# Super-Resolution through StyleGAN Regularized Latent Search: A Realism-Fidelity Trade-off

Marzieh Gheisari
École Normale Supérieure
Paris, France
gheisari@bio.ens.psl.eu

Auguste Genovesio
École Normale Supérieure
Paris, France
auguste.genovesio@ens.psl.eu

## Abstract

*This paper addresses the problem of super-resolution: constructing a highly resolved (HR) image from a low resolved (LR) one. Recent unsupervised approaches search the latent space of a StyleGAN pre-trained on HR images, for the image that best downscales to the input LR image. However, they tend to produce out-of-domain images and fail to accurately reconstruct HR images that are far from the original domain. Our contribution is twofold. Firstly, we introduce a new regularizer to constrain the search in the latent space, ensuring that the inverted code lies in the original image manifold. Secondly, we further enhanced the reconstruction through expanding the image prior around the optimal latent code. Our results show that the proposed approach recovers realistic high-quality images for large magnification factors. Furthermore, for low magnification factors, it can still reconstruct details that the generator could not have produced otherwise. Altogether, our approach achieves a good trade-off between fidelity and realism for the super-resolution task.*

Figure 1. Searching the latent space of a StyleGAN without proper constraints leads to an unrealistic HR image, represented by $\mathbf{w}_A$. However, RLS finds an optimal latent code $\mathbf{w}_B$ located in the dense regions of the image distribution. $\text{RLS}^+$ further improves the results by modulating the generator's weights within a small $\ell_1$-norm ball centered on $\mathbf{w}_B$, resulting in $\mathbf{w}_{B+}$, an image that accurately matches the input LR image when degraded.

## 1. Introduction

Super-resolution aims to reconstruct an unknown High Resolution (HR) image $\mathbf{x} \in \mathbb{R}^{n \times n}$ from a Low Resolution (LR) image $\mathbf{y} \in \mathbb{R}^{m \times m}$, related to one another by a down-sampling process described by $\mathbf{y} = \mathrm{D}(\mathbf{x}) + \delta$ with $\mathrm{D} : \mathbb{R}^{n \times n} \to \mathbb{R}^{m \times m}$ a down-sampling non-invertible forward operator and $\delta$ an independent noise with distribution $p_\delta$. As with many data generation tasks, super-resolution has largely benefited in recent years from the advent of Generative models. Two main research trends have emerged for this task: GAN-based methods and prior-guided methods.

GAN-based methods [10, 22, 24], learn a direct coupling of HR and LR images based on coupled image translation. Recently, GCFSR [7] introduced a generative and controllable face sup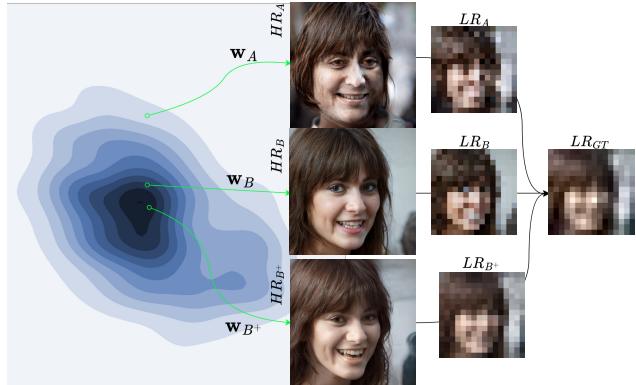er-resolution model that does not rely on ad-ditional priors and has been shown to reconstruct faithful images. However, GAN-based methods have limitations in that they train the SR generator from scratch using a combined objective function consisting of a fidelity term and an adversarial loss. This approach requires the generator to capture both the natural image characteristics and the fidelity to the ground-truth, which can result in limitations when approximating the natural image manifold. As a result, GAN-based methods often produce artifacts and unnatural textures.

