

Figure 1: System overview

four quarters of the phone screen. The front camera records a video clip containing the reflection of these light patterns from the user's face. These light patterns are not only used during video recording, but also help reconstruct 3D structure of the face and detect replay attacks. The recorded video then passes through a preprocessing module where first face region is extracted and aligned in each adjacent video frame. This is followed by an inverse gamma calibration operation applied to each frame to ensure linear camera response. Finally, the video is filtered by constructing its Gaussian Pyramid [13], where each frame is smoothed and subsampled to remove noise. After preprocessing, a temporal correlation between the passcode in the video frames and the one generated by the *Random Light Passcode Generator* is checked. If a high correlation is verified, the filtered video frames along with the random light passcode are fed into an *Image Recovery* module. The goal of this module is to recover the four stereo images corresponding to the four light sources, by utilizing the linearity of the camera response. The recovered stereo images are then used to compute face surface normals under unknown lighting conditions using a variant of photometric stereo technique [20]. A generalized human normal map template and its 2D wired mesh connecting the facial landmarks are used to compute these normals accurately. A 3D face is finally reconstructed from the surface normals by using a quadratic normal integration method [34]. Once the 3D structure is reconstructed, it is passed on to a liveness detection decision model. Here, a Siamese neural network [24] is trained to extract depth features from a known sample human face depth map and the reconstructed candidate 3D face. These feature vectors are then compared via L1 distance and a sigmoid activation function to give a similarity score for the two feature vectors. The decision model declares the 3D face as a real human face if this score is above a threshold and detects a spoofing attack otherwise.

4 FACEREVELIO ATTACK MODEL

Attacks to face authentication techniques can be classified into static and dynamic attacks. In a 2D static attack, a still object such as a photograph or mask is used, such that the face recognition algorithms would not be able to differentiate these presented objects from an actual face. Dynamic attacks aim at spoofing systems where some form of user action is required like making an expression or a gesture. In these attacks, a video of the user is replayed

performing the requested action. These videos can easily be forged by merging user's public photos with its facial characteristics. Adversaries can also launch a 3D static attack by using 3D models of the face. However, this requires advanced 3D printing capabilities which requires high cost. Similarly, 3D dynamic attacks involving building a 3D model in virtual settings, are impractical as described in [37].

In this paper, our goal is to prevent adversaries from spoofing face authentication systems with 2D static and dynamic attacks. We assume that an attacker has access to high-quality images of the legitimate user's face. We also assume that the adversary can record a video of the user while using *FaceRevelio*. In this case, the recorded video will capture the light patterns' reflections from the user's face. The attacker prepares these videos beforehand and launches an offline attack on our system by displaying them on a laptop screen/monitor. An adversary can possibly conduct an online attack if they have access to high-speed cameras, powerful computers, and a specialized screen with a fast refreshing rate such that it can capture and recognize the random passcode displayed on the screen on each use of the system, forge appropriate face responses depending on the passcode, and present the forged responses to the system. However, because of these difficult requirements for conducting such an attack, we believe that 2D attacks with photos/videos are still the major threat and the main focus of our paper.

5 FACEREVELIO SYSTEM DESIGN

5.1 Light Passcode Generator

To apply photometric stereo, we need to generate four images of the face illuminated under various light sources, from different directions. In order to simulate these light sources using the phone screen, we divide the screen into four quarters where each quarter is assumed a light source. During the video recording, each of these quarters is illuminated alternately in four equal intervals, while the other three quarters are dark. Figure 2 shows how the screen changes with different patterns during the four intervals and an example of the 3D reconstruction of the face using these patterns.

Random Passcode Generator: It could be argued that using these basic light patterns, the system would be prone to replay attacks. Keeping this in mind, we consider the idea of illuminating all the four quarters together for a certain period and changing the

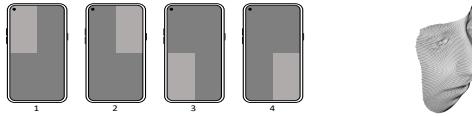


Figure 2: An example of 3D reconstruction using four basic light patterns displayed on four quarters of the screen.

screen lighting randomly at each time instance and each quarter to a random value drawn from a continuous range between -1 and 1 . Now, each quarter of the screen is illuminated simultaneously with a random pixel value, simulating four light sources. Based on this, we define a *light passcode* as a collection of four random light patterns displayed in the four quarters. In the rest of the paper, we will use *passcode* as a short-term for *light passcode*.

