arXiv:1908.08239v1 [cs.CV] 22 Aug 2019

# Progressive Face Super-Resolution via Attention to Facial Landmark

Deokyun Kim*
deokyunkim@kaist.ac.kr

Minseon Kim*
minseonkim@kaist.ac.kr

Gihyun Kwon*
cyclomon@kaist.ac.kr

Dae-Shik Kim
daeshik@kaist.ac.kr

School of Electrical Engineering,
Korea Advanced Institute of Science
and Technology,
Republic of Korea

## Abstract

Face Super-Resolution (SR) is a subfield of the SR domain that specifically targets the reconstruction of face images. The main challenge of face SR is to restore essential facial features without distortion. We propose a novel face SR method that generates photo-realistic $8\times$ super-resolved face images with fully retained facial details. To that end, we adopt a progressive training method, which allows stable training by splitting the network into successive steps, each producing output with a progressively higher resolution. We also propose a novel facial attention loss and apply it at each step to focus on restoring facial attributes in greater details by multiplying the pixel difference and heatmap values. Lastly, we propose a compressed version of the state-of-the-art face alignment network (FAN) for landmark heatmap extraction. With the proposed FAN, we can extract the heatmaps suitable for face SR and also reduce the overall training time. Experimental results verify that our method outperforms state-of-the-art methods in both qualitative and quantitative measurements, especially in perceptual quality.

## 1 Introduction

Face Super-Resolution (SR) is a domain-specific SR which aims to reconstruct High Resolution (HR) face images from Low Resolution (LR) face images while restoring facial details. When enlarging the LR face images to high-resolution images, the HR images suffer from face distortion. The finer details of faces disappear incurring misperception of facial attributes on faces. In an attempt to address this problem, the previous studies [13, 26] embedded additional facial attribute vectors into network feature maps to reflect facial attributes in super-resolved face images. These approaches require prior information for face SR; however, the additional information is difficult to obtain in the wild. Other studies incorporate facial landmark information by employing auxiliary networks such as face alignment

* The authors have equally contributed.
Github: https://github.com/DeokyunKim/Progressive-Face-Super-Resolution

network [1, 25], and prior estimation network [5]. However, these approaches tend to concentrate on the localization of facial landmarks without sufficient consideration of the facial attributes in the areas around landmarks.

Different from the previous works, we propose a face SR method which restores original facial details more precisely by giving strong constraints to the landmark areas. To stably generate photo-realistic $8\times$ upscaled images, we adopt a progressive training method [2, 11, 19, 20] which grows both generator and discriminator progressively. We also introduce a novel facial attention loss which makes our SR network to restore the accurate facial details. The attention loss is applied in both the intermediate and the last step of our progressive training.

Constraining the outputs by applying the attention loss at each step, the output images of each step reflect more accurate facial details. To obtain the attention loss, we extract the heatmaps from the pre-trained Face Alignment Network (FAN). The extracted heatmaps are used as weights of the pixel difference of the adjacent areas to the landmarks. Instead of using the state-of-the-art FAN [4], we suggest a compressed network of FAN, called distilled FAN, which is trained by a hint-based method [17]. The distilled FAN delivers comparable performance to the original FAN while being much more compact. With our approach, we obtain SR-oriented landmark heatmaps as well as significantly reduce the overall training time. Therefore, our method generates super-resolved face images which successfully reflect accurate details of facial components.

For the evaluation, we measure the performance of our method on both aligned and unaligned face images from CelebA [14] and AFLW [10] datasets. To compare the quality of our results, we calculate the conventional measurements of the average Peak Signal to Noise Ratio (PSNR), Structural SIMilarity (SSIM) [22], and Multi-Scale Structural SIMilarity (MS-SSIM) [21]. By conducting an ablation study, we verify that the proposed loss function is beneficial to super-resolving LR face images; we demonstrate the superiority of our method by comparing the results with those of previous studies. We further conduct Mean-Opinion-Score (MOS) [12] test to measure the perceptual quality. The experimental results show that our network successfully generates high-fidelity face images, accurately preserving the original features around the facial landmarks. In summary, our contributions are as follows:

1. To the best of our knowledge, progressive training method is used in natural image SR, but this is the first method which leverages the progressive training method for face SR. We give constraints to each step of the SR network and generate high-quality face images reflecting details of facial components.

2. Facial attention loss makes the SR network learn to restore facial details with the method of focusing on the adjacent area of facial landmarks, which is verified by our super-resolved results.

