# AgeTransGAN for Facial Age Transformation with Rectified Performance Metrics

Gee-Sern Hsu, Rui-Cang Xie, Zhi-Ting Chen, and Yu-Hong Lin

National Taiwan University of Science and Technology, Taipei, Taiwan {jison,m10703430,m10803432,m10903430}@mail.ntust.edu.tw

This document has the following contents:

- 1. More details about Cross-Age Face (CAF) Dataset
- 2. Determination of face verification threshold
- 3. Network settings
- 4. More experiments: performance on other datasets, different loss settings, user study, and more qualitative comparisons

(Line Number refers to the main paper)

## 1 More Details about Cross-Age Face (CAF) Dataset

(Line 417. More detail of the dataset is in supplementary document.)

The Cross-Age Face (CAF) dataset consists of 4000 images of 520 subjects with age between 0 to 94. Each face has a ground-truth age, and each individual has images in at least 5 age groups across  $G_{10}0 \sim G_{10}9^1$ . The numbers of subjects in  $G_{10}0, G_{10}1, ..., G_{10}9$  are 341, 364, 312, 399, 469, 515, 435, 296, 195, and 67, respectively. 120 (162/134/68/28/8) subjects have images across 5 (6/7/8/9/10) age groups. Several CAF samples are shown in Figure 1.

To make the CAF dataset, we hired 5 web-search workers who did not know each other. Each was asked to search for celebrities, politicians and athletes with photos that had specifications of ages. Each's search results were reaffirmed by the other four, and we verified the overall collection. We will follow the CVF or other necessary regulations to release this dataset.

# 2 Determination of Face Verification Threshold

#### (Line 421. See supplementary document for more details.)

We use the CAF dataset to rectify the face verification confidence measures made by the Face++ APIs [7] and the ArcFace [4]. We made 9,253 intra pairs and 300,000 inter pairs out of the dataset, and used Face++ APIs to measure the similarity confidence of each pair. For example, when defining the similarity confidence threshold for verifying the faces in  $G_{10}0$ , age  $0 \sim 2$ , we selected a desired FAR for the inter pairs, and obtained the corresponding similarity

<sup>&</sup>lt;sup>1</sup> The following 10 age intervals: 0–2, 3–6, 7–9, 10–14, 15–19, 20–29, 30–39, 40–49, 50–69 and  $\geq$ 70 years are labeled as  $G_{10}0, G_{10}1, ..., G_{10}9$ , respectively.



Fig. 1: CAF dataset samples, number under the image is the ground-truth age for each image

confidence threshold. The similarity confidence threshold for  $FAR=10^{-4}$  is 61.8, as shown in Table 2 in the main paper. Following the same way for other age groups, we obtained the corresponding CAF-rectified thresholds, 68.9, 72.7, ..., 65.2, compared with the common threshold (76.5) in the parentheses. All are shown in Table 2 in the main paper.

Table 1 shows the comparison of three performance metrics: 1) Face++ APIs with the common threshold (76.5), 2) Face++ APIs with CAF-rectified thresholds, and 3) the ArcFace [4] rectified by the CAF in the same manner as we did for Face++ APIs. When extracting the similarity by using ArcFace, we computing the cosine similarity between the facial features extracted by ArcFace. For all three metrics, the verification rates drop substantially when handling the faces for the youngest groups,  $G_{10}0$  and  $G_{10}1$ , and the drop is also clear for the most senior  $G_{10}9$ . However, the CAF-rectified Face++ APIs outperform the other two with clear robustness across all age groups.

# 3 Network Settings

(Line 189. See Supplementary document for details on network settings) Table 2, Table 3 and Table 4 show the network settings of the encoder  $G_{en}$ , the decoder  $G_{de}$  and the discriminator  $D_p$ , respectively.

Table 1: Comparison in face verification rates on the CAF by using 3 performance metrics: 1) Face++ APIs with common threshold (76.5), 2) Face++ APIs with rectified thresholds, and 3) the ArcFace with rectified thresholds.

