

## Volume Deep Face: A 3D Face Descriptor for Face Authentication System

Suttipat Srisuk<sup>1\*</sup>, Damrongsak Arunyagool<sup>1</sup>, Kitchanut Ruamboon<sup>1</sup>, Pantre Kompitaya<sup>2\*</sup>, Nattapong Jundang<sup>3\*</sup>

<sup>1</sup>Faculty of Engineering, Nakhon Phanom University, Thailand

<sup>2</sup>Thatphanom College, Nakhon Phanom University, Thailand

<sup>3</sup>Faculty of Engineering and Technology, Mahanakorn University of Technology, Thailand

\*Corresponding Author E-mail: srisuk.s@gmail.com, pantre@npu.ac.th, jnattapo@mut.ac.th

Received: 6/09/25, Revised: 5/10/25, Accepted: 6/10/25

### Abstract

In this paper, we introduce the Volume Deep Face (VDF), a novel face representation proposed for the face authentication system. VDF provides a fast and compact representation of faces using deep learning, enabling one to encode more distinctive features. Using our proposed method, images can be generated to form a 3D VDF representation or a 2D face descriptor (2DFD). The 3D VDF is created from multiple images in the training set, while the 2DFD is generated from a single image during the testing phase. The matching confidence is evaluated using our new volume matching. Our face authentication system is verified with extensive experiments in the XM2VTS database.

**Keywords:** Face Authentication, Face Descriptor, Convolutional Neural Networks, Volume Deep Face

### 1. Introduction

Image representation plays a crucial role in the face authentication system [1, 2, 3, 4, 7, 13, 14, 15, 18]. They provide rich information in which the transformed space is more distinctive than the image space. Due to the similar structure of face images, it is very hard to recognize the faces directly on image space. In addition, the gray level of human face images normally changes from time to time, resulting in a degradation of the performance of the face authentication system. Eigenfaces and Fisherfaces are among the most common methods for face representation in which the dimensionality of the data is reduced while minimizing the variance within the class [1, 2]. The local binary pattern (LBP) was also a popular method proposed for face representation with the ability to enhance the gray level of face images [3]. LBP measures the difference between the central pixel and its neighbors, and then encodes the changes as a binary number. LBP enhances robustness against illumination changes. Eigenfaces, Fisherfaces, and LBP can be categorized as hand-crafted feature extraction methods. The core strength of the convolutional neural network lies in its ability to extract useful features using the convolution operator [6, 7, 9, 10, 11, 12, 16, 17, 20]. CNN provides a dual mode that works both as a classifier and as a feature extractor. In this paper, we propose a novel face representation, termed a volume-deep face (VDF), which is constructed from multiple 2D face descriptors. We also propose a robust method for face and eye detec-

tions. The original contributions of this paper are: 1) It generates the face representation in both 2D and 3D forms, which provides the capability of constructing distinctive features. 2) It proposes robust 3D volume matching for face authentication.

### 2. The Proposed Face Authentication System

#### 2.1 Overview of the System

An overview of our proposed face authentication system is shown in Fig. 1. In the training phase, we first detect the face and eye using two-stage face and eye detections, 2SYOLO. The eye positions are then used to align the face in which the degree of rotation of the two eyes is zero. The face is then cropped to the normalized size  $130 \times 200$ . The aligned and cropped faces are then fed to the CNNs where VDF is generated. In the test phase, a similar 2SYOLO process is performed to obtain an aligned and cropped face. 2DFD is created from a single image for testing. The volume matching is used to measure the similarity between VDF and 2DFD, resulting in true or false identity. This unified pipeline ensures that the same preprocessing steps are applied consistently across both training and testing, reducing variability and increasing robustness. Furthermore, the framework can be efficiently extended to real-time applications, making it suitable for deployment in practical security systems.

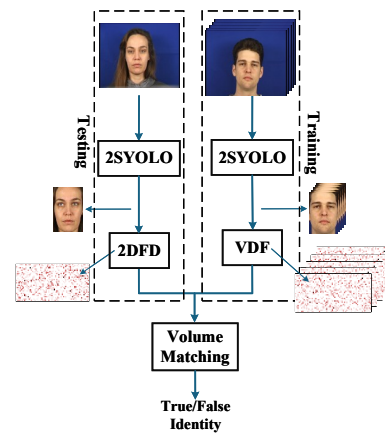


Fig. 1. The framework has two phases: (1) training, where a 3D VDF tensor is constructed; and (2) testing, where a 2DFD is generated and matched against the stored VDF for face authentication.

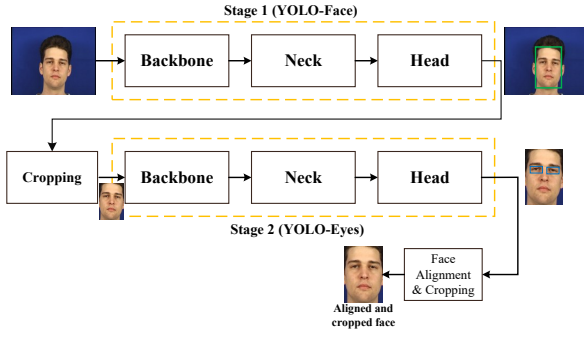


Fig. 2. Architecture of our two-stage YOLO-based approach for face and eye detection.