3. We compress the state-of-the-art FAN into a smaller network using hint-based method. With the distilled FAN, we are able to extract meaningful landmark heatmaps which are more suitable for a face SR task and reduce the overall training time.

## 2    Related work

Utilizing facial information, such as facial attributes and spatial configuration of facial components, is the key factor in face SR. Yu et al. [26] interweave multiple spatial transformer networks to satisfy the requirement of face alignment as well as embeds facial attribute vec-
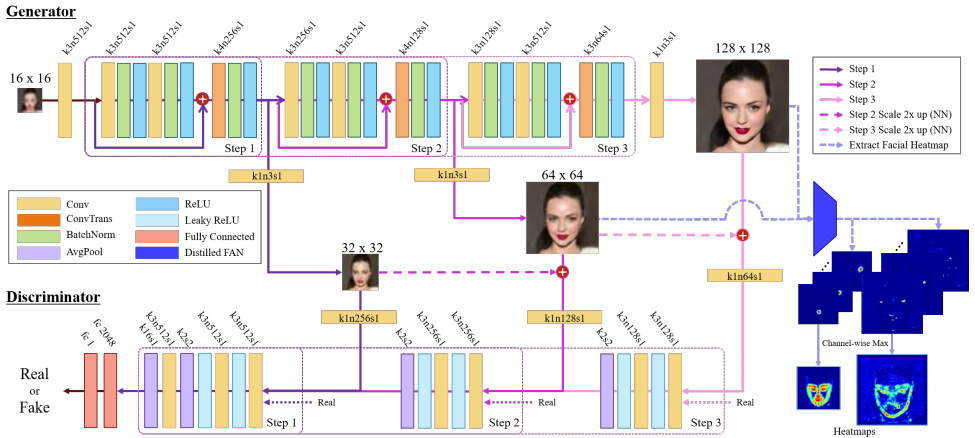
Figure 1: Our network architecture overview. *(k : kernel size, n : output channel, s : stride)*

tors to lower the ambiguity in facial attributes. Lee et al. [13] fuses the information of both image domain and attribute domain in order to reflect facial attributes in super-resolved images. These methods preserve facial attributes indicated by facial attribute vectors. However, attribute vectors are not only difficult to acquire in the wild but also limited to describing partial facial attributes.

While preserving the facial landmarks location, Chen et al. [5] propose the face shape prior estimation network, which provides a solution for accurate geometry estimation obtained from coarse HR face images. Yu et al. [25] estimate the spatial position of their facial components to preserve the original spatial structure in upscaled face images. In addition, Bulat et al. [1] propose a heatmap loss to localize landmarks in super-resolved images so as to upscale input face images 4×. Although these methods preserve the spatial configuration of facial components, they fail to fully reflect accurate facial attributes. In contrast to the previous works, our method carefully considers the facial attributes around the landmarks to restore the facial details without prior information as well as preserves landmark location.

# 3 Approach

In this section, we describe our methods for the enhanced face SR. To generate the high-fidelity super-resolved face images that reflect the facial attributes of target face images, three main approaches are used: progressive training, facial attention loss, and distillation of Face Alignment Network (FAN).

## 3.1 Progressive Face SR Network

The overview of our network architecture is shown in Figure 1. To incorporate the adversarial loss, our architecture is composed of the generator network, which is our face SR network, and the discriminator network. To train the generator and the discriminator stably, we construct both the networks which consist of layers stacked by steps. Our generator network consists of three residual blocks [6] with batch normalization layers (BatchNorm), transpose

convolution layers (ConvTrans), and Rectifier Linear Unit (ReLU) as an activation function. The discriminator network has a corresponding architecture to the generator network, which is comprised of convolution layers (Conv), average pooling layers (AvgPool), and Leaky ReLU. To improve the discriminator performance, we calculate the standard deviation of input batch, then replicate the value into a one-dimensional feature map, and concatenate it to the end of the discriminator [19]. We use additional convolution layers in each step in order to convert the intermediate feature maps into RGB images, and vice versa.

