

Simultaneous Face Hallucination and Translation for Thermal to Visible Face Verification using Axial-GAN

Rakhil Immidisetti¹Shuowen Hu²Vishal M. Patel¹¹Johns Hopkins University ²DEVCOM Army Research Laboratory

rimmidil@jhu.edu, shuowen.hu.civ@mail.mil, vpatel136@jhu.edu

Abstract

Existing thermal-to-visible face verification approaches expect the thermal and visible face images to be of similar resolution. This is unlikely in real-world long-range surveillance systems since humans are distant from the cameras. To address this issue, we introduce the task of thermal-to-visible face verification from low-resolution thermal images. Furthermore, we propose Axial-Generative Adversarial Network (Axial-GAN) to synthesize high-resolution visible images for matching. In the proposed approach we augment the GAN framework with axial-attention layers which leverage the recent advances in transformers for modelling long-range dependencies. We demonstrate the effectiveness of the proposed method by evaluating on two different thermal-visible face datasets. When compared to related state-of-the-art works, our results show significant improvements in both image quality and face verification performance, and are also much more efficient.

1. Introduction

In practical scenarios such as low-light or night-time conditions, one has to use thermal cameras for surveillance in order to detect and recognize faces. The acquired thermal images of faces in such scenarios have to be matched with existing biometric datasets that contain visible face images. Significant progress has been made by several works [5–7, 10, 11, 14, 40] to address the thermal-to-visible cross-spectrum face recognition problem. But existing works expect the thermal and visible face images to be of similar resolution. This is unlikely in real-world surveillance systems as humans are further away from cameras, thereby the region occupied by a face is much less when compared to an image in the visible face dataset. We illustrate the described issue in Figure 1. To address this, we introduce the task of

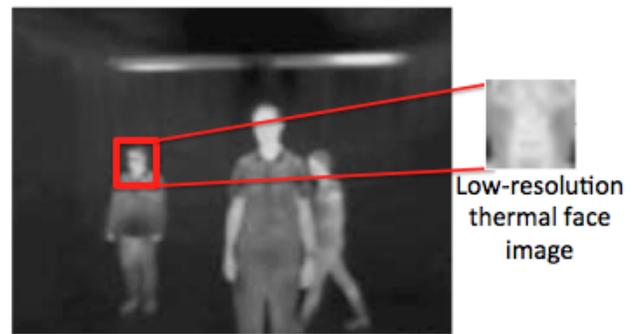


Figure 1. A typical thermal image [1]. Note that the captured images are of very low-resolution. In order to perform cross-modal face recognition, one needs to synthesize a high-resolution visible face image from a low-resolution thermal face image.

matching low-resolution (LR) thermal face images against high-resolution (HR) visible face images.

The large domain discrepancy between the thermal and visible images and the low resolution of the thermal images makes the introduced task quite challenging. To tackle it, we propose a hybrid network that augments an image-conditional generative adversarial network (GAN) [8] with axial-attention [31] layers. The generator synthesizes face images in the visible domain, which are then matched against a gallery of visible images using an off-the-shelf face matching algorithm. Using self-attention-based models [25, 30, 31] allows capturing the structural patterns of the face effectively, which is essential for tasks such as face verification. However, stand-alone self-attention models require large-scale datasets for training. Therefore, we develop a hybrid network that makes use of both convolutions and self-attention layers to efficiently capture the local and global information, respectively. Additionally, augmenting our network with self-attention avoids the use of several stacked convolutional layers for modelling global dependencies. This makes our network extremely parameter efficient without any reduction in performance. Although Di *et al.* [5] proposed a similar hybrid network, it doesn't uti-

