

Real-World Super-Resolution of Face-Images from Surveillance Cameras

Andreas Aakerberg¹ Kamal Nasrollahi^{1,2} Thomas B. Moeslund¹

¹Visual Analysis and Perception, Aalborg University, Rendsburggade 14 Aalborg, Denmark

²Research Department, Milestone Systems A/S, Milestone Systems, Brøndby, Denmark

{anaa, tbm}@create.aau.dk kna@milestone.dk

Abstract

Most existing face image Super-Resolution (SR) methods assume that the Low-Resolution (LR) images were artificially downsampled from High-Resolution (HR) images with bicubic interpolation. This operation changes the natural image characteristics and reduces noise. Hence, SR methods trained on such data most often fail to produce good results when applied to real LR images. To solve this problem, we propose a novel framework for generation of realistic LR/HR training pairs. Our framework estimates realistic blur kernels, noise distributions, and JPEG compression artifacts to generate LR images with similar image characteristics as the ones in the source domain. This allows us to train a SR model using high quality face images as Ground-Truth (GT). For better perceptual quality we use a Generative Adversarial Network (GAN) based SR model where we have exchanged the commonly used VGG-loss [24] with LPIPS-loss [52]. Experimental results on both real and artificially corrupted face images show that our method results in more detailed reconstructions with less noise compared to existing State-of-the-Art (SoTA) methods. In addition, we show that the traditional non-reference Image Quality Assessment (IQA) methods fail to capture this improvement and demonstrate that the more recent NIMA metric [16] correlates better with human perception via Mean Opinion Rank (MOR).

1. Introduction

Face Super-Resolution (SR) is a special case of SR which aims to restore High-Resolution (HR) face images from their Low-Resolution (LR) counterparts. This is useful in many different applications such as video surveillance and face enhancement. Current State-of-the-Art (SoTA) face SR methods based on Convolutional Neural Networks (CNNs) are able to reconstruct images with photo-realistic appearance from artificially generated LR images. However, these methods often assume that the LR images were downsampled with bicubic interpolation, and therefore fail



Figure 1: $\times 4$ SR of a real low-quality face image (100×128 pixels) from the Chokepoint DB [48]. Our method enhances details and removes noise while the ESRGAN [45] amplifies the corruptions.

to produce good results when applied to real-world LR images. This is mostly due to the fact that the downsampling operation with bicubic downscaling changes the natural image characteristics and reduces the amount of artifacts. Hence, when using algorithms trained with supervised learning on such artificial LR/HR image pairs, the reconstructed images usually contains strong artifacts due to the domain gap.

This paper is about SR of real low-resolution, noisy, and corrupted images, also known as Real-World Super-Resolution (RWSR). We apply our proposed method to face images, but the method is also applicable to other image domains. To create a SR model that is robust against the corruptions found in real images, we create a degradation framework that can produce LR images that have the same image characteristic as the images that we want to super-resolve, *i.e.* the source domain images. By creating LR images from clean high-quality images, *i.e.* the target domain, allows us to train a SR model that learns to super-resolve images with similar characteristics. This approach is inspired by the work of Ji *et al.* [22] who propose to perform RWSR via kernel estimation and noise injection. However, we observe that their framework for image degradation is not ideal for SR of LR face images from surveillance cam-

eras, as these are often also corrupted by compression artifacts. Hence, we extend the degradation framework from [22] to include JPEG compression artifacts. We use the ESRGAN [45] model, which is one of the SoTA models for perceptual quality, as our backbone SR model. However, we find that the combination of loss functions for the ESRGAN is not ideal for optimal perceptual quality. To this end, we exchange the VGG-loss [24] with PatchGAN [53] loss for the discriminator similar to [22]. Inspired by Jo *et al.* [23], we additionally exchange the VGG-loss [24] with Learned Perceptual Image Patch Similarity (LPIPS) loss [52] for better perceptual quality. Different from existing models for face SR [7, 12, 8], we do not restrict our model to only work for face images of fixed input sizes, which makes our model more useful in practice. To the best of our knowledge, we are the first to propose a method for SR of real LR face images of arbitrary sizes.

We evaluate our method on two different datasets. To enable comparison of the SR performance against Ground-Truth (GT) reference images, we artificially corrupt high-quality face images from Flickr-Faces-HQ Dataset (FFHQ) [25] and report quantitative results using conventional Image Quality Assessment (IQA) methods and the most recent methods for assessment of the perceptual quality. For evaluation on real LR face image from surveillance cameras we use the Chokepoint DB [48]. In this case, as no GT image is available, we report the results using Mean Opinion Rank (MOR) and several non-reference based IQA methods. In both cases we show the effectiveness of our method via quantitative and qualitative evaluations. Furthermore, our evaluations show that most existing non-reference based IQA methods correlate poorly with human perception, while the recent Neural Image Assessment (NIMA) [16] metric provides a good correlation with human judgment as proven with MOR.

In summary, our contributions are:

- We propose a novel framework for generation of LR/HR training pairs that includes the most common image degradation types in real-world face images. Our framework includes blur kernel estimation, noise injection and compression artifacts.
- We also propose an improved ESRGAN [45] based SR model with PatchGAN [53] and LPIPS loss [52] for better perceptual quality, and show the benefit on real LR face images from the Chokepoint DB [48] and artificially corrupted face images from the FFHQ DB [25].
- Quantitatively, we evaluate our method using the most popular non-reference based IQA methods, and find only the recent NIMA [16] metric to correlate with human judgment via MOR.

