# A STUDY OF MULTIMODAL PERSON VERIFICATION USING AUDIO-VISUAL-THERMAL DATA

*Madina Abdrakhmanova, Saniya Abushakimova, Yerbolat Khassanov,*
*and Huseyin Atakan Varol*

Institute of Smart Systems and Artificial Intelligence,
Nazarbayev University, Nur-Sultan City, Kazakhstan

## ABSTRACT

In this paper, we study an approach to multimodal person verification using audio, visual, and thermal modalities. The combination of audio and visual modalities has already been shown to be effective for robust person verification. From this perspective, we investigate the impact of further increasing the number of modalities by supplementing thermal images. In particular, we implemented unimodal, bimodal, and trimodal verification systems using the state-of-the-art deep learning architectures and compared their performance under clean and noisy conditions. We also compared two popular fusion approaches based on simple score averaging and soft attention mechanism. The experiment conducted on the SpeakingFaces dataset demonstrates the superiority of the trimodal verification system over both unimodal and bimodal systems. To enable the reproducibility of the experiment and facilitate research into multimodal person verification, we make our code, pretrained models and preprocessed dataset freely available in our GitHub repository[1].

***Index Terms***— Person verification, multimodal, audio-visual-thermal, data augmentation, fusion

## 1. INTRODUCTION

Person verification is the general task of verifying person's identity using various biometric characteristics. The two widely used biometric features are voice (i.e., audio) and face (i.e., visual), corresponding to the speaker and the face verification tasks, respectively. Recently, it has been observed that the combination of audio and visual modalities positively impacts the person verification task, yielding more accurate and robust systems [1, 2]. Presumably, increasing the number of modalities with the use of thermal images might boost the performance of a person verification system.

To date, most of the existing works on person verification focus on unimodal approaches utilizing various biometric features, such as voice, face, fingerprints, iris, and gait [3]. Remarkable progress has been achieved in the use

of voice [4, 5, 6] and face [7, 8] biometrics. However, the performance of unimodal verification systems degrades substantially under challenging conditions [1, 2]. In response, several works resort to multimodal approaches that introduce effective audio-visual fusion strategies. For example, [9, 10] combined the scores generated by unimodal verification systems. In contrast, [1, 2] explored fusion strategies using attention-based deep neural networks [11], and found that the proposed multimodal verification systems outperformed their unimodal counterparts and were more robust to noise. To further facilitate research into multimodal person verification, several multimedia datasets have been introduced [9, 12, 13], and challenges have been organized [14].

Recently, consumer electronics manufacturers have started to equip smartphones with thermal cameras [15, 16], bringing to market mobile devices that incorporate the triplet of audio, visual, and thermal modalities. A thermal camera presents new opportunities for a wide range of applications. In particular, it enables a more nuanced analysis of images in suboptimal physical environments (e.g., nighttime) and reduces the effects of variation in background, clothing, accessories, and appearance (e.g., makeup and contact lenses). Additionally, a high-resolution thermal camera can provide more granular association of temperature values with facial features. Its combination with other modalities can potentially strengthen multimodal person verification systems, by providing complementary information and alleviating the respective drawbacks of other modalities. Thermal images have been successfully deployed for speech recognition [17], lip reading [18], and person recognition using body images [19].

In this work, we investigate whether the addition of thermal images as a third modality can improve the performance (accuracy and robustness) of audio-visual person verification systems. To the best of our knowledge, this is the first attempt of audio-visual-thermal person verification. In particular, we adapted the SpeakingFaces dataset [20] to the trimodal person verification task and shared the preprocessed data with the community. We also developed unimodal, bimodal (i.e., audio-visual) and trimodal (i.e., audio-visual-thermal) verification systems using the state-of-the-art deep learning architectures, and assessed their performance under two data con-

---

[1] https://github.com/IS2AI/trimodal_person_verification

ditions:

1) *Clean:* Train, validation, and test sets did not contain any augmentations or noise.

2) *Noisy:* Each of the three sets contained 30% noise.

We also compared two popular fusion strategies based on simple score averaging and attention mechanism. Experimental results show that under both conditions the trimodal verification system outperforms the bimodal system by 18% and 20% relative equal error rates (EERs), respectively.

The rest of the paper is organized as follows: Section 2 details the specifications of the SpeakingFaces dataset used in our experiments. Section 3 describes the architecture of the developed person verification systems. Section 4 presents the experimental setup. Section 5 discusses the obtained results. Section 6 concludes the paper and highlights directions for future work.