	$G_{10}0$	$G_{10}1$	$G_{10}2$	$G_{10}3$	$G_{10}4$	$G_{10}5$	$G_{10}6$	$G_{10}7$	$G_{10}8$	$G_{10}9$
Face++ $(76.5)$	12.3	26.5	60.0	80.5	90.4	91.7	89.5	79.1	54.5	26.4
Face++	54.3	62.1	82.1	86.4	88.3	92.4	90.2	91.4	89.5	85.5
ArcFace	28.6	42.3	55.4	78.1	85.1	94.3	92.1	82.6	44.2	22.1

Table 2: Encoder architecture, k denotes kernel size , c is input channel number and s is stride.

	$G_{en}$			$B_{id}$		Bag			
Layer name	k, c, s	Dim.	Layer name	k, c, s	Dim.	Layer name	k, c, s	Dim.	
Input	-	$(3+N) \times 1024^2$		-			-		
Conv1	$1 \times 1,  32,  1$	$32 \times 1024^{2}$	Input	-	$512 \times 16^{2}$	Input	-	$512 \times 16^{2}$	
ResBlock1	$\begin{bmatrix} 3 \times 3, 32, 1 \\ 3 \times 3, 64, 2 \\ 1 \times 1, 64, 2 \end{bmatrix}$	$64\times512^2$	$ResBlock1_{id}$	$\begin{bmatrix} 3 \times 3, 512, 1 \\ 3 \times 3, 512, 1 \\ 1 \times 1, 512, 1 \end{bmatrix}$	$512 \times 16^2$	$ResBlock1_{ag}$	$\begin{bmatrix} 3 \times 3, 512, 1 \\ 3 \times 3, 512, 2 \\ 1 \times 1, 512, 2 \end{bmatrix}$	$512\times8^2$	
ResBlock2	$\begin{bmatrix} 3\times 3, 64, 1\\ 3\times 3, 128, 2\\ 1\times 1, 128, 2 \end{bmatrix}$	$128\times 256^2$	$ResBlock2_{id}$	$\begin{bmatrix} 3\times 3, 512, 1\\ 3\times 3, 512, 1\\ 1\times 1, 512, 1 \end{bmatrix}$	$512 \times 16^2$	$ResBlock2_{ag}$	$\begin{bmatrix} 3\times 3, 512, 1\\ 3\times 3, 512, 2\\ 1\times 1, 512, 2 \end{bmatrix}$	$512 \times 4^2$	
ResBlock3	$\begin{bmatrix} 3\times 3, 128, 1\\ 3\times 3, 256, 2\\ 1\times 1, 256, 2 \end{bmatrix}$	$256\times 128^2$				$Conv1_{ag}$	$3\times3,512,1$	$512 \times 4^2$	
ResBlock4	$\begin{bmatrix} 3 \times 3, 256, 1 \\ 3 \times 3, 512, 2 \\ 1 \times 1, 512, 2 \end{bmatrix}$	$512\times 64^2$				$AvgPool1_{ag}$	$4 \times 4, -, 4$	512	
ResBlock5	$\begin{bmatrix} 3 \times 3, 512, 1 \\ 3 \times 3, 512, 2 \\ 1 \times 1, 512, 2 \end{bmatrix}$	$512\times 32^2$				$FC_{ag} \times 8$	_	512	
ResBlock6	$\begin{bmatrix} 3 \times 3, 512, 1 \\ 3 \times 3, 512, 2 \\ 1 \times 1, 512, 2 \end{bmatrix}$	$512\times 16^2$							

 Table 3: Decoder Architecture. ToRGB shows the output size of the generated images from different layer.