2. Related Work

Recent advancements within deep-learning have proven very successful for use within super-resolution, and models of this type often achieve SoTA results. The first deep-learning based method for super-resolution was proposed by Dong *et al.* [15] who successfully trained a CNN to learn a non-linear mapping from LR to HR images. Later proposals relied on deeper networks and residual learning [27, 33], recursive learning [28], multi-path learning [21], and different loss functions [29] to reduce the reconstruction error between the super-resolved image and the GT image. However, while these methods yield high Peak Signal-to-Noise Ratio (PSNR) values, they tend to produce over-smoothed images which lack high-frequency details. To overcome this, Ledig *et al.* [32] proposed to use Generative Adversarial Networks (GANs) for SR with the SRGAN, to achieve realistic looking images according to human perception. The ESRGAN [45] further improves the SRGAN [32] by several changes to the discriminator and generator. The LR images needed for training the aforementioned deep-learning based super-resolution models are typically created by downsampling HR images with an ideal downscaling kernel, typically bicubic downscaling. However, the images generated by this kernel do not necessarily match real SR images. Additionally, in the downscaling process, important natural image characteristics, such as image sensor noise is removed, which the super-resolution algorithms are then prevented from learning. This results in poor reconstruction results and unwanted artifacts when a real-world noisy LR image is super-resolved [35].

Real-World Super-Resolution One way to address the lack of a proper imaging model for RWSR, is to create datasets that consist of real LR/HR image pairs captured using two cameras with different focal lengths [9, 43, 47]. However, this method is cumbersome and has inherent problems with the alignment of the image pairs. To overcome the problem of missing real-world training data, Shocher *et al.* [2] propose a zero-shot approach where a small CNN is trained at test time on LR/HR pairs extracted from the LR image itself. Soh *et al.* [42] extend the work of [2] by using meta-transfer learning phase to exploit information from an external dataset. Gu *et al.* [20] train a kernel estimator and corrector CNNs under the assumption that the downscaling kernel belongs to a certain family of Gaussian filters and uses the estimated kernel as input to a super-resolution model. To super-resolve LR images with arbitrary blur kernels, Zhang *et al.* [50] propose a deep plug-and-play framework which takes advantage of existing blind deblurring methods for blur kernel estimation. Bell-Kligler *et al.* [4] trains a GAN to estimate blur kernels from LR images and combines it with the ZSSR

SR model [2]. Fritsche *et al.* [17] train a GAN to introduce natural image characteristics to images downsampled with bicubic downscaling, which is then used to train a super-resolution for improved performance on real-world images. Zhang *et al.* [49] propose an iterative network for SR of blurry, noisy images for different scaling factors by leveraging both learning and model-based methods. Most recently Ji *et al.* [22] propose a degradation framework for the creation of LRHR image pairs for training. The degradation framework estimates blur kernels and noise distributions from real LR images in the source domain which are used to degrade HR images in the target domain. This enables training of a GAN based SR model which is shown to perform better on real LR images. However, a key limitation of this method is that it does not address the compression artifacts often found in real-world images.

Face Super-Resolution Face SR is a SR technique specialized for reconstruction of face images. One of the first methods for face SR was proposed by Baker and Kanade [3]. This method reconstructed face details by searching for the most optimal mapping between LR and HR patches. More recent work relies on deep learning based methods with CNNs and GANs. Dahl *et al.* [13] use pixel recursive learning with two CNNs to synthesize realistic hair and skin details. Chen *et al.* [11] combine face SR and face alignment to achieve previously unseen PSNR values. By searching the latent space of a generative model for images that downscale correctly, Menon *et al.* [37] are able to create face images of high resolution and perceptual quality. However, the problem with this approach is that the generated faces are often far from the true identity of the actual person, as illustrated in Figure 2. Additionally, none of the above mentioned methods are robust against noise or other corruptions in the input images [19].

There are very few publications available in the literature which address the problem of RWSR of face-images [19]. Furthermore, the few existing face RWSR methods are only compatible with LR images that have been squared to 16×16 pixels, meaning that the reconstructed image will be only 64×64 or 128×128 pixels depending on the scaling factor [7, 12, 8]. Hence, these models cannot perform true SR directly on the LR images. This means that the actual usefulness of the existing face SR models is limited. On the contrary, our work presents one possible solution for $\times 4$ RWSR of face images of arbitrary sizes, which we evaluate on real LR face images from surveillance cameras without any prior re-scaling.

3. The Proposed Framework

This section describes our two-step framework for RWSR. The first step aims to generate LR images from clean HR images in the target domain Y , such that these

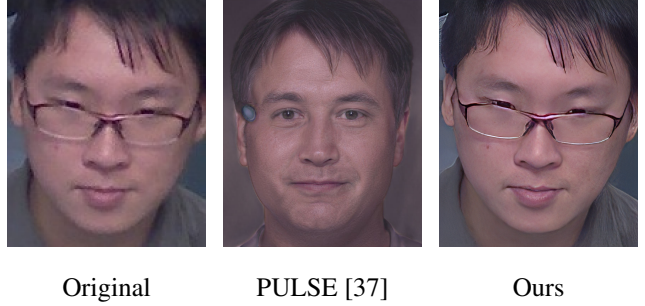


Figure 2: An example of SR of a real low-quality face image from the Chokepoint DB [48], where it can be seen that the PULSE [37] method changes the identity of the person, while our method preserves the identity and enhances details.

have similar image characteristics as the ones in the source domain X . The second step involves training a SR model on the constructed paired data, and optimizing for perceptual quality.

