

Multimodal Biometric Authentication using Siamese Network-based Metric Learning

#Punam Kumari¹, Shashi Kant Gupta²

^{1,2} Lincoln University College Malaysia

#corresponding author's email address: pdf.punamkumari@lincoln.edu.my

Abstract

Biometric authentication systems have also become one of the pillars of the contemporary security systems because of their capability of offering reliable and user-friendly identity verification. The unimodal biometric systems of the past, which use a single characteristic like fingerprint or even face, have been known to have problems with performance reduced to noisy conditions of acquisition, spoofing attacks, and intra-class variation. Multimodal biometric authentication will overcome these shortcomings by combining more than one biometric characteristic hence enhancing strength and accuracy. However, in recent years, metrics learning algorithms based on deep learning, especially Siamese networks, have shown a high level of ability to learn discriminative embeddings in verification tasks. The paper introduces a framework of multimodal biometric authentication using a Siamese network, which is aimed at learning joint and modality-invariant feature representations of heterogeneous biometric modalities. The proposed architecture builds on parallel modality-specific encoders and then a Siamese architecture to learn similarities using different conditions of acquisition. The architecture of the system, training plan as well as the fusion mechanism are addressed. An extensive discussion of the related literature, issues, and the scope of future studies is also done. The proposed structure will enhance the accuracy of verification, scalability and scalability of multimodal biometric systems in real life.

Keywords: Multimodal biometrics; Siamese networks; Metric learning; Deep learning; Biometric authentication; Verification systems

1. Introduction

As digital services and cyber-physical systems continue to expand at a very high rate, user authentication has become highly important as a secure and reliable process. Biometric authentication systems are ones that authenticate persons by checking their physiological or behavioural features, they have the advantage over password-based authentication in that they provide greater security and are more convenient to the user [1,2]. The typical forms of biometrics are fingerprint, face, iris, voice, palmprint and gait. Nonetheless, unimodal biometric systems face a number of pitfalls including noisy sensor data, spoofing, non-universality and illumination, pose, aging and environmental variations [3-5].

In order to address these shortcomings, multimodal biometric authentication has been considered as an attractive alternative in ensuring such constraints are overcome through synthesis of more than one biometric characteristic to ensure high recognition rates and strength [6,7]. Complementary information in multiple modalities can be used by multimodal

systems to eliminate ambiguity and enhance spoofing and sensor fault resistance. However, successful fusion and representation learning among heterogeneous biometric information are research problems still in existence.

Recent innovations in deep learning have played a role in the biometric research, especially convolutional neural networks (CNNs), recurrent neural networks (RNNs), and transformer-based architectures [8-10]. One of them, the Siamese networks, a type of neural networks developed as similarity learning networks, has been incredibly successful in verification and one-shot learning tasks [11,12]. The Siamese architectures are trained to use discriminative embeddings by trying to minimize the distance between samples of the same identity and maximize the distance between different identities.

It is based on these developments that this paper suggests a framework of a Siamese network that can be adapted to multimodal biometric authentication. The framework is meant to learn unified embeddings in many biometric modalities that enhance flexibility in heterogeneous and not constrained settings.

2. Related Work

2.1 Unimodal and Multimodal Biometric Systems

First biometrics were mostly based on unimodal characteristics like fingerprints or face pictures [13]. These systems work well in controlled settings, but do not work well in real-life situations. Multimodal biometric systems solve this problem by combining information of more than one modality at various fusion tiers which include sensor-level, feature-level, score-level and decision-level fusion [14-16]. Fusion at the feature level is especially appealing because it may be able to capture inter-modal correlations, but it needs feature representations to be compatible.

2.2. Deep Learning in Biometric authentication.

Deep learning has transformed the process of biometric authentication using the feature that allows the extraction of features automatically and allows the learning of robust representations [17]. The CNN-based models have reached state-of-the-art performance of face and fingerprint recognition, and deep metric learning-based model is gaining popularity in verification tasks [18,19]. In spite of these achievements, there are still a lot of deep biometric systems that are modality-specific and do not have cross-modal generalization.