In Step 1, each network employs one block for training and learns to upscale images $2\times$. These $2\times$ upscaled outputs from the generator go through the corresponding part of the discriminator, and the outputs are then compared with target images. In Step 2, the $2\times$ upscaled outputs from Step 1 are upscaled $2\times$ again by nearest-neighbor interpolation, and then the interpolated outputs are added to $4\times$ upscaled outputs from Step 2. This process is expressed as follows: $(1-\alpha) * f(G^{N-1}(I)) + \alpha * G^N(I)$, where $G$ is our SR network, $f$ is nearest neighbor (NN) interpolation, $I$ and $N \in \{2,3\}$ denote input images and number of step, respectively. A weight scale $\alpha$ increases linearly from zero to one. The upscaled outputs are compared to the corresponding target images through the discriminator. The same procedure above is implemented in Step 3 ($8\times$). The method allow the network learn super-resolving face images with different loss in each step effectively and stably.

## 3.2   Facial Attention Loss

We propose the *facial attention loss* to restore the attributes of the adjacent area to the facial landmarks. This *facial attention loss* makes the face SR network focus on the facial details around the predicted landmark area by element-wise multiplying landmark attention heatmaps $M^*$, and the $L1$ distance between the upscaled images and the corresponding target images. To achieve this, we employ facial landmark heatmaps which contain landmark location information. The *facial attention loss* is defined as:

$$L_{attention} = \frac{1}{r^2 WH} \sum_{x=1}^{rW} \sum_{y=1}^{rH} (M^*_{x,y} \cdot |I^{HR}_{x,y} - G(I^{LR})_{x,y}|) \tag{1}$$

where $G$ is face SR network, $I^{HR}$ and $I^{LR}$ are target face images and input LR images, respectively. The landmark attention heatmap $M^*$ is channel-wise max values of the target heatmap $M$ generated from target face images. To compensate for the variance between the landmarks, the heatmap $M$ is min-max normalized into [0,1]. The heatmap $M$ has the dimension of $N \times rW \times rH$, where $N$ is number of landmarks, $W$ and $H$ are width and height of the input image. The upscale factor $r$ is set to be 4 and 8 in Step 2 and Step 3, respectively. To give attention at images with enough information, we adjust facial attention loss at upscaled outputs with size of $64\times64$ and $128\times128$.

## 3.3   Distilled Face Alignment Network

The state-of-the-art FAN [4] predicts the location of all landmarks including occluded landmarks exploiting multiple-scale feature maps from four-stacked Hourglass architecture which consists of encoder-decoder and skip-layer[15]. This approach predicts landmark locations based on heatmap values, which are highly concentrated around the landmark points. However, for face SR, we aim to give attention to the overall facial landmark area in the images except for the occluded landmark area.
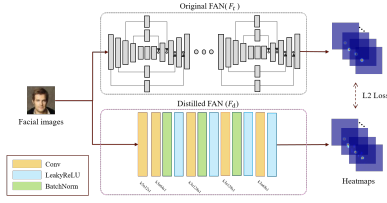
Figure 2: Distilled face alignment network hint based training

To generate heatmaps suitable for giving attention to accurate facial landmark area, we construct the network with neither encoder-decoder architecture nor skip-layer so as to predict landmarks based on single-scale feature maps. Also, in order to reduce overall training time and achieve comparable performance to state-of-the-art FAN, we compress the FAN into the network shown in Figure 2 based on a hint-based training method [11]. We train the distilled FAN to minimize $\sum (F_d (I) - F_t (I))^2$, where $F_d$ and $F_t$ represents our distilled FAN and original FAN, $I$ denotes input face images. This approach has two advantages: it provides face SR-oriented heatmaps, and it reduces the overall training time from $\sim 3$ days to $\sim 1$ days in our experiments.

## 3.4 Overall Training Loss

**MSE loss** We use the pixel-wise Mean-Square-Error (MSE) loss to minimize the distance between the HR target image and the super-resolved image.

$$L_{pixel} = \frac{1}{r^2 WH} \sum_{x=1}^{rW} \sum_{y=1}^{rH} \left( I_{x,y}^{HR} - G\left(I^{LR}\right)_{x,y} \right)^2 \tag{2}$$