3.1. Novel Image Degradation

Traditional approaches for SR assumes that a LR image I_{LR} is the result of a downscaling operation of the corresponding HR image I_{HR} using some kernel k and scaling factor s , namely:

$$I_{LR} = (I_{HR} * k) \downarrow_s \quad (1)$$

However, real LR images from cameras are influenced by multiple other factors that degrade the image as well. The RealSR [22] framework tries to address this issue by considering realistic noise distributions and blur kernels in the downscaling process. However, we observe that real images from surveillance cameras are often also degraded with compression artifacts, which makes the RealSR framework perform poorly on such images. To this end, we extend the degradation framework from [22] to include JPEG compression artifacts in addition to estimation of realistic noise distributions and blur kernels. Thus, we extend the basic SR formulation from Equation 3.1, and assume that the following image degradation model was used to create I_{LR} .

$$I_{LR} = (I_{HR} * k) \downarrow_s + n + c \quad (2)$$

where k , s , n , and c denotes the blur kernel, scaling factor, noise, and compression artifacts, respectively. I_{HR} is unknown together with k , n , and c . In our degradation framework, we estimate the kernel and noise directly from the images in the source domain X . We build a pool of the estimated kernels and noise patches which is used to generate corrupted LR images from clean HR images and finally JPEG compress the images, in order to create image pairs for training the SR model.



Figure 3: Comparison with SoTA methods for SR of a small face image (56×72 pixels) from the Chokepoint DB [48]. As visible, our method hallucinates more realistic face details than the existing methods.

3.2. Blur Kernel Estimation

For estimation of realistic blur kernels, we adopt the KernelGAN method by Bell-Kligler *et al.* [4]. This method estimates an image specific SR kernel k_i using an unsupervised approach. More specifically, a GAN is trained to down-scale the input image in a way that best preserves the image patch distributions across scales. We estimate realistic blur kernels from training images in X to form a pool of kernels that can be used to degrade the HR images in Y .

Downsampling To create the downsampled image I_D we randomly choose a blur kernels k_i from the pool of estimated kernels and perform cross-correlation with images in Y . More formally the process is described as:

$$I_D = (Y_n * k_i) \downarrow_s, i \in \{1, 2 \dots m\} \quad (3)$$

where I_D is the downsampled image, Y_n is a HR image, k_i refers to a kernel from the degradation pool $\{k_1, k_2, \dots k_m\}$ and s is the scaling factor.

3.3. Noise Estimation

For degradation with realistic image noise, we adopt the method from [10] to extract noise patches from the source images X . Here the assumption is that an approximate noise patch can be obtained from a noisy image by extracting an area with weak background and then subtracting the mean. We define two patches p_i and q_j^i . We obtain p_i by a sliding window approach across images in X , and similarly for q_j^i by scanning p_i . p_i is considered a smooth patch if the following constraints are met:

$$|Mean(q_j^i) - Mean(p_i)| \leq \mu \cdot Mean(p_i) \quad (4)$$

and

$$|Var(q_j^i) - Var(p_i)| \leq \gamma \cdot Var(p_i) \quad (5)$$

where $Mean$ and Var denotes the mean and variance respectively, and μ and γ are scaling factors. Different from

[10] we add an additional constraint to ensure that saturated patches are not extracted:

$$Var(p_i) \geq \phi \quad (6)$$

where ϕ denotes a minimum variance threshold. If all constraints are satisfied, p_i will be considered a smooth patch. We then create a pool of noise patches n_i by subtracting the mean value from all valid p_i .

Degradation with Noise We degrade the LR images by injecting real noise patches from the noise pool. For better regularization of the SR model we randomly pick a noise patch from the noise pool and inject it to the LR image during training. The downsampled and noisy LR image I_N is created as follows:

$$I_N = I_D + n_i, i \in \{1, 2 \dots l\} \quad (7)$$

where I_D is a downsampled image, and n_i is a noise patch from the noise pool $\{n_1, n_2, \dots n_l\}$

3.4. Degradation with Compression artifacts

Finally, we introduce compression artifacts to the LR training images to close the domain gap between these and the real JPEG compressed LR images in the source domain X . As there are no way of determining the compression strength of existing JPEG images we empirically compare images from X to similar images with different JPEG compression strengths applied and find that a compression strength of 30 results in similar compression artifacts.

3.5. Backbone Model

We base our SR model on the ESRGAN [45], which is one of the SoTA networks for perceptual SR with $\times 4$ upscaling, and train it on the paired LR and HR images generated with our degradation framework. Different from the SRGAN [32], the ESRGAN uses Residual-in-Residual Dense Blocks (RRDBs) in the generator network and the discriminator predicts the relative realness instead of an absolute value. Additionally, the ESRGAN removes the

batch normalization layers used in SRGAN.

