

Article

Person Identification Using Temporal Analysis of Facial Blood Flow

Maria Raia ¹, Thomas Stogiannopoulos ¹ , Nikolaos Mitianoudis ^{1,*}  and Nikolaos V. Boulgouris ² 

¹ University Campus Kimmeria, Democritus University of Thrace, 67100 Xanthi, Greece; tstogian@ee.duth.gr (T.S.)

² College of Engineering, Design and Physical Sciences, Brunel University London, Uxbridge UB8 3PH, UK; nikolaos.boulgouris@brunel.ac.uk

* Correspondence: nmitiano@ee.duth.gr; Tel.: +30-2541079572

Abstract: Biometrics play an important role in modern access control and security systems. The need of novel biometrics to complement traditional biometrics has been at the forefront of research. The Facial Blood Flow (FBF) biometric trait, recently proposed by our team, is a spatio-temporal representation of facial blood flow, constructed using motion magnification from facial areas where skin is visible. Due to its design and construction, the FBF does not need information from the eyes, nose, or mouth, and, therefore, it yields a versatile biometric of great potential. In this work, we evaluate the effectiveness of novel temporal partitioning and Fast Fourier Transform-based features that capture the temporal evolution of facial blood flow. These new features, along with a “time-distributed” Convolutional Neural Network-based deep learning architecture, are experimentally shown to increase the performance of FBF-based person identification compared to our previous efforts. This study provides further evidence of FBF’s potential for use in biometric identification.

Keywords: biometrics; motion magnification; facial blood flow



Citation: Raia, M.; Stogiannopoulos, T.; Mitianoudis, N.; Boulgouris N.V. Person Identification Using Temporal Analysis of Facial Blood Flow. *Electronics* **2024**, *13*, 4499. <https://doi.org/10.3390/electronics13224499>

Academic Editors: José Carlos Bregieiro Ribeiro, Rolando Miragaia and Sandra V.B. Jardim

Received: 15 October 2024

Revised: 11 November 2024

Accepted: 13 November 2024

Published: 15 November 2024



Copyright: © 2024 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

In recent decades, the importance of biometric identification has been increasing steadily. The biometric characteristics of people and their use in recognition systems are now popular for the authentication of individuals [1]. Most mobile devices employ traditional biometric traits, including fingerprints and face, for the verification of the identity of their user. The face and the visual facial characteristics have always been used for identity verification with great success [2–5]. A constraint in common face recognition systems is that they may be sensitive to facial expression variations [4], especially when faces are captured in the wild. Fingerprints have also been widely used in person identification in current mobile phones and police forensics [6–8], but their operation relies on the proximity of the sensing probe to the acquired fingerprint.

The research community always seeks novel biometric technologies in order to complement traditional biometric systems and increase their accuracy [9]. Apart from traditional biometric traits, which have been thoroughly tested and deployed, new biometric modalities have also been devised and explored, such as palm [10], vein [11], ear [12], eye blinking [13], and EEG [14]. These biometrics are either deployed in stand-alone identification applications, or are used as part of multi-modal biometric systems [1]. In such multi-modal cases, the novel biometric traits are usually complementary to traditional biometrics, aiming to enhance end-to-end system performance.

Video amplification was first introduced by Wu et al. [15] and aimed at amplifying, and making visible to the human eye, small motions in videos captured using an ordinary camera. In [15,16], it was demonstrated that, by applying video amplification on facial image sequences, it is possible to visualize the flow of facial blood. In [17], a complex steerable pyramid decomposition was employed and motion amplification was applied

solely on the phase component of the decomposition. This enabled motion amplification to focus on the edge information that is prevalent in the phase component. Due to its ability to reveal imperceptible information, motion magnification has been applied to various fields, ranging from computer vision and biomedical imaging to civil and mechanical engineering [18–21].

The widespread use of masks in public places during the COVID-19 pandemic opened another possible application domain for motion magnification. Facial masks cover most facial areas that are normally used by efficient face recognition systems. In a situation where part of the face is covered, it would be most useful to devise a facial biometric trait that focuses on the less textured facial areas and does not rely on the visibility of conventional facial landmarks, such as eyes, nose, mouth, and eyebrows. In addition, any new facial biometric should be robust to facial expression variations. For the above reasons, we will not be comparing traditional face biometrics with the proposed facial blood flow biometric trait.

Facial blood flow has seen a lot of medical applications recently [22–24]. Inspired by the application of video magnification to the visualization of subtle facial motion, we embarked on exploring the use of facial blood flow patterns as a biometric trait. To that end, we sought evidence that suggests that facial blood flow can act as a biometric. In [25], Buddharaju et al. explored the possibility of performing face identification using thermal infrared cameras. In their study, they support, based on medical evidence and their conducted experiments, that the facial network of veins and arteries can have great variability among individuals, and thus it can act as a distinguishing trait. In addition, they showed that the extraction of such facial traits can have repeatability in individuals over multiple days and, as such, they can serve as a biometric.

In our previous work, we investigated the use of facial blood flow (FBF) as a potential biometric trait. The method we developed uses commercial RGB cameras (24 bit color cameras with HD resolution (1920×1080 pixels) and 25–30 fps) and performs image processing and motion magnification to detect and visualize facial blood flow patterns that can be used for biometric identification. In [26], we proposed a baseline method that used facial blood flow (FBF) for person identification. In that method, FBF was extracted from a person's face using a commercial RGB camera. We used motion amplification [15] to enhance and reveal the actual blood flow in common RGB video streams. Our method in [26] was a contactless method that did not utilize any traditional facial features, i.e., eyes, nose, mouth, eyebrows. Instead, it extracted small facial areas that are not commonly obstructed by facial hair, and it used the motion-amplified video to extract spatio-temporal blood flow information. That approach was shown to work well as a distinctive biometric [26]. Unlike the work in [25], our method does not need a high-grade infrared camera to capture the vein–artery structure of the face. Instead, our approach uses a low-cost commercial RGB camera to capture three facial areas and use image processing to construct the FBF. In [27], we improved our baseline method by taking into account the temporal evolution of FBF within a period, which was not explicitly considered in [26]. The importance of temporal evolution was stressed before in the literature [28,29]. Furthermore, we adopted a deep Convolutional Neural Network (CNN) architecture, which improved the accuracy of the baseline system.

In this paper, we propose an improved new framework, based on several additional features that exploit the temporal evolution of FBF, i.e., temporal partitioning and FFT-based features. Further, we examine a number of different deep learning architectures for the classification task. Our ablation study over different features and network architectures yields that the final proposed system achieves efficient person identification. The proposed system features a “time-distributed” CNN architecture to capture the temporal evolution of FBF.

