



UNIVERSITÀ
DEGLI STUDI
DI UDINE

Università degli studi di Udine

Low-cost CNN for Automatic Violence Recognition on Embedded System

Original

Availability:

This version is available <http://hdl.handle.net/11390/1220689> since 2022-06-30T10:50:45Z

Publisher:

Published

DOI:10.1109/ACCESS.2022.3155123

Terms of use:

The institutional repository of the University of Udine (<http://air.uniud.it>) is provided by ARIC services. The aim is to enable open access to all the world.

Publisher copyright

(Article begins on next page)

Received February 7, 2022, accepted February 24, 2022, date of publication February 28, 2022, date of current version March 10, 2022.

Digital Object Identifier 10.1109/ACCESS.2022.3155123

Low-Cost CNN for Automatic Violence Recognition on Embedded System

JOELTON CEZAR VIEIRA¹, ANDREZA SARTORI^{1,2}, STÉFANO FRIZZO STEFENON^{3,4},
FÁBIO LUIS PEREZ¹, GABRIEL SCHNEIDER DE JESUS^{1,2},
AND VALDERI REIS QUIETINHO LEITHARDT^{5,6}, (Member, IEEE)

¹Department of Telecom., Electrical and Mechanical Engineering, Regional University of Blumenau (FURB), Rua São Paulo 3250, 89030-000 Blumenau, Brazil

²Department of Information Systems and Computing, Regional University of Blumenau (FURB), Rua Antônio da Veiga 140, 89030-903 Blumenau, Brazil

³Fondazione Bruno Kessler, Via Sommarive 18, 38123 Povo, Trento, Italy

⁴Computer Science and Artificial Intelligence, University of Udine, Via delle Scienze 206, 33100 Udine, Italy

⁵COPELABS, Lusófona University of Humanities and Technologies, Campo Grande 376, 1749-024 Lisboa, Portugal

⁶VALORIZA, Research Center for Endogenous Resources Valorization, Instituto Politécnico de Portalegre, 7300-555 Portalegre, Portugal

Corresponding author: Andreza Sartori (asartori@furb.br)

This work was supported by the National Funds through the Fundação para a Ciência e a Tecnologia, I.P. (Portuguese Foundation for Science and Technology) by the Project “VALORIZA—Research Centre for Endogenous Resource Valorization” under Grant UIDB/05064/2020 and Grant UIDB/04111/2020, and in part by the Instituto Lusófono de Investigação e Desenvolvimento (ILIND) under Project COFAC/ILIND/COPELABS/3/2020.

ABSTRACT Due to the increasing number of violence cases, there is a high demand for efficient monitoring systems, however, these systems can be susceptible to failure. Therefore, this work proposes the analysis and application of low-cost Convolutional Neural Networks (CNNs) techniques to automatically recognize and classify suspicious events. Thus, it is possible to alert and assist the monitoring process with a reduced deployment cost. For this purpose, a dataset with violence and non-violence actions in scenes of crowded and non-crowded environments was assembled. The mobile CNNs architectures were adapted and obtained a classification accuracy of up to 92.05%, with a low number of parameters. To demonstrate the models' validity, a prototype was developed by using an embedded Raspberry Pi platform, able to execute a model in real-time with 4 frames-per-second of speed. In addition, a warning system was developed to recognize pre-fight behavior and anticipate violent acts, alerting security to potential situations.

INDEX TERMS Neural networks, artificial neural networks, image processing, image classification.

I. INTRODUCTION

There is a growing interest in intelligent surveillance systems due to major concerns about global security and the need for effective monitoring of public places, such as airports, railway stations, malls, sports stadiums, tourist venues, etc. Indeed, the number of cameras installed in urban areas is increasing progressively to promote order and safety [1].

Currently, most public monitoring systems are performed through security cameras, mainly in areas with a large flow of people. In fact, monitoring systems are basically composed of several cameras positioned at strategic locations. The problem with these systems consists of having only one human agent responsible for tracking the video input from many cameras simultaneously [2]. This can lead to errors when identifying suspicious events, often caused by inattention or fatigue. As such, monitoring systems are used as a storage

The associate editor coordinating the review of this manuscript and approving it for publication was Byoung Wook Choi.

system for possible lawsuits, rather than an incident prevention system as it should be [3].

In an automatic surveillance system, computers continually process a video input and scrutinize each frame to find a suspicious event, in which they can immediately report to the supervisor for their attention [4]. For this reason, the use of Computer Vision and Machine Learning techniques applied to automatic recognition of violence can help monitor agents significantly.

The research in Computer Vision, specifically in action recognition, has mainly focused on detecting simple actions, such as walking or sports activities. The detection and recognition of fights or aggressive behavior, in general, has been comparatively less studied [2]. In practice, the automatic recognition of aggressive and irregular actions can be extremely useful for various video surveillance scenarios outside urban areas, such as in prisons, psychiatric hospitals, and even in monitoring daily activities and detecting falls of elderly people in homes [5]. Neural network models are

increasingly being applied in embedded systems [6]–[10] to improve the ability to analyze data from different equipment [11]–[14].

Thus, this work proposes an automatic violence recognition system in real-time, able to distinguish between acts of violence, such as fights, vandalism, shooting with firearms, and non-violence actions, such as walking, running, hugging, kissing, exercising, and celebrating. For this, several videos were assembled from a diverse range of public datasets, containing varied lighting conditions, scale, movement, number of people, and objects. Moreover, the study used and compared different models of Mobile Convolutional Neural Networks. To avoid overfitting and improve the model generalization, traditional data augmentation techniques were applied. Then, the models were implemented and evaluated on a Raspberry Pi 4 embedded platform.

To summarize, the contributions of this work are the following:

- A quantitative and qualitative analysis of state-of-the-art Mobile CNN architectures for the binary recognition of violence actions;
- A dataset for the violence acts recognition was set up, through a manual selection and combination of public datasets. This dataset has several human actions and activities divided into two classes: violence and non-violence. In addition, contains scenes of violence in crowded and non-crowded environments;
- It was possible to achieve an accuracy of up to 92.05%, with models containing 2.26 million parameters;
- The trained models are able to recognize various actions of violence, such as punching, kicking, fighting, attacking, destroying, aiming and firing guns, wrestling, and boxing.
- A prototype of a low-cost, intelligent monitoring system was developed on a Raspberry Pi embedded platform, able to run a mobile CNN model with a processing speed of up to 4 frames-per-second.
- A novel approach was created with the warning system, which is able to recognize pre-fight behavior and alert security to take appropriate action.

This paper is structured as follows: Section 2 presents the related work, with a brief description of existing automatic violence recognition works and the methods used. Section 3 corresponds to the development of the dataset, the selected public datasets, the distribution of training and testing set, and a video duration histogram. In Section 4 is described the adapted versions of the mobile CNNs with the preprocessing steps. Section 5 presents the experiments, with a comparison of the developed models' results, analyzing the error, accuracy, and number of parameters. Section 6 details the development, operation, and comparison of the prototype. Section 7 presents a novel approach with the warning system. The last section is devoted to the final discussion and conclusions.

II. RELATED WORK

Recently, Deep Learning techniques, such as Convolutional Neural Networks (CNN or ConvNet), have shown excellent results in image and video classification [15]–[18]. In different challenges and datasets, these structures have been performing much better than previous proposals [3]. In fact, there are three main advantages of using CNN models in intelligent monitoring systems. First, they are less affected by noise in the data. Second, they achieve higher accuracy than other methods, even sometimes greater than the human eye. Lastly, they have the ability to classify people into different orientations and postures. Moreover, they also do not require a hand-crafted extractor for encoding features [19], as was performed before the introduction of Deep Learning [20].

For example, [21] developed a violence detection system in movie scene, in which applied various elements of Deep Learning. First, a CNN was trained, then a two-stream CNN was used to extract both static and optical flow motion features. Finally, a Long Short-Term Memory (LSTM) [22] was applied to extract the long-term temporal features [23]. Complimentary motion and audio information were also extracted.

Also, [24] presented a model using 3D Convolutional Neural Networks directly on raw input data that automatically extracts the features. The 3D ConvNets extracts static features such as traditional CNN 2D, as well as adds the temporal dimension, thus allowing the extraction of motion features. Model evaluation is performed on the Hockey dataset [2]. Similarly, [20] proposed a Deep Learning model based on 3D CNN, adapting bottleneck units and the DenseNet architecture to promote efficient extraction of spatiotemporal features. Currently, [25] proposes a method that extracts the video spatiotemporal features through a convolutional neural network and combines them with the trajectory features in order to detect violence in video.