Loss Functions While traditional supervised SR models are trained with pixel loss to minimize the Mean Squared Error (MSE) between the reconstructed HR image and the GT image, we rely on loss functions that maximize the perceptual quality. The original ESRGAN [45] model uses several different loss functions during training. More specifically, the generator uses adversarial loss \mathcal{L}_{adv} [18] in combination with VGG perceptual loss \mathcal{L}_{vgg} [24] and pixel loss \mathcal{L}_{pix} , while the discriminator use VGG-128 [41] loss \mathcal{L}_{vgg} . However, we find that this combination of loss functions is not ideal for high perceptual quality. Following the work of [22], we first exchange the VGG-128 [41] discriminator loss with a PatchGAN discriminator from [53] to reduce the amount of artifacts in the reconstructed images. Different from the VGG loss, the PatchGAN loss \mathcal{L}_{patch} has a fully convolutional structure, and only penalizes structure differences at the scale of patches, to determine if an image is real or fake. For optimization of the generator, the loss from all patches are averaged and fed back to the generator. Continuing this track, we seek to also replace the VGG-loss in the generator. Inspired by [23], we find that using the LPIPS perceptual loss \mathcal{L}_{lips} [52] results in less noise and richer textures compared to using VGG-loss for the generator. This is mainly because the VGG network is trained for image classification, while LPIPS is trained to score image patches based on human perceptual similarity judgements. The LPIPS perceptual loss is formulated as:

$$\mathcal{L}_{lips} = \sum_k \tau^k (\phi^k(I_{gen}) - \phi^k(I_{gt})) \quad (8)$$

where I_{gen} is a generated image, I_{gt} is the corresponding GT image, ϕ is a feature extractor, τ is a transformation from embeddings to a scalar LPIPS score. The score is computed from k layers and averaged. In our implementation of LPIPS we use the pre-trained AlexNet model provided by the authors. In total, our full training loss for the generator is as follows:

$$\mathcal{L}_{generator} = \lambda_{pix} \cdot \mathcal{L}_{pix} + \lambda_{adv} \cdot \mathcal{L}_{adv} + \lambda_{lips} \cdot \mathcal{L}_{lips} \quad (9)$$

where λ_{pix} , λ_{adv} and λ_{lips} are scaling parameters.

3.6. Datasets

This section describes the datasets used for training and testing. For our experiments on real LR face images from surveillance cameras we use the Chokepoint Dataset [48] as our source domain images X . This dataset contains images of 29 different persons captured with three cameras in a real-world surveillance setting. All images have a resolution of 800×600 . We use a face detection algorithm to extract the faces from the images, and randomly split the

dataset, to obtain 72,282 images for training and 3,805 images for testing. The average resolution of the cropped faces is $\approx 92 \times 92$. We only use the Chokepoint training images to estimate realistic blur kernels and noise distributions for our degradation framework, and not for direct training of our SR model.

For the target domain of high-quality face images Y , we combine 571 face images from the SiblingsDB [44], 8,040 face images from the Radboud Faces Database [30] and 5,000 randomly selected face images from FFHQ database [25] for a total of 13,611 images. Both the SiblingsDB and Raboud Face Database contains portrait face images professionally captured in a studio setting with controlled lighting. The face images from the FFHQ are more diverse in appearance, and ethnicity of the subjects. We augment all images in the target domain by downsampling by 25, 50 and 75% with bicubic downscaling to obtain a more diverse dataset. We then apply our degradation framework described in Section 3.1 on the images in Y to obtain LR/HR image pairs for training of our SR model.

For evaluation on artificially corrupted faces images, we use the first 1,000 images from the FFHQ dataset. To generate LR/GT images we introduce three kinds of corruptions, namely, downsampling, sensor noise, and compression artifacts. For downsampling, we randomly choose a kernel from our blur kernel pool. For modeling of sensor noise we follow the protocol from [34] and use pixel-wise independent Gaussian noise, with zero mean and a standard deviation of 8 pixels. For compression artifacts, we convert the images to JPEG using a compression strength of 30.

3.7. Evaluation Metrics

Real-World Images Due to the nature of RWSR, no GT reference image exists, which makes it impossible to compare the different methods using traditional SR IQA methods *e.g.* PSNR and Structural Similarity index (SSIM). To this end, we follow the non-reference based IQA evaluation protocol from the NTIRE2020 RWSR challenge [1]. In particular, we assess the image quality using NIQE [39], BRISQUE [38], PIQE [40], NQRM [36] and PI [5], where PI is a weighted score computed as $\frac{1}{2}((10 - NQRM) + NIQE)$. However, these methods are known to correlate poorly with human ratings [1]. To address this issue, we supplement our evaluation protocol with MOR and NIMA [16], where NIMA is a learned metric based on human opinion scores, which can quantify image quality with high correlation to human perception. We use the pre-trained model for rating of the technical image quality. For the MOR, we ask the participants to rank overall image quality of the SR results. To simplify the ranking, we only include the predictions of the top-5 methods based on NIMA scores. To avoid bias, the order of the methods are randomly shuffled. We average