**Perceptual loss** A perceptual loss [7] is proposed to prevent generating blurry and unrealistic face images, and to obtain more realistic HR images. The loss over the pre-trained VGG16 [18] features at a given layer $i$ is defined as:

$$L_{feat/i} = \frac{1}{W_i H_i} \sum_{x=1}^{W_i} \sum_{y=1}^{H_i} \left( \phi_i \left(I^{HR}\right)_{x,y} - \phi_i \left(G\left(I^{LR}\right)\right)_{x,y} \right)^2 \tag{3}$$

where $\phi_i$ denotes the feature map obtained after the last convolutional layer of the $i-$th block.
**Adversarial loss** We use the WGAN Loss [5] to stabilize the training process. In WGAN, the loss function is defined as the Wasserstein distance between the distribution of target $I^{HR} \sim P_r$ and those of the generated images $\tilde{I} \sim P_g$. For further improvement of training stability, we apply the Gradient Penalty term proposed in WGAN-GP [23], which enforces the Lipschitz -1 condition of the discriminator. $\hat{I}$ is a randomly sampled image among the samples from $P_r$ and $P_g$. Therefore, the loss function is as follows:

$$L_{WGAN} = \mathbb{E}_{I^{HR} \sim P_r}[D(I^{HR})] - \mathbb{E}_{\tilde{I} \sim P_g}[D(\tilde{I})] + \lambda \mathbb{E}_{\hat{I} \sim P_{\hat{I}}}[||\nabla_{\hat{I}} D(\hat{I})_2 - 1||^2] \tag{4}$$

**Heatmap Loss** As proposed by [1], the heatmap loss improves the structural consistency of face images by minimizing the distance between the heatmaps of both generated images and target ones. The heatmap loss function is described as:

$$L_{heatmap} = \frac{1}{r^2 NWH} \sum_{n=1}^{N} \sum_{x=1}^{rW} \sum_{y=1}^{rH} (M_{x,y}^n - \tilde{M}_{x,y}^n)^2 \tag{5}$$

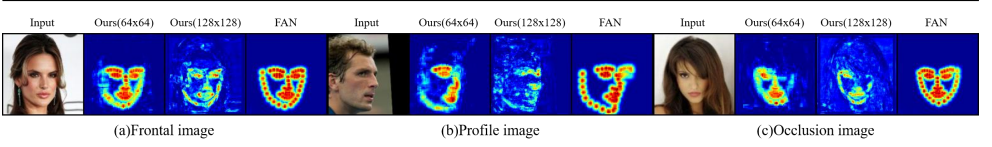where $N$ is the number of landmarks, $M$ and $\tilde{M}$ are calculated as $M = F_d(I^{HR})$ and $\tilde{M} = F_d(G\left(I^{LR}\right))$.

(a)Frontal image          (b)Profile image          (c)Occlusion image

Figure 3: Distilled FAN results(Ours) comparison with FAN results [4].



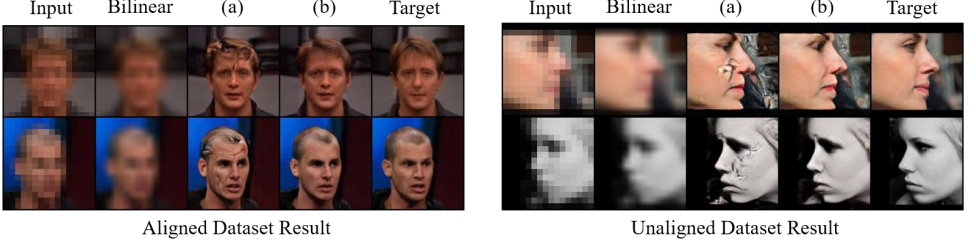Aligned Dataset Result          Unaligned Dataset Result

Figure 4:  Our image results (a) with the original FAN, (b) with the distilled FAN.

**Overall training loss** Since the landmark losses are applied to Step 2 & 3, we intend not to include the $L_{heatmap}$ and $L_{attention}$ in Step 1. The final loss term is shown as:

$$L_{Ours} = \alpha L_{pixel} + \beta L_{feat} + \gamma L_{WGAN}, \textit{ at step 1}$$
$$L_{Ours} = \alpha L_{pixel} + \beta L_{feat} + \gamma L_{WGAN} + \lambda L_{heatmap} + \eta L_{attention}, \textit{ at step 2 & step3}$$

(6)

where $\alpha$, $\beta$, $\gamma$, $\lambda$ and $\eta$ are the corresponding weights.