While the aforementioned works successfully recognize and classify acts of violence with good accuracy, all require significant processing power, as usually are employing high-end systems with multiple GPUs or TPUs [26].

The application of Deep Neural Networks in embedded systems or with limited computational capacity has been considerably less studied. Applications such as face recognition, gender detection [27], and emotion recognition [28] are some of the real-time models developed implementing the Raspberry Pi embedded platform. However, no applications or models were found on the violence classification and recognition by only utilizing the embedded system's capabilities.

Although typically the training phase requires more processing power than running the model [29], computationally limited systems such as small computers and embedded systems still face major difficulties in reproducing such models efficiently [30]. For this reason, this work conducts an extensive evaluation of Deep Neural Networks recently developed for embedded platforms and mobile applications, namely mobile CNN architectures [31]. With this study, it was developed a low-cost intelligent monitoring system on a Raspberry



FIGURE 1. Samples from the Violent-Flow dataset [33].

Pi embedded platform able to recognize pre-fight behavior and alert security to take appropriate action.

III. DATASET

To successfully train any deep CNN model, many video samples are needed [32]. However, there are only a few public video datasets available, specifically for violence recognition. In addition, they have an insufficient number of videos or are not applicable for this work. For this reason, action recognition videos were gathered in public datasets, in which were manually selected the classes and videos to be used.

The dataset developed is divided into training and testing set with two classes: violence and nonviolence. The violence class contains violent actions, such as punching, kicking, fighting, attacking, destroying, aiming, and firing guns, wrestling, and boxing. The nonviolence class contains typical actions for target places, such as malls, airports, subways and public parks, actions such as walking, jogging, running, sitting, hugging, kissing, walking the dog, exercising, biking, and celebrating. Thus, for this work, four public datasets were selected and combined: Violent-Flow¹ [33], UCF-101² [34], HMDB³ [35], and Moments in Time⁴ [36].

The Violent-Flow dataset [33] is the only selected dataset exclusively for violence recognition. It is a real-world dataset, that is, recorded videos are real acts of violence, mostly in crowds, taken by surveillance cameras or smartphones. The dataset contains 246 labeled training videos and 44 test videos. Figure 1 shows some samples of violence and nonviolence from the dataset. It is possible to see that, indeed, most actions are from crowded places.

The UCF-101 [34] is an action dataset for recognition and classification collected from YouTube. In total, contains 101 classes of actions and activities, and more than 13,320 videos. For this work, the classes of videos selected were: “punch” for the violence set, “walking the dog” and “bik-



FIGURE 2. Samples from the UCF-101 dataset [34].

ing” for the nonviolence set. Therefore, resulting in 417 training videos and 88 testing videos. Figure 2 presents some examples from the extracted classes for this work.

The HMDB [35] is a large dataset to recognize actions with 51 classes. Each contains at least 101 videos, for a total of 6,766 videos extracted from digital movies and YouTube videos. For this work were selected the classes “hit”, “kick”, “punch” and “gun” for the violence set; and “hug”, “kiss”, “run”, “walk” and, “sit” for the nonviolence set. Thus, resulting in 612 training videos and 107 testing videos. Figure 3 shows a few samples from the selected classes of the dataset.

Lastly, the Moments in Time [36] is a research project dedicated to building a dataset for recognizing and understanding video actions. Currently, the dataset includes a collection of approximately one million 3-second videos, corresponding to 339 different classes, involving people, animals, objects, or natural phenomena.

For the violence set the classes used were “aiming”, “attacking”, “boxing”, “destroying”, “fighting”, “hitting”, “kicking”, “punching”, “shooting”, and “wrestling”. For nonviolence were used “bicycling”, “celebrating”, “exercising”, “hugging”, “jogging” and “running”, “kissing”, “sitting”, and “walking”.

Due to the diversity of scenes in this dataset, even within the same class, a careful manual selection of videos was necessary, which resulted in 994 training videos and 162 testing videos. In Figure 4 is displayed some frames from the dataset.

Table 1 shows the contribution of videos from each dataset to the final combined dataset. The dataset, assembled for this work, contains a total of 2670 videos, of which 2269 videos belong to the training set and 401 to the testing set. From these videos, the violence class has a total of 1193 videos, of which 1014 are assigned to the training set and 179 to the testing set. The nonviolence class contains 1477 videos in which 1255 are assigned to the training set and 222 to the

¹<https://www.openu.ac.il/home/hassner/data/violentflows/>

²<https://www.crcv.ucf.edu/data/UCF101/>

³<https://serre-lab.clps.brown.edu/resource/hmdb-a-large-human-motion-database/>

⁴<http://moments.csail.mit.edu/>



FIGURE 3. Samples from the HMDB dataset [35].



FIGURE 4. Samples from the Moments in Time dataset [36].

TABLE 1. Distribution of videos gathered from public datasets.

| | Dataset | Violence | Nonviolence | Total |
|---------------------|-----------------|-------------|-------------|-------------|
| Train | Violent-Flow | 123 | 123 | 246 |
| | UCF-101 | 146 | 271 | 417 |
| | HMDB | 250 | 362 | 612 |
| | Moments in Time | 495 | 499 | 994 |
| Training Set | | 1014 | 1255 | 2269 |
| Test | Violent-Flow | 22 | 22 | 44 |
| | UCF-101 | 40 | 48 | 88 |
| | HMDB | 45 | 62 | 107 |
| | Moments in Time | 72 | 90 | 162 |
| Testing Set | | 179 | 222 | 401 |
| Total | | 1193 | 1477 | 2670 |

testing set. The difference between the number of train and test videos is due to the greater number of nonviolence classes in the action recognition datasets.

The total video length for the training set is 128 minutes, and 22 minutes for the testing set. The average length of videos is 3.35 seconds. The total duration length distribution of videos is presented in Figure 5. It is possible to observe that most of the videos have a duration lower than 3 seconds. This

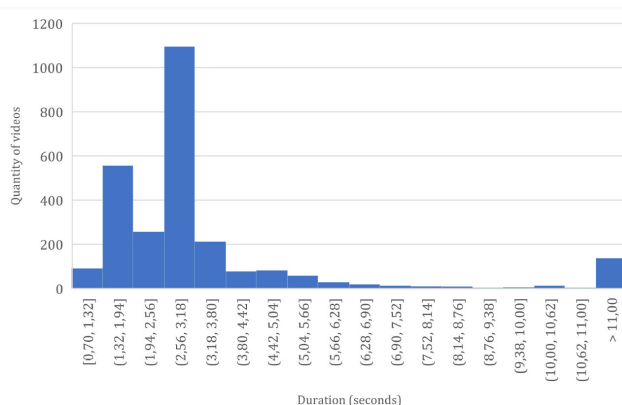


FIGURE 5. Histogram of video duration.

is a common characteristic video of violence, in which the actions happen very quickly. Instances higher than 4 seconds usually correspond to the nonviolence set, with actions such as walking and running.

Furthermore, the final dataset provides great diversity in actions, with large variations in camera position, appearance and pose of objects and people, scale, viewpoint, cluttered backgrounds, different lighting conditions, wide variations in motion, video quality, and occlusion. This allows the models to learn and predict the most diverse actions and activities characterized between violent and nonviolent.

IV. PROPOSED METHOD

CNNs have become ubiquitous in Computer Vision since the popularization of the AlexNet architecture [37]. The general trend has been to develop deeper and more complex networks to achieve greater accuracy. However, these improvements do not necessarily make architectures more efficient in terms of number of parameters, model size and processing speed [38]. In many applications, such as robotics, autonomous cars, augmented reality, and security. The recognition task needs to be employed at a specific time on a restricted computing platform, due to time constraints or space allocation.

Recently, there are attempts to develop CNN architectures for recognition and detection tasks focused on application in devices with limited computational power, such as smartphones and embedded systems. In this work, it was adapted the most popular of these architectures to the binary classification for violence recognition.

In the following subsections, a description of the used mobile CNN architectures is presented, with a brief introduction about the main method used by them. It is also shown the preprocessing steps, the network hyperparameters, and the final number of parameters for each architecture.