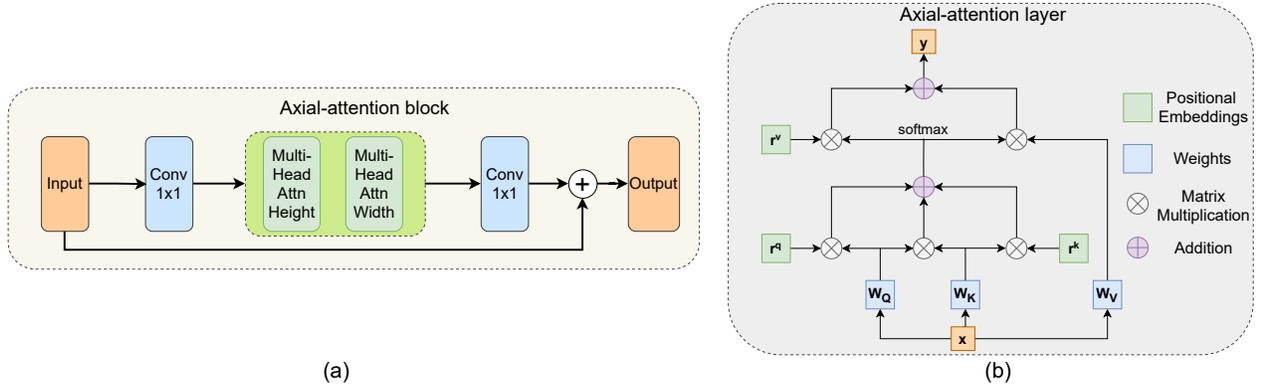


Figure 2. (a) The residual axial-attention block used in Axial-GAN. (b) Axial-attention layer, which is the basic building block of both height and width multi-head attention modules in the axial-attention block.

lize positional information and multi-head design that are essential for capturing spatial structures and a mixture of features, which we incorporate into our network. To the best of our knowledge, this is one of the first works to propose a transformer-based GAN for face translation and face hallucination.

We evaluate our approach on the ARL-VTF dataset [24] and the polarimetric thermal face recognition dataset [6]. We compare the performance of our approach with state-of-the-art methods in thermal-to-visible synthesis and also face hallucination. Our results show significant performance improvements in both image quality and face verification. Furthermore, an ablation study is conducted to demonstrate the effectiveness of axial-attention. Code is available at <https://github.com/sam575/axial-gan>.

2. Related Work

2.1. Thermal-to-visible face recognition

Several approaches have been proposed to address thermal-to-visible cross-spectrum face recognition, which can be mainly divided into two categories: feature-based and synthesis-based. Feature-based methods seek to find a common latent subspace where corresponding face images in each spectrum are closer in terms of some distance metric. Initial works used kernel prototype similarities [19], partial least squares [4, 13] and coupled neural networks [26] on hand-crafted features such as SIFT and HOG. Whereas, recent works leverage deep networks to extract domain-invariant features [7, 10, 11, 14] or disentangled features [35].

Synthesis-based methods have the advantage that they can leverage the recent advances in visible spectrum face recognition for matching the synthesized visible face images. Consequently, Riggan *et al.* [27] used features from both global and local regions, and developed a region-specific cross-spectrum mapping for estimating visible images. Zhang *et al.* [39, 40] and Di *et al.* [5, 6] leverage GANs to enhance the perceptual quality of the synthesized

images.

2.2. Face Hallucination

Face hallucination is a domain-specific image super-resolution problem aimed at enhancing the resolution of a LR face image to generate the corresponding HR face image [16]. Consequently, most of the works exploit face-specific information such as attributes, landmarks, parsing maps, and identity for effective reconstruction of HR face images [2, 3, 41]. This additional information is obtained either by human labelling or by using existing pre-trained models. Furthermore, to generate realistic faces many works tend to employ GANs [2, 3, 36, 37].

2.3. Transformers

Transformers, introduced by [30] leverage multi-head self-attention layers to compute pairwise correspondence between tokens to learn highly expressive features across long sequences. Recently, self-attention has been applied to many computer vision tasks such as classification, detection and segmentation [25, 29, 31, 33]. In contrast to non-local block models [5, 33, 38], transformer-based models [25, 29, 31] use relative positional encodings and multi-head design, which is essential as they capture spatial structures in an image and a mixture of affinities, respectively.

3. Proposed Method

In this section, we discuss the details of the proposed Axial-GAN for thermal-to-visible synthesis from LR thermal images. In particular, we first give an overview of axial-attention [12, 31], then discuss the proposed axial-attention-based generator and discriminator networks and finally address the objective functions and implementation details.

3.1. Axial-attention overview

Axial-attention [12, 31] factorizes 2D self-attention into two steps that apply 1D self-attention in height-axis and width-axis sequentially. Each step computes pairwise

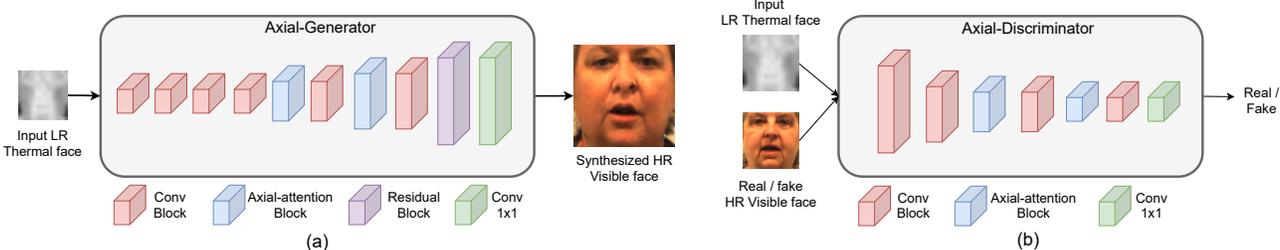


Figure 3. The proposed generator (a) and discriminator (b) augmented with axial-attention layers.