# 4  Experiments

## 4.1  Implementation details

**Datasets** In our experiments, we use two different datasets: aligned dataset and unaligned one. The aligned CelebA dataset [14] is used to test how accurately the facial details can be restored. The unaligned CelebA and AFLW [10] datasets are used to verify the applicability of our face SR network in real world. The aligned face images are cropped into square. The face images of the unaligned dataset are cropped based on the bounding box areas. The cropped images are resized into $128 \times 128$ pixels to be used as targets of Step 3, and bilinearly downsampled into $64 \times 64$ pixels as targets of Step 2, $32 \times 32$ pixels as targets of Step 1, and $16 \times 16$ pixels as LR inputs. We use all 162,770 images as a training set, and 19,867 images as a test set from aligned CelebA dataset. For the unaligned dataset, we use 80,000 cropped images from unaligned CelebA, and 20,000 from AFLW for training. As a test set, 5,000 images from CelebA and 4,384 images from AFLW are used.

**Training details** We implement our face SR network using PyTorch [16]. We train our networks using the Adam optimizer [9] with a learning rate of $1 \times 10^{-3}$, and the mini-batch size of 16. The training iteration of each step is set by hyperparameter. In our model, we train our model 50K, 50K and 100K iterations, empirically. Running totally 200K iterations takes ~1 day on single Titan X GPU. In addition, we train the distilled FAN using the Adam optimizer with a learning rate of $1 \times 10^{-4}$, mini-batch size of 16, and 100K iterations.

## 4.2 Distilled FAN Results

In this section, we compare our distilled FAN to the original FAN [4]. To verify how similarly the distilled FAN predicts landmark location compared to the original FAN, we use Normalized Mean Error (NME) metric [4, 27] between the predicted landmark locations from the distilled and the original FAN. The NME is calculated as $NME = \frac{1}{N} \sum_{k=1}^{N} \frac{\|g_k - p_k\|_2}{d}$, where $g$ denotes the landmark from original FAN, $p$ is the corresponding prediction from the distilled model, and $d$ is the facial image size. The NME evaluation results and the number of parameters are shown in Table 1. The results show that our distilled FAN predicts facial landmarks with comparable performance, and it has much fewer parameters than the original FAN.

As shown in Figure 3(b) and (c), the output heatmaps of the original FAN have high values around the landmark points even in the occluded area, but the output heatmaps of distilled FAN have relatively low values in the occluded landmark area. The heatmaps of our distilled FAN are suitable for facial attention weights.

Figure 4(a) shows the problem of using the heatmap from the original FAN as attention weights. There is no significant distortion in facial attributes, but it leads to some artifacts because the attention is applied only to the distinct points of facial landmarks. As shown in Table 1, the distilled FAN improves SR performance with a small number of parameters.

| | Parameters (ratio) | NME | | Aligned PSNR | SSIM | MS-SSIM | NME | | Unaligned PSNR | SSIM | MS-SSIM |
|---|---|---|---|---|---|---|---|---|---|---|---|
| FAN [4] (original) | 23,820,176 (100%) | (64x64) - | (128x128) - | 22.29 | 0.670 | 0.895 | (64x64) - | (128x128) - | 22.51 | 0.667 | 0.886 |
| **Ours** | 321,412 (1.35%) | 0.239% | 0.987% | **22.66** | **0.685** | **0.902** | 0.830% | 2.2643% | **22.96** | **0.695** | **0.897** |

Table 1: Parameters, NME evaluation, PSNR, SSIM, and MS-SSIM comparison results

## 4.3 Ablation Study

In order to observe the effects of each element of our method, we conduct an ablation study using the conventional measurements of the average PSNR, SSIM, and MS-SSIM; Minimizing the MSE maximize the PSNR, which is commonly used to evaluate the SR results. Since PSNR is defined based only on pixel-wise differences, the value of PSNR has limitation to represent perceptually relevant differences [12]. Therefore, we further measure SSIM and MS-SSIM.

**Effects of loss functions** We conduct three experiments to estimate the effect of perceptual loss, heatmap loss, and proposed *facial attention loss* on the aligned dataset and the unaligned dataset. Figure 5 shows the results of using different loss functions on both aligned and unaligned datasets. The result images without perceptual loss have severely deteriorated texture of face images. Moreover, the result images without heatmap loss have unclear shapes and distortion around the eyes and mouths. As the *facial attention loss* uses the landmark heatmaps as guidance, it is helpful for face images to restore facial details.

As shown in Table 2, using our *facial attention loss* shows the highest value in PSNR, SSIM and MS-SSIM. These results verify that our *facial attention loss* is helpful to generate more structurally meaningful face images.