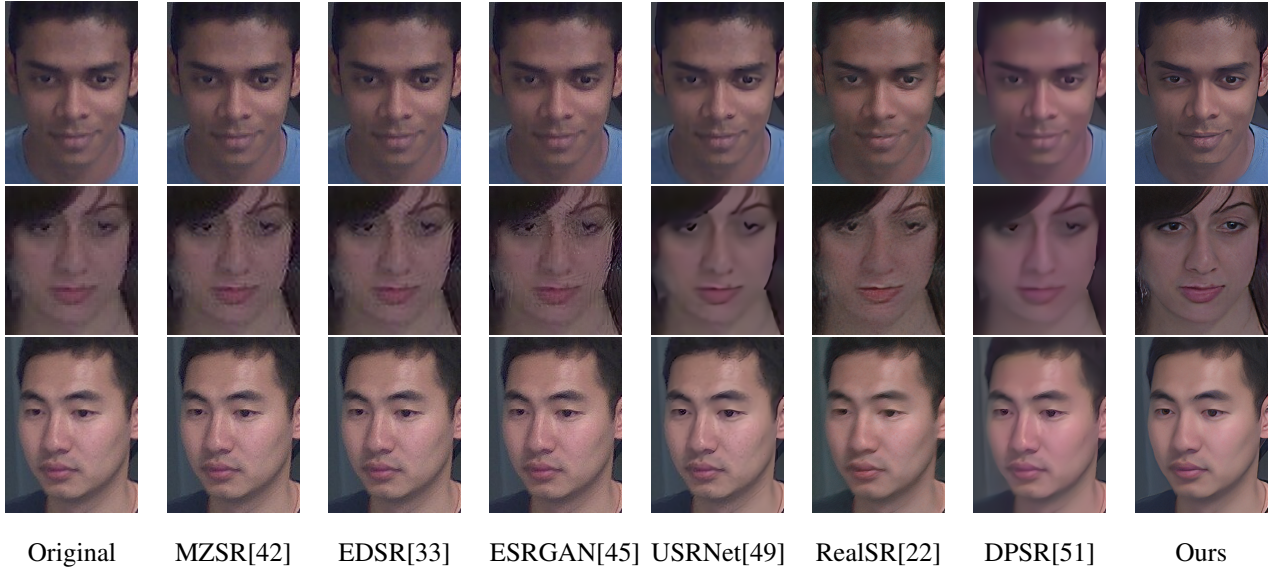


Figure 4: Comparison with SoTA methods for $\times 4$ SR of real low-quality face images from the Chokepoint DB [48]. As visible, our method generates superior reconstructions over the existing methods for different faces.



Figure 5: Comparison with SoTA methods for $\times 4$ SR of artificially corrupted face images from the FFHQ [25] testset. As seen, our method hallucinates faces with richer detail and less artifacts compared to the existing methods.

the assigned rank of each method over all images and participants to compute the MOR.

Artificially Corrupted Images For our experiments on artificially corrupted images we evaluate the performance using three conventional IQA methods, PSNR, SSIM, and the later Multi Scale Structural Similarity index (MS-SSIM) [46]. However, these metrics focus more on signal fidelity

rather than perceptual quality [6]. As our method is optimized towards perceptual quality, we also include three of the most recent full-reference metrics targeting perceptual quality, namely Normalized Laplacian Pyramid Distance (NLPD) [31], LPIPS [52], and Deep Image Structure and Texture Similarity (DISTS) [14].

Method	NIQE ↓	BRISQUE ↓	PIQE ↓	NRQM ↑	PI ↓	NIMA ↑	MOR ↓
Bicubic [26]	5.77	56.77	86.28	3.09	6.34	3.92	-
MZSR [42]	7.36	50.09	77.63	3.75	6.81	3.97	-
EDSR [33]	5.43	50.63	81.97	3.82	5.81	4.08	-
ESRGAN [45]	3.75	19.35	19.20	7.08	3.34	4.34	4.72
USRNet [49]	6.10	59.13	87.70	3.19	6.46	4.75	3.11
RealSR [22]	3.50	17.20	9.11	5.45	4.00	4.93	3.39
DPSR [51]	5.58	55.52	60.99	3.38	6.10	5.15	2.71
Ours	4.56	19.07	14.61	7.62	3.47	5.92	1.43

Table 1: Quantitative results on the Chokepoint testset. ↑ and ↓ indicate whether higher or lower values are desired, respectively. Our model scores lower on the traditional IQA metrics while being superior on the more recent NIMA metric and MOR which indicate that the traditional IQA metrics are not ideal for evaluation of perceptual quality.

Method	PSNR ↑	SSIM ↑	MS-SSIM ↑	NLPD ↓	LPIPS ↓	DISTS ↓
Bicubic [26]	28.39	0.79	0.88	0.32	0.52	0.20
MZSR [42]	29.56	0.78	0.89	0.29	0.43	0.18
EDSR [33]	28.27	0.78	0.88	0.33	0.50	0.19
ESRGAN [45]	28.09	0.77	0.88	0.34	0.40	0.19
USRNet [49]	28.53	0.80	0.89	0.32	0.53	0.21
RealSR [22]	29.14	0.79	0.90	0.29	0.29	0.18
DPSR [51]	27.45	0.79	0.88	0.33	0.51	0.25
Ours	30.20	0.79	0.91	0.28	0.25	0.16

Table 2: Quantitative results on the FFHQ testset. ↑ and ↓ indicate whether higher or lower values are desired, respectively.

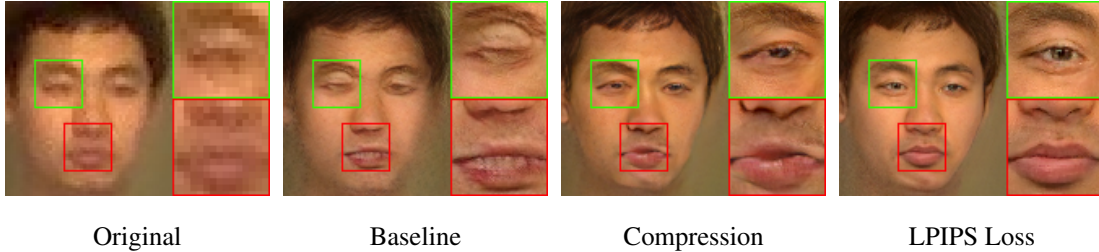


Figure 6: Ablation study of the effect of including compression artifacts in the degradation framework and exchanging the VGG-loss with LPIPS-loss for the generator in the SR model, compared to the baseline and the original LR image (56×64 pixels)

4. Experiments and Results

Implementation Details We perform all our experiments with a scaling factor $s = 4$. For our SR model we jointly train the generator and discriminator for 400K iterations with a batch size of 16. We initialize the weights from the PSNR optimized RRDB model from [45]. We use LR patches of size 32×32 , and empirically set λ_{pix} , λ_{adv} and $\lambda_{l_{lips}}$ to 0.01, 0.005 and 0.001 respectively. For noise estimation we set p_i to match the LR patch size and q_j^i to 8. Similar to [10] we set μ and γ to 0.1 and 0.25 respectively. We empirically set the minimum variance threshold ϕ to 0.5. For degradation with compression artifacts we JPEG

compress the LR training images with strength of 30 during training with a probability of 0.9 for better regularization of the SR model.