A. PREPROCESSING & NETWORK HYPERPARAMETERS

The videos applied on the mobile CNNs were converted into a series of images, with dimensions $227 \times 227 \times 3$ or $224 \times 224 \times 3$, depending on the CNN architecture. These dimensions are considered ideal for features representation and extraction during training, as large images require high

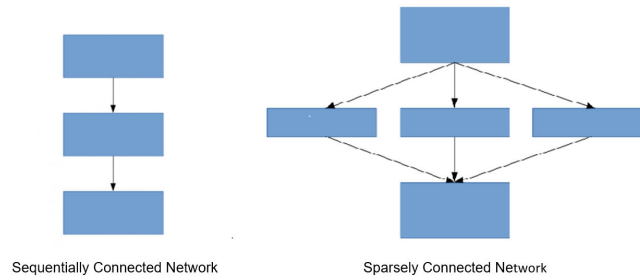


FIGURE 6. Types of neural network connections.

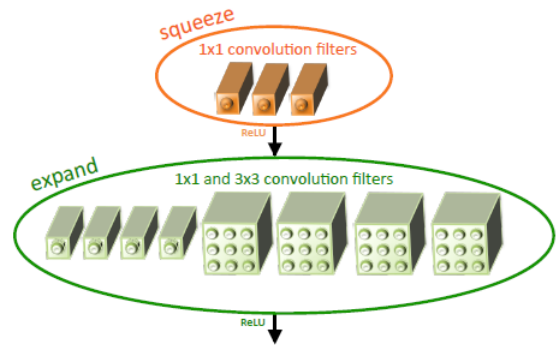


FIGURE 8. The squeezed module [50].

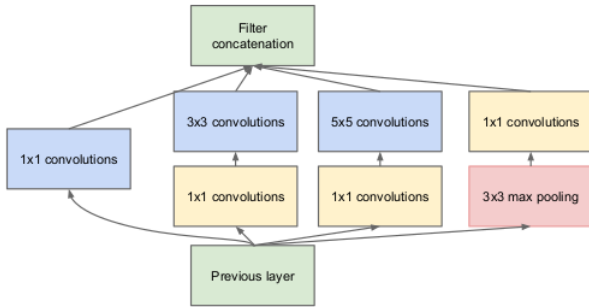


FIGURE 7. The inception module [43].

computational power, and in small images, the CNN may not be able to extract significant information. In addition, in-place/on-the-fly data augmentation techniques were used to better generalize the dataset [39]. Thus, each image in a batch is transformed by a series of random operations in real-time during the training process, among them:

- Mirroring;
- Approximation (up to 20%);
- Rotation (up to 20%);
- Width and height offset (up to 20%).

In addition, this work used the Adam optimization algorithm instead of the traditional SGD optimizer, responsible for minimizing the objective function [40], [41]. Adam’s algorithm uses adaptive learning methods to find individual learning rates for each parameter and enable ANNs to train faster. The learning rate lr , which represents the weight update step, has been set at $lr = 1 \cdot 10^{-5}$. Due to memory capabilities, the batch size was set in 32 samples per iteration.

B. MOBILE ARCHITECTURES

A strategy to reduce parameters and operations of mobile CNN architectures is to adapt deep sparse convolutional structures with small filter sizes. That is, instead of the traditional method of adding layer after layer in sequential order, these architectures add layers in sequential and parallel combination, normally structured in modules [42].

The final architecture is composed of merging these modules. This method is presented in Figure 6 and an example of a sparsely connected module seen in the popular Inception architecture [43] is shown in Figure 7.

The Inception module introduces a combination of layers, with 1×1 , 3×3 , and 5×5 convolution layers and max-pooling

layer sparsely connected. The output from these layers are concatenated into a single output vector, forming an input for the next step [44]. This technique was then expanded in the mobile architectures, mainly because it permits a CNN to capture more details and features in diverse scales, with a reduced number of parameters (compared to a sequentially connected network) [45].

1) SQUEEZENET

SqueezeNet [46] is a architecture that achieves AlexNet accuracy level [47], yet with fewer parameters. This architecture uses a combination of modules with the idea of squeezing and expanding layers. In the squeeze layer, there are only 1×1 convolutional filters, whereas in the expand layer there is a combination of 1×1 and 3×3 filters. This high use of 1×1 convolutional filters reduces considerably the number of parameters. After all, 1×1 filters use significantly fewer parameters than larger convolutional filters. For example, a 1×1 filter has 9 times fewer parameters than a 3×3 convolution filter.

Adapted for this work, the SqueezeNet model presents only 735.94 thousand parameters, which is a number far lower than the AlexNet original architecture, which has 62.3 million parameters [48]. Figure 8 shows the SqueezeNet module, with the squeeze and expand layers. Also, the architecture uses the ReLU activation function [49].

2) MOBILENET-V1 & MOBILENET-V2

The MobileNet architecture [51] is primarily based on depthwise separable convolution, in which factors a traditional convolution into a depthwise convolution followed by a pointwise convolution. In other words, a spatial convolution is performed independently for each channel, then by a 1×1 convolution across all channels. This approach was found to be easier than the normal 3D convolution [52]. Thus, the MobileNet model adapted for this work has a total number of 3.23 million parameters.

In the second version of the MobileNet architecture [53], an inverted residual sparse structure was introduced, consisting of 1×1 convolution, depthwise separable convolution, and the use of a linear function. Adapted for this work, the

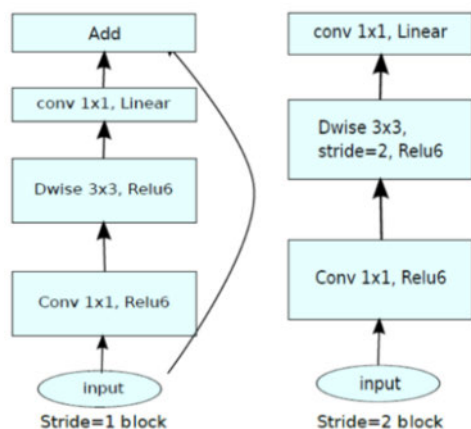


FIGURE 9. The MobileNet-v2 modules [55].

MobileNet-v2 model has 2.26 million parameters, a considerably lower number of parameters than its first version.

Figure 9 shows the MobileNet-v2 module. The module presents a residual cell (has a residual/identity connection introduced by [54]) with stride of 1, and a resizing cell with a stride of 2. From Figure 9, “conv” is a normal convolution, “dwise” is a depthwise separable convolution, “Relu6” is a ReLu activation function with a magnitude limitation, and “Linear” is the use of the linear function.

3) NASNET MOBILE

The NASNet [56] architecture searches the best combinations of convolution layers, first in a smaller dataset, then expanded the configuration to a larger one. The architecture is composed of a number of sparsely connected convolutions layers, normal and separable, with different filters sizes (1 × 1, 3 × 3, 5 × 5, and 7 × 7). The authors also developed a mobile version, which was adapted for this work and has 4.27 million parameters.

Figure 10 presents the NASNet module used. The NASNet architecture is composed of normal cells and reduction cells. Normal cells are convolution layers that return feature maps, and reduction cells resize the features maps by a factor. From Figure 10 “sep” means depthwise separable convolutions, “identity” are the residual/identity connections, “avg” are layers of average pooling, and “max” are layers of max pooling.

V. EXPERIMENTS

In this section, computational performance and results are compared and analyzed for each developed model. The error, accuracy and number of parameters were compared, as well as a diagram to verify the effectiveness of the models.