**Effects of progressive training** To verify that the progressive training method is helpful to super-resolve face images, we train our face SR network without the progressive training

Input  Bilinear  (a)  (b)  (c)  (d)  (e)  Target



Ablation results on aligned dataset

Input  Bilinear  (a)  (b)  (c)  (d)  (e)  Target
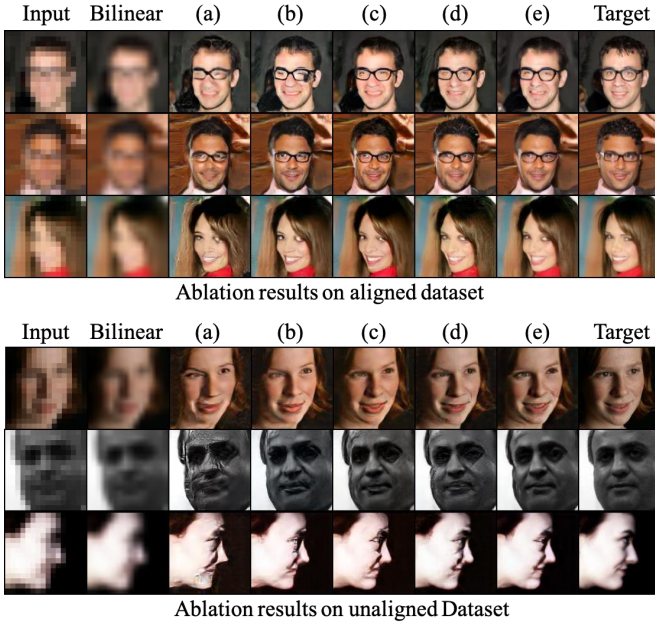


Ablation results on unaligned Dataset

Figure 5: Ablation study results on aligned and unaligned datasets. (a) $L_{pixel} + L_{WGAN}$ (b) $L_{pixel} + L_{WGAN} + L_{feat}$ (c) $L_{pixel} + L_{WGAN} + L_{feat} + L_{hm}$ (d) $L_{Ours}$-no progressive (e) $L_{Ours}$.

| Method | Aligned | | | Unaligned | | |
|---|---|---|---|---|---|---|
| | PSNR | SSIM | MS-SSIM | PSNR | SSIM | MS-SSIM |
| $L_{pixel} + L_{WGAN}$ | 21.62 | 0.616 | 0.873 | 21.62 | 0.626 | 0.863 |
| $L_{pixel} + L_{WGAN} + L_{feat}$ | 21.89 | 0.649 | 0.887 | 22.26 | 0.663 | 0.884 |
| $L_{pixel} + L_{WGAN} + L_{feat} + L_{hm}$ | 21.95 | 0.650 | 0.890 | 22.53 | 0.679 | 0.892 |
| $L_{Ours}$-no progressive | 22.24 | 0.660 | 0.893 | 22.83 | 0.680 | 0.895 |
| $L_{Ours}$ | **22.66** | **0.685** | **0.902** | **22.96** | **0.695** | **0.897** |

Table 2: PSNR, SSIM and MS-SSIM values for the ablation study results on aligned and unaligned datasets.

method. The qualitative comparison of outputs is shown in Figure 5(d). There is some degradation of facial details in super-resolved images from our non-progressively trained network, which is trained by minimizing $L_{ours}$. As shown in Table 2, the measurement values are also increased by progressive training method.

## 4.4   Comparison with State-of-the-Art

We compare our face SR method to the state-of-the-art SR methods both quantitatively and qualitatively. VDSR [8] employs a pixel-wise $L2$ loss in training. FSRNet and FSRGAN [5] employ a facial parsing map in training. URDGN [24] employs a spatial convolution layer and a discriminator.

Figure 6 provides results of various models, and Table 3 presents the quantitative comparisons on the test set. The images from VDSR achieve the highest PSNR, but the results
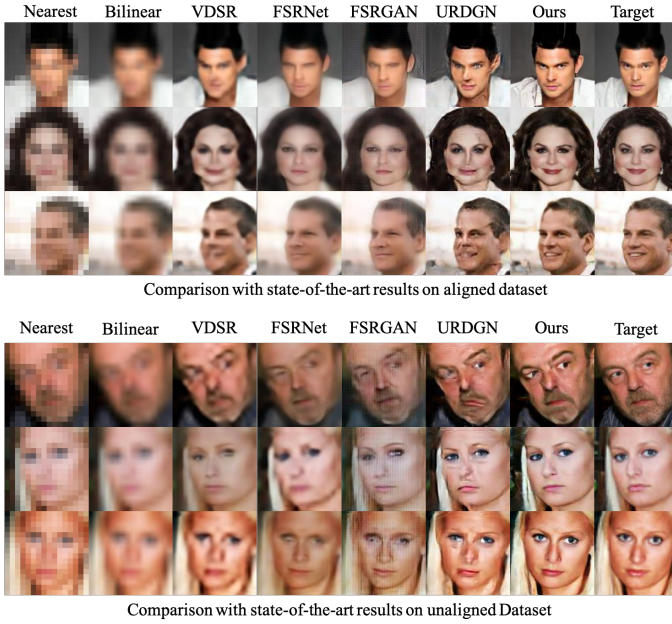
Comparison with state-of-the-art results on aligned dataset