2. A Facial Blood Flow Biometric System

This section describes the individual subsystems of the proposed facial blood flow biometric system. The proposed system is depicted in Figure 1. In the following subsections, each individual module of the system is described in detail.

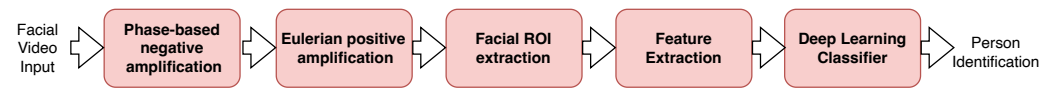


Figure 1. The proposed person identification system based on facial blood flow (FBF).

2.1. Video Capture

The proposed system uses an RGB camera in order to capture the facial blood flow (FBF) of a person. FBF is a periodical phenomenon with the period equal to the subject's heart rate. The common human pulse rate at rest varies between 60 to 120 bpm (beats per minute), which translates to frequencies between 1 Hz and 2 Hz. Thus, a commercial RGB camera with 30 fps (i.e., 30 Hz) video capture will be sufficient to capture several frames from each cycle of facial blood flow. Although video resolution is less important, modern commercial cameras offer HD quality video capture, which is most suitable for our target application. More details regarding the data collection process are provided in Section 3.

2.2. Facial Blood Flow Motion Amplification

The next step in the system is to perform motion amplification in the input video. Motion amplification is a remarkable algorithm that amplifies subtle movements in an RGB video. The algorithm detects and amplifies motion, and adds amplified versions of the motion to the original video in order to enhance it visually [15]. Motion amplification is the means by which facial blood flow (FBF) is calculated. Therefore, motion amplification is extremely important in our system. In other biometrics (conventional face recognition, gait recognition, or iris recognition), due to the fact that their respective features are directly visible, motion amplification is not required.

In [15], Wu et al. proposed an Eulerian motion amplification approach, where each image frame from the input video is decomposed into a multi-scale Laplacian pyramid and motion amplification is applied to each scale before reconstruction. In order to amplify movements that exist around object boundaries, i.e., image edges, Wadhwa et al. [17] proposed a phase-based motion amplification scheme, where each frame from the input video is decomposed into a complex steerable pyramid, where the amplification is applied to the phase component of the decomposition only. This approach exploits the well-known fact that edge information in images is best represented by the phase information in the Fourier domain [30]. This approach achieves superior motion amplification around object edges and boundaries, as opposed to the Eulerian motion amplification, which is more efficient in enhancing motion that exists within the texture of an object.

The motion amplification algorithm will not only amplify subtle motion taking place within the object, but will also amplify the subject's global motion. Hence, if the subject is moving as a whole, its global movement will also be amplified, producing output that is not useful for our application. Therefore, an important requirement for the efficient function of motion amplification is that the captured subject must be as still as possible, exhibiting minimal global movement. For that reason, in our application we capture human subjects who have been advised to stay as motionless as they can. In practice, however, it is unavoidable that there will be some slight global movement during capturing. Due to that, if motion amplification was used directly on the original video, the output would not be particularly informative regarding subtle facial movements. To avoid such problems, we apply pre-processing in order to filter out unavoidable small global motions before motion amplification is applied.

As a first processing step, we use the phase-based motion amplification [17] with a negative amplification factor α , in the range of $[-1, 0)$, in order to reverse any unwanted

global motion. The actual value of α was determined experimentally. As explained earlier, the phase-based motion amplification is more sensitive to edge movements that are due to unwanted small global movements, which are observed when the subject is moving as a whole. Therefore, applying negative amplification using the phase-based motion amplification algorithm will stabilize the subject within the video as much as possible without affecting the subtle movements in the facial texture area, which is where the facial blood flow can be observed.

Subsequently, we use the Eulerian video magnification method by Wu et al. [15] to amplify the facial blood flow. The Eulerian video magnification method tends to magnify the image value oscillations, thus it has proven [15] to be ideal for amplifying blood circulation. In contrast to [15], where a new (motion-enhanced) video sequence was created, i.e., the amplified motion is added to the original video, in order to create a highlighted video, our focus is only on the extraction of the amplified facial blood flow. Consequently, we extract the amplified motion, without adding it to the original video, as in [15]. In order to keep computational complexity low, we apply motion amplification to a grayscale version of the input video. The result of this procedure is the calculated facial blood flow of a subject, represented in the form of an image sequence $Q(x, y, t)$, where x, y denote the image coordinates and t the frame index. Here, we should stress that without motion magnification, the signal $Q(x, y, t)$ is zero, therefore there are no extracted features to model. Finally, in our approach we only use grayscale data. The reason is that different color planes would produce different apparent motions, which would lead to inaccurate facial blood flow calculation. Therefore, using aggregate information in the form of grayscale pixel intensities is a safer option.

2.3. Extraction of Facial Regions of Interest

To verify the validity of FBF as a biometric, we rely on facial areas that do not contain facial landmarks, which are used in conventional face biometry. Traditional facial recognition biometric approaches use the entire face, including the eyes, nose, and mouth, the shape of which are specific to each individual and, therefore, have significant discriminatory capacity. Instead, in our approach, we consider the use of facial areas that do not include discriminatory landmarks, but are certain to involve substantial blood flow that can be revealed. With that in mind, we focused on three selected areas of interest: (a) the lower forehead area, (b) the skin area below the left eye, and (c) the skin area below the right eye. These areas feature mainly facial skin without landmarks, and are ideal for FBF acquisition. This selection of facial areas is supported by the findings in the work of Buddharaju et al. [25], which is an exploration of the facial areas containing veins and arteries that can be used for person identification. As it can be seen, the forehead and the two left and right cheek areas contain multiple veins and arteries, a fact that makes them particularly suitable for the observation and representation of facial blood flow. This finding is also supported in the work by von Arx et al. [31], where a complete atlas of the facial blood vessels is presented. Another advantage of our chosen facial areas is that they are always visible or can easily become visible and easy to capture. The special situation where only the forehead area is visible is examined as a separate experiment. Such a situation arises when faces are partly covered by a mask.

To detect these regions, we use the Zhu–Ramanan face detector [32] to isolate the facial area. Then, we fit an Active Appearance Model (AAM) [33] with 68 facial features using the implementation of Bulat and Tzimiropoulos [34]. Using the facial features, inferred by the AAM, we select an area of 71×201 pixels above the eyes and eyebrows to extract the forehead area. The left and right facial areas are determined by positioning two 71×101 rectangles below the eyes.