affinities, thereby learning a rich set of associative features across that particular axis. Combining them sequentially allows in capturing the full global information. The factorization helps in reducing the computational complexity and allows to capture long-range dependencies on larger regions, which is infeasible for 2D self-attention. Additionally, Wang *et al.* [31] incorporates relative positional bias [28] in key, query and value as illustrated in Figure 2 (b), thereby effectively associating information with positional awareness. Formally, the position-sensitive axial-attention layer along the width-axis is as follows:

$$y_{ij} = \sum_{w=1}^W \text{softmax}(q_{ij}^T k_{iw} + q_{ij}^T r_{iw}^q + k_{iw}^T r_{iw}^k)(v_{iw} + r_{iw}^v), \quad (1)$$

where y_{ij} denotes the value at position (i, j) in the output feature map. Here, q, k, v denote query, key and value, respectively, which are computed from input feature map. r^q, r^k, r^v denote the relative positional encoding for query, key and value, respectively. In our work, we adopt the axial-attention blocks as illustrated in Figure 2 (a), which comprises of the multi-head position-sensitive axial-attention layers, to augment the GAN.

3.2. Generator

The generator, given a LR thermal image as input, synthesizes a HR visible image. Inspired by the works in super-resolution [34], a progressive upsampling framework is used in the proposed generator as shown in Figure 3 (a). This framework is efficient when compared to the pre-upsampling framework and also avoids the noise amplification caused by an upsampled input image. Initially, we use only convolutional layers to learn local features such as edges, which are difficult to model by content-based mechanisms such as self-attention. In the later stages, we augment the network with axial-attention blocks to model the global context. Additionally, in order to improve the stability of training, we use spectral normalization [21] for all convolutional layers except for those in the axial-attention block and the final output layer. Specifically, our generator architecture consists of the following components:

C64 - C128 - C256 - C512 - D256 - A256 - C256 - D128 - A128 - C128 - D64 - R64 - F

where Ck, Dk, Ak and Rk denote 3×3 Convolution-BatchNorm-ReLU layer, 3×3 Deconvolution-BatchNorm-ReLU layer, axial-attention block and residual block [9], respectively, with k filters. F denotes a 1×1 convolutional layer with Tanh as activation function, which produces a three channel output.

3.3. Discriminator

We use a patch-based discriminator [15] augmented with axial-attention layers as shown in Figure 3 (b), which is trained alternatively with the generator. The input to the discriminator is the concatenation of up-sampled thermal image and either real or fake visible image. Similar to the generator, spectral normalization is used for improving the training stability. The discriminator consists of the following components:

C64 - C128 - A128 - C256 - A256 - C512 - F

where Ck and Ak denote a 4×4 Convolution-LeakyReLU layer and axial-attention block, respectively, with k filters. F denotes a 1×1 convolutional layer, which produces a single channel output.

3.4. Objective function

The training dataset is given as a set of pairs $\{(x_i, y_i)\}$, where x_i is a LR thermal image and y_i is the corresponding HR visible image. We minimize the hinge version of adversarial loss [20] for training the generator G and discriminator D :

$$\begin{aligned} L_D &= -\mathbb{E}_{x,y}[\min(0, -1 + D(x, y))] \\ &\quad -\mathbb{E}_x[\min(0, -1 - D(x, G(x)))] \quad (2) \\ L_G &= -\mathbb{E}_x[D(x, G(x))]. \end{aligned}$$

The overall loss function for the generator is defined as follows:

$$L = L_G + \lambda_H L_H + \lambda_P L_P + \lambda_{FM} L_{FM}, \quad (3)$$

where L_G is the adversarial loss for generator in Eq. 2, L_H is the Huber loss in Eq. 4, L_P is the perceptual loss in Eq. 5, L_{FM} is the discriminator-based feature matching loss in Eq. 6. $\lambda_H, \lambda_P, \lambda_{FM}$ are the weights for Huber loss, perceptual loss and discriminator-based feature matching loss, respectively.