Comparison with state-of-the-art results on unaligned Dataset

Figure 6: Qualitative comparison with aligned and unaligned datasets

| | Aligned | | | | Unaligned | | | |
|---|---|---|---|---|---|---|---|---|
| | PSNR | SSIM | MS-SSIM | MOS | PSNR | SSIM | MS-SSIM | MOS |
| Bilinear | 20.75 | 0.574 | 0.782 | 1.72 | 21.86 | 0.624 | 0.795 | 1.52 |
| URDGN [24] | 21.99 | 0.622 | 0.875 | 2.55 | 23.01 | 0.643 | 0.874 | 2.45 |
| FSRGAN [5] | 22.27 | 0.601 | 0.841 | 2.46 | 20.95 | 0.515 | 0.741 | 2.28 |
| FSRNet [5] | 22.62 | 0.641 | 0.847 | 2.34 | 21.19 | 0.607 | 0.760 | 2.19 |
| VDSR [8] | **22.94** | 0.652 | 0.880 | 2.10 | **23.70** | 0.682 | 0.882 | 1.89 |
| Ours | 22.66 | **0.685** | **0.902** | **3.73** | 22.96 | **0.695** | **0.897** | **3.73** |

Table 3: PSNR, SSIM and MS-SSIM values for the baseline experimental results on aligned and unaligned datasets.

are significantly blurred. As we can see, the results of FSRNet and FSRGAN have realistic features in facial details, but they have artifacts and partially blurred facial components. The URDGN produces relatively clear images but generates distorted face images. The results show that our method outperforms other methods especially on SSIM and MS-SSIM, and generates photo-realistic face images with restoring accurate facial attributes. More image results are shown in the Supplementary Materials.

We also conduct a MOS test to quantify image quality based on human vision. We asked 26 raters to assign a score from 1 (bad quality) to 5 (excellent quality) to all the super-resolved images and the high-resolution target images. The raters were calibrated on the 20 images of Nearest Neighbor (score 1) and HR (score 5) before the main test. The raters rated 8 versions of each image on aligned and unaligned dataset: Nearest neighbor, Bilinear, URDGN, FSRNet, FSRGAN, VDSR, Ours, and HR images (GT). Each rater rated randomly presented 240 images with each dataset (total 480 images). More details are explained in the
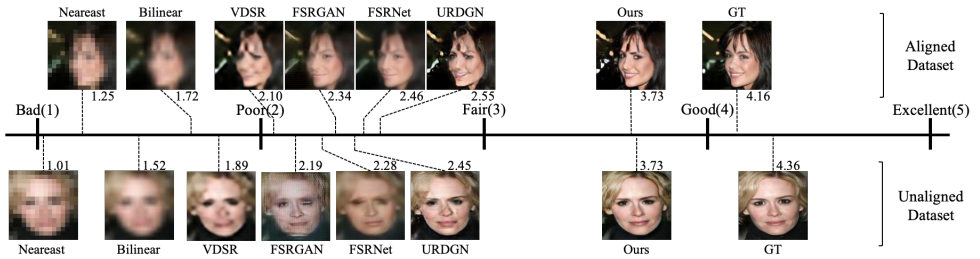
Figure 7: MOS result with aligned and unaligned datasets

Supplementary Materials. In Figure 7, our method shows overwhelming performance on MOS test, which indicates that our results are perceptually superior to other results.