The Figure 11 displays the accuracy and error per epoch for the training set. The graph shows that the accuracy increases progressively while the error decreases. The MobileNets and NASNet have similar outputs, and the SqueezeNet is slightly less effective. In addition, Figure 12 shows the accuracy and

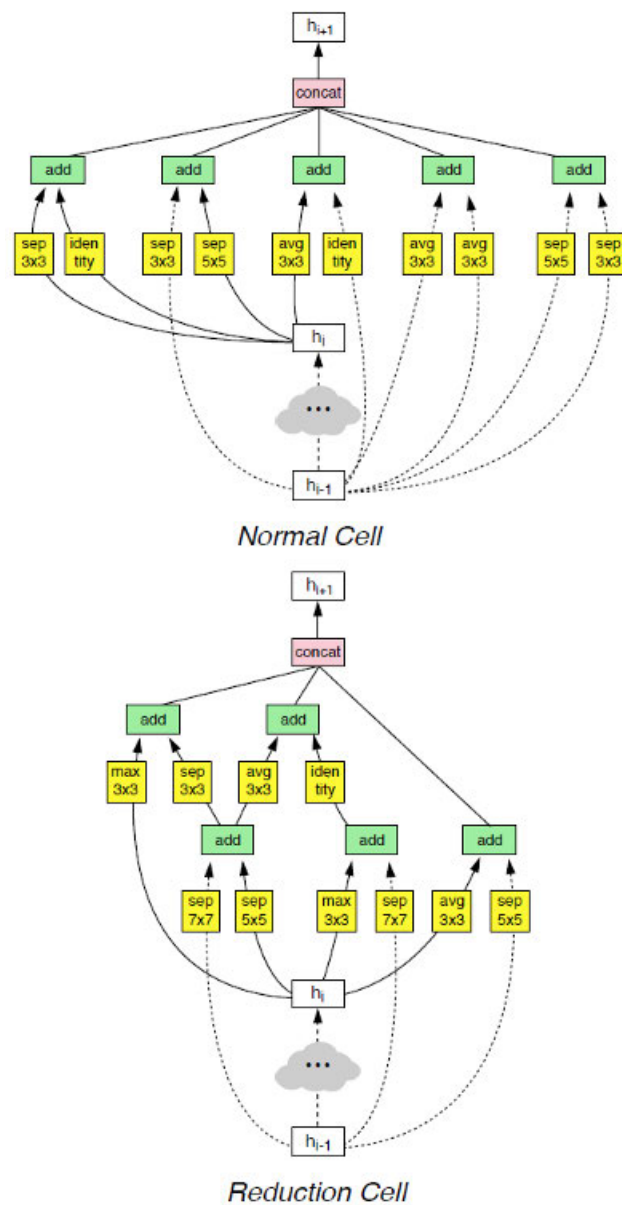


FIGURE 10. The NASNet modules [56].

error per epoch for the testing set. Again, the accuracy and the error of the MobileNets and NASNet models are similar, with slightly worse results from SqueezeNet.

It is also noticeable the presence of overfitting soon after the fifth epoch in the error graph from Figure 12. The errors increase as the models have already learned all patterns from the training set, including noise and outliers, and can no longer generalize new samples from the testing set. This graph also shows how quickly the models can learn all patterns from the dataset and correctly predict new samples.

Therefore, from Figure 12, it is possible to identify in which epoch the best result is obtained, that is, the epoch with the highest accuracy and lowest error. For the SqueezeNet and MobileNet-v1 architectures this happened in the fourth

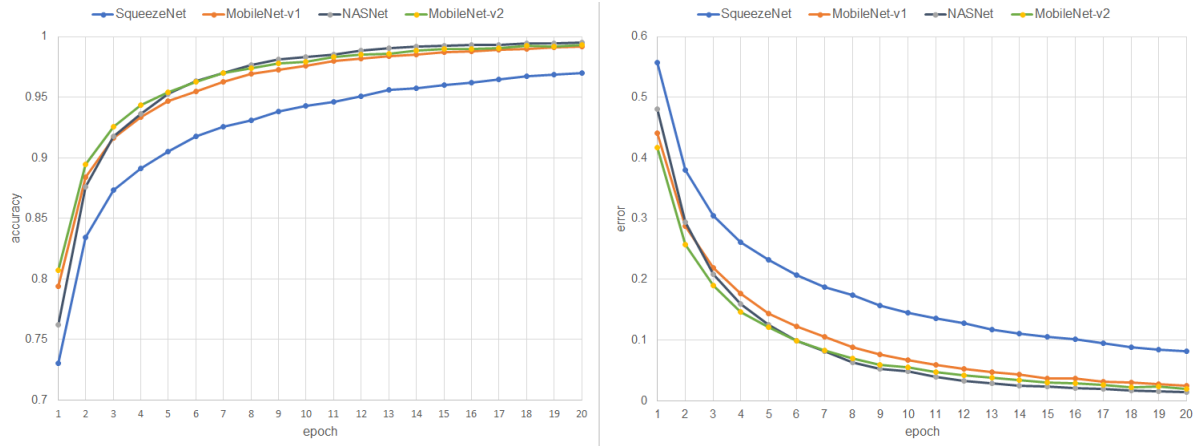


FIGURE 11. Accuracy (left) and error (right) for the training set.

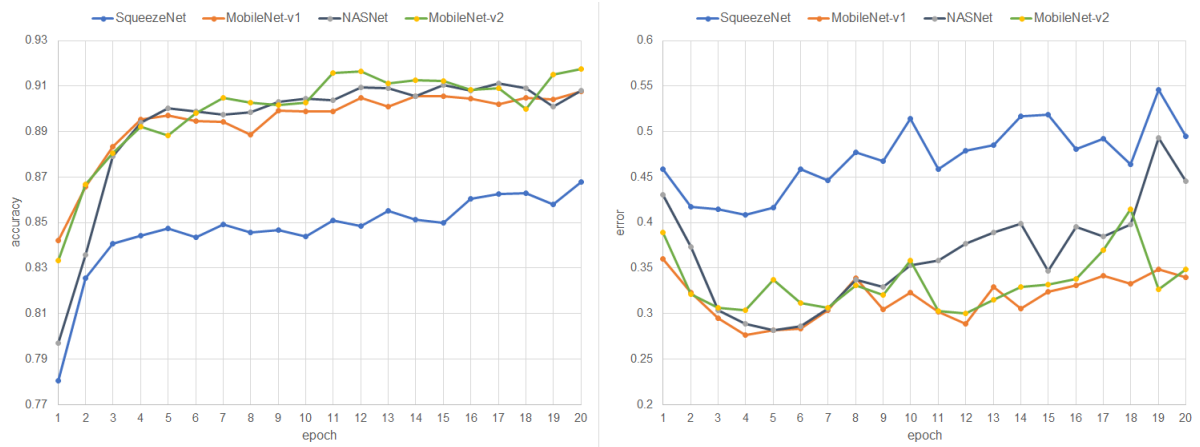


FIGURE 12. Accuracy (left) and error (right) for the testing set.

epoch. SqueezeNet achieved an accuracy of 0.8702 and an error of 0.3486, and MobileNet-v1 achieved an accuracy of 0.9048 and an error of 0.2888. For the Mobile NASNet it occurred in the fifth epoch, with an accuracy of 0.9001 and an error of 0.2813.

Lastly, for the MobileNet-v2 it happened in the twelfth epoch, with an accuracy of 0.9163 and an error of 0.2997. Thus, the best error result belongs to the MobileNet-v1 with an error of 0.2888 and the best accuracy belongs to the MobileNet-v2 with an accuracy of 0.9163. Although, as observed earlier, MobileNet-v1, MobileNet-v2 and NASNet show very close results, SqueezeNet performed slightly worse.

Figure 13 shows the total number of parameters from each adapted architecture for this work. It is possible to observe that the model with the lowest number of parameters belongs to the SqueezeNet with only 735.94 thousand parameters, three times lower than the second model with the lowest number of parameters, MobileNet-v2 with 2.26 million parameters. The MobileNet-v1 and NASNet have 3.23 and 4.27 million parameters, respectively.

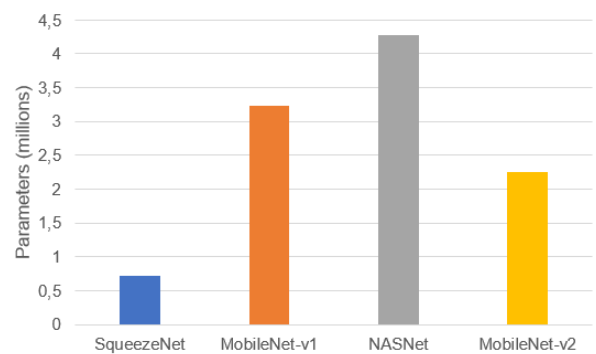


FIGURE 13. Comparative total number of parameters for each architecture.