## 2. AUDIO-VISUAL-THERMAL DATASET

In this work, we utilized the SpeakingFaces dataset [20] to train, validate, and test the person verification systems. SpeakingFaces is a publicly available multimodal dataset comprised of audio, visual, and thermal data streams. The dataset consists of 13,711 recordings of imperative sentences[2] spoken by 142 speakers from different backgrounds (i.e., around 100 utterances per speaker). The audio and visual streams were recorded using a webcam (Logitech C920 Pro HD), while the thermal stream was captured with a thermal camera (FLIR T540). The data were acquired approximately one meter away from the person.

To prepare the dataset for the person verification task, we first detected bounding boxes of facial regions in visual images using RetinaFace [21]. Next, we mapped the detected bounding boxes onto the thermal images by aligning the visual and thermal streams. Additionally, we manually filtered out from the final dataset those utterances where faces were not detected or where facial features were obscured based on previously documented artifacts. This left 13,036 recordings in the final version. We also downsampled the audio recordings to 16 kHz. The samples of visual and thermal facial image pairs are shown in Figure 1 and the specifications of the preprocessed dataset are listed in Table 1.

We prepared the train, validation, and test sets following the format of VoxCeleb [12]. In particular, the train list contains the paths to the recordings and the corresponding subject identifiers. The validation and test lists consist of randomly generated positive and negative pairs. For each subject, the same number of positive and negative pairs were selected. In total, the numbers of pairs in the validation and test sets are 38,000 and 46,200, respectively.

---

[2]Verbal commands given to virtual assistants and other smart devices, such as 'turn off the lights', 'play the next song', and so on.



**Fig. 1**. The samples of visual and thermal facial image pairs in the preprocessed SpeakingFaces dataset.

**Table 1**. Preprocessed SpeakingFaces dataset specification.

| Category | Train | Valid | Test | Total |
|---|---|---|---|---|
| # Speakers | 100 | 20 | 22 | 142 |
| # Utterances | 9,307 | 1,843 | 1,886 | 13,036 |
| Total duration | 10.4 hr | 2.2 hr | 2.3 hr | 14.9 hr |

To imitate more challenging scenarios, we deployed a noisy version of the train, validation, and test sets. The noisy data were constructed using Algorithm 1 presented in [2]. In particular, for the visual stream, we applied vertical and horizontal motion blur, to imitate the motion of a person in front of the camera, and the Gaussian blur, to imitate other noises. For the thermal stream, we applied only the Gaussian blur. For the audio stream, we added three kinds of noises from the Musan dataset [22]. In addition, for all streams, we used random noise sampled from the standard normal distribution (i.e., $\mathcal{N}(0, 1)$) to imitate a broken data stream.

## 3. SYSTEM ARCHITECTURE

We followed [23] to implement the multimodal person verification system. The system architecture consists of two main modules: encoder and fusion. The encoder module $Encoder(\cdot)$ is based on the ResNet34 [24] network. Specifically, for image input (i.e., visual and thermal), we used a variation of ResNet34, in which the number of channels in each residual block is halved in order to reduce computational cost. For audio, we used another variation of ResNet34, in which average-pooling was replaced with self-attentive pooling (SAP) [25] to aggregate frame-level features into utterance-level representation.

The encoder takes in a raw feature $x_i$ and outputs the corresponding embedding vector representation $e_i \in \mathbb{R}^{512}$:

$$e_i = Encoder(x_i) \qquad (1)$$

where $i \in \{a, v, t\}$ is used to represent the stream source (i.e., audio, visual, and thermal). We trained a separate encoder module for each data stream.

As a fusion module, we tried two different approaches: 1) embedding-level and 2) score-level. In the embedding-level fusion, we implemented a soft attention mechanism similar to that in [1], where the fusion module can pay attention to a salient modality among audio $e_a$, visual $e_v$, and thermal $e_t$ representations. Specifically, it first computes the attention score $\hat{\alpha}_{\{a,v,t\}} \in \mathbb{R}^3$ as follows:

$$\hat{\alpha}_{\{a,v,t\}} = \mathbf{W}[e_a, e_v, e_t] + \mathbf{b} \quad (2)$$

where the weight matrix $\mathbf{W} \in \mathbb{R}^{3 \times 1536}$ and the bias vector $\mathbf{b} \in \mathbb{R}^3$ are learnable parameters. Next, the fused person embedding vector $e_p$ is produced by the weighted sum:

$$e_p = \sum_{i \in \{a,v,t\}} \alpha_i e_i \ , \text{ where } \ \alpha_i = \frac{\exp(\hat{\alpha}_i)}{\sum_{k \in \{a,v,t\}} \exp(\hat{\alpha}_k)} \quad (3)$$

The embedding-level fusion for the audio-visual bimodal systems was designed in a similar manner, where two modalities are used instead of three. Note that the attention-based fusion module is jointly trained with the encoder modules in an end-to-end fashion.