Moreover, the super-resolution problem is ill-posed as, for a non-invertible forward operator $\mathrm{D}$ with $m < n$, there are infinitely many HR images that match a given LR image. Thus reconstruction procedure must be further constrained by prior information to better define the objective and lead to a stable solution. In addition, GAN-based methods often rely on a specific degradation model during training, which
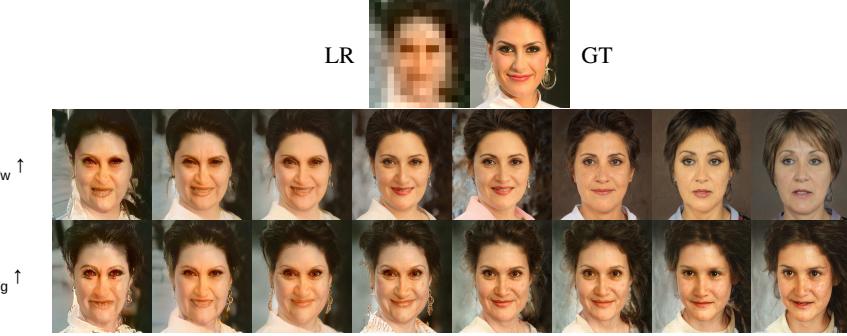
Figure 2. Impact of regularizer parameters $\lambda_w$ and $\lambda_g$ on fidelity and realism. The x-axis represents increasing values of $\lambda_w$ (respectively $\lambda_g$) from left to right, with $\lambda_g = 0$ (respectively $\lambda_w = 0$) and $\lambda_c = 0$.

can restrict their ability to handle the true degradation that can be encountered in real-world applications.

Prior-guided methods [3, 15, 23, 27] can be considered as blind restoration techniques, as they can adapt to the given problem without requiring re-training. These methods fall into two main categories: posterior sampling-based and optimization-based methods. DDRM [14] is a posterior sampling method that uses a pre-trained denoising diffusion generative model to gradually denoise a sample to the desired output, conditioned on the LR input image.

Optimization-based methods, learn the distribution of HR images in an unsupervised fashion using a GAN, and then search the latent space of this trained GAN to find the HR image that, once down-sampled, is the closest to the LR image. This idea was first introduced by Bora et al [3], and subsequent work by PULSE [18] and BRGM [17] improved upon it using StyleGAN.

To keep the search within the image manifold, BRGM assumed that the intermediate latent space $\mathcal{W}^+$ followed a standard Gaussian distribution, $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$, where $d$ is the dimensionality of the latent space. On the other hand, PULSE and a few other studies utilized an invertible transformation of $\mathcal{W}^+$ that included a leaky rectified linear unit (ReLU) [6] followed by an affine whitening transformation to ensure that transformed latent vectors approximately followed the standard Gaussian distribution. Sampled vectors are then constrained to lie around a hypersphere with radius $\sqrt{d}$ hypothesizing that most of the mass of a high-dimensional Gaussian distribution is located at or near $\sqrt{d}\mathcal{S}^{d-1}$, where $\mathcal{S}^{d-1}$ is the $d$-dimensional unit hypersphere. Constraining samples to lie in dense area of the StyleGAN style distribution resulted in increased realism of the generated images.

Although the above approach showed major improvements over previous work, it also presents three important caveats that in practice led to image artifacts. First, as we will show later, transforming the intermediate latent space this way does not lead to an accurate standard Gaussian distribution, and prevents proper regularization based on this hypothesis. Second, we argue that a search strictly limited

to the spherical surface $\mathcal{S}^{d-1}$ restricts access to the whole variety of images a StyleGAN can generate, thus preventing a close reconstruction of the HR image to be reached. Third, there is an inherent trade-off between realism and fidelity: the generator encoding the image prior can produce images from the learnt domain, however, there is no guarantee a pre-trained generator can produce a specific image we aim to reconstruct.

In this work, firstly, we operate a Regularized Latent Search (RLS) for a latent code located in "healthy" regions of the latent space. In this way, the system is constrained to produce images that belong to the original image domain StyleGAN was trained on. To do so, we take advantage of normalizing flow to Gaussianize the latent style sample distribution and show that it leads to a much closer standard Gaussian distribution. We then use this revertible transformation to regularize the search in $\mathcal{W}^+$ such that it remains in a high-density area of the style vector distribution.

Secondly, we mitigate the realism-fidelity trade-off issue by slightly modulating the generator's weights within a small $\ell_1$-norm ball centered to the previously identified latent code. To this end, we perform a small number of additional iterations to fine-tune the generator's training. In this way, it becomes possible to faithfully reconstruct a slightly out-of-domain HR image by simply increasing the generator's domain on demand. We then show experimentally, that the latter produces reconstructed images that are not only realistic but also more faithful to the original HR image. We also show that the approach is robust to noise and other image corruptions.