Layer name	k, c, s	Dim.	ToRGB
Input	-	$512 \times 16^2$	-
Conv1	$3\times3, 512, -$	$512 \times 16^2$	$3 \times 16^2$
ConvTrans1	$3 \times 3, 512, 2$	$512 \times 32^2$	$3 \times 32^2$
00002	$3 \times 3, 312, 1$		
ConvTrans2 Conv3	$3 \times 3, 512, 2$ $3 \times 3, 512, 1$	$256\times 64^2$	$3\times 64^2$
	0 / 0,012,1		
ConvTrans3 Conv4	$3 \times 3, 512, 2$ $3 \times 3, 256, 1$	$256\times 128^2$	$3\times 128^2$
ConvTrans4 Conv5	$\begin{array}{c} 3\times 3, 256, 2 \\ 3\times 3, 128, 1 \end{array}$	$128\times 256^2$	$3 \times 256^2$
ConvTrans5 Conv6	$\begin{array}{c} 3\times3,128,2\\ 3\times3,64,1 \end{array}$	$64\times512^2$	$3 \times 512^2$
ConvTrans6 Conv6	$\begin{array}{c} 3\times3,64,2\\ 3\times3,32,1 \end{array}$	$32\times 1024^2$	$3\times 1024^2$

Dim Layer nam Dim k, c, s k, c, sk, c, s Laver nar  $3 \times 1024^2$  $32 \times 1024$ Input  $1 \times 1, 32$  $3 \times 3, 32, 1$  $3 \times 3, 64, 2$ × 3, 32, × 3, 32, 3, 32,  $64 \times 512^2$ ResBlock1 $64 \times 512^2$ ResBlock  $64 \times 512$  $64 \times 512^{4}$  $3 \times 3, 64, 2$ ResBlock1 $3 \times 3, 64, 2$ ResBlock1 $3 \times 3, 64, 2$ × 1,64,  $\times 1,64,$  $\times 1,64,$  $\times 1, 64, 2$  $3 \times 3, 64, 3 \times 3, 128, 1 \times 1, 128,$  $3 \times 3, 64,$   $3 \times 3, 128,$   $1 \times 1, 128,$  $3 \times 3, 64, 1$  $3 \times 3, 128, 2$  $\times 3, 64, 1$   $\times 3, 128, 2$ ResBlock2  $128 \times 256$ ResBlock2  $128 \times 256$ ResBlock2  $128 \times 256$ ResBlock2  $128 \times 256^{2}$  $3 \times 3.128$ .  $\times 3.128.$  $\times$  3.128.  $256 \times 128$ ResBlock3 3 × 3, 256, 2 ResBlock3  $3 \times 3.256$  $256 \times 128^{2}$ ResBlock3 × 3.256.5  $256 \times 128^{2}$ 3.256.256 3. 256. ResBlock4 3,512,2  $512 \times 64^{2}$ ResBlock4512  $512 \times 64^{2}$ ResBlock4 3,512,  $512 \times 64^2$ 3 512 ResBlock5  $512 \times 32^{2}$ ResBlock  $512 \times 32^{2}$ 3, 512,  $\frac{1 \times 1,512,2}{3 \times 3,512,1}$  $3 \times 3,512,2$ ResBlock6  $512 \times 16^{\circ}$  $512 \times 16^2$ ResBlock6 × 3, 512, 5 < 1, 512.  $\times 1, 512, \\ \times 3, 512,$  $3 \times 3, 512, 2$  $512 \times 8^2$ ResBlock7 ResBlock8  $3 \times 3,512,2$  $512 \times 4^2$  $\times$  1, 512, 2 Minib  $513 \times 4^2$  $512 \times 4^2$  $3 \times 3.512.1$ 

Table 4: Discriminator Architecture, k, c, s are defined as in Table 2

### 4 More Experiments

#### 4.1 More Datasets

(Line 373. The experiments on the MORPH [21] and CACD [2] in the supplementary document. Line 430. See Supplementary Materials for more information about the CAF and MIVIA datasets.)

The MORPH Album-2[12] MORPH is one of the largest publicly available longitudinal face database with mugshot images, and it includes the meta data for race, gender, date of birth, and date of acquisition. It contains 55,134 images of 13,000 individuals with ages from 16 to 77 years, captured in controlled conditions. The largest age gap for the same individual in the MORPH is 5 years.

**CACD**[3] The CACD offers 163,446 face images of 2k celebrities captured in the wild with age ranging from 16 to 62. Is one of the largest publicly available cross-age face dataset. The largest age gap for the same individual in the CACD is 10 years. However, it contains lots of mislabeled and mismatched data.