# 5   Conclusion

We propose a novel face SR method which fully reflects facial details. To achieve this, we adopt progressive training method to generate photo-realistic face images and learn restoration of facial details with different guidance in each step. In addition, we propose a new facial attention loss which gives large weights to facial features in the adjacent area of landmarks. Therefore, the facial details are well expressed in super-resolved images. However, the original FAN produces landmark heatmaps including occluded landmark area, which results in degradation of super-resolution performance. Therefore, we suggest distillation of face alignment network to produce more suitable heatmaps for the SR. Besides, our distilled face alignment network has relatively light-weight architecture, so the overall training time is reduced from $\sim 3$ days to $\sim 1$ day. Our experiments demonstrate that our proposed method restore more accurate facial details. In particular, our method produces high-quality face images which are perceptually similar to the real images. As a summary, the proposed method allows our face SR network to super-resolve face images with more precise facial details.

We give attention to specific areas on the faces and propose a method to obtain the heatmaps suitable for face SR. If a better method is developed to obtain the heatmaps which well represent facial landmark areas, we will be able to achieve even better performance through our proposed method. Since our approach restores lost information by focusing on specific areas, we will further be able to restore the desired information by applying our mechanism to any task, such as super-resolution on the medical image, satellite image, and microscopic image, which requires restoration of lost information using super-resolution.

# 6   Acknowledgements

# References

[1] Georgios Tzimiropoulos Adrian Bulat. Super-fan: Integrated facial landmark localization and super-resolution of real-world low resolution faces in arbitrary poses with gans. In *CVPR*, 2018.

[2] Namhyuk Ahn, Byungkon Kang, and Kyung-Ah Sohn. Image super-resolution via progressive cascading residual network. In *CVPR Workshops*, 2018.

[3] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *ICML*, 2017.

[4] Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In *ICCV*, 2017.

[5] Yu Chen, Ying Tai, Xiaoming Liu, Chunhua Shen, and Jian Yang. Fsrnet: End-to-end learning face super-resolution with facial priors. In *CVPR*, 2018.

[6] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.

[7] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016.

[8] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *CVPR*, 2016.

[9] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.

[10] Martin Koestinger, Paul Wohlhart, Peter M Roth, and Horst Bischof. Annotated facial landmarks in the wild: A large-scale, real-world database for facial landmark localization. In *ICCV workshop*, 2011.

[11] Wei-Sheng Lai, Jia-Bin Huang, Narendra Ahuja, and Ming-Hsuan Yang. Deep laplacian pyramid networks for fast and accurate super-resolution. In *CVPR*, 2017.

[12] Christian Ledig, Lucas Theis, Ferenc Huszar, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew P. Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, and Wenzhe Shi. Photo-realistic single image super-resolution using a generative adversarial network. In *CVPR*, 2017.

[13] Cheng-Han Lee, Kaipeng Zhang, Hu-Cheng Lee, Chia-Wen Cheng, and Winston Hsu. Attribute augmented convolutional neural network for face hallucination. In *CVPR Workshops*, 2018.

[14] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *ICCV*, 2015.

[15] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *ECCV*, 2016.

[16] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NeurIPS workshop*, 2017.

[17] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. In *ICLR*, 2015.

[18] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.

[19] Samuli Laine Jaakko Lehtinen Tero Karras, Timo Aila. Progressive growing of gans for improved quality, stability, and variation. In *ICLR*, 2018.

[20] Yifan Wang, Federico Perazzi, Brian McWilliams, Alexander Sorkine-Hornung, Olga Sorkine-Hornung, and Christopher Schroers. A fully progressive approach to single-image super-resolution. In *CVPR Workshops*, 2018.

[21] Zhou Wang, Eero P. Simoncelli, and Alan C. Bovik. Multi-scale structural similarity for image quality assessment. In *IEEE Asilomar Conference on Signals, Systems and Computers*, 2003.

[22] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions Image Processing*, 2004.

[23] Xiang Wei, Boqing Gong, Zixia Liu, Wei Lu, and Liqiang Wang. Improving the improved training of wasserstein gans: A consistency term and its dual effect. In *ICLR*, 2018.

[24] Xin Yu and Fatih Murat Porikli. Ultra-resolving face images by discriminative generative networks. In *ECCV*, 2016.

[25] Xin Yu, Basura Fernando, Bernard Ghanem, Fatih Porikli, and Richard Hartley. Face super-resolution guided by facial component heatmaps. In *ECCV*, 2018.

[26] Xin Yu, Basura Fernando, Richard Hartley, and Fatih Porikli. Super-resolving very low-resolution face images with supplementary attributes. In *CVPR*, 2018.

[27] Xiangxin Zhu and Deva Ramanan. Face detection, pose estimation, and landmark localization in the wild. In *CVPR*, 2012.