We use Huber loss between the target visible image and the synthesized visible image:

$$L_H = \begin{cases} \mathbb{E}_{x,y}[0.5 * (G(x) - y)^2], & \text{if } |G(x) - y| < 1 \\ \mathbb{E}_{x,y}[|G(x) - y| - 0.5], & \text{otherwise.} \end{cases} \quad (4)$$

To generate visually pleasing results, perceptual loss [17] is minimized, which is computed using features extracted from an off-the-shelf pre-trained VGG-19 network. The features from the initial layers help in generating high-frequency details, whereas the deeper layers help in enhancing the discriminative details. The perceptual loss is formally stated below in which F_i denotes the i -th layer of the VGG-19 network:

$$L_p = E_{x,y}[\|F_i(y) - F_i(G(x))\|_1]. \quad (5)$$

Additionally, discriminator-based feature matching loss [32] is used to improve the training stability of GAN. Here, the i -th layer of discriminator D is denoted as D_i :

$$L_{FM} = E_{x,y}[\|D_i(x, y) - D_i(x, G(x))\|_1]. \quad (6)$$

3.5. Implementation

The entire network is trained on a single Nvidia 12 GB GPU. We choose $\lambda_H = 100$ for the Huber loss, $\lambda_P = 10$ for the perceptual loss and $\lambda_{FM} = 10$ for the discriminator-based feature matching loss. Adam [18] is used as the optimization algorithm with a learning rate of 0.0002 and the batch size is set to 32. The second convolutional layer at each scale in VGG-19 is used for the perceptual loss. The output features after each scale in the discriminator are used for the discriminator-based feature matching loss.

The input to the generator is a 16×16 thermal image downsampled from the HR thermal image. The generator synthesizes the corresponding 128×128 visible face image. Downsampling or upsampling is performed using the MATLAB bicubic kernel function. Moreover, the training dataset is augmented with random horizontal flips.

4. Experimental Results

We evaluate the performance of our method using image quality and face verification metrics. The quality of the synthesized visible images is evaluated using the peak signal-to-noise ratio (PSNR) and structural similarity (SSIM) index [42]. The face verification performance is evaluated using the area under the curve (AUC) of receiver operating characteristic (ROC) and equal error rate (EER). Face verification scores are computed using cosine similarity between features extracted from ‘maxp_5_3’ layer of a pre-trained VGG-Face model [23]. The image quality metrics are evaluated on all images in the test set, whereas the face verification performance is evaluated on the protocols as described in each of the dataset section.

4.1. Datasets and Protocols

We use the ARL-VTF dataset [24] and the polarimetric thermal face recognition dataset [6] for evaluating our method. We use only the baseline and expression images and ignore the pose images from these datasets in our evaluation, as it is extremely challenging to synthesize images from off-pose low-resolution images.

ARL-VTF dataset. This is the largest dataset of paired conventional thermal and visible images, containing over 500,000 images from 395 subjects. The cropped face images are obtained using the provided bounding box annotations. The development and test split consist of 295 and 100 subjects, respectively. For evaluating our method we follow the provided protocols, which prescribe different combinations of the gallery (G_VB0-, G_VB0+) and probe sets (P_TB0-, P_TE0-, P_TB0+). Here, ‘‘G’’ and ‘‘P’’ denote the gallery and probe sets, respectively. Visible and thermal spectrum data are represented as ‘‘V’’ and ‘‘T’’, respectively. ‘‘B’’ and ‘‘E’’ denote the baseline and expression sequences, respectively. ‘‘0’’ represents the images of subjects who do not possess glasses, while ‘‘-’’ and ‘‘+’’ represent the images of subjects who have their glasses removed or worn, respectively. For example, G_VB0- is the set of visible images in the gallery where no subjects are wearing glasses.

Extended ARL multi-modal polarimetric thermal face recognition dataset. This dataset contains a total of 5419 polarimetric thermal and visible image pairs corresponding to 121 subjects. A polarimetric thermal image consists of three Stokes images as its three channels: S0, S1, and S2, where S0 represents the conventional intensity thermal image, S1 represents the horizontal and vertical polarization-state information and S2 represents the diagonal polarization-state information. We follow the pre-processing steps outlined in [6] for obtaining the cropped face images. The subjects are randomly divided into train, validation, and test sets containing 71, 25, and 25 subjects respectively. Reported results are evaluated over three random splits. For evaluating thermal-to-visible face verification, the gallery set is formed using a random baseline visible image from each subject in the test set. The remaining disjoint polarimetric thermal images form the probe set.