MIVIA Age Dataset[5] The MIVIA Age Dataset is composed of 575,073 images of more than 9.000 identities, taken at different ages. The images are extracted from the VGGFace2 dataset and annotated with age labels by means of a knowledge distillation technique [5], making the dataset very heterogeneous in terms of face size, illumination conditions, facial pose, gender and ethnicity.

#### 4.2 More Details about Training

(Line 385. See supplementary document for more details about data preprocessing, other training and testing settings.) As for the preprocessing, all faces were aligned by the Face Alignment Network (FAN) [2] and cropped to  $256^2$  for the MORPH and CACD, and  $1024^2$  for the FFHQ-Aging. For the 5-fold subjectindependent cross validation on MORPH, each fold has 2,586 subjects, with 4,467, 3,030, 2,205 and 639 face images for  $G_40 \sim G_44$ , respectively. For CACD, each fold has 400 subjects, with 10,079, 8,635, 7,964 and 6,011 face images for the 4 age groups. The training and testing data split on FFHQ-Aging follows the same in [8].

We chose the Adam optimizer to train G and D at learning rate  $2e^{-4}$  on an Nvidia RTX Titan GPU. The batch size is 8 for training  $256^2$  images, and 4 for training  $1024^2$  images. For each iteration, we updated  $G = [G_{en}, G_{de}]$  while keeping  $D_p$  as was from the previous iteration, then updated  $D_p$  while keeping G as was from the previous iteration, and kept on.

The training on MORPH (50, 196 images,  $256^2$ ) took 36 hours, CACD (130, 756 images,  $256^2$ ) took 4 days, and FFHQ-Aging (60, 000 images,  $1024^2$ ) took 4 weeks. At runtime, a  $256^2$  image takes 11 sec to generate, and a  $1024^2$  image takes 18 sec.

#### 4.3 Performance Comparisons on MORPH and CACD

(Line 373. The experiments on the MORPH [21] and CACD [2] in the supplementary document.)

Table 5 presents the comparison with state-of-the-art approaches, including the GLCA-GAN [9], the WL-GAN [10] and the CPA-GAN [13], on the MORPH and CACD by using the Face++ APIs as the evaluation tool without any rectification. The performance of all approaches are duplicated from their papers as their codes are unavailable. Table 5 can be summarized as follows:

- The AgeTransGAN outperforms the WL-GAN on the CACD, but is slightly outperformed on the MORPH.
- To show the capability of handling bi-directional age transformation, the AgeTransGAN also presents the performance on age regression. Almost all other approaches in Table 5 require an additional model for handling age regression.

Table 5: Comparison with SOTA approaches on MORPH & CACD, using Face++ APIs [7]. Best two in each category shown in **boldface**. Progression transfers Group- $G_40$  to others, regression transfers Group- $G_43$  to others

	MORPH							CACD				
	Progression			Regression			Progression			Regression		
Age group	31-40	41-50	50+	41-50	31-40	30-	31-40	41-50	50+	41-50	31-40	30-
Raw Age data	38.87	48.03	58.29	48.03	38.87	27.93	38.92	46.95	53.75	46.95	38.92	30.96
	Mean Error / Verification Rate (%)											
GLCA-GAN	<b>0.23</b> /97.66	3.61/96.67	8.61/91.85	-	-	-	1.72/97.72	2.07/94.18	2.85/92.29	-	-	-
WL-GAN	0.13/100	0.19/100	0.68/98.26	-	-	-	0.37/99.76	0.88/98.74	0.66/98.44	-	-	-
CPA-GAN	0.75/100	0.87/100	1.75/99.98	-	-	-	1.60/100	1.08/100	0.30/99.88	-	-	-
AgeTransGAN	0.36/ 100	0.65/100	0.56/100	3.77/100	2.39/100	0.58/99.39	0.32/100	0.64/100	0.43/100	3.14/100	3.46/100	1.06/97.52



Fig. 2: Qualitative comparisons for age progression (Group-0 to 3) samples made by the AgeTransGAN and SOTA approaches



Fig. 3: Age progression and regression samples on Morph and CACD made by the AgeTransGAN.