The exact positions of the left and right facial areas are calculated by exploiting the positions of key points on the eyes and eyebrows. In Figure 2d, we depict the eyebrows' highest and rightmost points using blue and green points, respectively. Using this information we can infer the point where the forehead commences above the eyebrows. The x -coordinate of

the eyebrows' highest point and the y -coordinate of the eyebrows' rightmost point are used to form the starting point of the forehead, which is shown as a yellow cross in Figure 2d. The width and height of the forehead rectangle are fixed, thus extending the rectangle over both eyebrows. In a similar fashion, we exploit the leftmost eyebrows' point as a reference point (green point in Figure 2e), and proceed down by 120 pixels so as to detect the bottom left corner of the rectangle for the left facial area. Equally, using the rightmost eyebrow point, we can set a rectangle containing the right facial area (see Figure 2f). Figure 2 shows an example of the rectangle positioning procedure. The sizes of the rectangles were set based on experimentation. An alternative approach would be to deploy image registration algorithms to fully automate rectangle positioning. Nonetheless, in order to keep the setup as simple as possible, and not introduce distortions due to image transformation, we kept these sizes constant and asked the participants to sit at fixed distance from the camera. This ensured that minimal registration problems would appear. The AAM for the facial structure can provide sufficient information to tackle the face registration problems that would arise if the algorithm is applied in less controlled environments [35]. Snapshots from the three extracted videos from the three areas are depicted in Figure 3.

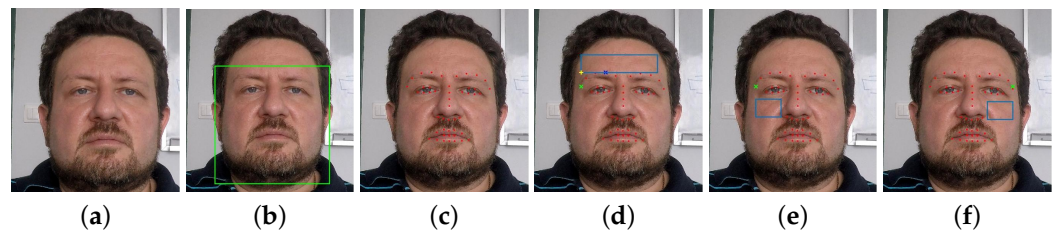


Figure 2. (a) Original face image. (b) Face detection using [32]. (c) Active Appearance Model (AAM) fit using [34]. The three control points (two from the AAM and another inferred from the other two) are highlighted. (d) Detection of the forehead region using two control points. (e,f) Detection of the left and right facial regions using the left and right control points, respectively. (The subject in this figure has agreed to have his image included in the paper for demonstration purposes.)



(a) Forehead



(b) Left area



(c) Right area

Figure 3. Snapshots from the three extracted areas for the subject in Figure 2. It is clear that these areas do not contain any traditional facial biometric traits.

2.4. Temporal and Spectral Feature Extraction

The next step is to extract efficient feature structures that will enable robust person identification. In order to achieve faster and more practical identification, we limit the duration of the clips that will be presented to the deep learning classifier to $T = 1$ s, so as to capture at least one cycle of blood circulation. Nonetheless, as it is not always easy to discern the start, end, or other stages within the facial blood cycle, synchronization

issues may occur. This highlights the importance of extracting features that are robust to phase variations within the observed periodic phenomenon. In [26,27], we proposed two baseline methodologies, which we extend here and show that they can lead to increased performance.

2.4.1. Temporal Features

In [26], we proposed to average all frames in each video clip, creating the average frame that was used for training and recognition. Although that approach provided robustness to phase variations, it ignored the temporal evolution of the FBF within the clip. Thus, in this paper, we propose a temporal partitioning, i.e., divide each clip into five sub-clips of equal duration T_s , and calculate the average FBF for each sub-clip separately. Assuming that we have a video sequence $\mathcal{Q}(x, y, t)$, where x, y are the spatial co-ordinates and t represents the temporal evolution, the proposed features $F_{temp}(x, y, k)$ can be described as follows (see Figure 4):

$$F_{temp}(x, y, k) = \frac{1}{T_s} \sum_{t=(k-1)T_s}^{kT_s} \mathcal{Q}(x, y, t), \quad 1 \leq k \leq 5 \quad (1)$$

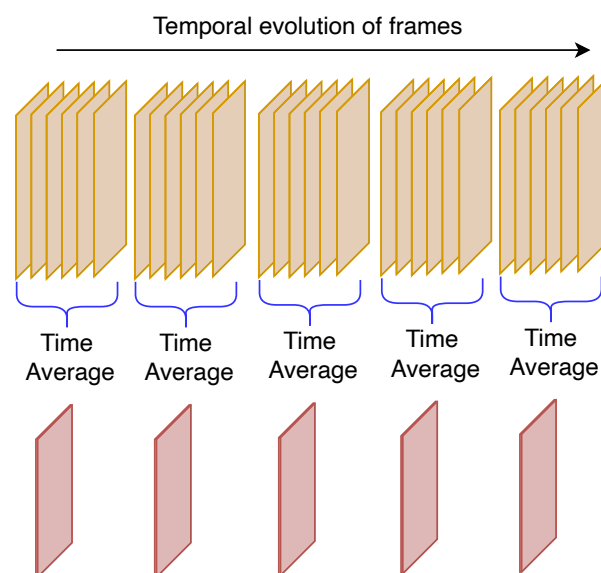


Figure 4. Division of a video clip into five sub-clips. Temporal averages are calculated for each sub-clip.

The average image and the temporal averages for the FBF biometric, extracted from the forehead region of two subjects, are depicted in Figure 5 and Figure 6 respectively. These templates clearly show the different FBF patterns in different individuals and, therefore, they demonstrate the discriminatory capacity of the FBF biometric.

2.4.2. Transform-Domain Phase-Invariant Features

In [27], we proposed to use the 1D-Discrete-Cosine Transform (DCT) along the temporal axis of the sequence $\mathcal{Q}(x, y, t)$, i.e.,

$$F_{DCT}(x, y, k) = \sum_{t=0}^T \mathcal{Q}(x, y, t) \cos \left[\frac{\pi k}{T} \left(t + \frac{1}{2} \right) \right], \quad 0 \leq k \leq T \quad (2)$$

The above DCT features provided a viable solution, since the DCT is real-valued and approximately phase-invariant.