In the score-level fusion, we simply averaged the scores $s_i \in [0, 2]$ produced by the unimodal verification systems. For example, for the trimodal system, the final score $s_{final}$ is computed as follows:

$$s_{final} = \frac{\sum_{i \in \{a,v,t\}} s_i}{3} \quad (4)$$

Likewise, for the bimodal system, the final score is computed by averaging the scores from two unimodal systems.

## 4. EXPERIMENTAL SETUP

All the person verification models were trained on a single V100 GPU running on the NVIDIA DGX-2 server using the training set. All hyper-parameters were tuned using the validation set, and the settings for each modality were tuned separately. The best performing model on the validation set was evaluated using the test set.

At the training stage, the input features from the audio stream were generated by randomly extracting a two-second segment from each recording, and then, transforming it into 40-dimensional Mel filterbank features. For the visual and thermal streams, we extracted the first frame of a recording. Model parameters were optimized using the angular prototypical [23] loss function. Each batch contained all the 100 speakers of the train set with nine utterances per speaker. As a regularization, we applied weight decay tuned for each model separately. We trained each model three times with a different seed number and report the mean and standard deviation of the results. The exact system implementation details are provided in our GitHub repository[1].

At the evaluation stage, from each test recording, we extracted ten two-second audio segments at regular intervals,

and ten equidistantly spaced image frames. In the multimodal setting with attention, for a given utterance, each of the ten audio segments was paired with an image frame from either or both modalities. We computed a similarity score between all possible combinations (10×10 = 100) from each pair of utterances present in the evaluation sets. As a similarity score, we used the Euclidean distance, and the mean of the 100 similarities was used as the final score. The performance of the models was evaluated using the EER metric.

## 5. EXPERIMENTAL RESULTS

We evaluated the performance of unimodal and multimodal person verification systems under two data conditions: 1) *clean* and 2) *noisy*. The noisy versions of the train, validation, and test sets were generated based on the procedure described in Section 2.

### 5.1. Unimodal Person Verification

The results of unimodal person verification experiment are given in the first part of Table 2. The results show that when train and evaluation sets are clean, the best EER performance on the test set is achieved by the visual modality (4.09%) followed by the audio (9.29%) and then the thermal (10.58%). In the noisy condition, the performance on the test set degrades by 28%, 7%, and 17% relative EER for the audio, visual, and thermal systems, respectively. According to these results, the visual system achieves better accuracy and is more robust to corrupted data.

We also analysed the verification errors made by the unimodal systems (see Figure 2). We observed that the number of overlapping errors between different modalities is lower than the errors made by a single modality. This indicates that these modalities posses strong complementary properties, and therefore, the multimodal systems, which can effectively combine them, have good potential. Future work should focus on analysing these errors in great detail to identify the weaknesses and strengths of each modality.
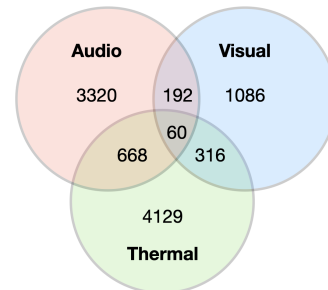


**Fig. 2**. The verification error statistics of unimodal systems.

| Modality | Fusion | clean | | noisy | |
|---|---|---|---|---|---|
| | | valid | test | valid | test |
| Audio | - | $10.82 \pm 0.39$ | $9.29 \pm 0.15$ | $14.89 \pm 0.60$ | $11.83 \pm 0.62$ |
| Visual | - | $4.04 \pm 0.89$ | $4.09 \pm 0.86$ | $5.32 \pm 0.63$ | $4.38 \pm 0.29$ |
| Thermal | - | $10.30 \pm 0.74$ | $10.58 \pm 1.28$ | $11.06 \pm 0.20$ | $12.39 \pm 1.06$ |
| Audio-Visual | Attention | $4.42 \pm 0.85$ | $3.69 \pm 0.44$ | $5.03 \pm 0.95$ | $4.05 \pm 0.61$ |
| | Score averaging | $2.39 \pm 0.42$ | $2.20 \pm 0.58$ | $3.70 \pm 0.35$ | $2.65 \pm 0.34$ |
| Audio-Visual-Thermal | Attention | $4.28 \pm 0.38$ | $4.17 \pm 1.19$ | $4.07 \pm 1.03$ | $3.65 \pm 1.06$ |
| | Score averaging | $\mathbf{2.26 \pm 0.10}$ | $\mathbf{1.80 \pm 0.31}$ | $\mathbf{2.98 \pm 0.20}$ | $\mathbf{2.13 \pm 0.51}$ |