4.1. Comparison with State-of-the-Art

We did not find any other $\times 4$ face image specific RWSR methods in the literature. Instead, we compare our method to bicubic upscaling, as well as with different groups of SoTA super-resolution methods including two generic SR models (ESRGAN [45], EDSR [33]), one SR method for arbitrary blur kernels (DPSR [51]), three real-world SR models (MZSR [42], USRNet [49], and RealSR [22]).

For a fair comparison, we adjust the competing models for optimal performance. For MZSR [42], which is an unsupervised method, we enable back-projection with 10 iterations and set a noise level of 0.5. For DPSR [51], we use the pre-trained DPSRGAN model with settings for real-world images. With USRNet [49] we set the noise value to 15 for best results. The results for the RealSR [22], is based on our re-implementation of the framework as the training code was not available. We adapt the RealSR method to our face data for a fair comparison. For ESRGAN, we use the pre-trained weights provided by the authors to better illustrate the difference from our method.

Real-World Images In this experiment we evaluate the SR performance on LR face images from the Chokepoint testset. Quantitative results can be seen in Table 1. Qualitative results for multiple images are shown in Figure 4 while a close-up view of facial components can be seen in Figure 3. Our method clearly outperforms the other methods in terms of perceptual quality. However, while the traditional non-reference IQA methods (NIQE [39], BRISQUE [38], PIQE [40] and NQRM [36]) fails to capture this, scores from the more recent NIMA [16] method correlates well with human perception, which is also backed by our MOR rankings. This shows that the traditional IQA metrics are not ideal for judgement of the perceptual quality.

Artificially Corrupted Images This experiment evaluate the SR performance on artificially corrupted images from the FFHQ testset. We show quantitative results of all methods in Table 2. Qualitative results for multiple images are shown in Figure 5. Our method produces sharp and detailed images with few artifacts which closely resembles the GT images, which is also reflected in the quantitative results. Most noteworthy are the DISTS results, which are very correlated with human perception of image quality. The results show that the reconstructed images produced by our method is superior in comparison to the other methods.

4.2. Ablation Study

We evaluate the effect of our proposed method for realistic image degradation and our improved ESRGAN based SR model in the same setting as described in Section 4.1. A qualitative comparison can be seen in Figure 6.

Baseline Here, we use kernel estimation and noise injection to generate training data for the ESRGAN with patch discriminator, similar to [22]. This SR model is fine-tuned to our face image dataset, and serves as our baseline. The resulting HR images contain unpleasing noise and lack detail.

Compression Artifacts In this setting, we add JPEG

compression artifacts to the LR images during training of the baseline model. This results in more noise-free reconstructions compared to the baseline.

LPIPS loss Here, we use the LPIPS loss function for the generator instead of VGG-loss combined with the addition of compression artifacts. When the baseline model is re-trained under these settings the resulting reconstructions becomes sharper with better texture and details.

4.3. Failure Cases

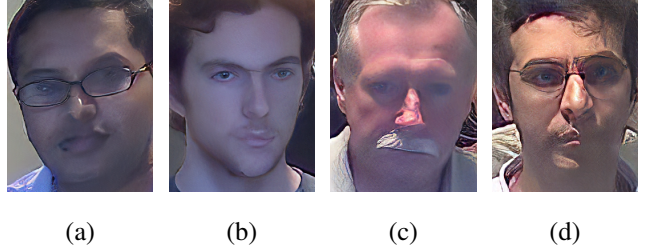


Figure 7: Examples of failure cases. Figure (a) and (b) illustrate cases where only parts of the image is super-resolved. Figure (c) shows a case where almost no high-frequency details are restored. Figure (d) shows a case where unrealistic facial features are introduced.

While our method produces reconstructed faces of better visual quality than the compared SoTA methods, it does not solve the problem RWSR of face images. Figure 7 shows several failure cases of our method. These occur when the input image is severely corrupted *e.g.* by motion blur or harsh lighting, or when out-of-focus. In these cases, our method might only super-resolve some parts of the face, *e.g.* a single eye, or even hallucinate unrealistic facial features.

5. Conclusion

In this paper, we have presented a novel framework for RWSR, which we have evaluated on low-quality face images from surveillance cameras, and artificially corrupted face images. Our method shows SoTA performance in both cases, which is achieved by using LPIPS-loss and making the SR model robust against the most common degradation types present in real LR images. Moreover, our model is the first to perform SR on real LR face images of arbitrary sizes, which makes it useful for practical applications. In the future, even better reconstructions could possibly be obtained by including more image degradation types in the framework *e.g.* chromatic aberration.