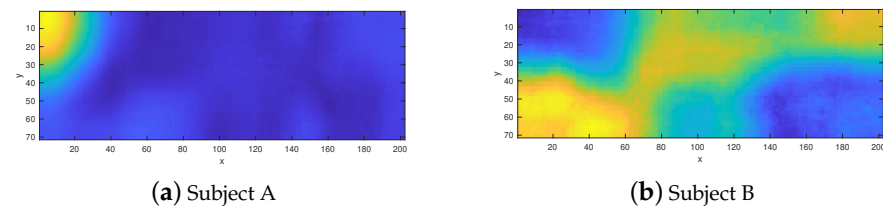


Figure 5. Temporal features for the FBF biometric extracted from the forehead region. The averaged image template [26] is shown for for subject A and B. Lighter colors represent greater values while darker colors represent smaller values (best seen in color).

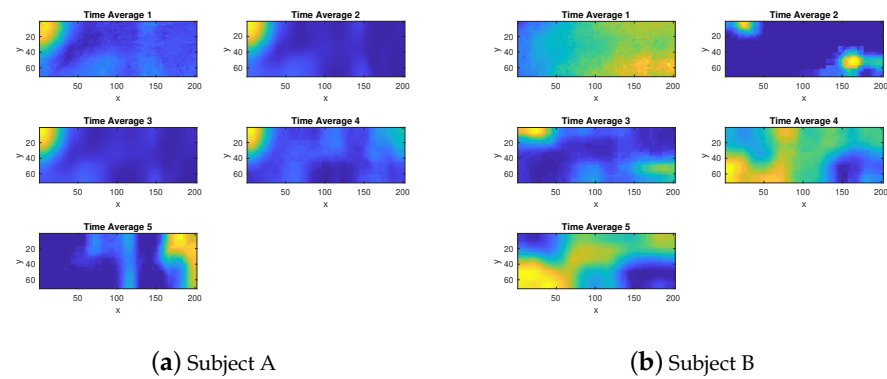


Figure 6. Temporal features for the FBF biometric extracted from the forehead region. The proposed temporal averages of (1) is shown for subject A and B. Lighter colors represent greater values while darker colors represent smaller values (best seen in color).

In the present work, however, we propose to use the abs-magnitude of the Fast-Fourier Transform (FFT), since the FFT is always phase-invariant. In addition, because of the magnitude FFT's even symmetry, the negative frequency components of the transformation can be dropped, yielding smaller-sized features. The proposed feature is given by:

$$F_{FFT}(x, y, k) = \left| \sum_{t=0}^T Q(x, y, t) e^{-j \frac{2\pi t k}{T}} \right|, \quad 0 \leq k \leq T \quad (3)$$

As in the case of DCT features, the proposed FFT features are calculated by applying the FFT transform along the temporal dimension of the input sequence. The two transform-domain features are depicted in Figure 7 for the forehead region. In either case, the extracted features are 3D matrices with the first two dimensions representing the image spatial coordinates (pixel location), and the third dimension representing the transform-domain content of each pixel (pixel value).

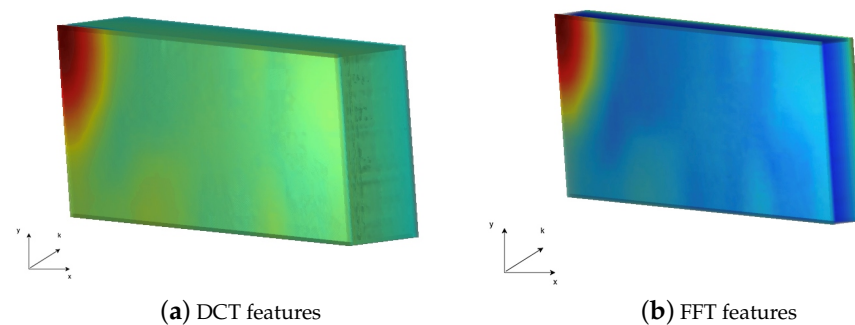


Figure 7. Frequency-domain features for the FBF biometric extracted from the forehead region of subject A. (a) DCT features, (b) FFT features calculated using (3).

2.5. Deep Learning Architectures for Classification

The next step is to use a deep learning architecture in order to perform person identification. As detailed in the previous subsection, the extracted feature data are in the form of 3D cubes, which represent spatio-temporal features (from temporal partitioning) or spatio-frequency features (DCT or FFT features for each pixel location). In either case, the input is a 3D tensor of size $M \times N \times K$, where $M \times N$ is the image resolution and K is the number of possible values for each feature.

Many deep learning architectures were investigated to model the evolution in time or frequency of the input features, including 3D CNN [36] and Long Short-Term Memory (LSTM) modules [37]. Nonetheless, in our experiments the most successful was an architecture that included “time-distributed” 2D CNN modules [38]. In other words, it is a CNN structure, consisting of 2D filters that are applied on the spatial co-ordinates x, y and remain unchanged over the z -axis, which describes time or frequency, depending on whether the spatio-temporal or spatio-frequency 3D feature is used. The proposed architecture is depicted in Figure 8. As seen, it consists of a smaller VGG16-like [39] structure of four “time-distributed” 2D convolutional layers, followed by two fully connected layers for classification. More specifically, the first “time-distributed” 2D CNN contains $16 \times 2 \times 2$ filters, and the second $32 \times 2 \times 2$ filters, which are then followed by a 2×2 max pooling layer. Consequently, there are two more CNN layers of $32 \times 2 \times 2$ filters each and a 2×2 max pooling layer. ReLUs are used after each CNN layer and the Dropout regularization method [40] using a parameter value of 0.2 is used after each max pooling layer. The features are flattened and presented to two fully connected layers (FCN or Dense), the first with 64 nodes and the second, being the output layer, with 13 nodes, equal to the number of people in the dataset. After the first FCN, Batch Normalization (BN) [41] is used along with ReLU and the Dropout regularization method with a parameter value of 0.3. The output of the last layer is presented to the Softmax activation function. The size of the input is not fixed in order to accommodate the features from the three areas of interest, which are of different size. The parameter values for Dropout, and the number and size of filters at each layer, were determined based on experimentation.