#### 4.4 Baseline Performance without $L_{px}$ and $L_{pl}$

# (Line 480. A comparison of the baselines with and without these losses is given in the supplementary document )

The pixel-wise attribute loss  $L_{px}$  can better maintain the color tone of the input, and the perceptual path length regularization  $L_{pl}$  can better preserve the gender and racial background of the input. Figure 4 shows a qualitative comparison of the baseline (B/L) model and those without either loss.

#### 4.5 Comparison of Margin in Triplet loss

(Line 373. The margin  $m_t$  in is experimentally determined as 0.5 out of a comparison study reported in the supplementary document.)

Table 6 shows the face verification rates and target age generation performance in EAM for different margins  $m_t$  in the triplet loss, Equation (3) in the main paper. When  $m_t = 0.1$ , the identity preservation appears slightly better than the other two on the cost of poorer target age generation. When  $m_t = 1$ , the



Fig. 4: Qualitative comparison of the baseline (B/L) model and the B/L without  $L_{px}$  or without  $L_{pl}$ 

Table 6: Performance for triplet loss margin settings with rectified and common thresholds for verification, and used for our age estimator. Best one in each category shown in **boldface**.

Age group	$0-2(G_{10}0)$	$3-6(G_{10}1)$	$7-9(G_{10}2)$	$10-14(G_{10}3)$	$15-19(G_{10}4)$	$30-39(G_{10}6)$	$40-49(G_{10}7)$	$50-69(G_{10}8)$	$70+(G_{10}9)$		
Face Verification Rate (%), Rectified Threshold (Common Threshold)											
$m_t = 0.1$	82.7(34.6)	97.4(88.4)	100(98.3)	100(99.3)	99.5(100)	99.5(99.5)	99.5(98.1)	100(96.6)	100(94.9)		
$m_t = 0.5$	80.3(10.3)	96.5(86.3)	95.9(85.4)	95.8(86.7)	100(100)	100(100)	100(97.7)	97.6(85.7)	96.8(84.7)		
$m_t = 1$	76.2(22.7)	95.1(85.1)	93.7(83.8)	94.4(8.1)	90.2(92.3)	82.5(81.7)	85.9(82.5)	90.6(85.2)	90.4(84.2)		
				EAM, Our	s / Mean Erro	or					
Raw data	1.5/-	4.9/-	8.6/-	12.8/-	18.9/-	31.9/-	43.9/-	57.2/-	68.9/-		
$m_t = 0.1$	1.8/0.3	4.0/0.9	10.1/1.5	13.6/0.8	20.6/1.7	31.1/0.8	41.5/2.4	53.9/3.3	68.2/0.7		
$m_t = 0.5$	1.1/0.4	4.5/0.4	8.8/0.2	13.5/0.7	18.7/0.2	32.3/0.4	41.7/2.2	55.5/1.7	68.4/0.5		
$m_t = 1$	1.7/0.2	4.4/0.5	7.8/0.8	13.2/0.4	17.6/1.3	31.4/0.5	42.5/1.4	55.0/0.2	68.8/0.1		

target age generation appears slightly better, but on the cost of poorer identity preservation. Therefore, we chose  $m_t = 0.5$ .

#### 4.6 User Study

We conducted a user study to justify our approach and offer an additional comparison to other approaches. But due to page limit, we can only report it in this supplementary document. Three metrics were considered: 1) Target age generation of the generated images, 2) Identity preservation of the source by the generated images, and 3) The overall identity preservation and target age generation. We hired 30 workers to participate, and each was randomly given 80 pairs of (source, synthetic) images from the FFHQ-aging testing set for each metrics. The results are given in Table 7. The AgeTransGAN, LATS and DLFS demonstrate similar performance for identity preservation along, and for the target age generation along, although the AgeTransGAN slightly outperforms others. However, for the overall performance evaluation, in which each user was asked to justify the combined performance for identity preservation and target age generation, the AgeTransGAN outperforms others with a clear margin.