Lastly, Figure 14 shows the accuracy versus the total number of parameters for each adapted architecture. The SqueezeNet has the lowest number of parameters from all the models, however have a lower accuracy rate. The MobileNet-v1, MobileNet-v2, and NASNet show high accuracy but demand more processing power. Therefore, analyzing the

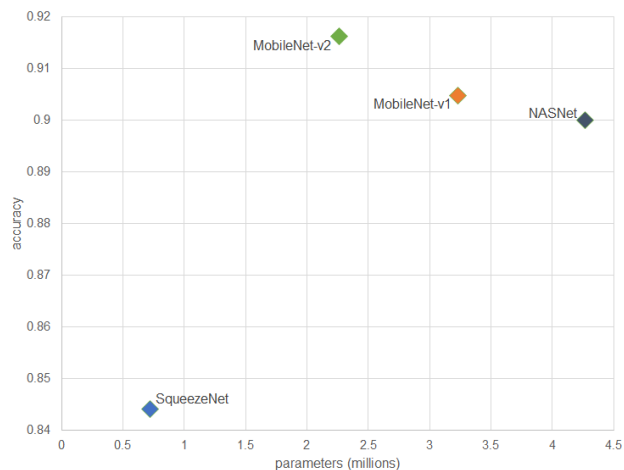


FIGURE 14. Scatter plot of accuracy vs. number of parameters.

TABLE 2. Development environment of raspberry Pi 4 and desktop.

| | Raspberry Pi 4 | Laptop |
|----------------|--------------------------------|------------------------------------|
| OS Distributor | Debian | Ubuntu |
| Description | Debian GNU/Linux 11 (bullseye) | Ubuntu 18.04.6 LTS (Bionic Beaver) |
| Release | 11 | 18.04 |
| Kernel Version | 5.10.92-v8+ | 5.4.0-84-generic |

graphic, it is possible to observe that the MobileNet-v2 presents the best cost-benefit as it has a high accuracy while keeping a low number of parameters. However, limited systems will still be able to run the SqueezeNet model more effectively, due to its reduced number of parameters.

VI. PROTOTYPE

SqueezeNet, MobileNet-v1, MobileNet-v2, and NASNet Mobile are architectures specially developed focused on mobile and embedded applications. Therefore, to prove if they are truly able to run a mobile application for complex tasks, such as violent recognition, the Raspberry Pi 4 embedded platform was selected as a prototype. Table 2 presents the development environment. The framework used for both Laptop with GPU and Raspberry Pi 4 was TensorFlow 2.6.0. Moreover, for the Laptop with GPU, was used CUDA 11.4 and cuDNN 8.2.4.

This section presents the operation process of the prototype, with experiments on the average runtime and the FPS (frames-per-second) that the prototype can reach while running the models. Comparative tables of the same models running on a more complex platform are also presented.

The prototype for this work can run a CNN model, read, and transform a video signal from a file or a USB webcam to an input signal, which allows the prediction to be made by the model. The result of this prediction is then displayed over the screen along with the respective label. The transformations necessary to feed the CNN model consist of transforming the video signal into a series of images, resizing, and normalizing

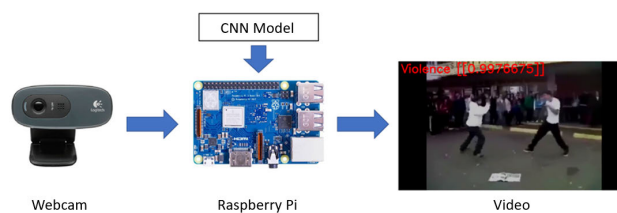


FIGURE 15. The prototype operation process.

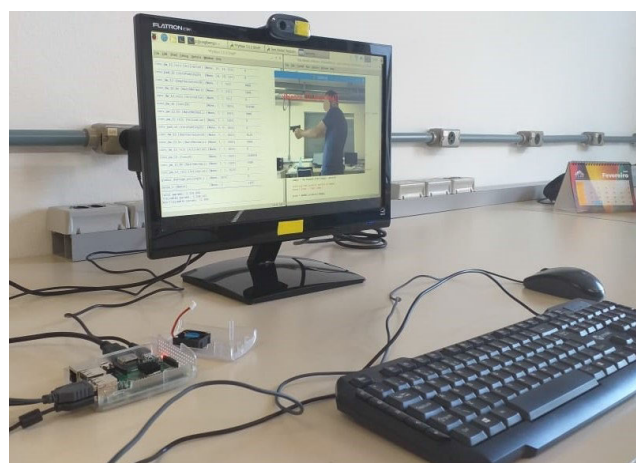


FIGURE 16. The prototype running a model.

TABLE 3. Comparative number of parameters, average runtime and FPS on raspberry Pi 4.

| Architecture | Parameters | Runtime (s) | FPS |
|--------------|------------|-------------|------|
| SqueezeNet | 735,937 | 98.25 | 4.19 |
| MobileNet-v1 | 3,229,889 | 113.60 | 3.74 |
| MobileNet-v2 | 2,259,265 | 103.05 | 4.05 |
| NASNet | 4,270,773 | 153.86 | 3.02 |

the images. Figure 15 presents an example of the prototype operation process and the Figure 16 shows the prototype running a model.

After the attempt to run the developed mobile models, Table 3 shows the adapted architecture, with its total number of parameters, the average runtime, and the average FPS rate. These values are obtained based on the average value to predict 30 seconds of lifestreaming with 10 repetitions.

From Table 3, it is possible to observe that the SqueezeNet architecture has the shortest average runtime. This result is expected, as the SqueezeNet architecture has the smallest number of parameters and makes high use of small convolutional filters as a technique for a lower computational cost.

However, the MobileNet-v1 and MobileNet-v2 architectures have a close average runtime, even though they have a higher number of parameters (compared to non-mobile CNNs). Nevertheless, the MobileNet architectures are much deeper, with many normal and depth-separable convolution layers, capable of extracting more features and information, and presenting superior accuracy.

TABLE 4. Comparative number of parameters, average runtime and FPS on desktop.

| Architecture | Parameters | Runtime (s) | FPS |
|--------------|------------|-------------|-------|
| SqueezeNet | 735,937 | 1.46 | 24.22 |
| MobileNet-v1 | 3,229,889 | 1.32 | 20.43 |
| MobileNet-v2 | 2,259,265 | 1.32 | 22.76 |
| NASNet | 4,270,773 | 2.24 | 14.70 |

The Raspberry Pi 4 platform was not efficient to run the NASNet Mobile architecture. This is due to NASNet is an extremely complex architecture, with several layers sparsely connected and various filter sizes, which makes it difficult to run even on more complex platforms.

To compare the influence of a more powerful setting, the same models were executed on a laptop with a dedicated GPU with the following specifications: Core i7 7700HQ processor, 16GB RAM, NVIDIA GTX 1050Ti GPU. The results are presented in Table 4.

According to Table 4, the influence of a more powerful setting with a dedicated GPU is clear. Models running on the laptop perform on average 98% better than those running on the Raspberry Pi platform. The use of GPUs improves performance due to parallel operations, where the multiple cores available in GPUs allow running multiple operations at the same time.

VII. WARNING SYSTEM

The work of violence recognition is typically a binary classification problem, that is, it has only two classes: violence and nonviolence. In these cases, the function for calculating the probability of belonging to a class is performed by the sigmoid function. The output values can only assume values between 0 and 1. Therefore, a threshold is established to define to which class an input belongs. For comparison purposes, Table 5 presents the classification results using different algorithms. In this case, the runtime is the average time to predict 1024 images with batch size of 32 and 10 repetitions.

In this application it can be observed that the results of accuracy and F1-score make the application possible, highlighting MobileNet-v2 that obtained a result of 92% accuracy and F1-score of 0.92 for the classification of situations of violence. Table 6 presents a benchmarking with well-established models based on deep neural networks. As can be seen, the use of models that require greater computational effort does not significantly improve the accuracy results in this application. In presenting the comparison between the models in Table 5 and Table 6 the time considered is the total time, which comprises the sum of the training time plus the execution time on the embedded system.