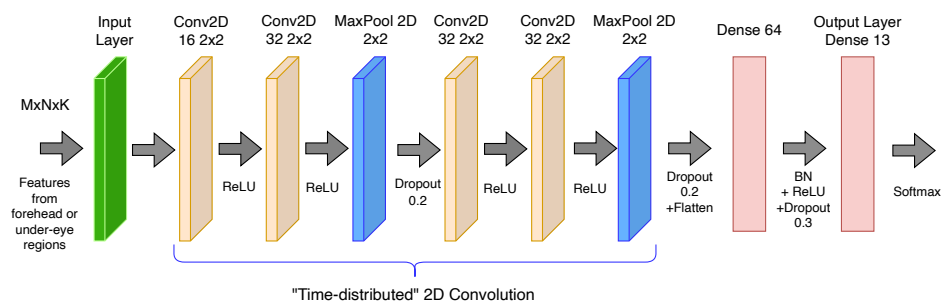


Figure 8. The proposed CNN structure with ‘time’-distributed 2D convolutions used in the convolutional layers of the network. Conv2D refers to a “time-distributed” 2D convolutional layer, ReLU and Softmax refer to the corresponding activation functions, Dropout refers to Dropout regularization [40], BN refers to Batch Normalization. The numbers of filters and the sizes of the filters are indicated at the top of the respective level.

The system should also combine features from three different facial areas in order to perform person identification. In [27], we determined that the features should be processed separately by the convolutive architecture, proposed in Figure 8, without the Dense layers (pipeline) and the final features should be concatenated and be presented to a single fully connected (Dense) stage for classification (Ensemble 1) (see Figure 9a). The Dense Layer consists of 64 weights. Figure 9b depicts a second method (Ensemble 2), where the features from each facial area are independently processed by the proposed convolutive architecture and a separate Dense stage, before being concatenated to the final FCN stage. Each of these

Dense stages has 64 nodes. The output dense layer has 13 nodes, equal to the number of subjects in the dataset.

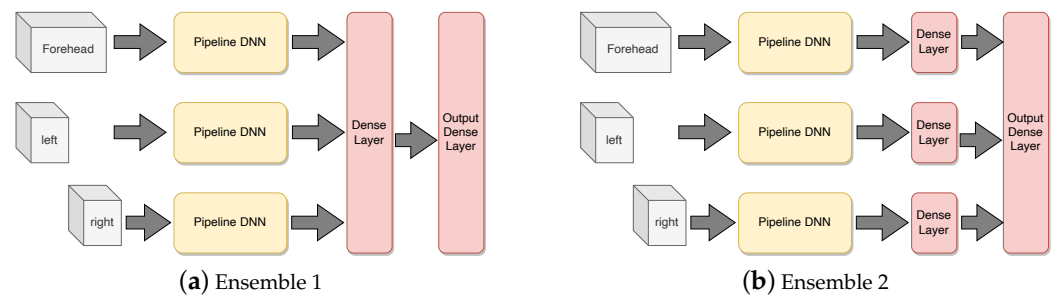


Figure 9. Two ensemble methods to combine the features from the three facial regions of interest. Pipeline refers to the architecture of Figure 8 without the fully connected (Dense) layers.

3. Experiments

3.1. Dataset Implementation

To test the efficiency of the proposed facial blood flow (FBF) biometric, we created a new dataset. A GoPro Hero 4 Black camera was used for video recording at 1920×1080 resolution and 30 frames per second. Facial image sequences from a total of 13 subjects were recorded in a room with natural light. This was needed because the motion magnification algorithm tends to magnify subtle, nearly invisible, movements and oscillations produced by the flickering of artificial light, which would result in contaminated output. Subjects were asked to sit on a chair at a fixed distance (~ 1 m) from the camera and to remain as motionless as possible during the recording. This guidance was meant to prevent the occurrence of registration problems due to incidental body movement. The camera was positioned at the same height for all subjects, almost at eyes level, in order to minimize registration issues. Eighteen (18) recordings were captured from each subject, each lasting 20 s. These recordings were captured for each subject over two different days. This produced a total of 360 s of recording per subject, which was used for training. A week later, three additional recordings per subject were captured on a single day. The additional recordings amounted to a total of 60 s per subject, and were used for system testing. Conducting the recordings over three different days strengthens our confidence in the temporal robustness of the proposed biometric trait and in its appropriateness for practical use. All recordings were segmented into 1-second clips, yielding 360 training clips and 60 validation clips for each subject. The compiled data can be made available upon request. An example of the captured video frames can be seen in Figure 2. Video capturing has been performed in a controlled environment, ensuring that no artificial light was used and that the subjects were instructed to remain as still as possible. Although this may seem as a limitation of the method, it is a realistic constraint to have a controlled environment for identity verification, e.g., in a police station or at border control.

The Movement Motion Attenuation and the Facial Blood Flow Amplification algorithms were implemented using MATLAB code, provided by [17] and [15], respectively. For the motion attenuation module, the value of $\alpha = -0.75$ was used, as it was experimentally seen to yield satisfactory results. The amplification factor was set to $\alpha = 120$, while the frequency range of amplification was set to 0.83–1 Hz, as suggested in [15]. The face isolation and AAM fitting stages were based on the method in [34]. The proposed deep architectures were developed in Python using TensorFlow on a PC with an NVidia RTX A6000 GPU, running Ubuntu Linux 20.04 (the code is available at https://github.com/mitia98/FBF_person_id, accessed on 12 November 2024). Stochastic Gradient Descent (SGD) was used to train the deep network, using categorical cross-entropy as a cost function, whereas the learning rate was set to 0.01 and momentum was set to 0.9. The network was trained for 30 epochs using a batch size of 64 samples.

3.2. Results of Ablation Study

As it was shown in [26], using only forehead features for person identification yields inferior performance in comparison to using the combination of all three proposed facial areas. In the present work, we re-evaluate the use of the forehead area in conjunction with our new FBF representations and the architecture shown in Figure 8. Results are shown in Table 1. As seen, similar to [26], the forehead area alone still exhibits moderate biometric performance.

Table 1. Identification accuracy for the validation dataset for the three different features and various scenarios. Salt-n-pepper augmentation on FFT features and the combination of the three areas with Ensemble 2 yields the best performance. Values in bold denote the best performance.

Features	Forehead Only	All Areas		
		Ensemble 1	Ensemble 2	
			No Augment.	With Augment.
Average Image [26]	68.76%	77.53%	79.72%	79.79%
DCT Features [27]	73.19%	82.61%	84.66%	84.78%
Temporal Partitioning	72.54%	84.51%	85.56%	85.6%
FFT Features	74.03%	85.02%	87.20%	89.04%

Subsequently, we evaluated the combination of the three facial areas, using the two combination approaches, shown in Figure 9. In Table 1, we present the identification accuracy achieved by the two combination methods. Clearly, the Ensemble 2 architecture yields the best results regardless of the feature used. Therefore, it seems that the most efficient strategy is to process each feature through individual FCNs before combining features in the output (final) layer. In addition, it is clear from Table 1 that the inclusion of the temporal element surely improves performance. All architectures that incorporate the temporal element (DCT features, temporal partitioning and FFT features) perform better than using the average image.