## 2.2 Face Detection and Alignment

The performance of a face authentication system can be improved by precisely locating and aligning facial components [19]. In particular, aligning human eyes across all face images ensures that key facial components are consistently positioned, enabling more reliable and accurate face verification. We introduce 2SYOLO, a two-stage YOLO-based approach for face and eye detection. YOLO provides real-time speed, easy fine-tuning and training, and clean two-stage (face→eyes) alignment; the pipeline does not depend on one specific face detector. Therefore, our method is not limited to YOLO; other detectors such as RetinaFace [23] or MTCNN [24] can also be integrated into the same pipeline. In the first stage, the face bounding boxes are located using YOLO-Face. In the second stage, eyes are localized within each detected face bounding box using YOLO-Eyes. Based on the detected eye positions, the face is then horizontally aligned to ensure that facial components are consistently positioned across all entire dataset. YOLO-Face is trained on face images, while YOLO-Eyes is trained specifically on eye region images to accurately detect eye positions within the detected face bounding boxes. Our two-stage YOLO-based approach is illustrated in Fig. 2 which has the following parts [8]: a backbone, a neck and a head. The backbone is responsible for feature extraction and handles most of the network computation. The neck aggregates features from various stages of the backbone, enhancing the network's ability to detect objects at multiple scale. Finally, the head generates the model's predictions based on the combined features.

## 2.3 2D Face Descriptor with Deep Learning

Convolutional neural networks (CNNs) were developed for more than two decades beginning with the introduction of LeNet architecture [10, 11, 12, 18, 13]. CNNs are a special type of multi-layer neural networks in which convolution operations are used to extract distinctive features from input images. Therefore, raw images are not fed directly from one layer to the next, but are transformed by convolution operations before being forward. Extensive

empirical evidence shows that deeper CNNs architectures produced greater performance [10, 11]. CNNs are composed of multiple layers including feature extractors (e.g., convolution operation), dimensionality reduction (e.g., max pooling) and classification layers (e.g., fully connected and softmax layers) [6]. The cropped and aligned face images resulting from the previous section are color images with the size of  $W \times H$ , where  $W = 130$  and  $H = 200$  in our implementation. Therefore, input image is represented as a tensor  $I \in \mathbb{Z}^{W \times H \times C}$  and is fed to the CNNs with shape  $W \times H \times C$ , where  $W$  and  $H$  are width and height of an image and  $C$  is the number of channels, i.e.  $C = 3$  for RGB images. Normally, the outputs of CNNs are produced by softmax layer. One can define CNNs as a  $K$  class probability distribution where the output is maximized at index  $k^{th}$  if and only if  $k^{th}$  class is identical to the identified object. Let us define by  $v_j$  the target of CNNs classification where  $j$  is the index of class  $j^{th}$ . By minimizing [10, 11, 12]

$$\|\sigma(\mathbf{y})_j - v_j\|; \forall j \in K, \quad (1)$$

for  $K$  class problem, the CNNs is converged.  $\sigma(\mathbf{y})_j$  is a softmax probability and can be described by [11]

$$\sigma(\mathbf{y})_j = \frac{e^{y_j}}{\sum_{k=1}^K e^{y_k}}; j = 1, \dots, K. \quad (2)$$

$y_j$  is a fully connected layer where all neurons from the previous layer are connected with weights to the next layer, hence the term fully connected (FC). Let us define by  $FC(x)$ , the output of the fully connected layer with [10]

$$FC(x) = a(Wx + b) \in \mathbb{R}^m. \quad (3)$$

where  $a(\cdot)$  is an activation function.  $FC(x)$  is fully connected with  $n$  inputs from the previous neuron layer and  $m$  output dimensions, computed using a trainable weight matrix  $W \in \mathbb{R}^{m \times n}$  and bias vector  $b \in \mathbb{R}^m$ . Typically, the activation function  $a(\cdot)$  used in CNNs is a nonlinear rectified linear unit (ReLU) which is defined by [12]

$$a(u) = \max(0, u), u \in \mathbb{R}, \quad (4)$$

It is easy to prove that  $a'(u) = 1$  for  $u > 0$  and that  $a'(u) = 0$  for  $u < 0$ . ReLU has an advantage over traditionally sigmoid function in that it promotes faster convergence. This is primarily because the sigmoid function may lead to the vanishing gradient problem when the input  $u$  is far from zero [11], i.e.,  $u \rightarrow \infty$  or  $u \rightarrow -\infty$ . The pooling layer is used to reduce the dimensionality of feature maps which are generated from conv and ReLU layers. In this paper, max pooling is used for selecting maximum value within each spatial neighbor with kernel size of 2 and the stride of 2. This leads to downsampling the feature maps, which helps CNNs facilitates faster convergence in training procedures.

