

# An Automatic Face Attendance Checking System using Deep Facial Recognition Technique

Thuy Nguyen-Chinh, Thien Do-Tieu, Sy Nguyen-Tan, Phuong Le-Van-Hoang, Qui Nguyen-Van, Phu Nguyen-Tan

**Abstract**—Nowadays, as computers are powerful enough for implementing complex algorithms, there are numerous applications that people utilize computers to run. In which, facial recognition is one of the most active fields of applications. In fact, computers can not only automatically identify who a person is, but also operate 24/7, which human beings cannot endure. This leads to the replacement of people by computers in some repetitive and real-time applications.

In this work, we apply the facial recognition into an attendance checking system that uses faces of registered people to check their attendance. This system has a GUI which allows easy user-to-system interaction. The core of the system is a deep facial recognition technique, which has four stages (e.g., removing motion-blur frames, detecting faces, removing non-frontal-view faces, and recognizing). Particularly, in the recognition phase, we consider this stage as an open-set facial recognition problem, so the system is able to detect people who have not registered in the database before. Also, we boost the performance of the system by utilizing hardware resources of users' computers. Although the system is designed to run with a low-resolution webcam, its performance is reasonably accurate on our private dataset.

**Index Terms**—Face Attendance Checking, Facial Recognition, Deep Learning

## I. INTRODUCTION

Face recognition systems are applying widely in real life, such as: tracking, managing employees, finding information of celebrities, etc. There are many approaches to design a face recognition system, but these systems frequently are affected by light, non-frontal faces, resolution of cameras, etc, each method have many separable challenges. Overall, a face recognition has two main stages which are face detection and face recognition, yet we want to create a constraint on frontal faces for users, that lead our system to have three stages: face detection, face landmark detection and face recognition.

### A. Face detection

Face detection and alignment are essential to many face applications such as face recognition and facial expression analysis. However, the large visual variations of faces, such as occlusions, large pose variations and extreme lightings, impose great challenges for these tasks in real world applications.

This work is a final project in the course of "Artificial Intelligence in Control Engineering" Sep-Dec 2018, guided by Dr. Cuong Pham-Viet (email: pvcuong@hcmut.edu.vn), Faculty of Electrical and Electronics Engineering, HoChiMinh city University of Technology.

Authors are senior students in the Faculty of Electrical and Electronics Engineering, University of Technology, HoChiMinh city Vietnam National University (e-mail: {thuy.ng.ch, dotieuthien9997, tansyab1, hpcqt97, nvqui97, tanphu97.nguyen}@gmail.com). The software is open source and can be found in <https://github.com/AntiAegis/Face-Attendance-System>.

The cascade face detector proposed by Viola and Jones [1] utilizes Haar-Like features and AdaBoost to train cascaded classifiers, which achieve good performance with real-time efficiency. However, quite a few works [2, 3, 4] indicate that this detector may degrade significantly in real-world applications with larger visual variations of human faces even with more advanced features. Besides the cascade structure, [5, 6, 7] introduce deformable part models (DPM) for face detection and achieve remarkable performance. However, they need high computational expense and may usually require expensive annotation in the training stage. Recently, convolutional neural networks (CNNs) achieve remarkable progresses in a variety of computer vision tasks, especially face detection task. Li et al. [19] use cascaded CNNs for face detection, but it requires bounding box calibration from face detection with extra computational expense and ignores the inherent correlation between facial landmarks localization and bounding box regression. Face alignment also attracts extensive interests. Regression-based methods [12, 13, 16] and template fitting approaches [14, 15, 7] are two popular categories.

However, most of the available face detection and face alignment methods ignore the correlation between these two tasks. Though there exist several works attempt to jointly solve them, there are still limitations in these works. For example, Chen et al. [18] jointly show alignment and detection with random forest using features of pixel value difference. But, the handcraft features used limits its performance. From those previous experiments, we choose an new approach which integrate these two tasks using unified cascaded CNNs by multi-task learning called Multi-task Convolutional Network in section III-D.