For the passcode, a random light pattern P_j is generated for a quarter j . During a time interval t_s , P_j is the light pattern represented as a sequence of random numbers, between -1 and 1 , of length t_s . The light pattern represents what each pixel of the screen is set to in the quarter j . In order to account for the smartphone screen refreshing rate, we apply an ideal low pass filter with a frequency threshold of 3Hz to each of the four light patterns. Although current smartphone screens support a refreshing rate of 60Hz , there is a delay when the screen is gradually updated from top to bottom. As a result, when the frequency threshold is set to a higher value, the intensity within each quarter may not be consistent. Additionally, setting a higher frequency threshold would result in rapid changes in the screen intensity, making it uncomfortable for users' eyes. These filtered patterns are then normalized such that each pattern is zero-mean.

One problem in illuminating the four quarters together is that the recorded video has a mixture of reflections of the four light patterns from the face. To be able to recover the stereo images from the mixture of reflections, we guarantee independence when combining the light patterns into a passcode. On top of ensuring their independence, we also introduce orthogonality between these four patterns. We apply Gram-Schmidt [18] process to the four light patterns to get their orthogonal basis and use these as patterns. Orthogonality assures a good separation between the impact of the four patterns on the human face and hence helps in the recovery of stereo images. Using induction and the fact that Gram-Schmidt process is linear, we can prove that if each of the original patterns satisfies the frequency threshold of 3Hz , the resulting orthogonal patterns are also within 3Hz . Figure 3 shows an example of a passcode with four patterns and the FFT of these patterns before and after the application of Gram-Schmidt process. We can see that the FFT of the patterns generated after applying Gram-Schmidt to the filtered random sequence only has components below 3Hz . On a side note, the above process is analogous to code-division multiple access (CDMA) [39] used in radio communications. In our case, the face is analogous to the shared media, the camera is the receiver and our orthogonal patterns are like the codes in CDMA. The stereo images generated by each independent quarter are like the data bit sent by each user. The difference is that in our case, we design and use patterns of continuous values that satisfy a frequency bound requirement.

As a result of the above steps, we obtain four orthogonal zero-mean light patterns, forming a passcode. Each value in the passcode

is then multiplied with an amplitude of 60 and finally the passcode is added on top of a constant base pixel intensity value of 128 to be displayed on the screen. Here, note that the Section 5.5 describes how the passcodes are used to defend against replay video attacks.

5.2 Video Preprocessing and Filtering

After generating a random passcode, the corresponding light patterns are displayed on the smartphone screen. Meanwhile, a video of their reflections from a user's face is recorded using the front camera. From the recorded video, first, we locate and extract the face in each frame by identifying the facial landmarks (83 landmarks) using Face++ [2]. We then use these landmarks to align the face position in every adjacent frame to neutralize the impact of slight head movements and hand tremors.

Since our following algorithms focus on how the changes in lighting conditions affect the captured face images, we preprocess the recorded video by converting each frame from the color space to the HSV space [10]. Only the V component will be kept and the other two components are discarded since the V component reflects the brightness of an image. Then, each video frame represented by the V component is further processed using Gaussian pyramid [13] which is a standard technique used in signal processing to filter noise and achieve a smoother output. We used Gaussian pyramid to remove any inherent camera sensor noise. Additionally, pyramids reduce the size of the input video frames by decreasing the spatial sampling density while retaining the important features within the frame, which in turn reduces the system's processing time. We use two levels of pyramid and select the peak of the pyramid in the subsequent steps for video analysis.

5.3 Image Recovery

Recall that in photometric stereo, at least three stereo images with different single light sources are needed for computing the surface normals. However, what we obtained so far is a series of frames, in which the lighting on the face at any given time is a combined effect of all four light patterns on the screen. Therefore, we need to recover these stereo images for each quarter from the preprocessed video frames, which is different from the traditional way of directly collecting stereo images used for photometric stereo.