In the face authentication system, we need to create a face template and store it in ID card or in database. The face

template (face descriptor) will be used in face matching process where the matching may be 1-to-1 or 1-to-many. However, the output of CNNs produced by softmax layer is not appropriate for this case and can not be used as a face template. We proposed here to use the output from several FCs layer as a 2D face descriptor (2DFD). Small feature representation can be used for faster classification with the risk of weak identity. We need templates for face authentication where the templates of the same class must be similar while those of different individuals remain distinct. Based on our experiments, more and more FCs layers can be used to construct a 2DFD and helps us to improve the overall accuracy of face authentication system. The 2DFD serves as a face template for face matching (face authentication in our case). 2DFD is a general face template usable for 1:1 verification and 1:N identification with cosine or  $\ell_2$  matching; authentication is however the scope of our paper.

The architecture of our 2DFD, as shown in Fig. 3, is based on fully connected layers from CNNs where conv is the convolution operator, ReLU is a rectified linear unit and FC is the fully connected layer. The conv, ReLU and Max Pooling are repeated for several layers to produce more and more salient features before passing them to the next layer. The 2DFD is constructed by stacking the outputs of 16 FC layers. Each FC output (128 units) is split into two 64-element halves and reshaped into a  $2 \times 64$  slice; concatenating the 16 slices column-wise yields a compact  $32 \times 64$  2D face template with enhanced discriminability. From a deployment perspective, the  $32 \times 64$  template (2,048 elements) can be stored as 8-bit integers in roughly 2 KB per identity, enabling low-latency matching on resource-constrained devices. Moreover, the layer-wise stacking design aligns naturally with our Volume Deep Face (VDF) representation: multiple 2DFDs computed from different enrollment images can be concatenated along a third (depth) axis to form a 3D tensor, preserving both within-layer structure and across-layer diversity. This structural compatibility simplifies the transition from single-image (2DFD) verification to multi-image (VDF) aggregation without changing the downstream matching logic.

Let us define by  $\lambda_i^k(r, c) \in \mathbb{R}$  the 2D face template for class  $k^{th}$  of image  $i^{th}$  where  $r \in R$  and  $c \in C$ , i.e.  $R = 32$  and  $C = 64$ , respectively. In order to standardize the face template, we normalize it by  $\hat{\lambda}_i^k(r, c) = \frac{\lambda_i^k(r, c) - \min}{\max - \min}$  so that  $\hat{\lambda}_i^k(r, c) \in [0, 1]$ . For illustration proposes,  $\hat{\lambda}_i^k(r, c)$  should be scaled to gray level range by subtracting and multiplying it by 255, i.e.,  $\hat{\lambda}_i^k(r, c) := 255 - (\hat{\lambda}_i^k(r, c) * 255)$ . One can observe the difference between 1D (the  $1 \times 128$  1D descriptor) and 2D (the  $2 \times 64$  and  $32 \times 64$ ) descriptors that the 2D descriptor exhibit more structured patterns compared to the 1D vector. Hence, the between-class variances can be maximized while maintaining the within-class variances as low as possible. This face template  $\hat{\lambda}_i^k$  will be used to construct the volume deep face in the next section.

## 2.4 Volume Deep Face

In this section, we propose the volume deep face (VDF) which is constructed from 2DFD  $\hat{\lambda}_i^k(r, c)$ . Let us suppose that we generate multiple 2DFDs  $\hat{\lambda}_i^k(r, c), i = 1, \dots, N$  where  $N$  is the number of training images for class  $k^{th}$ . The VDF is shown in Fig. 4.

This VDF is a 3D representation inheriting from multiple 2DFDs. Therefore, more distinctive features can be achieved. The VDF can be regarded as a tensor-valued function and formulated as

$$\Lambda = \hat{\lambda}_1^k \oplus \hat{\lambda}_2^k \oplus \dots \oplus \hat{\lambda}_i^k \oplus \dots \oplus \hat{\lambda}_N^k, 1 \leq i \leq N \quad (5)$$

where  $\oplus$  is a concatenate operation. Therefore, each  $\hat{\lambda}_i^k$  is concatenated in the direction of index  $i$ .  $N$  is the number of 2DFD generated from FC layers of CNNs. In summary, given input images from class  $k^{th}$ , create multiple 2DFDs  $\hat{\lambda}_i^k$  from FC layer, then concatenating all features to form the VDF  $\Lambda$ . It should be noted that the VDF is generated from CNNs where the learnable weight  $W$  has been trained with error minimization. The details of generating VDF are summarized in the algorithm 1.