**Table 2**. EER (%) performance results (mean $\pm$ std) of person verification systems under two conditions: 1) *clean* and 2) *noisy*. In the *noisy* condition, the percentages of noisy samples in the train, validation, and test sets were set to 30%.

## 5.2. Multimodal Person Verification

The experimental results for our bimodal (audio-visual) and trimodal (audio-visual-thermal) verification models are presented in the second part of the Table 2. The models were constructed using two fusion methods, soft attention and score averaging, as mentioned in the previous sections. The latter approach provides superior performance for both bimodal and trimodal verification systems, which is consistent with the observations made in [2], but different from the findings in [1].

The experimental results show that the multimodal systems outperform the unimodal systems. The best bimodal system reduced EER on the test set by 46% (from 4.09 to 2.20) and 39% (from 4.38 to 2.65) relative to the visual system under the clean and noisy conditions, respectively. Similarly, the best trimodal system reduced EER on the test set by 56% (from 4.09 to 1.80) and 51% (from 4.38 to 2.13) under the clean and noisy conditions, respectively. Remarkably, these improvements were achieved by simply averaging the scores of the unimodal systems.

The trimodal system performed better than the bimodal under both clean and noisy conditions, improving EER on the test set by 18% (from 2.20 to 1.80) and 20% (from 2.65 to 2.13). Interestingly, the attention-based trimodal verification system achieved better EER under the noisy condition (4.17 versus 3.65). The analysis of the attention network parameters suggests that the mechanism learned to prioritize the visual stream when the streams are not corrupted. In contrast, the network focused on all the three when the data were noisy. Therefore, augmenting the train set with noise plays an important role in learning more robust features, and future work should study different augmentation methods for the multimodal person verification task. Unfortunately, combining three modalities using the soft attention mechanism is challenging due to the instability of the training process, and further studies should be conducted to explore other attention methods. Overall, the results highlight that the addition of thermal image data does indeed enhance the EER performance of multimodal person verification systems.

## 6. CONCLUSION

In this work, we explored multimodal learning for the person verification task using audio, visual and thermal data streams. We adapted the SpeakingFaces dataset to build the state-of-the-art deep learning-based models. Specifically, we developed unimodal, bimodal, and trimodal verification systems and compared their performance under clean and noisy data conditions. Among the unimodal systems, the visual modality achieved the best result, followed by the audio and then the thermal modality. For the multimodal systems, we compared two popular fusion strategies based on simple score averaging and soft attention mechanism. The former achieved a substantial performance gain for the trimodal system compared to the other systems. Specifically, the relative EER improvements over the visual-based system were 56% and 51%, and over the audio-visual system were 18% and 20%, under the clean and noisy data conditions, respectively.

In the future, we plan to thoroughly analyze the errors made by the different unimodal systems to identify the weaknesses and strengths of each modality. We also plan to implement more advanced attention mechanisms and study the impact of different data augmentation techniques. The robustness evaluation of multimodal verification systems under different noise rates is also left for future work.

## 7. REFERENCES

[1] Suwon Shon, Tae-Hyun Oh, and James Glass, "Noise-tolerant audio-visual online person verification using an attention-based neural network fusion," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 3995–3999.

[2] Yanmin Qian, Zhengyang Chen, and Shuai Wang, "Audio-visual deep neural network for robust person verification," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 1079–1092, 2021.

[3] Shervin Minaee, Amirali Abdolrashidi, Hang Su, Mo-

hammed Bennamoun, and David Zhang, "Biometrics recognition using deep learning: A survey," *arXiv preprint arXiv:1912.00271*, 2019.

[4] Najim Dehak, Patrick J Kenny, Réda Dehak, Pierre Dumouchel, and Pierre Ouellet, "Front-end factor analysis for speaker verification," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 4, pp. 788–798, 2010.

[5] David Snyder, Daniel Garcia-Romero, Gregory Sell, Daniel Povey, and Sanjeev Khudanpur, "X-vectors: Robust dnn embeddings for speaker recognition," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5329–5333.