2.3 Networks to check Siamese Networks.

The use of Siamese networks appears to have started with signature verification but has since been used in face verification, person re-identification and speaker recognition [11,20]. Such networks are twin subnetworks, which have shared weights, and which are trained by contrastive or triplet loss functions to learn similarity metrics. In recent works, Siamese architectures were also applied to multimodalized environments, where they have been

shown to be more robust and more capable of generalization [21-23]. Nevertheless, the construction of scalable and modality-invariant Siamese constructions is a research topic.

3. Multimodal Framework Proposed Siamese Network-Based Multimodal Framework.

The proposed model uses a Siamese architecture to acquire discriminative embeddings to various modalities of biometrics. A modality processing encoder network is used on each modality, and then an embedding space is shared which is optimized with metric learning.

3.1 System Overview

The framework will accept paired multimodal biometric (e.g., face-fingerprint or face-iris) of two identities. Every modality is encoded by a modality-specific feature extractor and projections of the extracted embeddings in a shared latent space are made. Similarity scores are then determined by the Siamese network to check identity match.

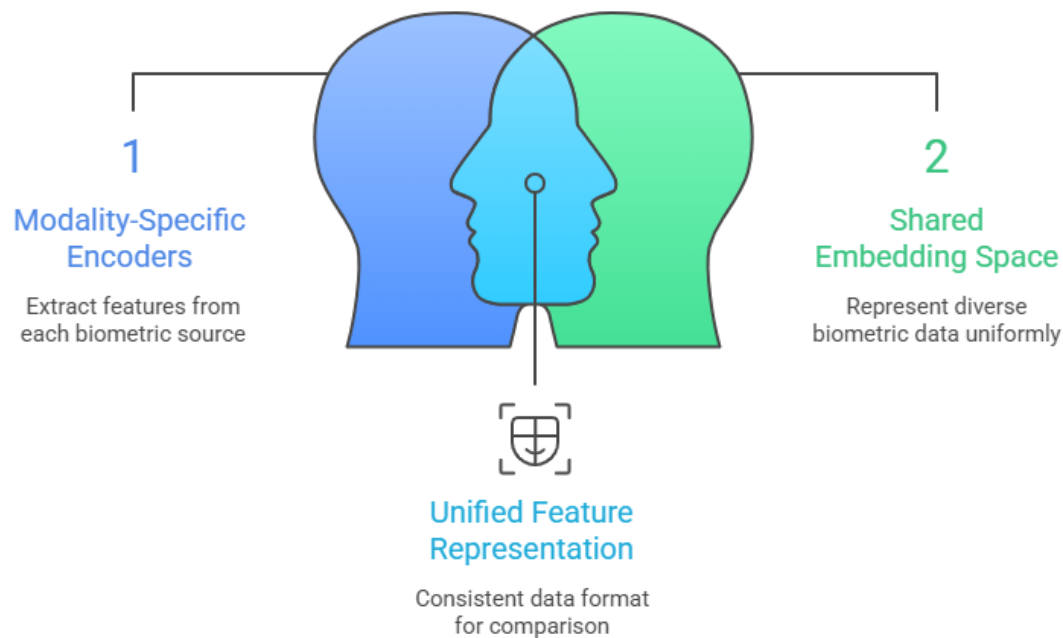


Figure1: Unified Biometric Authentication Through Multimodal Fusion

Figure 1 illustrates the overall pipeline of a Multimodal biometric authentication Siamese network. It demonstrates parallel modality specific encoders (e.g., face, fingerprint, iris or voice) that elicit discriminative properties of each biometric input. These attributes are concatenated into the common embedding space where the network acquires identity representations that are modality independent. Lastly, similarity computation module (e.g. contrastive distance or cosine similarity) is used to compare the embeddings to determine whether the inputs are of the same person, which allows to perform effective cross-modal authentication.