## 2. Prior work

### 2.1. Style-based Generative Adversarial Networks

StyleGAN models are well known for generating highly realistic images. The StyleGAN architecture consists of two sub-networks: a mapping network $G_m : \mathbb{R}^d \to \mathbb{R}^d$, and an $L$-layer synthesis network $G_s : \mathbb{R}^{L \times d} \to \mathbb{R}^n$. The mapping network maps a sample $\mathbf{z} \in \mathbb{R}^d$ from a standard normal

distribution to a vector $\mathbf{w} \in \mathcal{W}$. The $d \times L$ dimensional vector $\mathbf{w} \in \mathcal{W}^+$ containing $L$ copies of $\mathbf{w}$ is fed to the $L$-layer synthesis network $\mathsf{G}_s$. The $i$-th copy of $\mathbf{w}$ represents the input to the $i$-th layer of $\mathsf{G}_s$, which controls the $i$-th level of detail in the generated image. In addition to these, $\mathsf{G}_s$ also takes as input a collection of latent noise vectors $\boldsymbol{\eta}$ that control minor stochastic variations of the generated image at each scale.

The ability of a StyleGAN to control features of the generated image at different scales is partially due to this architecture, and partially due to the style-mixing regularization occurring during training [12, 13]. In addition to these basic characteristics, StyleGAN2 introduces path length regularization, which helps in reducing the representation error. Prior works demonstrated that performing latent search in the extended latent space $\mathcal{W}^+$ led to more accurate reconstructions but at the cost of a reduced editability [1,26]. In PULSE, GEOCROSS is introduced as a penalty term to force the embedding in the extended latent space $\mathcal{W}^+$ to remain close to the latent space $\mathcal{W}$, which in turn promotes the embedded object to be close to the range of the generator network $\mathsf{G}_s$ with the latent space $\mathcal{W}$.

## 2.2. Learning the prior with Normalizing Flow

Using a sequence of invertible mappings, a Normalizing Flow $\mathsf{F} : \mathbb{R}^d \to \mathbb{R}^d$ is a transformation of an unknown complex distribution into a simple probability distribution that is easy to sample from and whose density is easy to evaluate such as standard Gaussian [16].

Let $\mathbf{z} = \mathsf{F}(\mathbf{w})$ with probability density function $p(\mathbf{z})$. Using the change-of-variable formula, we can express the log-density of $\mathbf{w}$ by [20]:

$$\log p_\mathsf{F}(\mathbf{w}) = \log p(\mathbf{z}) + \log|\det J_\mathsf{F}(\mathbf{w})|, \quad \mathbf{z} = \mathsf{F}(\mathbf{w}) \quad (1)$$

where $J_\mathsf{F}(\mathbf{w})$ is the Jacobian of $\mathsf{F}$ evaluated at $\mathbf{w}$.

In practice the Jacobian determinant in Equation (1) should be easy to compute, so that the density $p_\mathsf{F}(\mathbf{w})$ can be evaluated. Furthermore, as a generative model, the invertibility of $\mathsf{F}$ allows new samples $\mathbf{w} = \mathsf{F}^{-1}(z)$ to be drawn through sampling from the base distribution.

In the literature, several flow models were proposed, such as Real Non-volume Preserving Flow (RealNVP) [5] and Masked Auto-regressive Flow (MAF) [21].

## 3. Proposed Method

In this work, we aim to find the latent code of a pre-trained StyleGAN that produces a realistic HR image that, when degraded, matches the input LR image accurately. The proposed method includes two steps to ensure that the results reside within the image domain and are close to the HR image. First, in order to leverage the StyleGAN prior, we restrict the optimization solution to remain on the
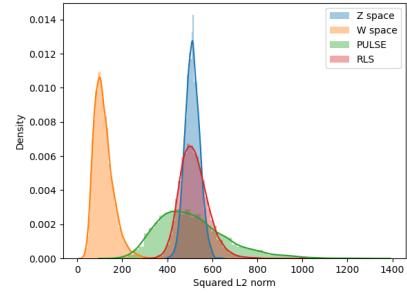


Figure 3. Evaluation of the gaussianization prior.

manifold. This restriction can be efficiently implemented by incorporating the image prior introduced in Sec. 3.1. However, especially in the case of slightly out-of-domain images, this approach yields less faithful reconstructions. Therefore, in a second step (Sec. 3.2), the generator is refined around this anchor point to recover the missing details without affecting the image prior.