6. Acknowledgments

This work was supported by Danmarks Frie Forskningsfond under grant number 8022-00360B, and the Milestone

References

- [1] A. Lugmayr et al. Ntire 2020 challenge on real-world image super-resolution: Methods and results. *CVPR Workshops*, 2020.
- [2] Michal Irani Assaf Shocher, Nadav Cohen. "zero-shot" super-resolution using deep internal learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [3] Simon Baker and Takeo Kanade. Limits on super-resolution and how to break them. In *2000 Conference on Computer Vision and Pattern Recognition (CVPR 2000)*, 13-15 June 2000, Hilton Head, SC, USA, pages 2372–2379. IEEE Computer Society, 2000.
- [4] Sefi Bell-Kligler, Assaf Shocher, and Michal Irani. Blind super-resolution kernel estimation using an internal-gan. In *Advances in Neural Information Processing Systems 32*, pages 284–293. Curran Associates, Inc., 2019.
- [5] Yochai Blau, Roey Mechrez, Radu Timofte, Tomer Michaeli, and Lihi Zelnik-Manor. The 2018 pirm challenge on perceptual image super-resolution. In Laura Leal-Taixé and Stefan Roth, editors, *Computer Vision – ECCV 2018 Workshops*, pages 334–355, Cham, 2019. Springer International Publishing.
- [6] Yochai Blau and Tomer Michaeli. The perception-distortion tradeoff. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 6228–6237. IEEE Computer Society, 2018.
- [7] Adrian Bulat and Georgios Tzimiropoulos. Super-fan: Integrated facial landmark localization and super-resolution of real-world low resolution faces in arbitrary poses with gans. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 109–117. IEEE Computer Society, 2018.
- [8] Adrian Bulat, Jing Yang, and Georgios Tzimiropoulos. To learn image super-resolution, use a GAN to learn how to do image degradation first. In Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss, editors, *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part VI*, volume 11210 of *Lecture Notes in Computer Science*, pages 187–202. Springer, 2018.
- [9] J. Cai, H. Zeng, H. Yong, Z. Cao, and L. Zhang. Toward real-world single image super-resolution: A new benchmark and a new model. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3086–3095, 2019.
- [10] Jingwen Chen, Jiawei Chen, Hongyang Chao, and Ming Yang. Image blind denoising with generative adversarial network based noise modeling. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 3155–3164. IEEE Computer Society, 2018.
- [11] Yu* Chen, Ying* Tai, Xiaoming Liu, Chunhua Shen, and Jian Yang. Fsrnet: End-to-end learning face super-resolution with facial priors. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [12] Zhiyi Cheng, Xiatian Zhu, and Shaogang Gong. Characteristic regularisation for super-resolving face images. In *IEEE Winter Conference on Applications of Computer Vision, WACV 2020, Snowmass Village, CO, USA, March 1-5, 2020*, pages 2424–2433. IEEE, 2020.
- [13] Ryan Dahl, Mohammad Norouzi, and Jonathon Shlens. Pixel recursive super resolution. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 5449–5458. IEEE Computer Society, 2017.
- [14] Keyan Ding, Kede Ma, Shiqi Wang, and Eero P. Simoncelli. Image quality assessment: Unifying structure and texture similarity. *CoRR*, abs/2004.07728, 2020.
- [15] Chao Dong, C.C. Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep convolutional networks. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 38(2):295–307, Feb 2016.
- [16] Hossein Talebi Esfandarani and Peyman Milanfar. NIMA: neural image assessment. *IEEE Trans. Image Process.*, 27(8):3998–4011, 2018.
- [17] Manuel Fritsche, Shuhang Gu, and Radu Timofte. Frequency separation for real-world super-resolution. In *IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, 2019.
- [18] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In Zoubin Ghahramani, Max Welling, Corinna Cortes, Neil D. Lawrence, and Kilian Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 2672–2680, 2014.
- [19] Klemen Grm, Martin Pernus, Leo Cluzel, Walter J. Scheirer, Simon Dobrsek, and Vitomir Struc. Face hallucination revisited: An exploratory study on dataset bias. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 2405–2413. Computer Vision Foundation / IEEE, 2019.
- [20] Jinjin Gu, Hannan Lu, Wangmeng Zuo, and Chao Dong. Blind super-resolution with iterative kernel correction. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [21] W. Han, S. Chang, D. Liu, M. Yu, M. Witbrock, and T. S. Huang. Image super-resolution via dual-state recurrent networks. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1654–1663, 2018.
- [22] Xiaozhong Ji, Yun Cao, Ying Tai, Chengjie Wang, Jilin Li, and Feiyue Huang. Real-world super-resolution via kernel estimation and noise injection. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2020.
- [23] Younghyun Jo, Sejong Yang, and Seon Joo Kim. Investigating loss functions for extreme super-resolution. In *2020 IEEE/CVF Conference on Computer Vision and Pat-*