Based on the theory that the intensities of incoherent lights add linearly [22], we propose to recover the stereo images by directly solving the equation, $G = WX$, where G is a $f \times n$ matrix representing the light intensity values received on each pixel in the recorded video frames, where f is the number of frames and n is the number of pixels in one frame. W represents the $f \times 4$ light patterns $[P_1; P_2; P_3; P_4]$ used while recording the video. $X (= [I_1; I_2; I_3; I_4])$ is a $4 \times n$ matrix representing the four stereo images that we aim to recover. This equation utilizes the fact that under a combined lighting condition, the light intensity received on a certain pixel is a weighted sum of four light intensities with a single light from each quarter.

However, we cannot directly use the above equation unless under the assumption that camera sensors can accurately capture light intensities and reflect the actual values. Problems, e.g. inaccurate image recovery, will arise if we ignore the possible effects of camera

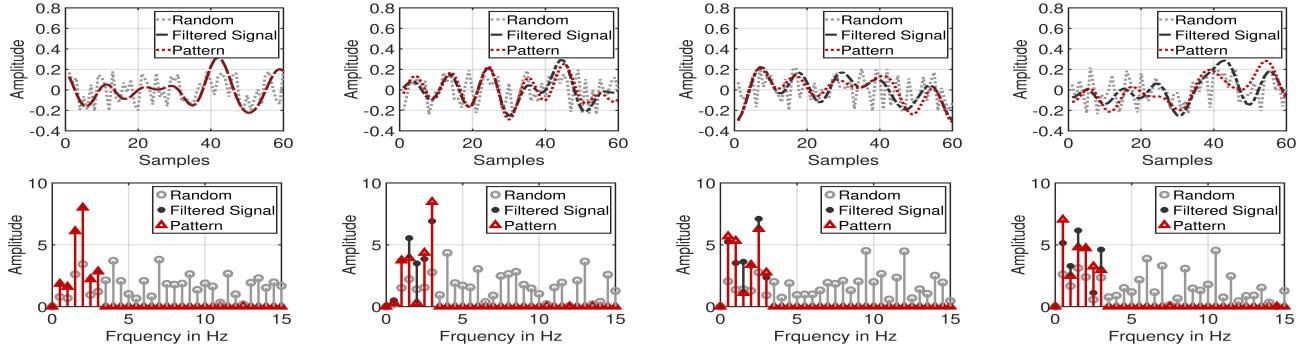


Figure 3: An example of a random passcode. The top row shows the four random patterns in the passcode before and after low-pass filtering and the final patterns after applying Gram-Schmidt process to the filtered pattern. The bottom row shows the FFT of these patterns before and after applying Gram-Schmidt process. The frequency bound still holds after applying Gram-Schmidt process.

parameters and sensitivity. Recently, smartphone camera APIs¹ started supporting manual camera mode which gives the user full control of the exposure parameters, i.e. aperture, shutter speed (exposure time) and sensitivity (ISO). In automatic mode, the camera continuously adjusts its ISO to compensate for lighting changes in the scene. In order to have a smoother camera response with changing light intensity, we use the camera in manual mode where its ISO is set to a fixed value.

Although the camera response curve is smooth after fixing the ISO, we still need to linearize the relationship between the image captured and the light from the screen to be able to use the equation for solving G . For this purpose, we dig deep into the mechanics of the camera sensor and image processing involved in generating the final output images. Cameras typically apply a series of operations on the raw camera sensor data to give us the final output images. These include linearization of sensor data, white balancing, demosaicing [6] and gamma calibration [7]. Gamma calibration is where non-linearity arises between the captured pixel values and the light intensity from the scene. In order to make use of linear relationship between these two, we apply an inverse of the gamma calibration, to the recorded video frames obtained from the camera. As a result, the resulting pixel values in the range between black and saturation level have a linear relationship with the actual light present in the scene. This relationship can be formulated as the linear model, $y = kx + b$, where b is the y-intercept introduced to account for the non-zero black level of the camera sensor. This inverse calibration is applied to each frame in the video preprocessing before face extraction. Now by generalizing the linear model to every frame, containing multiple pixels, we get

$$K = kG + B, \quad (4)$$

where K is the video frames that the camera actually captured for the duration of the passcode. By substituting the definition of G into Equation 4, we get