## 3.1. Regularized Latent Search

MAP estimation is a common approach in Bayesian statistics for estimating an unknown parameter $\mathbf{x}$ based on observed data $\mathbf{y}$. We can formulate the super-resolution problem in terms of MAP estimation [17]. The unknown parameter we are trying to estimate is the HR image $\mathbf{x}$, and the observed data is the LR image $\mathbf{y}$. For a given LR image $\mathbf{y}$, we wish to recover the HR image $\mathbf{x}$ as the MAP estimate of the conditional distribution $p_\mathsf{G}(\mathbf{x}|\mathbf{y})$:

$$\begin{aligned} &\arg\max_{\mathbf{x}} \log p_\mathsf{G}(\mathbf{x}|\mathbf{y}) \\ &= \arg\max_{\mathbf{x}}[\log p(\mathbf{y}|\mathbf{x}) + \log p_\mathsf{G}(\mathbf{x}) + \log p(\mathbf{y})] \end{aligned} \quad (2)$$

Since the marginal density $\log p(\mathbf{y})$ is constant we drop it. We also model the likelihood $p(\mathbf{y}|\mathbf{x})$ as a delta function $p_\delta(\mathbf{y} - \mathsf{D}(\mathbf{x}))$, where $\mathsf{D}$ is a degradation operator that maps HR images to LR images. We can then rewrite the MAP objective as:

$$\arg\max_{\mathbf{x}}[\log p_\delta(\mathbf{y} - \mathsf{D}(\mathbf{x})) + \log p_\mathsf{G}(\mathbf{x})] \quad (3)$$

where the first term is the likelihood term and the second term is the prior which describes the manifold of real HR images.

**Image prior** Let $\mathsf{G}_s$ be the synthesis network of a Style-GAN [13] pre-trained on the considered image domain. $\mathsf{G}_s$ takes as input $\mathbf{w}$, produced by the mapping network, and outputs an image *i.e.* $\mathbf{x} = \mathsf{G}_s(\mathbf{w})$ (Sec. 2.1), which is a deterministic transformation of $\mathbf{w}$ through a differentiable
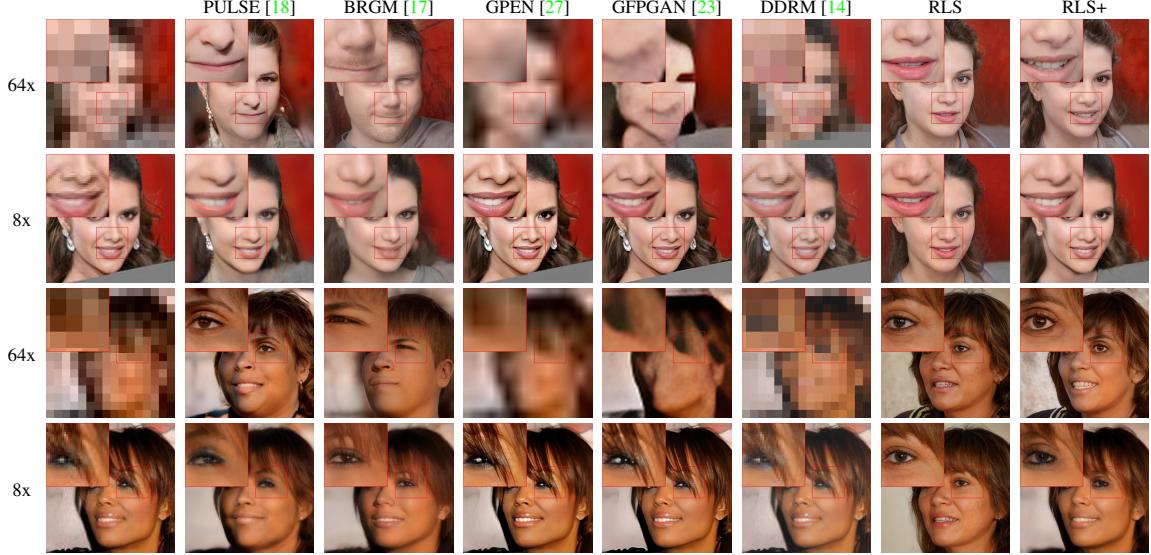
Figure 4. Comparison of the reconstructions of a high resolution face from CelebA. (Zoom-in for best view)

function $G_s$. A change of variables can be used for a non-invertible mapping [4], then the probability density function of $x$ can be obtained from the probability density function of $w$:

$$p_G(x) = p_w(w) \left| \frac{\partial w}{\partial x} \right|$$

where $w = G_s^{-1}(x)$ is the inverse transformation of $G_s$, and $\frac{\partial w}{\partial x}$ is the derivative of $w$ with respect to $x$.