After training the model, the outputs generated by the sigmoid function was carefully analyzed frame by frame. It was observed that the outputs close to the threshold correspond to cases in which actions and activities have common characteristics, such as the position of the bodies or the arms around someone. For example, nonviolence scenes that

TABLE 5. Comparison of algorithms.

| Model | Accur. | Prec. | Recall | F1-score | Time (s) |
|--------------|--------|-------|--------|----------|----------|
| SqueezeNet | 0.87 | 0.87 | 0.87 | 0.87 | 114.79 |
| MobileNet-v1 | 0.90 | 0.91 | 0.90 | 0.90 | 114.92 |
| MobileNet-v2 | 0.92 | 0.92 | 0.91 | 0.92 | 104.37 |
| NASNetMobile | 0.90 | 0.90 | 0.90 | 0.90 | 156.28 |

TABLE 6. Convolutional neural network benchmark.

| Model | Accur. | Prec. | Recall | F1-score | Time (s) |
|--------------|--------|-------|--------|----------|----------|
| VGG-16 | 0.91 | 0.92 | 0.91 | 0.91 | 1,276.85 |
| Inception V3 | 0.91 | 0.92 | 0.90 | 0.91 | 262.26 |
| ResNet50 | 0.92 | 0.92 | 0.92 | 0.92 | 483.2 |

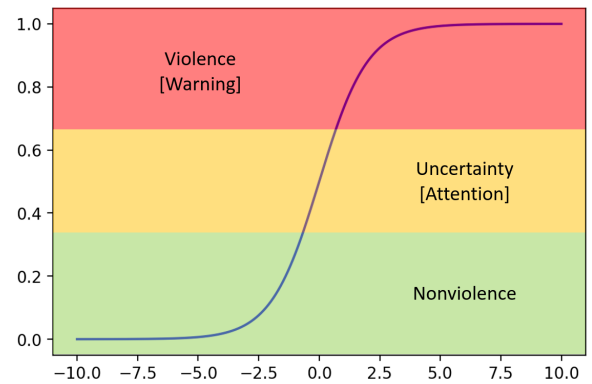


FIGURE 17. Warning system limits for the sigmoid function.

feature actions of “hugging” and “kissing” show many characteristics in common to violence scenes that feature actions of “fighting”, “attacking”, and “beating”. Precisely, it is in these cases that the algorithm can perform the prediction incorrectly.

To reduce this problem, a “third class” was created in the sigmoid function, that is, instead of using only a threshold, a zone of uncertainty was delimited. Thus, in cases where the algorithm is unable to correctly classify an input, the vigilant operator is alerted that the scene may or may not contain violence. The uncertainty zone was defined between the limits [0.35 – 0.65] as presented in Figure 17.

Another important observation is that, in some cases, the frames right before acts of violence, usually with quick and expressive actions and fast movement gestures, the algorithm tends to generate a prediction in the uncertainty zone. Thus, in cases where there is a gradual evolution of violence in the behavior of those involved, such as sudden movements and accusations to finally acts of physical attacks, the algorithm is able to draw the attention of the vigilant operator to be aware of a possible occurrence and take appropriate preventive action.

Figure 18 shows the output generated by the SqueezeNet model, executed on the Raspberry Pi embedded platform, where the label and the prediction result is displayed on the screen. It is possible to observe that nonviolence actions have a prediction result close to zero, violent actions have



FIGURE 18. Model output.



FIGURE 19. Model output of violence in non-crowded environments.



FIGURE 20. Model output of violence in crowded environments.

a prediction result close to one and cases of pre-fights, with quick movements expressing violence, have a prediction result around 0.5.

Moreover, Figure 19 and Figure 20 display the model output applied in crowded and non-crowded environments. The image samples were taken from the testing set and show that the trained models are able to successfully recognize violence in both environments with diverse backgrounds.

VIII. CONCLUSION

This work evaluated how mobile CNNs can perform the task of automatic violence recognition in a new dataset of 2670 videos containing scenes of violence in crowded and non-crowded environments, collected from various public datasets. The experiments have shown that high classification accuracy can be achieved using mobile architectures with a lower number of parameters and, it was able to achieve up to 92.05% of accuracy.

A low-cost prototype of an intelligent monitoring system was presented on a Raspberry Pi embedded platform and used to compare the performance of different developed mobile CNN models, which was able to run a real-time model on up to 4.19 FPS. Experimental results demonstrated that it is possible to achieve with mobile CNN models even on platforms with limited processing power, proving a higher efficiency of the models without high deployment costs.

After careful analysis of the models' output, it was noticed that the incorrect prediction of the models occurs when actions present similarities between the two classes of violence and non-violence. For example, hugging can have many characteristics in common with violent actions when the video is analyzed frame by frame. To try to minimize this behavior, a third class was developed on the response of models with warning function or safety monitors in cases where the algorithm is not able to predict correctly. Through this system, in some cases prior to acts of violence, usually with actions and gestures of fast and expressive movements, the algorithm is used to generate a response in this attention zone. Thus, when there is a gradual evolution of violence in the behavior of those involved, the algorithm is able to call the attention of monitoring agents in order to prevent occurrences.

As future work, we can focus on a study of the impact of using mobile CNN architectures on more powerful embedded platforms. Moreover, we aim to improve the dataset, with videos of actions of violence and nonviolence, above all in public places such as malls, airports, subways, parks, and sports stadiums, in order to improve the applicability and accuracy of the models. Another possibility is the detection and location of the violent occurrences in the videos. This could be accomplished by using BoVW (Bag of Visual Words) or by more advanced segmentation techniques.

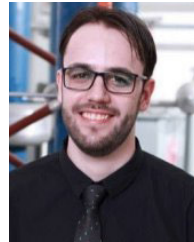
ACKNOWLEDGMENT

The authors would like to thank the Coordination for the Improvement of Higher Education Personnel (CAPES–Brazil), the PIBIC/FURB Program, and the University of Blumenau (FURB) for their support in this work. This work was supported by the National Funds through the Fundação para a Ciência e a Tecnologia, I.P. (Portuguese Foundation for Science and Technology) by the Project “VALORIZA—Research Centre for Endogenous Resource Valorization” under Grant UIDB/05064/2020 and Grant UIDB/04111/2020, and in part by the Instituto Lusófono de Investigação e Desenvolvimento (ILIND) under Project COFAC/ILIND/COPELABS/3/2020.