---

### Algorithm 1 Construction of volume deep face (VDF).

---

Let  $I_i^k \in \mathbb{Z}^{W \times H \times C}$  be the training image  $i^{th}$  of class  $k^{th}$ .  
**repeat** for each  $k \in K, i = 1, \dots, N$   
     compute conv and ReLU from  $I_i^k$   
     apply max pooling to ReLU  
      $FC(x)^j \leftarrow a(Wx + b) \in \mathbb{R}^m, 1 \leq j \leq 16$   
     stack  $FC(x)^j \forall j$  to generate  $\hat{\lambda}_i^k$   
     construct  $\Lambda^k$  for each class  $k^{th}$   
**until** all classes  $K$  have been generated

---

## 2.5 Volume Matching

Let  $\Lambda_A$  be the VDF of images formulated by the algorithm 1 and let  $\hat{\lambda}^t(r, c)$  denote the 2DFD of a test image, as described in the previous section. As our VDF can be generated from multiple images in the training set, we assume that during the testing phase, only a single image may be used per test. In such cases, images in the training set can be used to construct the VDFs, while in the test phase only 2DFD will be generated. We measure the similarity between the VDF and 2DFD by

$$\gamma = \frac{1}{L} \sum_{i=1}^N \sum_{r=1}^R \sum_{c=1}^C \left\| \hat{\lambda}_i^A(r, c) - \hat{\lambda}^t(r, c) \right\|_2, \forall \hat{\lambda}_i^A \in \Lambda_A \quad (6)$$

where  $L = R \times C \times N$  is the total number of elements in the VDF. If  $\hat{\lambda}^t(r, c)$  is identical to  $\hat{\lambda}_i^A(r, c)$ , then  $\gamma \rightarrow 0$  otherwise  $\gamma \rightarrow \infty$

## 3. Experimental Results

The images used in this article were collected from the standard XM2VTS face database [5] as shown in Fig. 5. XM2VTS targets verification with a standardized protocol and strong baselines; identification-oriented sets (e.g.

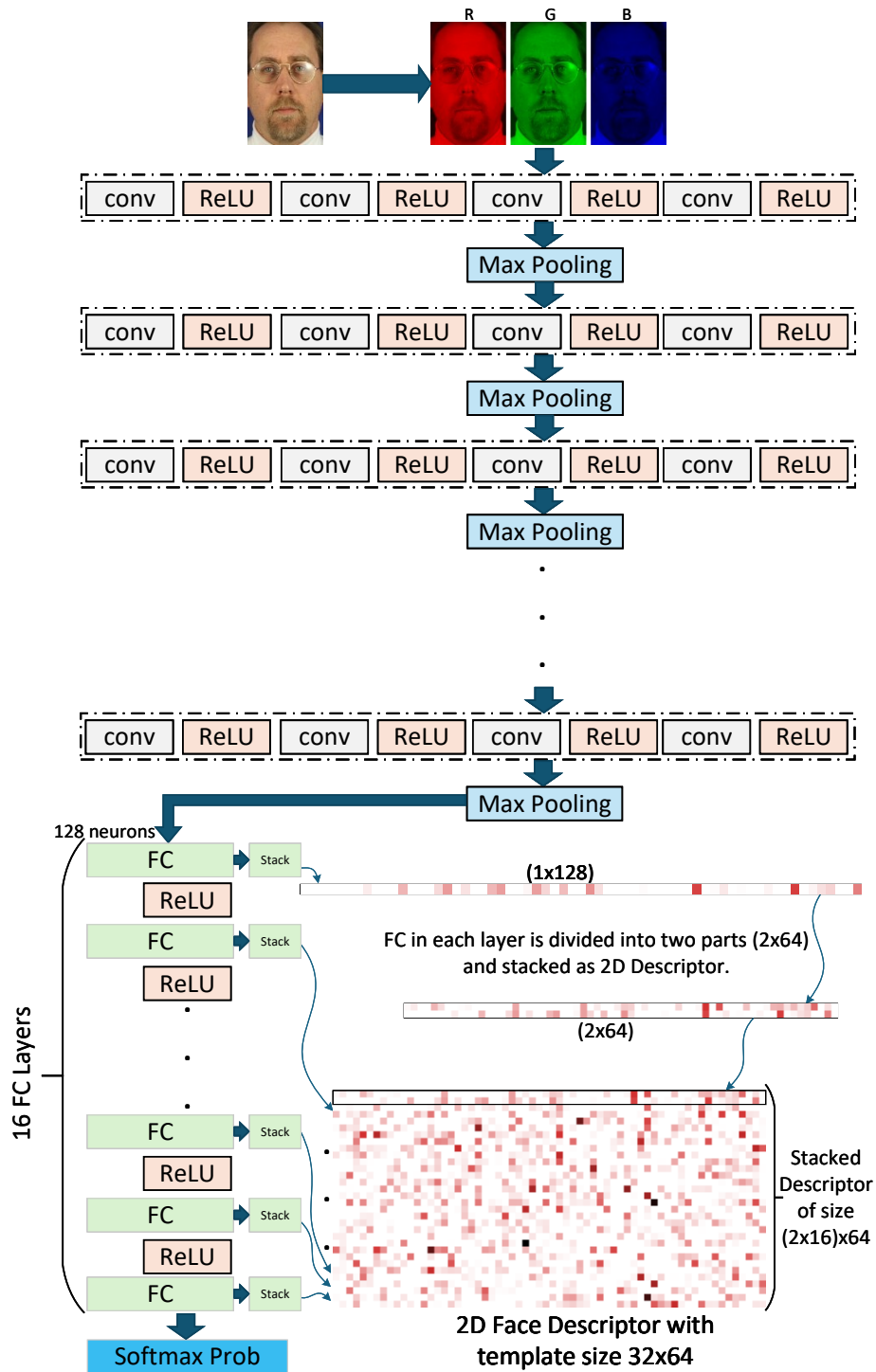


Fig. 3. The architecture of our deep CNNs for generating 2D face descriptor with compact template size  $32 \times 64$ . For visualization only,  $\hat{\lambda}$  may be linearly mapped to 8-bit grayscale (and optionally pseudo-colored); the color appearing in figure illustrations is artificially added to aid interpretation and is not used during matching.