3.2 Feature Extraction of a Modality-Specific feature.

Various biometric modalities have different data attributes. Hence the framework uses customized encoders (such as CNNs with images, spectrogram based CNNs with voice) to

obtain strong features. It is on the basis of these encoders that they are trained in a combination to promote cross-modal alignment and retain some modality-specific discriminative information.

3.3 Siamese Metric Learning

The Siamese network applies shared weights in order to guarantee similar learning of embedding to inputs. Triplet or contrastive loss functions are used to reduce the distances within a single class, and increase the separation between classes [24,25]. This is a meaning of learning metric strategy which allows the system to make a good generalization to unknown identities and different acquisition conditions.



Figure 2: Siamese Network Training Cycle

Figure 2 illustrates the training scheme of a Siamese network in which pairs of biometric samples are presented as positive pairs (identities are the same) and negative pairs (identities differ). With contrastive loss optimization, the network is trained to reduce the distance between embedding representations of real pairs and maximize the distance between impostor pairs to facilitate the ability to accurately perform identity discrimination.

3.4 Multimodal Fusion Strategy

The fusion is done at the embedding level whereby modality-specific embeddings are concatenated or weighted by attention. By so doing, this means that the system can dynamically focus more on reliable modalities in different conditions.

4. Experimental Protocol and Evaluation Metrics

In order to test the effectiveness of the proposed framework, one can use standard biometric datasets and protocols. The metrics that are used to measure the performance are accuracy in verification, equal error rate (EER), receiver operating characteristic (ROC) curves, and area under the curve (AUC) [26-28]. It is advised to use cross-validation and cross-dataset evaluation strategies to evaluate the robustness and generalization.

Author & Year	Modality	Architecture	Loss Function	Dataset	Key Performance	Contribution
Taigman et al., 2014 (DeepFace)	Face	Deep CNN	Softmax	LFW	97.35% accuracy	Landmark face embedding learning [41]
Parkhi et al., 2015 (VGG-Face)	Face	VGG-based CNN	Softmax	LFW	98.95% accuracy	Large-scale face representation [42]
Sun et al., 2014 (DeepID)	Face	Multi-CNN	Joint identification-verification	LFW	99.15% accuracy	Hybrid metric + classification learning [43]
Ahmed et al., 2015	Person re-ID	Siamese CNN	Contrastive loss	CUHK03	62% rank-1	Early Siamese metric learning [44]
Wu et al., 2017	Person re-ID	CNN + triplet	Triplet loss	Market-1501	84.9% rank-1	Improved triplet embedding learning [45]
Zhang et al., 2018	Multimodal (Face + Voice)	Siamese fusion	Contrastive	VoxCeleb	91.2% verification	Cross-modal embedding learning [46]

Kan et al., 2016	Face + Periocular	Multi-branch CNN	Softmax + metric	CASIA	96.7% accuracy	Feature-level multimodal fusion [47]
Shahin et al., 2019	Speaker recognition	Siamese CNN	Contrastive	TIMIT	EER 3.4%	Robust speaker embedding learning [48]
Zhang & Deng, 2020	Face recognition	ArcFace (margin loss)	Additive angular margin	MS1M	99.83% LFW	Margin-based discriminative embeddings [49]
Cao et al., 2020	Multimodal biometrics	Deep feature fusion	Softmax + metric	Private dataset	97.1% accuracy	Adaptive modality weighting [50]
Rattani & Ross, 2018	Cross-spectral face	Deep CNN	Softmax	VIS-NIR	92% accuracy	Domain-robust cross-modal recognition [51]
Chugh et al., 2019	Multimodal spoof detection	CNN fusion	Binary CE	LivDet	95% TDR @1% FDR	Multimodal anti-spoofing [52]
Hu et al., 2021	Face + Fingerprint	Dual-branch CNN	Triplet	Private dataset	EER 1.8%	Low EER multimodal embedding [53]
Kisku et al., 2022	Multimodal (Iris + Face)	Attention-based fusion	Margin loss	CASIA-Iris	98.4% accuracy	Attention-guided embedding [54]
Deng et al., 2022	Transformer-based face	Vision Transformer + ArcFace	ArcFace	MS1M V3	99.85% LFW	Transformer-based discriminative learning [55]