Finally, we experimented with data augmentation. We attempted three types of augmentation that seemed suitable for our problem: Gaussian noise, salt-n-pepper noise, and vertical flipping. Noise addition was considered because it was shown to improve network generalization and reduce overfitting [42]. Out of the three options, only the addition of salt-n-pepper noise to the final extracted features appeared to improve performance. Specifically, the saturation of 10% of each image pixel to equal percentage of black and white pixels seemed to yield performance improvement. This augmentation doubled the available training data. Results are shown in Table 1 in terms of classification accuracy.

As seen in Table 1, augmentation with salt-n-pepper noise improved the performance of all feature sets. Further, the proposed FFT features are consistently superior in all scenarios, reaching a recognition rate equal to 89.04%. The temporal partitioning features rank second in performance, followed by the DCT features [27] and the average image feature of [26]. The performance achieved by our new features and architecture is an improvement over the system in [27] and shows the potential of the FBF biometric trait for effective person identification. The evolution of the loss function and accuracy over 30 epochs for the proposed “time-distributed” VGG network with the FFT features and the Ensemble 2 strategy is depicted in Figure 10. The difference between the final training and validation losses is a sign of overfitting, which is mainly due to the small size of the training dataset. The confusion matrix for the proposed architecture and features is shown in Figure 11. It is clear that the system does not overfit and in addition the classification performance is robust since only minor mis-classifications are observed.

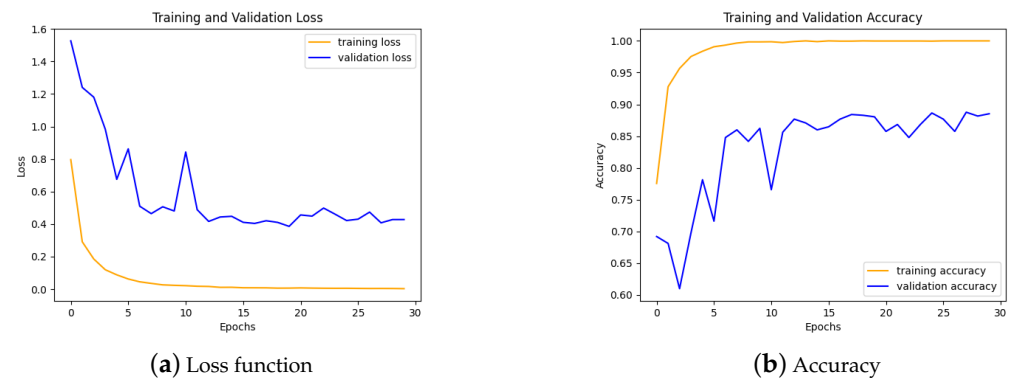


Figure 10. The evolution of the loss function and accuracy over 30 epochs for the proposed “time-distributed” VGG network with the FFT features and the Ensemble 2 strategy.

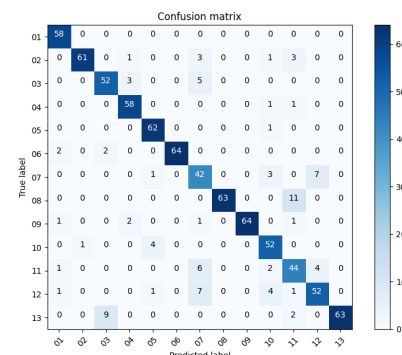


Figure 11. The confusion matrix for the proposed “time-distributed” VGG network with the FFT features and the Ensemble 2 strategy.

In Table 2, we present an ablation study with other deep architectures that were examined in our investigation. The study was performed using the FFT features from all regions of interest and the Ensemble 1 architecture in terms of accuracy and false positive rate. The study included 2D-CNN VGG, 3D-CNN VGG, the “time-distributed” VGG of Figure 8, stacked LSTM, and the “time-distributed” VGG with LSTM layers. The exact details of these architectures are not presented here, due to limited interest, but they are of similar size to the one presented in Figure 8, and therefore are a good basis for direct comparisons. From the results presented in Table 2, it is clear that the inclusion of LSTM modules did not improve performance. Instead, a more shallow “time-distributed” VGG is more efficient.

Table 2. Comparison between different Ensemble 1 architectures, in conjunction with FFT features extracted from all facial regions of interest. No data augmentation was applied. Values in bold denote the best performance. The proposed “time-distributed” VGG architecture yields the best performance.

Model	Accuracy	False Positive Rate
2D-CNN VGG version	71.01%	2.6%
3D-CNN VGG version	81.26%	1.64%
Time Distr. CNN VGG version	85.02%	1.27%
Stacked LSTM	78.05%	1.62%
Time Distr. CNN VGG version + LSTM	79.59%	1.51%

Table 3 presents comparisons in terms of computational complexity. Complexity is measured in terms of the required number of network parameters, the network size (in MBs), and the average inference time per batch. The compared systems are the aforementioned reference architectures, which use FFT features from all regions of interest,

and the Ensemble 1 architecture. In terms of parameters and size, the “time-distributed” CNN (VGG version), which achieved the highest identification performance, is the largest network, with 55.7 million parameters requiring 435.6 MB of storage. In terms of inference time, that architecture ranks in the middle. Although inference times were measured in a high-specification machine, they are very small and can comfortably support a practical real-time identification scenario. In particular, the per-sample inference time achieved by the best-performing “time-distributed” CNN network is equal to $0.742/64 = 11.6$ ms, where the total inference time was divided by the number of samples (64) in the batch. Such low inference times make our proposed application viable for use even on everyday computational platforms (e.g., common PCs).

Table 3. Comparison in terms of computational complexity, presenting the number of parameters, model size, and inference time (per batch) for each of the examined architectures. The proposed “time-distributed” VGG architecture has the greatest number of parameters and model size, nonetheless it does not require the greatest inference time per batch.