4.2. Results and Comparisons

We compare the performance of our method with pix2pix [15], Di *et al.* [5] and HiFaceGAN [36]. Di *et al.* [5] uses a self-attention [33] based CycleGAN [43] (SAGAN) for thermal-to-visible synthesis. We use our implementation for comparison and ignore the cyclic-consistency part *i.e.*, we do not train the GAN for visible-to-thermal synthesis for a fair comparison. HiFaceGAN is a recent work in face hallucination that uses content-adaptive convolutions



Input LR Thermal pix2pix [15] SAGAN [5] HiFaceGAN [36] Ours Axial-GAN HR GT HR Visible
 Figure 4. Synthesized visible images from different methods on the ARL-VTF dataset

Gallery	Method	P_TB0-		P_TE0-		P_TB0+	
		AUC (%)	EER (%)	AUC (%)	EER (%)	AUC (%)	EER (%)
G_VB0-	pix2pix	91	17.77	88.94	19.37	80.06	28.13
	SAGAN	92.29	15.3	90.78	17	79.18	28.38
	HiFaceGAN	91.29	17.13	89.42	18.49	82.82	26.21
	Axial-GAN (Ours)	94.4	12.38	92.71	14.86	84.62	24.67
	Ours w/o axial-attention	90.89	17.75	88.87	19.64	80.07	27.69
	Ours w self-attention	92.72	14.62	90.52	17.29	82.15	26.44
	Ours w HR Thermal	99.05	4.98	98.07	7.1	91.45	17.85
G_VB0+	pix2pix	86.32	22.43	83.7	24.58	91.15	17.89
	SAGAN	87.8	22.27	85.85	23.93	87.37	21.44
	HiFaceGAN	86.12	23.7	84.24	25.01	91.03	16.61
	Axial-GAN (Ours)	89.71	19.75	88.01	21.58	93.62	14.05
	Ours w/o axial-attention	86.73	22.43	84.37	24.3	89.96	17.63
	Ours w self-attention	87.89	22.59	85.96	23.75	91.48	16.5
	Ours w HR Thermal	96.53	10.21	94.72	13.75	98.53	6.67

Table 1. Face verification results corresponding to the ARL-VTF Dataset

Method	PSNR	SSIM	Resolution	Method	AUC	EER	PSNR	SSIM
pix2pix	16.243	0.549	24 × 24	pix2pix	90.66	17.66	16.75	0.55
SAGAN	17.67	0.61		SAGAN	92.26	15.72	17.66	0.60
HiFaceGAN	17.764	0.62		HiFaceGAN	91.32	17.09	18.21	0.63
Ours (Axial-GAN)	18.173	0.607		Ours (Axial-GAN)	95.52	11.44	18.56	0.63
Ours w/o axial-attention	17.067	0.577	8 × 8	pix2pix	75.55	31.81	15.85	0.54
Ours w self-attention	17.736	0.591		SAGAN	72.62	33.81	16.30	0.53
Ours w HR thermal	18.267	0.643		HiFaceGAN	75.71	31.61	16.36	0.57
				Ours (Axial-GAN)	78.79	28.39	16.57	0.55