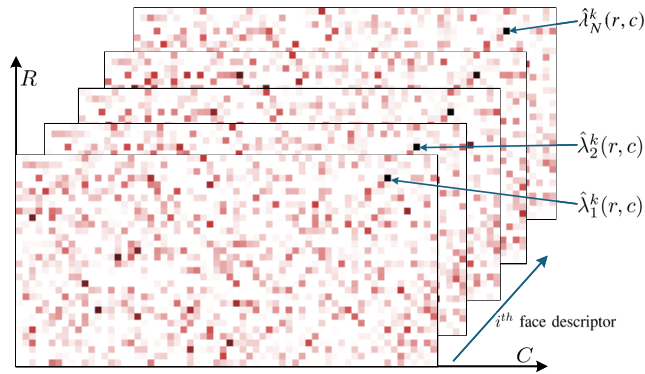


Fig. 4. The volume deep face constructed from multiple 2DFDs  $\hat{\lambda}_i^k(r, c)$  for  $i^{\text{th}}$  face descriptor at  $k^{\text{th}}$  class with size  $R \times C$ . By concatenating along the descriptor index  $i$ , these slices form a 3D tensor  $\Lambda^k \in \mathbb{R}^{R \times C \times N}$  that captures intra-class variability and enables robust volume matching.

VGGFace2 [21]) are outside the current scope, with cross-dataset verification planned. The database has 295 subjects, each of which was captured by 8 shots with 4 distinct sessions during 4 months resulting in a total of 2,360 images. The database was divided into three sets: training set, evaluation set, and test set [5].

The corresponding 2DFD for each cropped and aligned face image were generated and depicted in Fig. 6. Columns (a), (c) and (e) are examples of aligned and cropped face images, while columns (b), (d) and (f) represent the corresponding 2DFDs of (a), (c) and (e), respectively. It was cleared that the 2DFDs of the same individual are highly similar in which the within-class variance is minimized. Moreover, the 2DFDs of different classes show significant differences, which helps maximize the between class variance. We can conclude that the micro pattern on each element  $(r, c)$  of 2DFD  $\hat{\lambda}$  represents the 2D structure generated from the face image exhibiting more discriminative features for classification. All experiments were conducted on a workstation equipped with an RTX 5070Ti Super GPU (16 GB VRAM), an AMD Ryzen 7 9800X3D CPU (4.70 GHz, 8C/16T), and 64 GB RAM. Using CUDA programming for GPU acceleration, the proposed system achieves computation times below 16 ms per frame, enabling real-time face authentication.

### 3.1 Face Verification

The performance of our proposed method was evaluated based on the protocol described in [5]. False acceptance (FA) occurs when an imposter (false identity) is incorrectly accepted as a client (true identity). In contrast, false rejection (FR) occurs when the true identity has been rejected. The FA and FR were measured by  $FA = \frac{EI}{I} \times 100$  and  $FR = \frac{EC}{C} \times 100$ , where  $EI$  and  $EC$  represent the number of imposter acceptances and client rejections, respectively.  $I$  and  $C$  denote the total number of imposter and clients claims in the system. In this paper,  $I = 112000$  and  $C = 400$ . Table 1 shows the error rates of

our proposed method in comparison to the other approaches. We obtain the total error rate (TER) with 0.01161% which is the lowest error rate, where  $TER = FA + FR$ . Under the same XM2VTS protocol, prior baselines report TER  $\sim 1.48\text{--}11.36\%$ ; weighted DeepFace reaches 0.01429%, while ours attains 0.01161% ( $FR = 0$ ), an  $\approx 18.8\%$  reduction. It is worth noting that, with a very low error rate of our proposed method, implementing a face authentication system for airport check-in process is feasible.

Table 1. Error Rates on XM2VTS Database. (The bold values indicate the best performance.)

| Methods [5]           | Test Set       |        |                |
|-----------------------|----------------|--------|----------------|
|                       | FA (%)         | FR (%) | TER (%)        |
| UniS-ICPR2000         | 2.30           | 2.50   | 4.80           |
| IDIAP-Marcel          | 1.748          | 2.0    | 3.75           |
| IDIAP-Cardinaux       | 1.84           | 1.50   | 3.34           |
| MUT-UniS-STT          | 0.97           | 0.50   | 1.47           |
| UCL                   | 1.71           | 1.50   | 3.21           |
| TB                    | 5.61           | 5.75   | 11.36          |
| UniS-ECOC             | 0.86           | 0.75   | 1.61           |
| UniS-NC               | 0.48           | 1.00   | 1.48           |
| weighted DeepFace [7] | 0.01429        | 0.0    | 0.01429        |
| Our Proposed Method   | <b>0.01161</b> | 0.0    | <b>0.01161</b> |

### 4. Discussion and Future Work

The proposed face authentication framework demonstrates strong performance on the XM2VTS benchmark by integrating two major components: a compact 2D Face Descriptor (2DFD) and the higher-dimensional Volume Deep Face (VDF). While the VDF effectively aggregates discriminative features from multiple images, the 2DFD provides efficiency for single-image scenarios. Together, these representations enable a balance between accuracy and computational cost.