Model	No. of Parameters	Model Size (MB)	Inference Time per Batch (s)
2D-CNN VGG version	7M	54.7	0.733
3D-CNN VGG version	13.9M	109.3	0.736
Time Distr. CNN VGG version	55.7M	435.6	0.742
Stacked LSTM	14.89M	116.4	0.771
Time Distr. CNN VGG version + LSTM	14.03M	109.9	0.786

4. Conclusions

In this paper, we presented evidence of the validity of facial blood flow (FBF) as a biometric trait. We examined two new features: the FFT features and the temporal partitioning features, which yield improved performance over the performance achieved using the DCT features we proposed in our past work. In the present paper, we also proposed a new “time-distributed” CNN architecture that consolidates the performance gains using the new features. The presented algorithms and experiments demonstrate the effectiveness of FBF as a stand-alone or a complementary biometric trait and serve as a well-founded proof-of-concept. The main limitation of the proposed biometric is that it requires special capturing conditions for all subjects, i.e., no artificial light and no movement from the subject. In addition, in this study we have not verified the efficiency of FBF on a more ethnically diverse database.

In our future work, we will expand our experiments using larger datasets and multi-modal features. More specifically, we are planning to study the performance of the FBF biometric on a more diverse database, and also exploring the use of multi-spectral cameras. In future research, 3D modeling can be used for very accurate compensation of incidental head movements before blood flow calculation is performed.

Author Contributions: Conceptualization, N.V.B.; Methodology, N.M. and N.V.B.; Software, M.R. and T.S.; Validation, M.R. and T.S.; Formal analysis, N.M. and N.V.B.; Investigation, M.R.; Writing—original draft, M.R. and T.S.; Writing—review & editing, N.M. and N.V.B.; Visualization, M.R. and T.S.; Supervision, N.M. and N.V.B.; Project administration, N.M. All authors have read and agreed to the published version of the manuscript.

Funding: This research received no external funding.

Institutional Review Board Statement: Ethical review and approval were waived for this study due to the anonymization of the final data used by the proposed algorithm.

Informed Consent Statement: Written informed consent has been obtained from all participants in the study, including the patient(s) whose information is included in this paper. The consent form includes a full explanation of the nature and purpose of the study, as well as any potential risks and benefits of participation. The participants were informed that their participation is voluntary, and

that they have the right to withdraw from the study at any time without any negative consequences. The participants were also informed that their information will be kept confidential and anonymous, and that the data collected will only be used for the purposes of the research study.

Data Availability Statement: The raw data supporting the conclusions of this article will be made available by the authors on request.

Acknowledgments: We gratefully acknowledge the support of NVIDIA Corporation with the donation of the RTX A6000 GPU used for this research. We would also like to thank the 13 students at the Electrical and Computer Engineering Department at Democritus University of Thrace, who participated in our experiment.

Conflicts of Interest: The authors declare no conflicts of interest.

References

- Ross, A.A.; Nandakumar, K.; Jain, A.K. *Handbook of Multibiometrics*; Springer: New York, NY, USA, 2006; p. 220.
- Adjabi, I.; Ouahabi, A.; Benzaoui, A.; Taleb-Ahmed, A. Past, Present, and Future of Face Recognition: A Review. *Electronics* **2020**, *9*, 1188. [[CrossRef](#)]
- Taskiran, M.; Kahraman, N.; Erdem, C.E. Face recognition: Past, present and future (a review). *Digit. Signal Process.* **2020**, *106*, 102809. [[CrossRef](#)]
- Oloyede, M.O.; Hancke, G.P.; Myburgh, H.C. A review on face recognition systems: Recent approaches and challenges. *Multimed. Tools Appl.* **2020**, *79*, 27891–27922. [[CrossRef](#)]
- Li, L.; Mu, X.; Li, S.; Peng, H. A Review of Face Recognition Technology. *IEEE Access* **2020**, *8*, 139110–139120. [[CrossRef](#)]
- Yang, W.; Wang, S.; Hu, J.; Zheng, G.; Valli, C. Security and Accuracy of Fingerprint-Based Biometrics: A Review. *Symmetry* **2019**, *11*, 141. [[CrossRef](#)]
- Yao, Z.; Bars, J.M.L.; Charrier, C.; Rosenberger, C. Literature review of fingerprint quality assessment and its evaluation. *IET Biom.* **2016**, *5*, 243–251. [[CrossRef](#)]
- Kaushal, N.; Kaushal, P. Human Identification and Fingerprints: A Review. *J. Biom. Biostat.* **2011**, *2*, 123. [[CrossRef](#)]
- Ghayoumi, M. A review of multimodal biometric systems: Fusion methods and their applications. In Proceedings of the 2015 IEEE/ACIS 14th International Conference on Computer and Information Science (ICIS), Las Vegas, NV, USA, 28 June–1 July 2015; pp. 131–136.
- Liang, X.; Yang, J.; Lu, G.; Zhang, D. CompNet: Competitive Neural Network for Palmprint Recognition Using Learnable Gabor Kernels. *IEEE Signal Process. Lett.* **2021**, *28*, 1739–1743. [[CrossRef](#)]
- Kuzu, R.S.; Maiorana, E.; Campisi, P. Vein-Based Biometric Verification Using Densely-Connected Convolutional Autoencoder. *IEEE Signal Process. Lett.* **2020**, *27*, 1869–1873. [[CrossRef](#)]
- Nikose, S.; Meena, H.K. Ear-biometrics for human identification. In Proceedings of the 2020 Advanced Computing and Communication Technologies for High Performance Applications (ACCTHPA), Cochin, India, 2–4 July 2020; pp. 8–13.
- Abo-Zahhad, M.; Ahmed, S.M.; Abbas, S.N. A Novel Biometric Approach for Human Identification and Verification Using Eye Blinking Signal. *IEEE Signal Process. Lett.* **2015**, *22*, 876–880. [[CrossRef](#)]
- Alsumari, W.; Hussain, M.; AlShehri, L.; Aboalsamh, H. EEG-Based Person Identification and Authentication Using Deep Convolutional Neural Network. *Axioms* **2023**, *12*, 74. [[CrossRef](#)]
- Wu, H.Y.; Rubinstein, M.; Shih, E.; Gutttag, J.; Durand, F.; Freeman, W.T. Eulerian Video Magnification for Revealing Subtle Changes in the World. *ACM Trans. Graph.* **2012**, *31*, 65. [[CrossRef](#)]
- Rubinstein, M. Analysis and Visualization of Temporal Variations in Video. Ph.D. Thesis, Massachusetts Institute of Technology, Cambridge, MA, USA, 2013.
- Wadhwa, N.; Rubinstein, M.; Durand, F.; Freeman, W. Phase-based Video Motion Processing. *ACM Trans. Graph.* **2013**, *32*, 80. [[CrossRef](#)]
- Lin, Z.; Chen, X.; Chen, G. Cervical Pulse Wave Extraction Based On Video Motion Amplification. In Proceedings of the 2021 IEEE 2nd International Conference on Big Data, Artificial Intelligence and Internet of Things Engineering (ICBAIE), Nanchang, China, 26–28 March 2021; pp. 784–787.
- Śmieja, M.; Mamala, J.; Prażnowski, K.; Ciepliński, T.; Szumilas, L. Motion Magnification of Vibration Image in Estimation of Technical Object Condition-Review. *Sensors* **2021**, *21*, 6572. [[CrossRef](#)] [[PubMed](#)]
- Fioriti, V.; Roselli, I.; Tati, A.; Romano, S.; Canio, G.D. Motion Magnification Analysis for structural monitoring of ancient constructions. *Measurement* **2018**, *129*, 375–380. [[CrossRef](#)]
- Zhang, D.; Zhu, A.; Hou, W.; Liu, L.; Wang, Y. Vision-Based Structural Modal Identification Using Hybrid Motion Magnification. *Sensors* **2022**, *22*, 9287. [[CrossRef](#)]
- Zhang, Y.; Chen, S.; Ruan, Y.; Lin, J.; Li, C.; Li, C.; Xu, S.; Yan, Z.; Liu, X.; Miao, P.; et al. The Facial Skin Blood Flow Change of Stroke Patients with Facial Paralysis after Peripheral Magnetic Stimulation: A Pilot Study. *Brain Sci.* **2022**, *12*, 1271. [[CrossRef](#)]
- Yoshida, K.; Nishidate, I. Phase Velocity of Facial Blood Volume Oscillation at a Frequency of 0.1 Hz. *Front. Physiol.* **2021**, *12*, 627354. [[CrossRef](#)]