Table 2. Image quality results on ARL-VTF dataset

Table 3. Comparison of results for different resolutions on ARL-VTF dataset

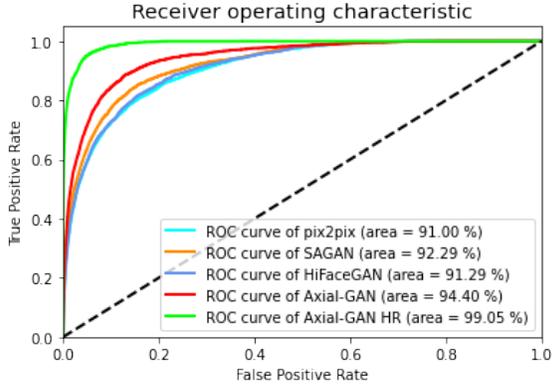


Figure 5. The ROC curve comparison on ARL-VTF dataset

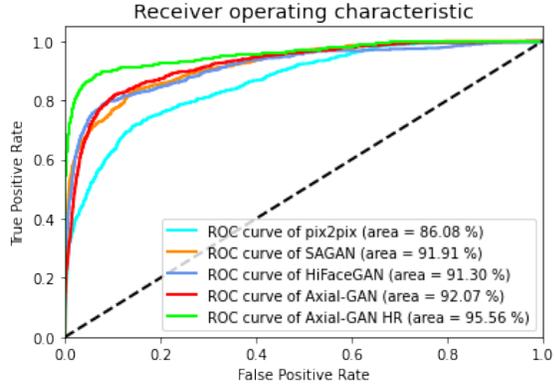
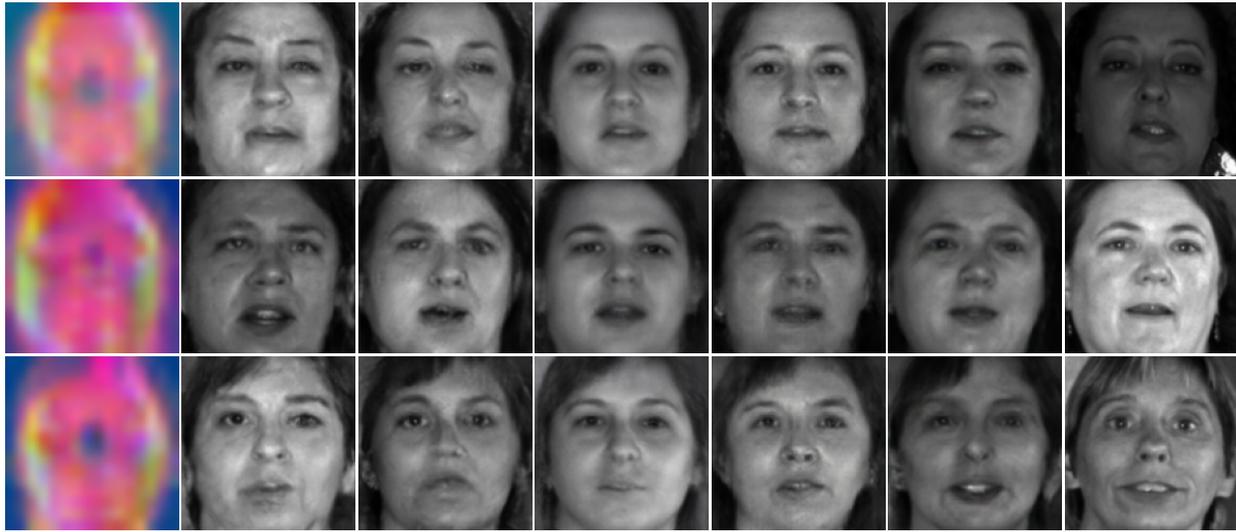


Figure 6. The ROC curve for polarimetric thermal dataset



Input LR Thermal pix2pix [15] SAGAN [5] HiFaceGAN [36] Ours Axial-GAN HR GT HR Visible

Figure 7. Synthesized visible images from different methods on the ARL polarimetric thermal face dataset

to extract features for semantic guidance [22] during replenishment. The competing methods use bicubic upsampled thermal images as input. Additionally, we compare with Axial-GAN that runs on HR thermal images to show the limitations of thermal-to-visible synthesis.

Table 1 shows the thermal-to-visible verification performance of different methods on the ARL-VTF dataset. Compared to the other state-of-the-art methods, our method performs better with higher AUC scores and lower EER scores across all protocols. We also show the ROC curve for G_VB0- vs. P_TB0- protocol in Figure 5. Additionally, Table 2 shows the image quality results, where our method outperforms the other methods based on PSNR but is comparable to HiFaceGAN based on SSIM. Similar results can also be observed in Figure 4, which shows the synthesized visible images for all methods. SAGAN has many artifacts in the synthesized images but the identity information is well retained. HiFaceGAN synthesizes images that are

smoother but the identity information is lost in this process. Our method comparatively synthesizes more realistic faces with well-defined face contours while preserving identity information. Furthermore, we conduct an additional study which shows the quantitative comparisons for different resolutions of input thermal image (see Table 3). Here, we report the average AUC and EER scores across all protocols. We obtain results that are consistent with the results for 16×16 resolution.

For the extended ARL multi-modal polarimetric face recognition dataset, the quantitative results are shown in Table 4 and the ROC curve for one of the splits is shown in Figure 6. When compared to the ARL-VTF dataset, the performance of all methods decrease for this dataset. This is mainly because of its dataset size, which is 100 times lesser than that of ARL-VTF. Additionally, there is a lot of variation in illumination which can be seen in the visible images of Figure 7. As a consequence, there is a

Method	AUC (%)	EER (%)	PSNR	SSIM	Parameters
pix2pix	81.573 ± 3.298	25.93 ± 2.733	16.488 ± 0.197	0.513 ± 0.021	41.8 M
SAGAN	84.377 ± 6.393	23.557 ± 6.699	17.329 ± 0.28	0.565 ± 0.016	7.9 M
HiFaceGAN	84.067 ± 6.801	22.947 ± 6.538	17.456 ± 0.176	0.583 ± 0.009	79.9 M
Ours (Axial-GAN)	85.557 ± 5.307	22.347 ± 5.7	17.739 ± 0.165	0.582 ± 0.012	4.1 M
Ours w HR Thermal	91.227 ± 4.17	15.53 ± 5.176	17.58 ± 0.05	0.588 ± 0.008	6.1 M

Table 4. Results for face verification and image quality on extended ARL multi-modal face recognition dataset

significant impact on the PSNR metric as thermal images fail to capture such variations. Our method performs better than the competing methods but the improvements are less when compared to the ARL-VTF dataset. HiFaceGAN performs relatively well when compared to its performance on the ARL-VTF dataset. This can potentially be attributed to the larger dataset size required by the attention-based models for observing reasonable performance boosts. Table 4 also shows the parameters in the generator for each of the methods. Our method has approximately 10×, 2×, and 20× lesser parameters than pix2pix, SAGAN, and HiFaceGAN, respectively. Furthermore, Figure 7 shows the qualitative comparisons of the competing methods. The qualitative observations for this dataset are consistent with that of the ARL-VTF dataset.