$$K = kWX + B. \quad (5)$$

¹Android supports manual camera mode starting from Android Lollipop 5.1

Finally, we use the least square method to solve

$$WX = \frac{1}{k}(K - B) \quad (6)$$

which can be written as

$$X = (W^T W)^{-1} W^T \left(\frac{1}{k}(K - B) \right) \quad (7)$$

Here, notice that B is a constant matrix and since each of the four patterns in the passcode W are zero-mean, the term $W^T B$ will be eliminated. Hence Equation 7 becomes:

$$X = (W^T W)^{-1} W^T \left(\frac{1}{k} K \right) \quad (8)$$

Note that this solution X will have an uncertainty of a scale factor. For any $\alpha > 0$, let $X' = \alpha X$, $k' = \frac{1}{\alpha} k$. X' , k' will also minimize the above function.

However, this will not have an impact on the reconstructed surface normals. Recall, that surface normals are computed by taking SVD of the stereo images. So, when X and X' are both factorized using SVD, the decompositions are

$$X = U\Sigma V^T, \quad (9)$$

$$X' = U(\alpha\Sigma)V^T. \quad (10)$$

The surface normal V^T will stay the same in these two cases. From the above observation, we can set $k = 1$ without any impact on the surface normals. Now, we can solve for X' by

$$X' = (W^T W)^{-1} W^T K \quad (11)$$

So far, we assumed that the only light present in the scene is due to the passcode displayed on the screen. However, we still need to consider the ambient light present in the scene as well as the base intensity value of the screen on top of which the passcode is added. To account for these other light sources, Equation 5 now becomes

$$K = kWX + B + C \quad (12)$$

where C is the constant light present in the scene. Again, since C is a constant, because of the orthogonal and zero-mean nature of our passcode, W , $W^T C$ will become 0. As a result, Equation 11 will give a solution for X even when ambient light is present.

Due to the inherent delay in the camera hardware, the recorded video may have some extra frames and the timestamps for each



Figure 4: The recovered stereo images corresponding to the four patterns in the passcode. The bottom row shows a binary representation to emphasize the differences in these stereo images.

video frame captured and the four patterns displayed on the screen at that point may differ. To ensure that we obtain a correct and fine alignment between these two, we first compute the average brightness of each frame and then apply a low pass filter on the average brightness across frames. The peaks and valleys in the average brightness are matched with those of the passcode and finally, DTW [11] is used to align the two series correctly. Once aligned, the result is the video frames which exactly represent the reflection of the passcode from the face. These video frames are then given as input to Equation 11 to recover the four stereo images as X . We define the average brightness of these video frames as the recorded passcode for later sections.

An example of the recovered four stereo images corresponding to every single light i.e. four patterns displayed in each quarter is shown in Figure 4. The top 4 images are the recovered stereo images. The bottom images are the binary representation of these stereo images such that in each image, a pixel value is 1 if the pixel in the corresponding stereo image is larger than the mean value of the same pixel in the other three stereo images. This binary representation is just to visually emphasize how different these stereo images are and how they represent the face illuminated from lighting in four different directions.

5.4 Photometric Stereo and 3D Reconstruction

The stereo images recovered from the least squared method approximate the facial images taken with four different point lights. Now, we can use these stereo images to compute the surface normals of the face as described in Section 2.

However, as mentioned earlier, these surface normals have an ambiguity of matrix A . We design an algorithm illustrated in Algorithm 1 to compute the normals without this uncertainty. We use a generalized template, N_t , for the surface normals of a human face and use this to solve for A . This template can be the surface normals of any human face recovered without any ambiguity like surface normals computed when the lighting is known. Note that obtaining this template is a one-time effort and the same normal template is used for all users. Along with the normal map, we also have a 2D wired triangulated mesh, M_t , connecting the facial landmarks (vertices), for this template. Now, when computing the normals of a user subject, we use the facial landmarks detected from an RGB image of the face to build a triangulated mesh of the face, M , using M_t as a reference for connecting the vertices and triangles. A representation of this mesh can be seen in Figure 5 (left). An affine transformation from the template mesh, M_t to M is then found independently for each corresponding pair of triangles in the two

ALGORITHM 1: Surface Normal Computation

Data: normal map template N_t , template mesh M_t , stacked four stereo images I and face RGB image R