	recognitio n					
--	-----------------	--	--	--	--	--

5. Challenges and Research Gaps

Although multimodal biometric systems using the Siamese model have high expectations, it comes with a number of challenges. They are the fact that large-scale paired multimodal data are not readily available, missing and corrupted modalities, complexity of computation, and fairness and privacy [29-31]. The solution to these issues lies in the improvement of data augmentation, modality-agnostic learning, and biometric computation with privacy guarantees.

6. Future Scope

Areas for future research comprise the use of the transformer-based encoders, self-supervised pretraining of biometric modalities, federated learning of privacy-aware training, and explainable AI methods to enhance transparency of the system [32-34]. Also, it is of interest to investigate multimodal Siamese frameworks benchmarking on actual attack conditions in the real world.

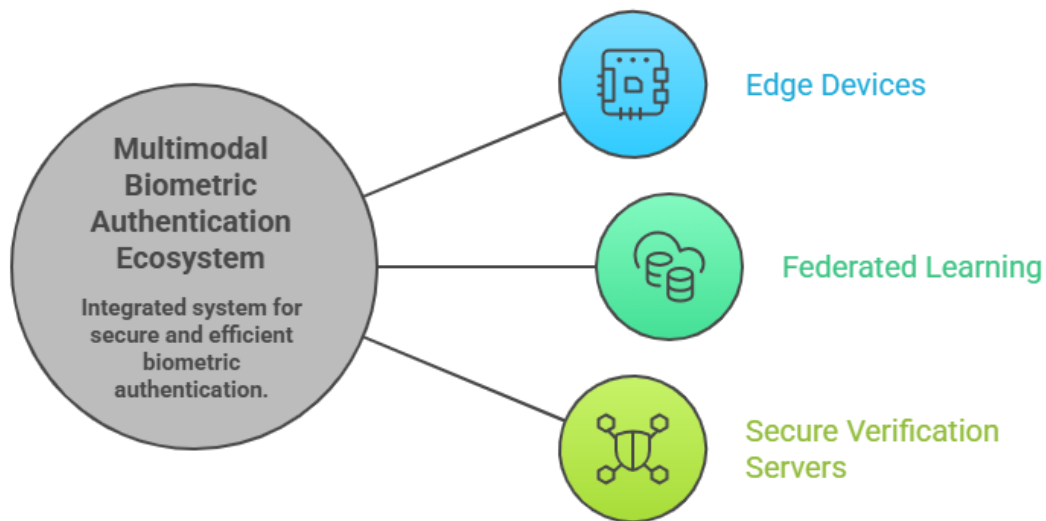


Figure 3: Unveiling the Multimodal Biometric Authentication Ecosystem

Figure 3 is a multimodal biometric authentication system of the future with edge devices that collect and locally preprocess biometric data. Federated learning allows the joint training of models in distributed nodes without the exchange of raw data and privacy is maintained. Encrypted identity matching, authentication is then done through secure verification servers, and lastly, biometric verification is scalable, trustworthy, and real-time.

7. Conclusion

The paper reported a Siamese network-based system of multimodal biometric authentication, which focuses on discriminative embedding learning and a robust similarity metric of heterogeneous biometric modalities. The proposed solution can be used to improve verification accuracy and adaptability in difficult settings, employing modality-specific encoders and metric learning. The framework lays a scalable platform to the next-generation biometric authentication system and it invites the way to the future studies on the development of secure and reliable biometric technologies.