Ablation study. In Table 1 and 2 we also show the effectiveness of axial-attention. When we remove the axial-attention blocks from our method, we observe a significant decrease in performance. This shows the importance of capturing global information using self-attention-based models. Additionally, we also replace the axial-attention blocks with self-attention layers [5, 33] in our method. As expected, this performs similarly to SAGAN and falls short when compared to our method. This shows that the positional information and multi-head design are essential in improving the performance of self-attention-based models.

5. Conclusion

We introduced the task of thermal-to-visible face verification from low-resolution thermal images to deal with lower resolution faces in surveillance systems. To address this task, we proposed Axial-GAN in which we augment the GAN framework with axial-attention layers. Axial-attention effectively captures long-range dependencies with high efficiency. Our quantitative and qualitative results on multiple thermal-visible face datasets show improvements when compared to previous related works. In future work, we would like to investigate the challenging task of face verification using off-pose LR thermal faces.

Acknowledgement

This work was supported by ARO grant W911NF-21-1-0135.

References

- [1] G. A. Bilodeau, A. T. P. L. St-Charles, and D. Riahi. Thermal-visible registration of human silhouettes: a similarity measure performance evaluation. *Infrared Physics & Technology*, 64:79–86, May 2014. 1
- [2] A. Bulat and G. Tzimiropoulos. Super-fan: Integrated facial landmark localization and super-resolution of real-world low resolution faces in arbitrary poses with gans. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 2
- [3] Y. Chen, Y. Tai, X. Liu, C. Shen, and J. Yang. Fsrnet: End-to-end learning face super-resolution with facial priors. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 2
- [4] J. Choi, S. Hu, S. S. Young, and L. S. Davis. Thermal to visible face recognition. In *Proc.SPIE*, pages 8371 – 8371 – 10, 2012. 2
- [5] X. Di, B. S. Riggan, S. Hu, N. J. Short, and V. M. Patel. Polarimetric thermal to visible face verification via self-attention guided synthesis. In *2019 International Conference on Biometrics (ICB)*, pages 1–8, 2019. 1, 2, 4, 5, 6, 7
- [6] X. Di, B. S. Riggan, S. Hu, N. J. Short, and V. M. Patel. Multi-scale thermal to visible face verification via attribute guided synthesis. *IEEE Transactions on Biometrics, Behavior, and Identity Science*, 3(2):266–280, 2021. 2, 4
- [7] C. N. Fondje, S. Hu, N. J. Short, and B. S. Riggan. Cross-domain identification for thermal-to-visible face recognition. In *2020 IEEE International Joint Conference on Biometrics (IJCB)*, pages 1–9, 2020. 1, 2
- [8] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, volume 27, 2014. 1
- [9] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 3
- [10] R. He, X. Wu, Z. Sun, and T. Tan. Learning invariant deep representation for nir-vis face recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2017. 1, 2
- [11] R. He, X. Wu, Z. Sun, and T. Tan. Wasserstein cnn: Learning invariant features for nir-vis face recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,