Result: surface normals S

```

1:  $V \leftarrow getFaceLandmarks(R)$ 
2:  $M \leftarrow buildMesh(V, M_t)$ 
3:  $\hat{S}, \hat{L}^T \leftarrow SVD(I)$ 
4:  $N'_t \leftarrow transform(N_t, M_t, M)$ 
5: Solve  $N'_t = A\hat{S}$  for  $A$ 
6:  $S' \leftarrow A\hat{S}$ 
7:  $M_s \leftarrow symmetrizeMesh(M)$ 
8:  $S \leftarrow transform(S, M, M_s)$ 
9:  $S \leftarrow adjustNormalValues(S)$ 
10:
11: function TRANSFORM( $Z, T_1, T_2$ )
12:   for each pair of triangles  $< t_1, t_2 > \in T_1, T_2$  do
13:      $a \leftarrow affineTransformation(t_1, t_2)$ 
14:      $Z_{out} \leftarrow warp(Z(t_1), a)$ 
15:      $Z(t_2) \leftarrow Z_{out}$ 
16: end function

```

meshes and applied to the matching piece in N_t . As a result, the transformed normal map template, N'_t now fits the face structure of the user. This transformed template can finally be used to find the unknown A , by solving $N'_t = A\hat{S}$ where \hat{S} are the approximate normals recovered from SVD, and obtain the surface normals, S' . The last step in normal map computation is to make the normal map symmetric. This is needed to reduce noise in the recovered stereo images and hence the surface normals. We first find the center axis of the 2D face mesh using landmarks on the face contour, nose tip and mouth. Once the center is found, each pairing landmarks like eyes, eyebrow corner etc. are adjusted such that they have equal distance to the center to get a symmetric mesh. After symmetrizing the mesh, we fit S into this symmetrized mesh. Now, we can easily apply inverse symmetry to the x component of S and symmetrize the values in y and z components of S . Note that by introducing symmetry, we might lose some tiny details of the facial features as all human faces are not symmetrical. However, since our goal is to distinguish the human face from their spoofing counterpart and not another human, the information retained in the surface normals is more than sufficient. Figure 5 (right) shows an example of the x , y and z components of a normal map generated from our algorithm.



Figure 5: Normal map calculation (left) shows 2D triangulated face mesh generated by using facial landmarks. (right) shows the X , Y and Z components of the normal map generated from Algorithm 1.

After we have successfully recovered the surface normals, we can reconstruct the 3D surface of the face from them. For 3D reconstruction, we follow the quadratic normal integration approach

described in [34]. The results of 3D reconstruction are shown in Figure 6. Side and top view are shown for each reconstructed model.

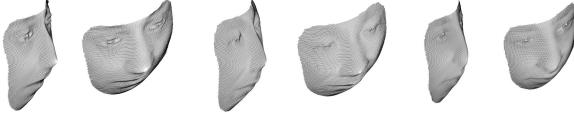


Figure 6: Examples of 3D reconstruction from human faces. Side and top views are shown.

5.5 Liveness Detection

FaceRevelio aims to provide security against two broad categories of spoofing attacks: 2D printed photograph attack and video replay attack.

2D Printed Photograph Attack: To defend against the 2D printed photograph attacks, we need to determine whether the reconstructed 3D face belongs to a real/live person or a printed photograph. Figure 7 shows examples of 3D reconstruction from a printed photograph using the approach described in the previous section. It is interesting to note here that the same general human face normal map template is used for computing the surface normals of a photograph. As a result, the overall structure of the reconstructed model looks similar to a human face. However, even when using this human normal map template, the freedom provided by solving for A is only up to 9 dimensions. Therefore, despite having a similar structure, the reconstruction from the 2D photograph lacks depth details in facial features, e.g. nose, mouth and eyes, as is clear in the examples in Figure 7.

Based on these observations, we employ a deep neural network to extract facial depth features from the 3D reconstruction and classify it as a human face or a spoofing attempt. We train a Siamese neural network adapted from [24] for this purpose. The Siamese network consists of two parallel neural networks whose architecture is the same, however, their inputs are different. One of these networks takes in a known depth map of a human face while the other is given the candidate depth map obtained after the 3D reconstruction. Therefore, the input to the Siamese network is a *pair of depth maps*. Both the neural networks in the Siamese network output a feature vector for their inputs. These feature vectors are then compared using L1 distance and a sigmoid activation function. The final output of the Siamese network is the probability of the candidate depth map being that of a real human face. If this output value is above a predefined threshold, τ_s , the system detects a real face. Otherwise,

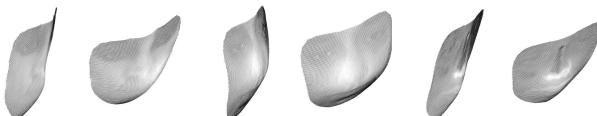


Figure 7: Examples of 3D reconstruction from 2D printed photographs.