24. Fu, G.; Zhou, X.; Wu, S.J.; Nikoo, H.M.; Panesar, D.; Zheng, P.P.; Oatley, K.; Lee, K. Discrete emotions discovered by contactless measurement of facial blood flows. *Cogn. Emot.* **2022**, *36*, 1429–1439. [[CrossRef](#)]
25. Buddhharaju, P.; Pavlidis, I.; Manohar, C. Face Recognition Beyond the Visible Spectrum. In *Advances in Biometrics: Sensors, Algorithms and Systems*; Ratha, N.K., Govindaraju, V., Eds.; Springer: London, UK, 2008; pp. 157–180.
26. Pistola, T.; Papadopoulos, A.; Mitianoudis, N.; Boulgouris, N.V. Biometric Identification using Facial Motion Amplification. In Proceedings of the IEEE International Conference on Image Processing (ICIP), Taipei, Taiwan, 22–25 September 2019.
27. Gkentsidis, K.; Pistola, T.; Mitianoudis, N.; Boulgouris, N.V. Deep Person Identification using Spatiotemporal Facial Motion Amplification. In Proceedings of the IEEE International Conference on Image Processing (ICIP), Abu Dhabi, United Arab Emirates, 25–28 October 2020.
28. Roccetti, M.; Gerla, M.; Palazzi, C.E.; Ferretti, S.; Pau, G. First Responders’ Crystal Ball: How to Scry the Emergency from a Remote Vehicle. In Proceedings of the 2007 IEEE International Performance, Computing, and Communications Conference, New Orleans, LA, USA, 11–13 April 2007; pp. 556–561.
29. Sun, Z.; Torrie, S.A.; Sumsion, A.W.; Lee, D.J. Self-Supervised Facial Motion Representation Learning via Contrastive Subclips. *Electronics* **2023**, *12*, 1369. [[CrossRef](#)]
30. Gonzalez, R.C.; Woods, R.E. *Digital Image Processing*; Prentice Hall: Upper Saddle River, NJ, USA, 2008.
31. von Arx, T.; Tamura, K.; Oba, Y.; Lozanoff, S. The Face—A Vascular Perspective. *Swiss Dent. J.-SSO-Sci. Clin. Top.* **2018**, *128*, 382–392. [[CrossRef](#)] [[PubMed](#)]
32. Zhu, X.; Ramanan, D. Face detection, pose estimation and landmark localization in the wild. In Proceedings of the Computer Vision and Pattern Recognition (CVPR), Providence, RI, USA, 16–21 June 2012.
33. Cootes, T.; Edwards, G.; Taylor, C. Active Appearance Models. In Proceedings of the European Conference on Computer Vision (ECCV), Freiburg, Germany, 2–6 June 1998.
34. Bulat, A.; Tzimiropoulos, G. Binarized convolutional landmark localizers for human pose estimation and face alignment with limited resources. In Proceedings of the Computer Vision (ICCV), Venice, Italy, 22–29 October 2017.
35. Liu, X. Generic Face Alignment using Boosted Appearance Model. In Proceedings of the 2007 IEEE Conference on Computer Vision and Pattern Recognition, Minneapolis, MN, USA, 17–22 June 2007; pp. 1–8.
36. Tran, D.; Bourdev, L.; Fergus, R.; Torresani, L.; Paluri, M. Learning spatiotemporal features with 3D convolutional networks. In Proceedings of the IEEE International Conference on Computer Vision (ICCV), Santiago, Chile, 7–13 December 2015.
37. Hochreiter, S.; Schmidhuber, J. Long short-term memory. *Neural Comput.* **1997**, *9*, 1735–1780. [[CrossRef](#)] [[PubMed](#)]
38. Koeppe, A.; Bamer, F.; Markert, B. An intelligent nonlinear meta element for elastoplastic continua: Deep learning using a new Time-distributed Residual U-Net architecture. *Comput. Methods Appl. Mech. Eng.* **2020**, *366*, 113088. [[CrossRef](#)]
39. Simonyan, K.; Zisserman, A. Very Deep Convolutional Networks for Large-Scale Image Recognition. *arXiv* **2014**, arXiv:1409.1556.
40. Srivastava, N.; Hinton, G.E.; Krizhevsky, A.; Sutskever, I.; Salakhutdinov, R. Dropout: A simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.* **2014**, *15*, 1929–1958.
41. Ioffe, S.; Szegedy, C. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. In Proceedings of the International Conference on Machine Learning, Lille, France, 6–11 July 2015.
42. Grandvalet, Y.; Canu, S.; Boucheron, S. Noise Injection: Theoretical Prospects. *Neural Comput.* **1997**, *9*, 1093–1108. [[CrossRef](#)]

Disclaimer/Publisher’s Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.