- 41(7):1761–1773, 2019. 1, 2
- [12] J. Ho, N. Kalchbrenner, D. Weissenborn, and T. Salimans. Axial attention in multidimensional transformers. *arXiv preprint arXiv:1912.12180*, 2019. 2
- [13] S. Hu, J. Choi, A. L. Chan, and W. R. Schwartz. Thermal-to-visible face recognition using partial least squares. *JOSA A*, 32(3):431–442, 2015. 2
- [14] S. M. Iranmanesh, A. Dabouei, H. Kazemi, and N. M. Nasrabadi. Deep cross polarimetric thermal-to-visible face recognition. In *2018 International Conference on Biometrics (ICB)*, pages 166–173, Feb 2018. 1, 2
- [15] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. *CVPR*, 2017. 3, 4, 5, 6
- [16] J. Jiang, C. Wang, X. Liu, and J. Ma. Deep learning-based face super-resolution: A survey. *arXiv preprint arXiv:2101.03749*, 2021. 2
- [17] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision*, pages 694–711, 2016. 4
- [18] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2014. 4
- [19] B. F. Klare and A. K. Jain. Heterogeneous face recognition using kernel prototype similarities. *IEEE transactions on pattern analysis and machine intelligence*, 35(6):1410–1422, 2013. 2
- [20] J. H. Lim and J. C. Ye. Geometric gan. *arXiv preprint arXiv:1705.02894*, 2017. 3
- [21] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. Spectral normalization for generative adversarial networks. In *International Conference on Learning Representations*, 2018. 3
- [22] T. Park, M.-Y. Liu, T.-C. Wang, and J.-Y. Zhu. Semantic image synthesis with spatially-adaptive normalization. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 6
- [23] O. M. Parkhi, A. Vedaldi, and A. Zisserman. Deep face recognition. In *BMVC*, 2015. 4
- [24] D. Poster, M. Thielke, R. Nguyen, S. Rajaraman, X. Di, C. N. Fondje, V. M. Patel, N. J. Short, B. S. Riggan, N. M. Nasrabadi, and S. Hu. A large-scale, time-synchronized visible and thermal face dataset. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2021. 2, 4
- [25] P. Ramachandran, N. Parmar, A. Vaswani, I. Bello, A. Levskaya, and J. Shlens. Stand-alone self-attention in vision models. In *Advances in Neural Information Processing Systems*, 2019. 1, 2
- [26] B. S. Riggan, N. J. Short, and S. Hu. Optimal feature learning and discriminative framework for polarimetric thermal to visible face recognition. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2016. 2
- [27] B. S. Riggan, N. J. Short, and S. Hu. Thermal to visible synthesis of face images using multiple regions. In *IEEE Winter Conference on Applications of Computer Vision*, 2018. 2
- [28] P. Shaw, J. Uszkoreit, and A. Vaswani. Self-attention with relative position representations. In *NAACL-HLT*, 2018. 3
- [29] J. M. J. Valanarasu, P. Oza, I. Hacihaliloglu, and V. M. Patel. Medical transformer: Gated axial-attention for medical image segmentation. *arXiv preprint arXiv:2102.10662*, 2021. 2
- [30] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. *NIPS’17*, page 6000–6010, 2017. 1, 2
- [31] H. Wang, Y. Zhu, B. Green, H. Adam, A. Yuille, and L.-C. Chen. Axial-deeplab: Stand-alone axial-attention for panoptic segmentation. In *European Conference on Computer Vision (ECCV)*, 2020. 1, 2, 3
- [32] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 4
- [33] X. Wang, R. Girshick, A. Gupta, and K. He. Non-local neural networks. In *IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018. 2, 4, 7
- [34] Z. Wang, J. Chen, and S. C. Hoi. Deep learning for image super-resolution: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 2020. 3
- [35] X. Wu, H. Huang, V. M. Patel, R. He, and Z. Sun. Disentangled variational representation for heterogeneous face recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 9005–9012, 2019. 2
- [36] L. Yang, S. Wang, S. Ma, W. Gao, C. Liu, P. Wang, and P. Ren. Hifacegan: Face renovation via collaborative suppression and replenishment. In *28th ACM International Conference on Multimedia*, pages 1551–1560, 2020. 2, 4, 5, 6
- [37] X. Yu, F. Porikli, B. Fernando, and R. Hartley. Hallucinating unaligned face images by multiscale transformative discriminative networks. *International Journal of Computer Vision*, 128(2):500–526, 2020. 2
- [38] H. Zhang, I. Goodfellow, D. Metaxas, and A. Odena. Self-attention generative adversarial networks. *arXiv preprint arXiv:1805.08318*, 2018. 2
- [39] H. Zhang, V. M. Patel, B. S. Riggan, and S. Hu. Generative adversarial network-based synthesis of visible faces from polarimetric thermal faces. In *2017 IEEE International Joint Conference on Biometrics (IJCB)*, pages 100–107, 2017. 2
- [40] H. Zhang, B. S. Riggan, S. Hu, N. J. Short, and V. M. Patel. Synthesis of high-quality visible faces from polarimetric thermal faces using generative adversarial networks. *International Journal of Computer Vision*, 127(6):845–862, 2019. 1, 2
- [41] K. Zhang, Z. Zhang, C.-W. Cheng, W. Hsu, Y. Qiao, W. Liu, and T. Zhang. Super-identity convolutional neural network for face hallucination. In *ECCV*, 2018. 2
- [42] Zhou Wang and A. C. Bovik. A universal image quality index. *IEEE Signal Processing Letters*, 9(3):81–84, 2002. 4
- [43] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *IEEE International Conference on Computer Vision (ICCV)*, 2017. 4