REFERENCES

- [1] P. Zhou, Q. Ding, H. Luo, and X. Hou, “Violence detection in surveillance video using low-level features,” *PLoS ONE*, vol. 13, no. 10, 2018, Art. no. 203668.
- [2] E. B. Nieves, O. D. Suarez, G. B. García, and R. Sukthakar, “Violence detection in video using computer vision techniques,” in *Computer Analysis of Images and Patterns*, vol. 6855. Berlin, Germany: Springer, 2011, pp. 332–339.
- [3] A. B. Mabrouk and E. Zagrouba, “Abnormal behavior recognition for intelligent video surveillance systems: A review,” *Expert Syst. Appl.*, vol. 91, pp. 480–491, Jan. 2018.
- [4] S. Vishwakarma and A. Agrawal, “A survey on activity recognition and behavior understanding in video surveillance,” *Vis. Comput.*, vol. 29, no. 10, pp. 983–1009, Oct. 2013.
- [5] O. P. Popoola and K. Wang, “Video-based abnormal human behavior recognition—A review,” *IEEE Trans. Syst., Man, Cybern. C, Appl. Rev.*, vol. 42, no. 6, pp. 865–878, Nov. 2012.
- [6] M. Rusci, D. Rossi, E. Farella, and L. Benini, “A sub-mW IoT-endnode for always-on visual monitoring and smart triggering,” *IEEE Internet Things J.*, vol. 4, no. 5, pp. 1284–1295, Oct. 2017.
- [7] S. F. Stefenon, C. Kasburg, R. Z. Freire, F. C. S. Ferreira, D. W. Bertol, and A. Nied, “Photovoltaic power forecasting using wavelet neuro-fuzzy for active solar trackers,” *J. Intell. Fuzzy Syst.*, vol. 40, no. 1, pp. 1083–1096, Jan. 2021.
- [8] G. Cerutti, R. Prasad, A. Brutti, and E. Farella, “Compact recurrent neural networks for acoustic event detection on low-energy low-complexity platforms,” *IEEE J. Sel. Topics Signal Process.*, vol. 14, no. 4, pp. 654–664, May 2020.
- [9] A. Cimatti and S. Tonetta, “Contracts-refinement proof system for component-based embedded systems,” *Sci. Comput. Program.*, vol. 97, pp. 333–348, Jan. 2015.
- [10] A. Cimatti, S. Mover, and S. Tonetta, “Quantifier-free encoding of invariants for hybrid systems,” *Formal Methods Syst. Des.*, vol. 45, no. 2, pp. 165–188, Oct. 2014.
- [11] F. Montagna, M. Buiatti, S. Benatti, D. Rossi, E. Farella, and L. Benini, “A machine learning approach for automated wide-range frequency tagging analysis in embedded neuromonitoring systems,” *Methods*, vol. 129, pp. 96–107, Oct. 2017.
- [12] S. F. Stefenon, C. Kasburg, A. Nied, A. C. R. Klaar, F. C. S. Ferreira, and N. W. Branco, “Hybrid deep learning for power generation forecasting in active solar trackers,” *IET Gener., Transmiss. Distrib.*, vol. 14, no. 23, pp. 5667–5674, Dec. 2020.
- [13] M. Rusci, D. Rossi, M. Lecca, M. Gottardi, E. Farella, and L. Benini, “An event-driven ultra-low-power smart visual sensor,” *IEEE Sensors J.*, vol. 16, no. 13, pp. 5344–5353, Jul. 2016.
- [14] N. F. S. Neto, S. F. Stefenon, L. H. Meyer, R. Bruns, A. Nied, L. O. Seman, G. V. Gonzalez, V. R. Q. Leithardt, and K.-C. Yow, “A study of multilayer perceptron networks applied to classification of ceramic insulators using ultrasound,” *Appl. Sci.*, vol. 11, no. 4, p. 1592, Feb. 2021.
- [15] S. F. Stefenon, M. P. Corso, A. Nied, F. L. Perez, K. Yow, G. V. Gonzalez, and V. R. Q. Leithardt, “Classification of insulators using neural network based on computer vision,” *IET Gener., Transmiss. Distrib.*, pp. 1–12, Dec. 2021.
- [16] D. Liu, H. Zhang, and P. Zhou, “Video-based facial expression recognition using graph convolutional networks,” in *Proc. 25th Int. Conf. Pattern Recognit. (ICPR)*, Milan, Italy, vol. 25, Jan. 2021, pp. 607–614.
- [17] N. Varshney and B. Bakariya, “Deep convolutional neural model for human activities recognition in a sequence of video by combining multiple CNN streams,” *Multimedia Tools Appl.*, pp. 1–13, Aug. 2021.
- [18] G. H. dos Santos, L. O. Seman, E. A. Bezerra, V. R. Q. Leithardt, A. S. Mendes, and S. F. Stefenon, “Static attitude determination using convolutional neural networks,” *Sensors*, vol. 21, no. 19, p. 6419, Sep. 2021.
- [19] M. P. Corso, F. L. Perez, S. F. Stefenon, K.-C. Yow, R. G. Ovejero, and V. R. Q. Leithardt, “Classification of contaminated insulators using k-nearest neighbors based on computer vision,” *Computers*, vol. 10, no. 9, p. 112, Sep. 2021.
- [20] J. Li, X. Jiang, T. Sun, and K. Xu, “Efficient violence detection using 3D convolutional neural networks,” in *Proc. 16th IEEE Int. Conf. Adv. Video Signal Based Surveill. (AVSS)*, Taipei, Taiwan, Sep. 2019, pp. 1–8.
- [21] Q. Dai, R.-W. Zhao, Z. Wu, X. Wang, Z. Gu, W. Wu, and Y.-G. Jiang, “Fudan-Huawei at MediaEval 2015: Detecting violent scenes and affective impact in movies with deep learning,” in *Proc. Multimedia Benchmark Workshop*, Wurzen, Germany, 2015, pp. 1–3.
- [22] S. F. Stefenon, R. Z. Freire, L. H. Meyer, M. P. Corso, A. Sartori, A. Nied, A. C. R. Klaar, and K.-C. Yow, “Fault detection in insulators based on ultrasonic signal processing using a hybrid deep learning technique,” *IET Sci., Meas. Technol.*, vol. 14, no. 10, pp. 953–961, Dec. 2020.
- [23] F. Fernandes, S. F. Stefenon, L. O. Seman, A. Nied, F. C. S. Ferreira, M. C. M. Subtil, A. C. R. Klaar, and V. R. Q. Leithardt, “Long short-term memory stacking model to predict the number of cases and deaths caused by COVID-19,” *J. Intell. Fuzzy Syst.*, pp. 1–14, Dec. 2021.
- [24] C. Ding, S. Fan, M. Zhu, W. Feng, and B. Jia, “Violence detection in video by using 3D convolutional neural networks,” in *Advances in Visual Computing*, vol. 8888. Cham, Germany: Springer, 2014, pp. 551–558.
- [25] P. Wang, P. Wang, and E. Fan, “Violence detection and face recognition based on deep learning,” *Pattern Recognit. Lett.*, vol. 142, pp. 20–24, Feb. 2021.
- [26] M. Sabokrou, M. Fayyaz, M. Fathy, Z. Moayed, and R. Klette, “Deep-anomaly: Fully convolutional neural network for fast anomaly detection in crowded scenes,” *Comput. Vis. Image Understand.*, vol. 172, pp. 88–97, Jul. 2018.
- [27] M. H. Gauswami and K. R. Trivedi, “Implementation of machine learning for gender detection using CNN on raspberry Pi platform,” in *Proc. 2nd Int. Conf. Inventive Syst. Control (ICISC)*, Coimbatore, India, vol. 2, Jan. 2018, pp. 608–613.
- [28] S. P. Suchitra and S. Tripathi, “Real-time emotion recognition from facial images using raspberry Pi II,” in *Proc. 3rd Int. Conf. Signal Process. Integr. Netw. (SPIN)*, Noida, India, vol. 3, Feb. 2016, pp. 666–670.
- [29] A. Medeiros, A. Sartori, S. F. Stefenon, L. H. Meyer, and A. Nied, “Comparison of artificial intelligence techniques to failure prediction in contaminated insulators based on leakage current,” *J. Intell. Fuzzy Syst.*, pp. 1–14, Dec. 2021.
- [30] V. Mazzia, A. Khaliq, F. Salvetti, and M. Chiaberge, “Real-time apple detection system using embedded systems with hardware accelerators: An edge AI application,” *IEEE Access*, vol. 8, pp. 9102–9114, 2020.
- [31] V. Leithardt, D. Santos, L. Silva, F. Viel, C. Zeferino, and J. Silva, “A solution for dynamic management of user profiles in IoT environments,” *IEEE Latin Amer. Trans.*, vol. 18, no. 7, pp. 1193–1199, Jul. 2020.
- [32] M. Dua, R. Singla, S. Raj, and A. Jangra, “Deep CNN models-based ensemble approach to driver drowsiness detection,” *Neural Comput. Appl.*, vol. 33, pp. 3155–3168, Apr. 2021.
- [33] T. Hassner, Y. Itcher, and O. Kliper-Gross, “Violent flows: Real-time detection of violent crowd behavior,” in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. Workshops*, Providence, RI, USA, vol. 1, Jun. 2012, pp. 1–6.
- [34] K. Soomro, A. Roshan Zamir, and M. Shah, “UCF101: A dataset of 101 human actions classes from videos in the wild,” 2012, *arXiv:1212.0402*.