### B. Landmark detection

The locations of the fiducial facial landmark points around facial components and facial contour capture the rigid and non-rigid facial deformations due to head movements and facial expressions. They are hence important for various facial analysis tasks. Many facial landmark detection algorithms have been developed to automatically detect those key points over the years. In this paper, we use dlib library which is a powerful source for face and facial landmark detection. We will discuss our implement in detail in section III-C.

### C. Face recognition

After face detection and alignment, those regions of face is extracted to get feature vectors. With conventional way, One of the most popular feature for face recognition is Gabor

feature. Tudor Barbu ?? uses Gabor transform to extract feature, and then using K-Nearest Neighbour (K-NN) based on clustering feature to predict identity of a face. This implement achieve quite impressed performance with accuracy of 90% on Yale Face Database B. In Opencv library which focuses on algorithms of Computer Vision introduces a method called Local Binary Patterns (LBP) based on Haar-Like feature. In term of speed, LBP has relly real-time efficiency, whereas it is not stable in term of accuracy, this method cannot face with arounded noise which is the reason why LBP and Haar-Like feature are rarely applied in practical systems. Because of limitations of conventional features, deep learning models gradually instead and get better and better. Yi Sun, Xiaogang Wang, Xiaoou Tang ?? build a deep model Deep hidden IDentity features (DeepID) which uses convolutional neural network to extract face feature. Advantage of this model is using a small dataset for training, that is consideraby good for systems which cannot collect a large dataset of users. However, to reach a high accuracy, DeepID model become really complex with many neural network branches for each person. There are 10 patches of face which contain interested information are chosen from each image, then they are scaled with three figures in RGB and gray chanel. Totally, model have 60 different networks to extact feature of an image, then feedforwarding feature into a classifier using Joint Bayesian. DeepID achieves excellent accuracy of 97.45% on dataset Labeled Faces in the Wild (LFW). In 2014, the authors of DeepID show DeepID2 which is a improved version of DeepID. In new version, interested regions of face algorithm is built to eliminate useless patches which cannot extract high-level feature. That work really helpfull affect accuracy, specifically there is a increase in accuracy at 99.63%. In 2015, Google Inc.?? use deep convolutional network Inception and triplet loss function in FaceNet mdoel to extract feature. Their outstanding work in this model is using hard triplet loss to separate feature for each person, so FaceNet feature is robust in both face verification and face recognition. The accuracy of 99.63% on LFW and 95.12% on Youtube Faces DB dataset is high enough to represent the perfection of model.

## II. PROPOSED SYSTEM

In this paper, we apply deep facial recognition techniques into the problem of face attendance checking. A system is built in order to manage appearances of students in a class, which is revealed in Fig. 1. As normally, a facial recognition system is organized as a pipeline of typical stages such as face detection, landmark detection, and face recognition. However, to ensure input frames for underlying algorithms are high quality, we append an early filter (the Blur detection stage) that are able to discard blur frames, which are caught by motions of people in front of a standard webcam. Then, clean frames are passed through the Face detection block to count the number of faces existing inside these frames. In our specific case, there is only one face per frame allowed to be processed, so frames that contain more than one faces are rejected by the block. Afterward, the Landmark detection is to localize statistic-salient points in the face in order to verify whether

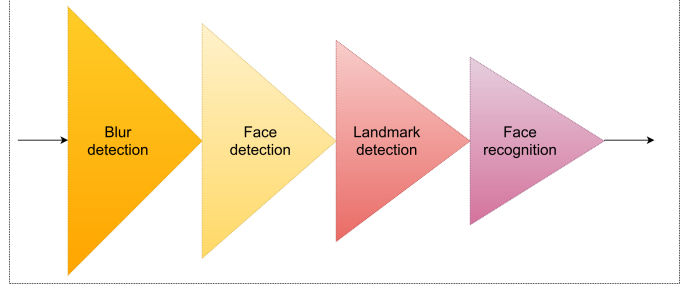


Fig. 1: The system pipeline. There are four stages, namely blur detection, face detection, landmark detection, and face recognition. These four blocks are in the descending order of size in the direction from input to output. This points out that our system are tougher to input frames from the camera when such frames passed through the system. Therefore, best frames are likely to be processed, which may improves the final recognition accuracy.

the face is in frontal view of the camera. This frontal-view check will help improve the accuracy of the face recognition algorithm. Finally, the Face recognition uses blur-clean, one-face, and frontal-view frames from previous stages to extract relevant features and then perform a classification task to indicate which category the input most likely belong to.