One important observation from our experiments is



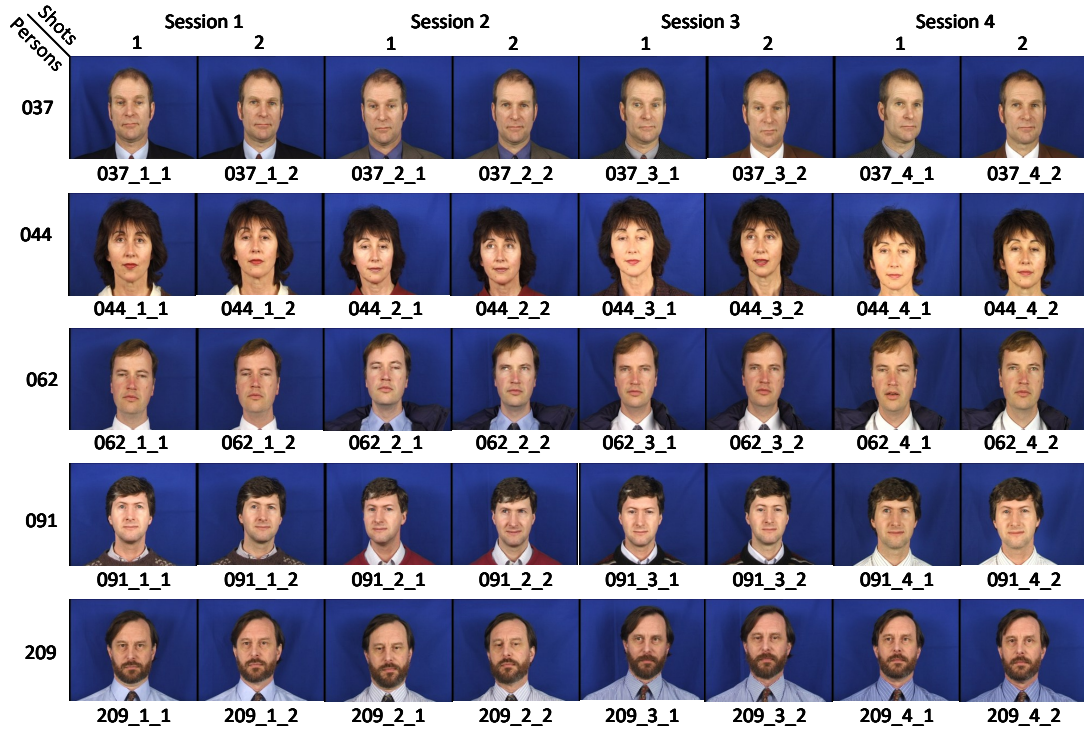


Fig. 5. Examples face images obtained from the XM2VTS database used in the face authentication experiments. For fair comparison, all methods in Table 1 were evaluated on identical XM2VTS inputs under the standard XM2VTS protocol and evaluation procedure [5].

that the use of eye-based alignment through the two-stage YOLO (2SYOLO) detector significantly reduces intra-class variation caused by pose and rotation. This step ensures that both the 2DFD and VDF descriptors are constructed on normalized facial regions, which improves stability and robustness in real-world conditions. Furthermore, the proposed volume matching strategy consistently produced low error rates even when the test set included natural variations. Although promising, there are several open research directions worth pursuing. First, extending the training process with large-scale, unconstrained face datasets such as VGGFace2 [21] or MS-Celeb-1M [22] may help generalize the learned representations to more diverse demographics. Second, the framework could be adapted for video-based face authentication where temporal information across frames may provide additional discriminative cues beyond static images. Third, incorporating lightweight neural architectures (e.g., MobileNet or ShuffleNet) would make the system deployable on edge devices, which is critical for real-time security applications such as airport check-in and access control systems. In summary, the proposed system establishes a strong baseline for compact and discriminative face representations. By addressing scalability, real-time deployment, and robustness to uncontrolled conditions, future research can further advance the applicability of VDF-based face authentication in practical biometric systems.