[6] Arsha Nagrani, Joon Son Chung, Weidi Xie, and Andrew Zisserman, "Voxceleb: Large-scale speaker verification in the wild," *Computer Speech & Language*, vol. 60, pp. 101027, 2020.

[7] Omkar M. Parkhi, Andrea Vedaldi, and Andrew Zisserman, "Deep face recognition," in *Proceedings of the British Machine Vision Conference (BMVC)*. 2015, pp. 41.1–41.12, BMVA Press.

[8] Florian Schroff, Dmitry Kalenichenko, and James Philbin, "Facenet: A unified embedding for face recognition and clustering," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2015, pp. 815–823, IEEE Computer Society.

[9] Gregory Sell, Kevin Duh, David Snyder, Dave Etter, and Daniel Garcia-Romero, "Audio-visual person recognition in multimedia data from the iarpa janus program," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 3031–3035.

[10] Jahangir Alam, Gilles Boulianne, Lukáš Burget, et al., "Analysis of ABC submission to NIST SRE 2019 CMN and VAST challenge," in *Odyssey: The Speaker and Language Recognition Workshop*, 2020, pp. 289–295.

[11] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio, "Neural machine translation by jointly learning to align and translate," in *International Conference on Learning Representations (ICLR)*, 2015.

[12] Arsha Nagrani, Joon Son Chung, and Andrew Zisserman, "Voxceleb: A large-scale speaker identification dataset," in *Interspeech*. 2017, pp. 2616–2620, ISCA.

[13] Joon Son Chung, Arsha Nagrani, and Andrew Zisserman, "Voxceleb2: Deep speaker recognition," in *Interspeech*. 2018, pp. 1086–1090, ISCA.

[14] Seyed Omid Sadjadi, Craig S Greenberg, Elliot Singer, Douglas A Reynolds, Lisa Mason, Jaime Hernandez-Cordero, et al., "The 2019 NIST audio-visual speaker recognition evaluation," *Proc. Speaker Odyssey*, 2020.

[15] Teledyne FLIR, "FLIR ONE Pro," https://www.flir.com/products/flir-one-pro/, Accessed: September 8, 2021.

[16] Caterpillar, "CAT S62 Pro," https://www.catphones.com/en-dk/cat-s62-pro-smartphone/, Accessed: September 8, 2021.

[17] Steven J. Anderson, Alvis Cheuk M. Fong, and Jie Tang, "Robust tri-modal automatic speech recognition for consumer applications," *IEEE Transactions on Consumer Electronics*, vol. 59, no. 2, 2013.

[18] Takeshi Saitoh and Ryosuke Konishi, "Lip reading using video and thermal images," in *SICE-ICASE International Joint Conference*. IEEE, 2006, pp. 5011–5015.

[19] Dat Tien Nguyen, Hyung Gil Hong, Ki-Wan Kim, and Kang Ryoung Park, "Person recognition system based on a combination of body images from visible light and thermal cameras," *Sensors*, vol. 17, no. 3, pp. 605, 2017.

[20] Madina Abdrakhmanova, Askat Kuzdeuov, Sheikh Jarju, Yerbolat Khassanov, Michael Lewis, and Huseyin Atakan Varol, "Speakingfaces: A large-scale multimodal dataset of voice commands with visual and thermal video streams," *Sensors*, vol. 21, no. 10, pp. 3465, 2021.

[21] Jiankang Deng, Jia Guo, Yuxiang Zhou, Jinke Yu, Irene Kotsia, and Stefanos Zafeiriou, "Retinaface: Single-stage dense face localisation in the wild," *arXiv preprint arXiv:1905.00641*, 2019.

[22] David Snyder, Guoguo Chen, and Daniel Povey, "Musan: A music, speech, and noise corpus," *arXiv preprint arXiv:1510.08484*, 2015.

[23] Joon Son Chung, Jaesung Huh, Seongkyu Mun, Minjae Lee, Hee-Soo Heo, Soyeon Choe, Chiheon Ham, Sunghwan Jung, Bong-Jin Lee, and Icksang Han, "In defence of metric learning for speaker recognition," in *Interspeech*. 2020, pp. 2977–2981, ISCA.

[24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016, pp. 770–778, IEEE.

[25] Weicheng Cai, Jinkun Chen, and Ming Li, "Exploring the encoding layer and loss function in end-to-end speaker and language recognition system," in *Odyssey*. 2018, pp. 74–81, ISCA.