In addition, to leverage the ease in use, we design a friendly graphic user interface (GUI) so that people who want to use the system to manage (teachers) or check (students) attendance can interact with the application without any specific knowledge. To make the system more robust, we carefully analyze the distribution outlier of features representing for registered accounts. Therefore, the algorithm has ability to detect people who have not registered in the application before, which is equivalent to the open-set problem in face recognition.

Our work is organized as follows. In the section III, stages of the proposed system are described clearly, including motion-blur detection, face detection, frontal-view detection, and face recognition. Then, section IV is for reporting some experimental results.

## III. IMPLEMENTATION

### A. Motion-blur detection

The first stage of this system is detecting blurred image and rejecting them out of next stage. We know that the blurred image means each pixel in the source image gets spread out and mixed into surrounding neighbour pixels. For our attendance checking system, the motion blur happens when an object (namely face or webcam) moves during the exposure. So as to detect whether an image is blurred, we use the 2D-FFT (2D-Fast Fourier Transform) method.

We will review about Fourier Transform of Images. To compute the Fourier transform of an image, you need to:

- Compute DFT of each row, in place.
- Compute DFT of each column, in place.

When a signal is discrete and periodic, we use the discrete Fourier transform, or DFT. Suppose our signal is  $a_n$  for  $n = 0 \dots N - 1$ , and  $a_n = a_{n+jN}$  for all  $n$  and  $j$ . The spectrum of  $a$  is:

$$A_k = \sum_{n=0}^{N-1} W_N^{kn} a_n \quad (1)$$

where

$$W_N = e^{-i\frac{2\pi}{N}}$$

and  $W_N^k$  for  $k = 0 \dots N-1$  are called the  $N$ th roots of unity. The sequence  $A_k$  is the discrete Fourier transform of the sequence  $a_n$ . Each is a sequence of  $N$  complex numbers.

The FFT is a fast algorithm for computing the DFT. If we take the 2-point DFT and 4-point DFT and generalize them to 8-point, 16-point, ...,  $2^m$ -point ( $n$  is an integer), we get the FFT algorithm.

There are several ways to write an FFT. For instance, let  $m$  be an integer and let  $N = 2^m$ . Suppose that  $x = [x_0, \dots, x_{N-1}]$  is an  $N$  dimensional complex vector. Let  $\omega = \exp(-\frac{2\pi i}{N})$ . Then the DFT,  $c = F_N(x)$  is given by

$$c_k = \frac{1}{N} \sum_{j=0}^{N-1} x_j \omega^{jk}. \quad (2)$$

Let  $n = N/2$ , let  $u$  and  $v$  be  $n$  dimensional vectors defined by

$$u_j = x_j + x_{j+n}, \quad j = 0, \dots, n-1 \quad (3)$$

$$v_j = (x_j - x_{j+n})\omega^j, \quad j = 0, \dots, n-1. \quad (4)$$

Then

$$c_{2j} = \frac{1}{2}(F_n(u))_j, \quad j = 0, \dots, n-1 \quad (5)$$

$$c_{2j+1} = \frac{1}{2}(F_n(v))_j, \quad j = 0, \dots, n-1. \quad (6)$$

To compute the DFT of an  $N$ -point sequence using equation (1) would take  $(N^2)$  multiplies and adds. The FFT algorithm computes the DFT using  $(N \log N)$  multiplies and adds.

Practical issues: We translate the picture so that pixel (0,0), which now contains frequency  $(\omega_x, \omega_y) = (0, 0)$ , moves to the center of the image. Then, we display pixel values proportional to  $\log(\text{magnitude})$  of each complex number. For color images, do the above to each of the three channels (R, G, and B) independently.

Apply to our system, firstly, we calculate FFT of image. Secondly, we will compute mean amplitude spectrum value of entire pixel in image and. Finally, the result of this operation is compared to an optimal threshold which distinguishes blurred and non-blurred image as accurate as possible. The image is called non-blurred if and only if its average value greater than the threshold value, and vice versa. After that, non-blurred images are applied to face detection stage of system.