## 5. Conclusions

We have proposed a novel face representation, the Volume Deep Face (VDF), for a face authentication system. The VDF is composed of multiple concatenated 2D face descriptors (2DFDs) in the tensor form allowing us to generate more discriminative features of face representations. The 2DFDs were derived from multiple fully connected layers of CNNs that has been pre-trained prior to the VDF construction process. Additionally, input of face images were aligned and cropped with our proposed two-stage YOLO based face and eye detection which helps us increasing the performance of the overall system. We obtained the total error rate with 0.01161% which is the lowest error reported to date on the XM2VTS database. Beyond the current evaluation, the proposed framework is also scalable and can be adapted for large-scale applications, such as airport check-in, smart surveillance, and secure access control, where both accuracy and efficiency are critical. Moreover, because the VDF can naturally extend to video sequences and multimodal biometric fusion, the framework opens new opportunities for developing real-time, privacy-aware, and user-friendly authentication systems that meet the demands of future intelligent security infrastructures. Future work will examine cross-database generalization under illumination and pose variations and include comprehensive ablations on the number of FC slices and volume-matching criteria.

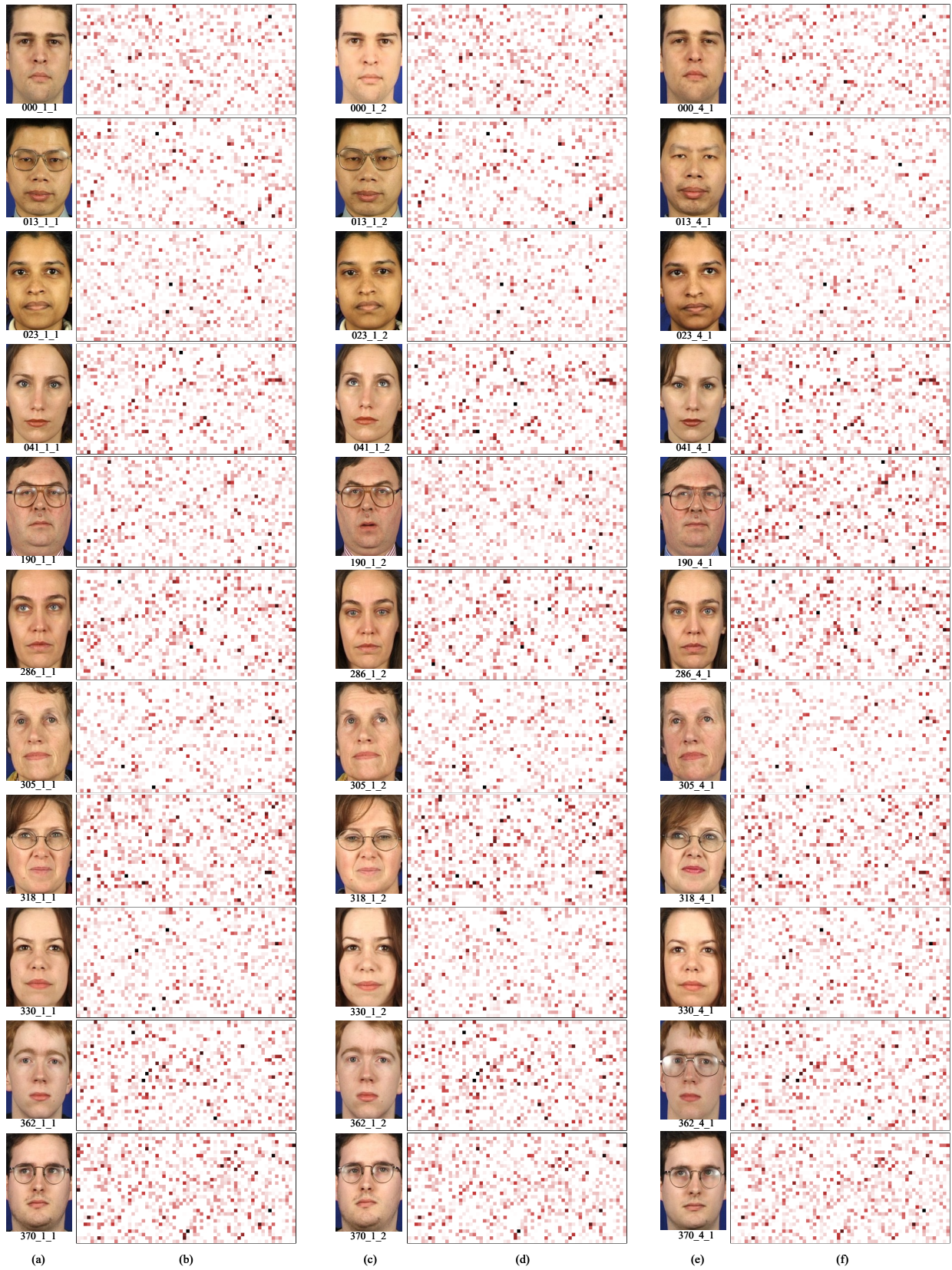


Fig. 6. The 2DFDs generated from the XM2VTS database: (a) face images from session 1 shot 1; (b) the corresponding 2DFDs of (a); (c) face images from session 1 shot 2; (d) the corresponding 2DFDs of (c); (e) face images from session 4 shot 1; (f) the corresponding 2DFDs of (e).