- [35] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre, "HMDB: A large video database for human motion recognition," in *Proc. Int. Conf. Comput. Vis.*, Barcelona, Spain, Nov. 2011, pp. 2556–2563.
- [36] M. Monfort, C. Vondrick, A. Oliva, A. Andonian, B. Zhou, K. Ramakrishnan, S. A. Bargal, T. Yan, L. Brown, Q. Fan, and D. Gutfreund, "Moments in time dataset: One million videos for event understanding," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 2, pp. 502–508, Feb. 2020.
- [37] H.-J. Yoo, "Deep convolution neural networks in computer vision: A review," *IEIE Trans. Smart Process. Comput.*, vol. 4, no. 1, pp. 35–43, Feb. 2015.
- [38] P. R. R. De Souza, K. J. Matteussi, A. D. S. Veith, B. F. Zanchetta, V. R. Q. Leithardt, A. L. Murciego, E. P. De Freitas, J. C. S. D. Anjos, and C. F. R. Geyer, "Boosting big data streaming applications in clouds with BurstFlow," *IEEE Access*, vol. 8, pp. 219124–219136, 2020.
- [39] S. C. Wong, A. Gatt, V. Stamatescu, and M. D. McDonnell, "Understanding data augmentation for classification: When to warp?" in *Proc. Int. Conf. Digit. Image Comput., Techn. Appl. (DICTA)*, Gold Coast, QLD, Australia, vol. 1, Nov. 2016, pp. 1–6.
- [40] S. F. Stefenon, L. O. Seman, N. F. S. Neto, L. H. Meyer, A. Nied, and K.-C. Yow, "Echo state network applied for classification of medium voltage insulators," *Int. J. Electr. Power Energy Syst.*, vol. 134, Jan. 2022, Art. no. 107336.
- [41] S. F. Stefenon, M. H. D. M. Ribeiro, A. Nied, K.-C. Yow, V. C. Mariani, L. D. S. Coelho, and L. O. Seman, "Time series forecasting using ensemble learning methods for emergency prevention in hydroelectric power plants with dam," *Electr. Power Syst. Res.*, vol. 202, Jan. 2022, Art. no. 107584.
- [42] D. Rajkumar, "Image classification using network inception-architecture & applications," *Int. J. Innov. Res. Sci. Eng. Technol.*, vol. 10, no. 1, pp. 329–333, 2021.
- [43] C. Szegegy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Boston, MA, USA, Jun. 2015, pp. 1–9.
- [44] C. Lin, G. Zhao, Z. Yang, A. Yin, X. Wang, L. Guo, H. Chen, Z. Ma, L. Zhao, H. Luo, T. Wang, B. Ding, X. Pang, and Q. Chen, "CIR-Net: Automatic classification of human chromosome based on inception-ResNet architecture," *IEEE/ACM Trans. Comput. Biol. Bioinf.*, early access, Jun. 18, 2020, doi: [10.1109/TCBB.2020.3003445](https://doi.org/10.1109/TCBB.2020.3003445).
- [45] M. G. D. Dionson and E. J. P. Bibangco, "Inception-V3 architecture in dermatoglyphics-based temperament classification," *Philippine Social Sci. J.*, vol. 3, no. 2, pp. 173–174, Nov. 2020.
- [46] F. Ucar and D. Korkmaz, "COVIDiagnosis-Net: Deep Bayes-SqueezeNet based diagnosis of the coronavirus disease 2019 (COVID-19) from X-ray images," *Med. Hypotheses*, vol. 140, Jul. 2020, Art. no. 109761.
- [47] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems*, vol. 25. Red Hook, NY, USA: Curran Associates, 2012, pp. 1–9.
- [48] Y. Xu, G. Yang, J. Luo, and J. He, "An electronic component recognition algorithm based on deep learning with a faster SqueezeNet," *Math. Problems Eng.*, vol. 2020, pp. 1–11, Oct. 2020.
- [49] S. F. Stefenon, L. O. Seman, C. S. F. Neto, A. Nied, D. M. Seganfredo, F. G. D. Luz, P. H. Sabino, J. T. González, and V. R. Q. Leithardt, "Electric field evaluation using the finite element method and proxy models for the design of stator slots in a permanent magnet synchronous motor," *Electronics*, vol. 9, no. 11, p. 1975, 2020.
- [50] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, "SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <0.5 MB model size," 2016, *arXiv:1602.07360*.
- [51] P. N. Srinivasu, J. G. SivaSai, M. F. Ijaz, A. K. Bhoi, W. Kim, and J. J. Kang, "Classification of skin disease using deep learning neural networks with MobileNet V2 and LSTM," *Sensors*, vol. 21, no. 8, p. 2852, 2021.
- [52] S. Phiphatphaisit and O. Surinta, "Food image classification with improved MobileNet architecture and data augmentation," in *Proc. 3rd Int. Conf. Inf. Sci. Syst.*, Cambridge, U.K., vol. 3, Mar. 2020, pp. 51–56.
- [53] U. Kulkarni, M. S. Meena, S. V. Gurlahosur, and G. Bhogar, "Quantization friendly MobileNet (QF-MobileNet) architecture for vision based applications on embedded platforms," *Neural Netw.*, vol. 136, pp. 28–39, Apr. 2021.
- [54] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Las Vegas, NV, USA, Jun. 2016, pp. 770–778.
- [55] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "MobileNetV2: Inverted residuals and linear bottlenecks," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Salt Lake City, UT, USA, Jun. 2018, pp. 4510–4520.
- [56] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Salt Lake City, UT, USA, Jun. 2018, pp. 8697–8710.



JOELTON CEZAR VIEIRA received the Graduate degree in electrical engineering from the University of Blumenau (FURB), Brazil, in 2016, with undergraduate thesis in low-cost home automation implemented on an Arduino embedded platform, and the master's degree in electrical engineering with an area of expertise in power electronics from FURB, in 2019, under the guidance of Prof. Dra. Andreza Sartori. During graduation, he participated in an exchange program with the Aschaffenburg University of Applied Sciences, Germany, where he took classes related to industrial engineering and computing.

He has been a Software Developer at WEG Electrical Equipment S.A., since 2020, where he works in the research and development of intelligent systems for power transformers. His current research interests include action and signals recognition, integrated and embedded systems, automation, industry 4.0, and pattern recognition in industrial systems, mainly working with signals, images, and videos.



ANDREZA SARTORI received the Graduate degree in information systems from the University Center of Brusque, in 2007, the master's degree in neteconomy: technology and information and knowledge management from the University of Trento, Italy, in 2011, and the Ph.D. degree in informatics from the International Doctoral School on Information and Communication Technology, University of Trento, in 2015. During her Ph.D. period, she worked with the Semantics & Knowledge Innovation Laboratory (SKIL Lab), Telecom Italia's Research Center, associated with EIT ICT Labs. She did part of her doctorate at Bogaziçi University, Turkey. Her Ph.D. thesis discusses the use of computer vision and machine learning techniques to analyze emotions in abstract paintings. She took part in the master's at Åbo Akademi University, Finland, with the Erasmus Program, in 2009.

She has been a full-time Professor with the Regional University of Blumenau (FURB), since 2016, where she teaches for the undergraduate in computer science and graduate programs in electrical engineering. Her research interests include automatic emotion recognition, computer vision, machine learning, and pattern recognition of industrial systems.

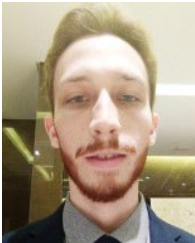


STÉFANO FRIZZO STEFENON received the B.E. and M.E. degrees in electrical engineering (power systems) from the Regional University of Blumenau, Brazil, in 2012 and 2015, respectively, and the Ph.D. degree in electrical engineering (artificial intelligence) from the State University of Santa Catarina, Brazil, in 2021. During his Ph.D. period, he developed a research project in the field of deep learning applied to computer vision with the Faculty of Engineering and Applied Science, University of Regina, Canada.

He is currently working at Fondazione Bruno Kessler, Italy, in co-foundation with the University of Udine in the field of computer science and artificial intelligence. His research interests include artificial intelligence for fault identification in electrical power systems, time series forecasting, deep learning, computer vision, and wavelet transform.



FÁBIO LUIS PEREZ received the B.S. degree in electrical engineering, the degree with a specialization in automation and control, the M.Sc. degree in production engineering, and the Ph.D. degree in electrical engineering from the Federal University of Santa Catarina, Brazil, in 1990, 1993, 2003, and 2015, respectively. He joined the Department of Telecommunication, Electrical and Mechanical Engineering, University of Blumenau, Santa Catarina, Brazil, in 1992. He is currently a Professor in electrical engineering and the Director of the Technological Area. His research interests include control systems, adaptive signal processing, image and speech processing, and artificial intelligence.



GABRIEL SCHNEIDER DE JESUS received the bachelor's degree in computer science from the University of Blumenau (FURB), Brazil. He is currently pursuing the master's degree with Lomonosov Moscow State University, Russia. During his graduation, he was a Mentor and a Researcher in computer science. He started his researcher career taking part on scientific research projects with the University of Blumenau, after receiving his Machine Learning Engineering Nanodegree Certificate from Udacity Inc., in 2018. His research interests include time series predictions applied to electrical engineering, machine learning applications to stock market forecasting, the usage of deep learning techniques to behavior detection, and violence recognition.



VALDERI REIS QUIETINHO LEITHARDT (Member, IEEE) received the Ph.D. degree in computer science from INF-UFRGS, Brazil, in 2015. He is currently an Adjunct Professor with the Polytechnic Institute of Portalegre and a Researcher integrated with the VALORIZA Research Group, School of Technology and Management (ESTG). He is also a Collaborating Researcher with the following research groups, such as COPELABS, Universidade Lusófona de Lisboa, Portugal, and the Expert Systems and Applications Laboratory, University of Salamanca, Spain. His research interests include distributed systems with a focus on data privacy, communication, and programming protocols, involving scenarios and applications for the Internet of Things, smart cities, big data, cloud computing, and blockchain.

• • •