### B. Face detection

In this paper, we have used Histogram of Oriented Gradients method for extracting features of the face and Linear Support Vector Machine (SVM) method for face detections.

The implementation of this method using sliding window technique with the different sizes of the windows. Using the sliding window technique we could complete the calculation of HOG features, applied to detect and differentiate the face and the false face recognition using the SVM technique. All of the pre-processing steps are automatically implemented

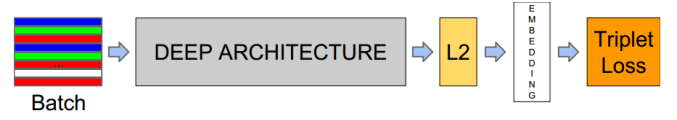


Fig. 2: Model structure. Our network consists of a batch input layer and a deep CNN followed by L2 normalization, which results in the face embedding. This is followed by the triplet loss during training.

before using Dlib library with the input of being given facial images and the output of localization the identified faces.

### C. Frontal-view detection

To check whether the shape of the faces has to be frontal, we implement these 3 following steps: Step 1.

- 1) Focusing on the center of the image. Only accept these faces that locate in the most central of the images.  $((120 - 360), (90, 360))$
- 2) Identify the skew of the image: calculate the coordinate of these eyes and the angle deviation between two eyes in the horizontal direction. If the angle deviation is larger than 10 degree, the image will be ignored.
- 3) Identify the rotation of the image: choose the point which is the midpoint of the right and the left eye. If the nose which is deviated from the selected point is greater than 10 pixels in the horizontal direction, the image is ignored.

These steps are implemented based on 5-point facial landmark technique with Dlib library instead of 68-point facial landmark in order to improve performance. If the image satisfies the condition, it will be accepted.

### D. Face recognition

In this stage, faces in raw images are detected and aligned by Multi-task CNN, we use convenient pre-trained FaceNet model to extract feature (in Figure 2) and then feedforward it to a SVM classifier for recognition.

1) *Multi-task Convolution Network*: The overall pipeline Multi-task CNN is shown in Figure 3. An image is initially resized to different scales to build an image pyramid, which is the input of the following three-stage cascaded framework with CNN architectures in Figure 4:

**Stage 1:** A fully convolutional network is exploited, called Proposal Network (P-Net), to obtain the candidate windows and their bounding box regression vectors. Then using the estimated bounding box regression vectors to calibrate the candidates. After that, employing non-maximum suppression (NMS) to merge highly overlapped candidates.

For each candidate window, P-CNN predict the offset between it and the nearest ground truth (i.e., the bounding boxes left top, height, and width). The learning objective is formulated as a regression problem, and the Euclidean loss is employed for each sample  $x_i$ :

$$L_i^{box} = \|y_i^{prediction} - y_i^{truth}\|_2^2 \quad (7)$$

The figure illustrates the architecture of the proposed face detection framework, showing three main components: P-Net, R-Net, and O2-Net. Each component takes an input image and processes it through a series of convolutional and pooling layers to produce three outputs: face classification, bounding box regression, and facial landmark localization.

- P-Net:** Takes an input of size  $12 \times 12 \times 3$ . It consists of two convolutional layers (Conv:  $3 \times 3$ ) and a max pooling layer (MP:  $3 \times 3$ ). The output is a stack of three feature maps: face classification (1x1x2), bounding box regression (1x1x2), and facial landmark localization (1x1x10).
- R-Net:** Takes an input of size  $24 \times 24 \times 3$ . It consists of two convolutional layers (Conv:  $3 \times 3$ ), a max pooling layer (MP:  $3 \times 3$ ), and a fully connected layer (fully connect:  $4 \times 48$ ). The output is a stack of three feature maps: face classification (1x1x2), bounding box regression (1x1x2), and facial landmark localization (1x1x10).
- O2-Net:** Takes an input of size  $48 \times 48 \times 3$ . It consists of four convolutional layers (Conv:  $3 \times 3$ , Conv:  $3 \times 3$ , Conv:  $3 \times 3$ , Conv:  $2 \times 2$ ) and a fully connected layer (fully connect:  $3 \times 128$ ). The output is a stack of three feature maps: face classification (1x1x2), bounding box regression (1x1x2), and facial landmark localization (1x1x10).