## References

- [1] D. Ying., "Exploring PCA-based feature representations of image pixels via CNN to enhance food image segmentation," 10.48550/arXiv.2411.01469., 2024.
- [2] P. N. Belhumeur, J. P. Hespanha and D. J. Kriegman, "Eigenfaces vs. Fisherfaces: recognition using class specific linear projection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 19, no. 7, pp. 711-720, July 1997.
- [3] T. Ahonen, A. Hadid and M. Pietikainen, "Face Description with Local Binary Patterns: Application to Face Recognition," in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 28, no. 12, pp. 2037-2041, Dec. 2006.
- [4] Y. Taigman, M. Yang, M. Ranzato, L. Wolf, "DeepFace: Closing the Gap to Human-Level Performance in Face Verification," *Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1701-1708, 2014.
- [5] K. Messer, J. Kittler, M. Sadeghi, S. Marcel, C. Marcel, S. Bengio, F. Cardinaux, C. Sanderson, J. Czyz, L. Vandendorpe, S. Srisuk, M. Petrou, W. Kurutach, A. Kadyrov, R. Paredes, B. Kepenekci, F. B. Tek, G. B. Akar, F. Deravi and N. Mavity, "Face Verification Competition on the XM2VTS database," in *Proc. of 4th Int. Conf. Audio and Video Based Biometric Person Authentication, Lecture Notes in Computer Science*, Vol. 2688, Springer-Verlag, pp. 964-974, 2003.
- [6] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, T. Darrell, "Caffe: Convolutional Architecture for Fast Feature Embedding," *arXiv preprint arXiv:1408.5093*, 2014.
- [7] S. Srisuk, et.al., "Robust Face Recognition based on weighted DeepFace," 2017 International Electrical Engineering Congress (IEEECON), Pattaya, Thailand, 2017, pp. 1-4, doi: 10.1109/IEEECON.2017.8075885.
- [8] C. Li, et. al., "YOLOv6: A Single-Stage Object Detection Framework for Industrial Applications," *Technical Reports*, 2022, <https://arxiv.org/abs/2209.02976>.
- [9] J. Yangqing, et. al, "Caffe: Convolutional Architecture for Fast Feature Embedding," *arXiv preprint arXiv:1408.5093*.
- [10] Y. LeCun, et. al, "Handwritten Digit Recognition with a Back-Propagation Network," *Advances in Neural Information Processing Systems*, vol. 2, 1989.
- [11] I. Goodfellow, Y. Bengio and A. Courville, "Deep Learning," MIT Press, 2016.
- [12] N. Buduma, N. Buduma and J. Papa, "Fundamentals of Deep Learning," "O'Reilly Media, Inc., 2022.
- [13] F. Schroff, D. Kalenichenko and J. Philbin, "FaceNet: A Unified Embedding for Face Recognition and Clustering," *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, pp. 815-823, 2015.
- [14] J. Deng, J. Guo, N. Xue and S. Zafeiriou, "ArcFace: Additive Angular Margin Loss for Deep Face Recognition," *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, pp. 4690-4699, 2019.
- [15] O. M. Parkhi, A. Vedaldi and A. Zisserman, "Deep Face Recognition," *Proc. British Machine Vision Conf. (BMVC)*, 2015.
- [16] I. Kemelmacher-Shlizerman and S. M. Seitz, "Face Reconstruction in the Wild," *Proc. IEEE Int. Conf. Computer Vision (ICCV)*, pp. 1746-1753, 2011.
- [17] K. He, X. Zhang, S. Ren and J. Sun, "Deep Residual Learning for Image Recognition," *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, pp. 770-778, 2016.
- [18] A. Krizhevsky, I. Sutskever and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," *Adv. Neural Information Processing Systems (NeurIPS)*, pp. 1097-1105, 2012.
- [19] J. Redmon, S. Divvala, R. Girshick and A. Farhadi, "You Only Look Once: Unified, Real-Time Object Detection," *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, pp. 779-788, 2016.
- [20] H. Wang, Y. Wang, Z. Zhou, X. Ji, D. Gong, J. Zhou and Z. Li, "CosFace: Large Margin Cosine Loss for Deep Face Recognition," *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, pp. 5265-5274, 2018.
- [21] Q. Cao, L. Shen, W. Xie, O. M. Parkhi and A. Zisserman, "VGGFace2: A Dataset for Recognising Faces across Pose and Age," *Proc. IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, pp. 67-74, 2018.
- [22] Y. Guo, L. Zhang, Y. Hu, X. He and J. Gao, "MS-Celeb-1M: A Dataset and Benchmark for Large-Scale Face Recognition," *Proc. European Conf. Computer Vision (ECCV)*, pp. 87-102, 2016.
- [23] J. Deng, J. Guo, Y. Zhou, J. Yu, I. Kotsia and S. Zafeiriou, "RetinaFace: Single-stage Dense Face Localisation in the Wild," in *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5202-5211, 2020.
- [24] K. Zhang, Z. Zhang, Z. Li and Y. Qiao, "Joint Face Detection and Alignment Using Multitask Cascaded Convolutional Networks," in *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1499-1503, 2016.