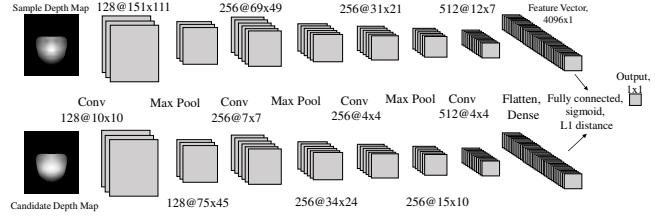


Figure 8: Architecture of the Siamese neural network. One of the twin neural networks takes a known human depth map as input while the other is passed the candidate 3D reconstruction.

a spoofing attempt is identified. Figure 8 shows the architecture of the Siamese network.

To elaborate the training process for our Siamese network, suppose we have N depth maps collected from human subjects and N depth maps from their photos/videos. From these depth maps, we obtain $N(N - 1)/2$ pairs of positive samples where both the depth maps in the pair are from human subjects. For the negative samples, we have N^2 pairs where one depth map is of a human subject while the other is from a photo/video. Since the total number of negative samples is larger than the positive samples, we randomly select $N(N - 1)/2$ samples from the negative pairs. These positive and negative samples are then used as input to train the Siamese network. Every time a subject tries to authenticate using *FaceRevelio*, the reconstructed 3D model along with a sample human depth map is fed as the input pair to the Siamese network. Here, note that the sample human depth map can be any depth map obtained from our reconstruction algorithm from a human subject in the training set and does not require registration by the test subject.

Since Siamese network uses the concept of one-shot learning [19] and takes pairs as input for training, the amount of data required for training is much smaller than traditional convolutional neural networks. Here, one may argue that why not train the model with the raw images/videos captured by the front camera, for the duration of passcode, instead of the depth map to decide if the subject is a human or not? Although intuitive, training such classifiers would require huge amounts of data; datasets for different environment settings, different light passcodes, different distances between the face and phone and various face orientations. In contrast, our image recovery module and approach for reconstructing the 3D surface of the face account for the ambient lighting and different orientations of the face before generating the depth map. Training a model using these depth maps ensures that input to our network is not impacted by the various ambient environment conditions; hence, much less data is required for training. Furthermore, models with video input are more complex with larger number of trainable parameters, resulting in higher storage and computation costs.

Video Replay Attacks: *FaceRevelio* has a two-fold approach for defending against video replay attacks. The first line of defense is to utilize the randomness of the passcode. When a human subject tries to authenticate via *FaceRevelio*, the passcode displayed on the screen is reflected from the face and captured by the camera. As a result, the average brightness of the video frames across time has a high correlation with the light incident upon the face i.e. the sum

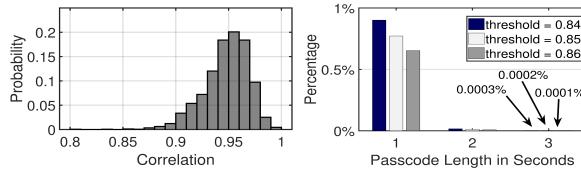


Figure 9: Video Replay Attack: (left) shows the distribution of correlation between recorded passcodes from human face and the original passcode. (right) shows the percentage of passcodes which have a correlation with another random passcode higher than a threshold for different thresholds.