The sections that follow describe how to create GUIs with PyQt5. This includes laying out the components, programming them to do specific things in response to user actions, and



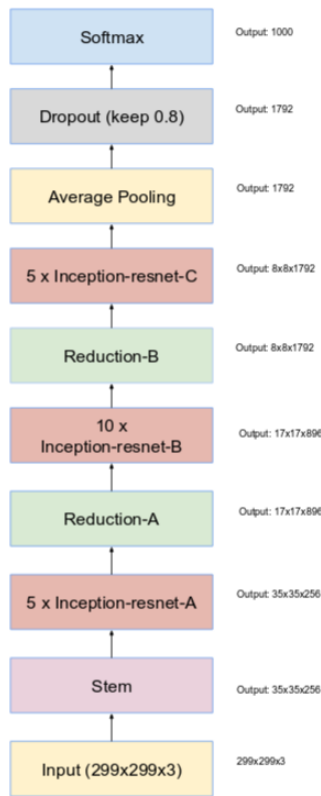


Fig. 6: Schema for Inception-ResNet-v1 and Inception-ResNet-v2 networks. This schema applies to both networks but the underlying components differ.

saving and launching the GUI; in other words, the mechanics of creating GUIs.

PyQt5 implements GUIs as figure windows containing various styles of uicontrol objects. You must program each object to perform the intended action when activated by the user of the GUI. In addition, you must be able to save graphical user interface development environment.

### GUI Development Environment

The process of implementing a GUI involves two basic tasks:

- Laying out the GUI components
- Programming the GUI components

GUIDE primarily is a set of layout tools. However, you must create a PY-file that contains code to handle the initialization and launching of the GUI. This PY-file provides a framework for the implementation of the callbacks - the functions that execute and launch your GUI.

### The Implementation of A GUI

#### Laying out the GUI components

While it is possible to write an PY-file that contains all the commands to lay out a GUI, it is easier to use GUIDE to lay out the components interactively and to generate one file that save the GUI:

UI-file -contains a complete description of the GUI figure and all of its children (uicontrols and axes), as well as the values of all object properties.

#### Programming the GUI components

Once you have the UI-file, you will use PyQt5 to extract the UI-file to PY-file to execute the program. PY-file -contains the functions that launch and control the GUI and the callbacks, which are defined as sub functions. This PY-file is referred to as the application PY-file in this documentation. Note that the application PY-file does not contain the code that lays out the uicontrols; this information is saved in the UI-file. The following diagram illustrates the parts of a GUI implementation.

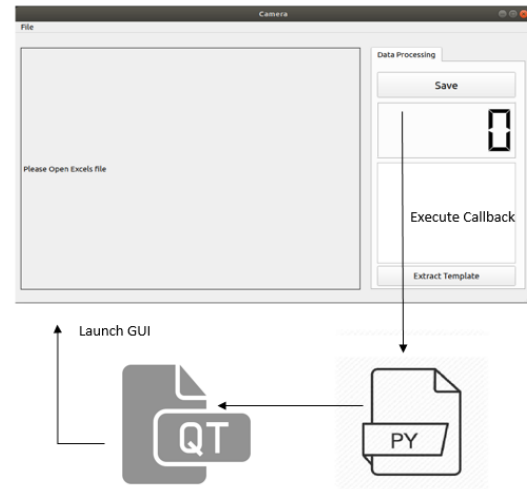


Fig. 7: Parts of GUI Implementation

The GUI are available from the QtDesigner shown in the figure below. The design is described briefly below. Subsequent sections show you how to use them.

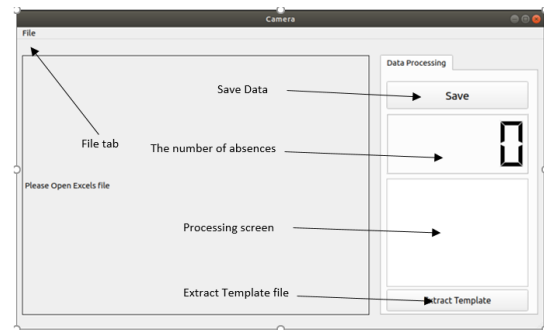


Fig. 8: GUI layout

### F. Attendance management

This is the final phase of Face Attendance Checking System. It was designated to mark the presence of one resulted from our algorithm in a file of excel format, namely xlsx extension. To be used by the system, the excel file must meet a stringent format made up of essential contents and be generated by the GUI.