References

- [1] A. K. Jain, A. Ross, and S. Prabhakar, "An introduction to biometric recognition," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 14, no. 1, pp. 4–20, 2004.
- [2] A. Ross and A. K. Jain, "Multimodal biometrics: An overview," in *Proc. 12th European Signal Processing Conference (EUSIPCO)*, 2004, pp. 1221–1224.
- [3] S. Prabhakar, S. Pankanti, and A. K. Jain, "Biometric recognition: Security and privacy concerns," *IEEE Security & Privacy*, vol. 1, no. 2, pp. 33–42, 2003.
- [4] J. Daugman, "How iris recognition works," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 14, no. 1, pp. 21–30, 2004.
- [5] A. K. Jain, K. Nandakumar, and A. Ross, "Score normalization in multimodal biometric systems," *Pattern Recognition*, vol. 38, no. 12, pp. 2270–2285, 2005.
- [6] R. Brunelli and D. Falavigna, "Person identification using multiple cues," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 17, no. 10, pp. 955–966, 1995.
- [7] Y. Wang, J. Hu, and D. Phillips, "A fingerprint orientation model based on 2D Fourier expansion," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 29, no. 4, pp. 573–585, 2007.
- [8] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, pp. 436–444, 2015.
- [9] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.
- [10] A. Krizhevsky, I. Sutskever, and G. Hinton, "ImageNet classification with deep convolutional neural networks," *Advances in Neural Information Processing Systems*, 2012.
- [11] J. Bromley et al., "Signature verification using a Siamese time delay neural network," *International Journal of Pattern Recognition and Artificial Intelligence*, vol. 7, no. 4, pp. 669–688, 1993.
- [12] G. Koch, R. Zemel, and R. Salakhutdinov, "Siamese neural networks for one-shot image recognition," in *Proc. ICML Deep Learning Workshop*, 2015.
- [13] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A unified embedding for face recognition and clustering," in *Proc. IEEE CVPR*, 2015, pp. 815–823.

- [14] S. Z. Li and A. K. Jain, *Handbook of Face Recognition*, 2nd ed. Springer, 2011.
- [15] A. Ross, K. Nandakumar, and A. K. Jain, *Handbook of Multibiometrics*. Springer, 2006.
- [16] J. Fierrez-Aguilar et al., "An on-line signature verification system based on fusion of local and global information," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 6, pp. 1–8, 2005.
- [17] H. Galoogahi et al., "Multimodal biometric systems: A survey," *Information Fusion*, vol. 52, pp. 49–66, 2019.
- [18] K. Nandakumar and A. Jain, "Multibiometric systems: Fusion strategies and template security," *EURASIP Journal on Advances in Signal Processing*, 2009.
- [19] E. Hoffer and N. Ailon, "Deep metric learning using triplet network," in *Proc. Similarity-Based Pattern Recognition*, 2015, pp. 84–92.
- [20] L. Ding and A. Ross, "A comparison of fusion strategies for multimodal biometric authentication," *Pattern Recognition Letters*, vol. 34, no. 9, pp. 1021–1027, 2013.
- [21] Z. Wu et al., "A comprehensive survey on deep learning-based biometric recognition," *Neurocomputing*, vol. 439, pp. 191–207, 2021.
- [22] A. J. Ma and P. C. Yuen, "Learning efficient multimodal biometric representations with deep learning," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 2, pp. 407–421, 2017.
- [23] S. Chopra, R. Hadsell, and Y. LeCun, "Learning a similarity metric discriminatively, with application to face verification," in *Proc. IEEE CVPR*, 2005, pp. 539–546.
- [24] W. Liu et al., "Large-margin softmax loss for convolutional neural networks," in *Proc. ICML*, 2016.
- [25] H. Wang et al., "CosFace: Large margin cosine loss for deep face recognition," in *Proc. IEEE CVPR*, 2018.
- [26] J. Daugman, "High confidence visual recognition of persons by a test of statistical independence," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 15, no. 11, pp. 1148–1161, 1993.
- [27] N. Damer et al., "ToF face recognition: A survey," *IEEE Transactions on Biometrics, Behavior, and Identity Science*, vol. 2, no. 4, pp. 1–20, 2020.
- [28] ISO/IEC 19795-1, *Biometric performance testing and reporting*, International Organization for Standardization, 2006.
- [29] P. Grother, M. Ngan, and K. Hanaoka, "Face recognition vendor test (FRVT)," *NIST Interagency Report*, 2018.