Now, we can express the image prior with respect to the latent variables $w$, which allows us to work with the more tractable latent space of the StyleGAN network, rather than the high-dimensional space of the high-resolution images:

$$\log p_G(G_s(w)) = \log p_w(w) + \log |\det J_{G_s}(w)|. \quad (4)$$

$J_{G_s}(w)$ is the Jacobian matrix of the mapping $G_s$ evaluated at $w$ that describes how small changes in the input $w$ result in changes in the output $x$.

In StyleGAN2, the authors introduced a new regularization term to encourage smoothness and disentanglement of the latent space. This regularization term is called the "path length regularization," and it is based on the notion of a "path" in the latent space. In particular, given two latent vectors $w_1$ and $w_2$, we can define a "path" in the latent space as a function $w(s)$ that smoothly interpolates between $w_1$ and $w_2$ as $s$ varies from $0$ to $1$. The authors of StyleGAN2 then introduce a penalty term on the length of this path and claim that this regularization term implies that the Jacobian determinant of the network $G_s$ is approximately constant for all $w$. Based on this property of StyleGAN2 the Jacobian determinant term in the above equation can be

dropped and the image prior can be expressed directly by the image prior $p_w(w)$, which is defined on $w \in \mathcal{W}^+$ by:

$$\log p_w(w) = \lambda_w \mathcal{P}_w + \lambda_g \mathcal{P}_{gaussian} + \lambda_c \mathcal{P}_{cross} \quad (5)$$

where:

- $\mathcal{P}_w$ is a prior that keeps $w$ in the area of high density in $\mathcal{W}^+$: $\mathcal{P}_w = \frac{1}{L} \sum_{i=1}^{L} \log p_F(w_i)$ , where $p_F(w)$ is estimated by a normalizing flow model $F$ explained in Sec. 2.2.

- $\mathcal{P}_{gaussian}$ is a gaussianization prior. Using the normalizing flow model $F$, a gaussianized latent vector $w_n$ can be obtained and a $L2$ regularization applied on it to keep it near the surface of the hypersphere: $\mathcal{P}_{gaussian} = -\frac{1}{L} \sum_{i=1}^{L} (||F(w)||_2 - \sqrt{d})^2$

- $\mathcal{P}_{cross}$ is a pairwise euclidean distance prior on $w = [w_1, \dots, w_L] \in \mathcal{W}^+$ that ensures $w \in \mathcal{W}^+$ remains close to the trained manifold in $\mathcal{W}$: $\mathcal{P}_{cross} = -\sum_{i=1}^{L-1} \sum_{j=i+1}^{L} ||w_i - w_j||_2^2$

**Optimization** In the likelihood term in Equation (3), we assume that the noise follows a Laplace distribution *i.e.* $\delta \sim Laplace(0, \lambda_l I)$, then the log-density of $\delta$ becomes: $\log p_\delta(\delta) = -||\delta||_1 - C$ for a constant $C$. With the parameters of $G_s$ denoted by $\theta$, the problem in Equation (3) can be recast as an optimization over $w$, leading to the final objective function:

$$\hat{w} = \arg \min_w ||y - D(G_s(w, \theta))||_1 - \log p_w(w) \quad (6)$$

| Scale | Method | FID↓ | KID$^{(\times 10^3)}$↓ | NIQE↓ | ID↑ | LPIPS↓ | PSNR↑ | MSSIM↑ |
|---|---|---|---|---|---|---|---|---|
| | PULSE [18] | 42.9331 | 30.2643 | 5.0957 | 0.6709 | 0.5197 | 19.5775 | 0.5430 |
| | BRGM [17] | 58.3559 | 47.5288 | 4.3817 | 0.6426 | 0.5412 | 18.8989 | 0.5300 |
| | GPEN [27] | 474.275 | 693.3078 | 14.6695 | 0.6623 | 0.6704 | 20.2558 | 0.6014 |
| 64x | GFPGAN [23] | 197.3977 | 181.3333 | 12.9577 | 0.7151 | 0.6431 | 19.6302 | 0.6047 |
| | DDRM [14] | 391.2105 | 538.9393 | 8.0546 | 0.6938 | 0.6924 | 18.6866 | 0.5296 |
| | RLS | 47.8888 | 30.8534 | 4.1032 | 0.7210 | 0.5037 | 17.9680 | 0.5183 |
| | RLS$^+$ | 36.3110 | 20.9327 | 4.1348 | 0.7335 | 0.4749 | 19.7064 | 0.5671 |
| | Bicubic | 86.9839 | 104.7387 | 9.9253 | 0.8100 | 0.5346 | 28.0067 | 0.8568 |
| | PULSE [18] | 34.5038 | 21.6472 | 5.9383 | 0.7511 | 0.4618 | 23.4985 | 0.7090 |
| | BRGM [17] | 38.0316 | 27.2339 | 7.6593 | 0.7634 | 0.4998 | 21.9977 | 0.6817 |
| 8x | GPEN [27] | 27.9026 | 19.4761 | 4.9814 | 0.8746 | 0.3217 | 26.3693 | 0.8472 |
| | GFPGAN [23] | 28.2971 | 18.5216 | 6.0168 | 0.8775 | 0.3323 | 27.1016 | 0.8512 |
| | DDRM [14] | 30.1999 | 25.1746 | 7.1486 | 0.8358 | 0.5386 | 23.4707 | 0.8566 |
| | RLS | 45.8778 | 29.6858 | 4.2293 | 0.7539 | 0.4738 | 18.4833 | 0.5677 |
| | RLS$^+$ | 27.6691 | 13.0044 | 4.7241 | 0.8152 | 0.3925 | 24.2802 | 0.7577 |

Table 1. Quantitative comparison on CelebA for 64x and 8x super-resolution. (The best and the second-best are emphasized by blue and red respectively.)

## 3.2. Boosting Reconstruction Fidelity

**Realism-Fidelity Trade-off** The impact of regularizer parameters $\lambda_w$ and $\lambda_g$ on the reconstruction quality is depicted in Figure 2. The experiment considers two variants, where only the corresponding prior is retained, and the others are disconnected. For instance, when examining the effect of $\lambda_w$ (i.e., as $\lambda_w$ increases), $\lambda_g$ and $\lambda_c$ are set to zero. The variant in which $\lambda_g$ is increased showcases the rationale behind $P_{gaussian}$, which maintains $w$ around the surface of a sphere, rather than precisely on it. This variant demonstrates that as $\lambda_g$ approaches infinity, the LR consistency (reconstruction fidelity) decreases, indicating that the exact surface of the sphere does not encompass the entire image distribution learned by StyleGAN.

It is also worth noting that Figure 2 shows that retaining only the prior $\mathcal{P}_w$ (for some parameters) can result in a reconstructed image that is closer to the ground truth, indicating the consequences of searching beyond regions near the hypersphere. Conversely, when $\lambda_w$ is large, the reconstruction quality decreases. Thus, we can employ $P_{gaussian}$ to strengthen the prior and ensure that the latent code resides in healthy regions.

As depicted in Figure 2, both variants of RLS yield excellent reconstruction outcomes for specific parameter values. This allows us to strike a balance between realism (with respect to the dataset StyleGAN was trained on) and reconstruction fidelity. However, this realism-fidelity trade-off can be further enhanced by RLS$^+$, which is explained in the following.

**RLS$^+$** Using a pre-trained StyleGAN as the image prior introduces a limitation in maintaining low reolution consistency. As with any GAN model, the representational capacity of StyleGAN is directly linked to the size and diversity of the dataset on which it was trained. In the case of faces, it is unrealistic to expect a given trained model to reconstruct all possible faces. It will typically excel at reconstructing common faces but may struggle to accurately reproduce rare details that were not present during training.