Fig. 9 depicts a new standard empty excel table generated by our GUI. After obtaining a new file, we should fill in the table with the desired data (Fig. 10). The most special things in this table are column ID and Total. ID is considered a primary key because the algorithm will mark the presence of a

uicontrols	function
<b>Open File Option</b>	We use this option to open the existing file Open option appears and you can choose to open the file.
<b>Start Camera</b>	Open camera to start checking process
<b>Stop Camera</b>	Stop camera
<b>Exit</b>	We use this option to close the program. we can follow the steps: Click on File tab >Exit, active file will be closed. When we close the file, we get the confirmation message to save the file or not or cancel the command.
<b>Save</b>	Store completely-checked student IDs on the document. Once program exits, this process will be automatically started so that you wont lose your work that has been completed if there is a power interruption or other system malfunction, the first time you Save, it will take you to the Save File dialog box.
<b>Extract Template</b>	We use this option to export the template in XLSX document. To Extract the file, we can follow the steps: Click Extract Template . And then we can export it as per our requirement.

TABLE I. USER GUIDE

DANH SÁCH SINH VIÊN				
TRÍ TUỆ NHÂN TẠO TRONG ĐIỀU KHIỂN				
1 = present blank = absent				
ID	Last Name	First Name	Group	Total

Fig. 9: New standard excel form

specific person via his ID. To help the host in easy attendance management, we designed the column Total with a view to showing the number of absences in all.

Fig. 11 depicts an excel file's content after a checking progress finished. The GUI will automatically insert the only one new day column between Group and Total ones and in the tail of previous checked day. Letter 1 will be marked as presence in a cell of this column accordant to an ID. After attendance checking process is completed, the Total column will display the number of absences of previous days and the current one. Smartly can it display as we specially assigned a size-dynamic sum function to each cell of this column.

DANH SÁCH SINH VIÊN				
TRÍ TUỆ NHÂN TẠO TRONG ĐIỀU KHIỂN				
1 = present blank = absent				
ID	Last Name	First Name	Group	Total
1511844	Lương Hữu Phú	Lộc	1	
1512221	Phạm Ngọc Khôi	Nguyễn	1	
1512396	Bùi Tấn	Phát	1	
1512534	Nguyễn Trọng	Phúc	1	

Fig. 10: Excel form contain pre-inputed data

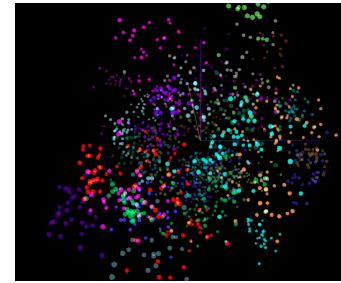
DANH SÁCH SINH VIÊN				
TRÍ TUỆ NHÂN TẠO TRONG ĐIỀU KHIỂN				
1 = present blank = absent				
ID	Last Name	First Name	Group	09/06/2018
1511844	Lương Hữu Phú	Lộc	1	1
1512221	Phạm Ngọc Khôi	Nguyễn	1	1
1512396	Bùi Tấn	Phát	1	1
1512534	Nguyễn Trọng	Phúc	1	1

Fig. 11: Form is under checking

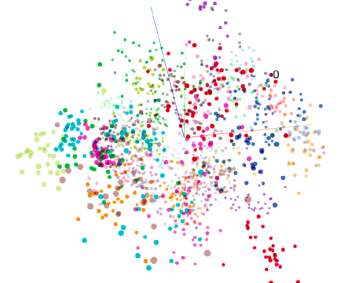
#### IV. EXPERIMENTAL RESULT

In this section, we first evaluate the effectiveness of the feature extracted from FaceNet. Then we will compare our system in different context such as: background, illumination, resolution of camera. Finally, we evaluate the computational efficiency of our system.