- tern Recognition, *CVPR Workshops 2020, Seattle, WA, USA, June 14-19, 2020*, pages 1705–1712. IEEE, 2020.
- [24] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II*, volume 9906 of *Lecture Notes in Computer Science*, pages 694–711. Springer, 2016.
- [25] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 4401–4410. Computer Vision Foundation / IEEE, 2019.
- [26] R. G. Keys. Cubic Convolution Interpolation for Digital Image Processing. *IEEE Transactions on Acoustics Speech and Signal Processing*, 29:1153–1160, Jan. 1981.
- [27] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR Oral)*, June 2016.
- [28] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Deeply-recursive convolutional network for image super-resolution. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR Oral)*, June 2016.
- [29] Wei-Sheng Lai, Jia-Bin Huang, Narendra Ahuja, and Ming-Hsuan Yang. Deep laplacian pyramid networks for fast and accurate super-resolution. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [30] Oliver Langner, Ron Dotsch, Gijsbert Bijlstra, Daniel H. J. Wigboldus, Skyler T. Hawk, and Ad van Knippenberg. Presentation and validation of the radboud faces database. *Cognition and Emotion*, 24(8):1377–1388, 2010.
- [31] Valero Laparra, Johannes Ballé, Alexander Berardino, and Eero P. Simoncelli. Perceptual image quality assessment using a normalized laplacian pyramid. In Huib de Ridder, Thrasyvoulos N. Pappas, and Bernice E. Rogowitz, editors, *Human Vision and Electronic Imaging, HVEI 2016, San Francisco, California, USA, February 14-18, 2016*, pages 1–6. Ingenta, 2016.
- [32] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi. Photo-realistic single image super-resolution using a generative adversarial network. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 105–114, 2017.
- [33] B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee. Enhanced deep residual networks for single image super-resolution. In *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 1132–1140, 2017.
- [34] Andreas Lugmayr, Martin Danelljan, and Radu Timofte. Un-supervised learning for real-world super-resolution. In *2019 IEEE/CVF International Conference on Computer Vision Workshops, ICCV Workshops 2019, Seoul, Korea (South), October 27-28, 2019*, pages 3408–3416. IEEE, 2019.
- [35] A. Lugmayr, M. Danelljan, R. Timofte, M. Fritsche, S. Gu, K. Purohit, P. Kandula, M. Suin, A. N. Rajagoapalan, N. H. Joon, Y. S. Won, G. Kim, D. Kwon, C. Hsu, C. Lin, Y. Huang, X. Sun, W. Lu, J. Li, X. Gao, S. Bell-Kligler, A. Shocher, and M. Irani. Aim 2019 challenge on real-world image super-resolution: Methods and results. In *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*, pages 3575–3583, 2019.
- [36] Chao Ma, Chih-Yuan Yang, Xiaokang Yang, and Ming-Hsuan Yang. Learning a no-reference quality metric for single-image super-resolution. *Comput. Vis. Image Underst.*, 158:1–16, 2017.
- [37] Sachit Menon, Alexandru Damian, Shijia Hu, Nikhil Ravi, and Cynthia Rudin. PULSE: self-supervised photo upsampling via latent space exploration of generative models. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 2434–2442. IEEE, 2020.
- [38] Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik. No-reference image quality assessment in the spatial domain. *IEEE Trans. Image Process.*, 21(12):4695–4708, 2012.
- [39] A. Mittal, R. Soundararajan, and A. C. Bovik. Making a “Completely Blind” Image Quality Analyzer. *IEEE Signal Processing Letters*, 20(3):209–212, Mar. 2013.
- [40] Venkatanath N., Praneeth D., Maruthi Chandrasekhar Bh., Sumohana S. Channappayya, and Swarup S. Medasani. Blind image quality evaluation using perception based features. In *Twenty First National Conference on Communications, NCC 2015, Mumbai, India, February 27 - March 1, 2015*, pages 1–6. IEEE, 2015.
- [41] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.
- [42] Jae Woong Soh, Sunwoo Cho, and Nam Ik Cho. Meta-transfer learning for zero-shot super-resolution. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020.
- [43] Hamid Vaezi Joze, Ilya Zharkov, Karlton Powell, Carl Ringler, Luming Liang, Andy Roulston, Moshe Lutz, and Vivek Pradeep. Imagepairs: Realistic super resolution dataset via beam splitter camera rig. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2020.
- [44] Tiago F. Vieira, Andrea Bottino, Aldo Laurentini, and Matteo De Simone. Detecting siblings in image pairs. *The Visual Computer*, 30(12):1333–1345, Dec 2014.
- [45] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In Laura Leal-Taixé and Stefan Roth, editors, *Computer Vision – ECCV 2018 Workshops*, pages 63–79, Cham, 2019. Springer International Publishing.
- [46] Zhou Wang, Eero P. Simoncelli, and Alan C. Bovik. Multi-scale structural similarity for image quality assessment. In *Proc. IEEE Asilomar Conf. on Signals, Systems, and Computers, (Asilomar)*, pages 1398–1402, 2003.
- [47] Pengxu Wei, Ziwei Xie, Hannan Lu, ZongYuan Zhan, Qixiang Ye, Wangmeng Zuo, and Liang Lin. Component divide-

and-conquer for real-world image super-resolution. In *Proceedings of the European Conference on Computer Vision*, 2020.

- [48] Yongkang Wong, Shaokang Chen, Sandra Mau, Conrad Sanderson, and Brian C. Lovell. Patch-based probabilistic image quality assessment for face selection and improved video-based face recognition. In *IEEE Biometrics Workshop, Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 81–88. IEEE, June 2011.
- [49] Kai Zhang, Luc Van Gool, and Radu Timofte. Deep unfolding network for image super-resolution. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3217–3226, 2020.
- [50] Kai Zhang, Wangmeng Zuo, and Lei Zhang. Deep plug-and-play super-resolution for arbitrary blur kernels. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 1671–1681. Computer Vision Foundation / IEEE, 2019.
- [51] Kai Zhang, Wangmeng Zuo, and Lei Zhang. Deep plug-and-play super-resolution for arbitrary blur kernels. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 1671–1681. Computer Vision Foundation / IEEE, 2019.
- [52] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 586–595. IEEE Computer Society, 2018.
- [53] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 2242–2251. IEEE Computer Society, 2017.