 28 April 2019.
- [35] Google. Google colaboratory, 2017. accessed 29 April 2019.