of the four patterns in the passcode displayed on the screen. Figure 9 (left) shows a distribution of the correlation between recorded passcodes and the original passcode for experiments conducted with humans. The correlation between the two passcodes is higher than 0.85 for more than 99.9% of the cases. An adversary may try to spoof our system by recording a video of a genuine user while using *FaceRevelio* and replay this video on a laptop screen or monitor in front of the phone later. In this case, the video frames captured by the camera will have the reflections of the passcode on the phone screen as well as the passcode present in the replay video. Since *FaceRevelio* chooses a totally random passcode each time as described in 5.1, the probability that the passcode displayed on the screen and the passcode in the video has a high correlation is extremely low. To give an idea, for a passcode duration of 3s, if we compare 300 million pairs of random passcodes, only 0.0003% of the pairs will have a correlation greater than 0.84. Figure 9 (right) shows the percentage of passcode with a correlation higher than threshold values 0.84, 0.85 and 0.86 for passcode lengths of 1, 2 and 3s. Hence, just by computing and setting a threshold on the correlation between the recorded passcode and the sum of passcode from the screen, the chances of detecting a replay attack are very high.

For the rare cases when the correlation is higher than the predefined threshold, our second line of defense comes into play. Similar to 2D photograph attack, video replay attacks can also be detected using the reconstructed 3D model. The reconstruction from the replayed video suffers from two main problems. First, it is hard for the adversary to accurately synchronize playing the attack video with the start of the passcode display on the smartphones. Second, even if the correlation passes the threshold, there will be some differences in the replayed passcode and *FaceRevelio*'s passcode. Because of this, the DTW matching will not match the recorded video frames with the displayed passcode very well. Hence, the four stereo images, X , obtained by solving equation 11 will not be representative of the subject's face being illuminated from four different lighting directions. As a result, the surface normals and 3D reconstruction from these wrong stereo images do not capture the 3D features of the face and is sufficient to identify a spoofing attempt.

6 EVALUATION

We describe the implementation and evaluation of our system in this section. We first describe the experiment settings and the data

collection details and then the performance of our system in different settings.

6.1 Implementation and Data Collection

We implemented a prototype for *FaceRevelio* on Samsung S10 which runs Android 9, with 10 MP front camera that supports Camera2 API. The videos collected for our authentication system have a resolution of 1280x960 and a frame rate of 30fps. For each experiment setting, we display the passcode patterns on the smartphone screen and record a video of the reflections from the user's face via the front camera. We use Face++ [2] for landmark detection and OpenCV in the image recovery and reconstruction modules of our system. Python libraries for TensorFlow [9] and Keras were used to train the neural network for liveness detection while TensorFlow Lite was used for inference on Android.

We evaluated *FaceRevelio* with 30 volunteers using our system for liveness detection. The volunteers included 19 males and 11 females with ages ranging from 18 to 60. These volunteers belonged to different ethnic backgrounds including Americans, Asians, Europeans and Africans. During the experiments, the volunteers were asked to hold the phone in front of their faces and press a button on the screen to start the liveness detection process. Once the button was clicked, the front camera started recording a video for the duration of the passcode. During all experiments, we collected a high-quality image of the user to test the performance of our system against photo attacks. For the video replay attack, we used the videos collected from real users and replayed them to the system.

We collected a total of 4900 videos from the 30 volunteers over a duration of 3 weeks. We evaluated the performance of our system in natural daylight as well as in completely dark environment (0 lux). For the daylight setting, all experiments were conducted during daytime however the light intensity varied (between 200 to 5000 lux) based on the weather conditions on the day and time of the experiment. Each volunteer performed 10 trials of liveness detection using our system for each of the two light settings. A random passcode of 1s duration was added on top of a gray background (grayscale intensity value of 128) for these trials. We also tested *FaceRevelio* with passcode durations of 2 and 3s in the two light settings. We also evaluated the impact of indoor lighting (~ 250 lux), the distance between the face and the smartphone screen and the orientation of the face, on the performance of our system. For these scenarios, we collected data from 10 volunteers with a passcode duration of 1s. These volunteers used the system 30 times for each scenario. In addition, we also explored whether using a background image affects *FaceRevelio*'s performance.

We used the Siamese neural network described in section 5.5 to test each user. We employed a leave-one-out method for each test user where we used the depth maps generated from the data collected from the remaining 29 users for training. From these 29 users' data, we used 80% of the data as the training set while the remaining 20% was used for validation. Hence, the test user's data remained unseen by the network during the training process. At inference time, the depth map from the test subject along with a sample human depth map, randomly selected from the human depth maps collected from the other 29 subjects, was given as the input pair to the Siamese network. The predefined threshold, τ_s ,