	01		
Methods	Identity preservation	Age accuracy	Overall better
Ours	84.5%	85.5%	32.3%
LATS[11]	83.5%	83.1%	20.2%
DLFS[6]	83.8%	85.3%	28.8%
SAM[1]	69.5%	84.5%	18.7%

Table 7: User study performance for 3 metrics

#### 4.7 Qualitative Comparisons

**Different Source Age Groups** In the main paper, we only report the performance with  $G_{10}5$  as the only source group. The performance of using ALL other age groups as the source groups to transfer to each specific age group is presented in Table 8, with sample images shown in Figure 5.

Table 8: Performance of using ALL age groups as source groups to be transferred to each specific age group

$0-2(G_{10}0)$	$3-6(G_{10}1)$	$7-9(G_{10}2)$	$10-14(G_{10}3)$	$15-19(G_{10}4)$	$20-29(G_{10}5)$	$30-39(G_{10}6)$	$40-49(G_{10}7)$	$50-69(G_{10}8)$	$70+(G_{10}9)$		
Verification Rate (%) with common threshold											
84.6	97.2	99.5	98.9	100	100	99.7	99.7	98.3	93.4		
Estimated Mean Age, Our estimator (Face++)											
1.4(14.4)	3.8(23.1)	7.4(25.2)	13.6(28.0)	18.3(24.1)	23.8(26.6)	33.1(40.6)	41.9(53.3)	53.0(65.4)	68.8(76.4)		

More Qualitative Comparisons with Ground-Truth Images in CAF Figure 6 to Figure 8 shows more samples of the generated images compared with the ground-truth images in the CAF. SAM [1] works fine for target age generation, but poorly for identity preservation. Both our approach and the DLFS [6] demonstrate better visual quality in identity preservation and target age generation. However, the DLFS frequently generates artifacts, especially on the younger groups, e.g.,  $G_{10}0 \sim G_{10}3$ . Besides, as emphasized in the main paper, both DLFS and LATS only offers the transfer to 6 age groups. The AgeTransGAN and SAM offer the transfer to all 10 age groups.



Fig. 5: Generated image samples by using different age groups from FFHQ-aging as the source images



Fig. 6: Qualitative comparison of the generated images with the ground truth in the CAF dataset.



Fig. 7: Qualitative comparison of the generated images with the ground truth in the CAF dataset.



Fig. 8: Qualitative comparison of the generated images with the ground truth in the CAF dataset.

# References

- Alaluf, Y., Patashnik, O., Cohen-Or, D.: Only a matter of style: Age transformation using a style-based regression model. ACM Transactions on Graphics (TOG) 40(4), 1–12 (2021)
- 2. Bulat, A., Tzimiropoulos, G.: How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In: ICCV (2017)
- 3. Chen, B.C., Chen, C.S., Hsu, W.H.: Face recognition and retrieval using cross-age reference coding with cross-age celebrity dataset. TMM (2015)
- Deng, J., Guo, J., Xue, N., Zafeiriou, S.: Arcface: Additive angular margin loss for deep face recognition. In: CVPR (2019)
- Greco, A., Saggese, A., Vento, M., Vigilante, V.: Effective training of convolutional neural networks for age estimation based on knowledge distillation. Neural Computing and Applications (2021)
- 6. He, S., Liao, W., Yang, M.Y., Song, Y.Z., Rosenhahn, B., Xiang, T.: Disentangled lifespan face synthesis. In: ICCV (2021)
- 7. Inc., M.: Face++ research toolkit. http://www.faceplusplus.com
- 8. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: CVPR (2019)
- Li, P., Hu, Y., Li, Q., He, R., Sun, Z.: Global and local consistent age generative adversarial networks. In: ICPR (2018)
- 10. Liu, Y., Li, Q., Sun, Z.: Attribute-aware face aging with wavelet-based generative adversarial networks. In: CVPR (2019)
- 11. Or-El, R., Sengupta, S., Fried, O., Shechtman, E., Kemelmacher-Shlizerman, I.: Lifespan age transformation synthesis. In: ECCV (2020)
- 12. Ricanek, K., Tesafaye, T.: Morph: A longitudinal image database of normal adult age-progression. In: FG (2006)
- Yang, H., Huang, D., Wang, Y., Jain, A.K.: Learning continuous face age progression: A pyramid of gans. TPAMI (2019)