##### A. Embeddings



(a)



(b)



(c)

Fig. 12: Embeddings with PCA visualization

Because we use pre-trained model FaceNet, we need to test specification of embeddings which is output of model. We use PCA ?? and t-SNE ?? to visualize embeddings in 3 dimension space in Figure 12 and Figure 13. In PCA visualization, embeddings of the same person close together, although they are not completely separable. In another of t-SNE, because t-SNE method include clustering stage, so embeddings totally belong to their classes. If embeddings of FaceNet model is

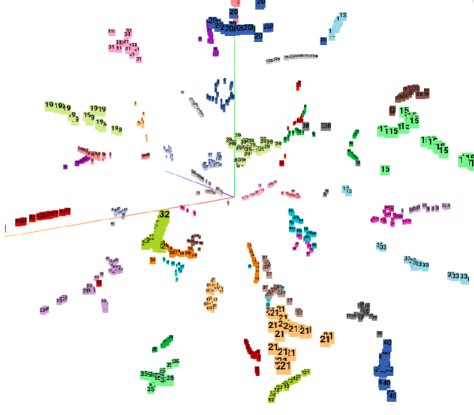


Fig. 13: Embeddings with t-SNE visualization

not contain high-level of specification, reducing dimension algorithms cannot show or cluster embeddings properly.

### B. Training

Training data is carefully collected with different views from  $-70^\circ$  to  $70^\circ$ . This work can improve accuracy in practical system, because it is difficult for users to keep their faces in a correct position. In training data include 52 identities and 1560 images totally.

The accuracy of SVM classifier is 99.36%, after training classifier, we train to get the best threshold. We divide threshold in range  $[0; 1]$ . As a result, threshold for 52 identities is 0.18825. Testing accuracy with threshold achieve 98.85% on testing subset. In practical environment, we test on 32 identities, three are 29 identities recognized easily and 3 identities who are not recognized continuously. In Figure 14, the (a),(b) are training images and (c),(d) are testing images, the effect of different illumination lead the probabilities of testing anchor are lower than threshold, so training data have to cover many real-life cases to create the best classifier.

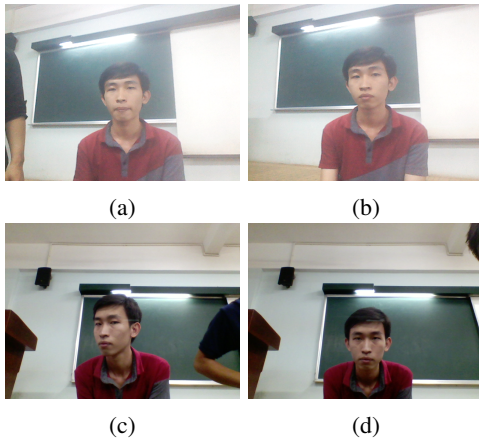


Fig. 14: Unrecognized identity, (a),(b) are training images, and (c),(d) are testing images in practical condition.

The system has a pipeline with four stages (e.g., motion-blur detection, face detection, landmark detection, and face recognition). Besides, the system is also integrated a friendly GUI, which allows users both teachers and students interact with it in an easy way. On our private dataset, the application perform accurate despite of the low-resolution webcam of typical laptops. This demonstrates that our underlying algorithm is effective to deal with this poor-quality input problem.

In the future, we will target to widen our dataset so that the dataset will be asymptotic to real applications. In addition, more algorithms will be considered to improve the ability of the algorithm to discriminate feature distributions of output classes.

### ACKNOWLEDGMENT

The authors would like to thank Dr. Pham Viet Cuong for providing documents as well as chance for us to do this work. Also, the authors would like to thank ...

### REFERENCES

- [1] B. Yang, J. Yan, Z. Lei, and S. Z. Li, "Aggregate channel features for multi-view face detection", IEEE International Joint Conference on Biometrics, 2014, pp. 1-8.
- [2] P. Viola and M. J. Jones, "Robust real-time face detection", International journal of computer vision, vol. 57, no. 2, pp. 137-154, 2004.

### V. CONCLUSION

In this work, we applied the deep facial recognition techniques to solve the problem of face attendance checking.