- [30] C. Rathgeb and A. Uhl, "A survey on biometric cryptosystems and cancelable biometrics," *EURASIP Journal on Information Security*, 2011.
- [31] M. Gomez-Barrero et al., "Biometric security: A survey," *IEEE Access*, vol. 6, pp. 30544–30557, 2018.
- [32] R. Ranjan et al., "L2-constrained softmax loss for discriminative face verification," *IEEE Signal Processing Letters*, vol. 24, no. 11, pp. 1579–1583, 2017.
- [33] Y. Sun et al., "Deep learning face representation by joint identification-verification," *Advances in Neural Information Processing Systems*, 2014.
- [34] Q. Cao et al., "VGGFace2: A dataset for recognising faces across pose and age," in *Proc. IEEE FG*, 2018.
- [35] S. Marcel, M. Nixon, and S. Li, *Handbook of Biometric Anti-Spoofing*. Springer, 2014.
- [36] R. S. Malhotra et al., "Privacy-preserving biometric authentication using deep learning," *IEEE Access*, vol. 8, pp. 123456–123468, 2020.
- [37] D. Yambay et al., "LivDet iris 2017: Iris liveness detection competition," in *Proc. IEEE IJCB*, 2017.
- [38] K. Patel, H. Han, and A. K. Jain, "Secure face unlock: Spoof detection," *IEEE Transactions on Information Forensics and Security*, vol. 14, no. 4, pp. 1–12, 2019.
- [39] M. S. Nixon and J. Carter, "Automatic recognition by gait," *Proceedings of the IEEE*, vol. 94, no. 11, pp. 2013–2024, 2006.
- [40] A. Ross, "An introduction to multibiometrics," in *Proc. IEEE ICIP*, 2010, pp. 1–4.
- [41] Y. Taigman et al., "DeepFace: Closing the gap to human-level performance in face verification," *Proc. CVPR*, 2014.
- [42] O. M. Parkhi et al., "Deep face recognition," *BMVC*, 2015.
- [43] Y. Sun et al., "Deep learning face representation by joint identification-verification," *NeurIPS*, 2014.
- [44] E. Ahmed et al., "An improved deep learning architecture for person re-identification," *Proc. CVPR*, 2015.
- [45] L. Wu et al., "Sampling matters in deep embedding learning," *Proc. ICCV*, 2017.
- [46] J. Zhang et al., "Multimodal deep metric learning for cross-modal biometric authentication," *IEEE Access*, 2018.
- [47] M. Kan et al., "Multi-view deep network for cross-modal face recognition," *IEEE Trans. Image Processing*, 2016.

- [48] I. Shahin et al., "Deep Siamese network for speaker verification," *Neural Computing and Applications*, 2019.
- [49] J. Deng et al., "ArcFace: Additive angular margin loss for deep face recognition," *Proc. CVPR*, 2019.
- [50] Y. Cao et al., "Deep multimodal biometric fusion using feature-level integration," *Pattern Recognition Letters*, 2020.
- [51] A. Rattani and A. Ross, "Cross-spectral face recognition," *IEEE Trans. Information Forensics and Security*, 2018.
- [52] T. Chugh et al., "Fingerprint spoof detection using deep learning," *IEEE TIFS*, 2019.
- [53] X. Hu et al., "Deep multimodal feature embedding for biometric verification," *Information Fusion*, 2021.
- [54] D. Kisku et al., "Attention-based multimodal biometric fusion," *IEEE Access*, 2022.
- [55] J. Deng et al., "Vision transformer for face recognition," *IEEE T-PAMI*, 2022.