To mitigate the aforementioned issue, we propose RLS$^+$, a refinement step that aims to enhance the reconstruction quality by further training the generator, the latent code, and the noise inputs simultaneously. In this way, the generator range can be expanded to recover an image that better matches the original HR image. In this step, we aim at modifying G as little as possible in order to faithfully reconstruct the HR image still without affecting the semantic prior learned previously that ensures the HR image will remain a face. Moreover, the reconstruction quality can be further improved by optimizing the noise inputs $\boldsymbol{\eta}$, which comprises high-frequency details. Finally, after obtaining the anchor point with Equation (6), the following optimization problem can be solved by initializing the generator with the pre-trained weights $\boldsymbol{\theta}$:

$$\min_{\mathbf{w},\boldsymbol{\theta},\boldsymbol{\eta}} ||\mathbf{y} - \mathsf{D}(\mathsf{G}_s(\mathbf{w},\boldsymbol{\theta},\boldsymbol{\eta}))||_1, \quad \mathbf{w} \in \mathcal{B}(\mathbf{w}_{ancor},r) \quad (7)$$

where $\mathcal{B}(\mathbf{w}_{anchor},r) = \{\mathbf{w} \in \mathbb{R}^d | \, ||\mathbf{w} - \mathbf{w}_{anchor}||_1 \leq r\}$ denote the set of points with $\ell_1$-norm bounded by radius $r$ from the anchor point. This set comprises latent codes that can generate realistic images with almost the same facial
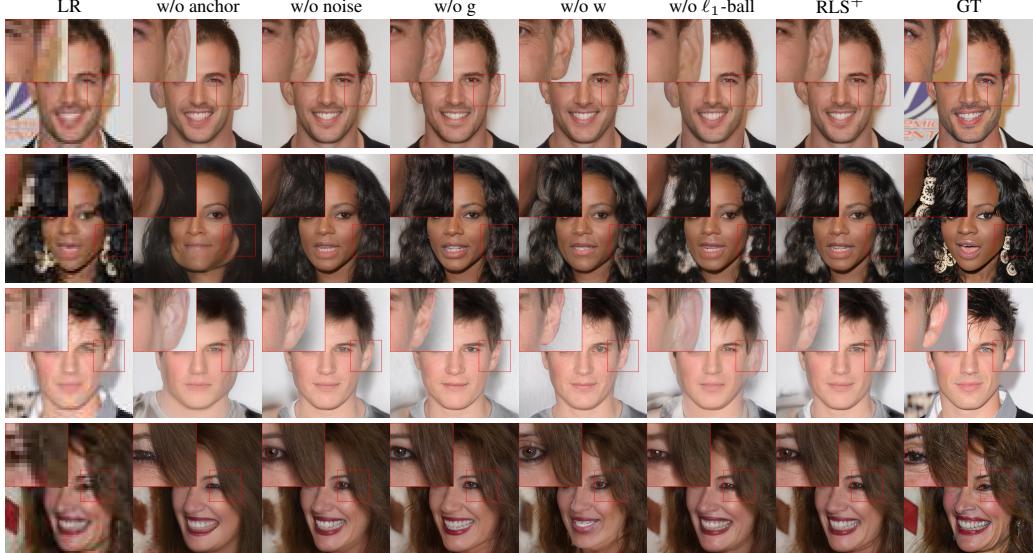
Figure 5. Ablation study. w/o anchor indicates optimization over $\mathbf{w}$, $\boldsymbol{\theta}$ and $\boldsymbol{\eta}$ using a mean value of random samples as initialization. The variants w/o noise, w/o g, and w/o w remove noise inputs, generator's weights, and latent vectors from the optimization, respectively. w/o $\ell_1$-ball removes the locality constraint during the fine-tuning. (16x)

attributes. Our algorithm attempts to explore the extended range of $\mathsf{G}_s$, to find the latent vector that best explains the input LR image, but we only allow solutions that lie within an $\ell_1$-norm ball centered at $\mathbf{w}_{anchor}$. Our empirical analysis shows that by keeping the radius of the ball relatively small, we can enhance both realism and identity-similarity.

Intuitively, allowing deviations from the anchor point increases the capability of the generator to produce the closest reconstruction of the target image. However, limiting its deviation, prevent over-fitting to unrealistic details. Obviously, as $r$ increases, $\mathcal{B}$ contains latent vectors that are further away from the anchor point. The last allows finding latent codes that are more expressive and diverse, but also further away from faces and thus producing unrealistic images. Altogether, $r$ offers control over a trade-off between realism and reconstruction error.

We increase the expressivity of the StyleGAN by further optimizing $\boldsymbol{\theta}$. However, to avoid over-expanding the range of the generator to non-realistic images, we employ early stopping to maintain the generative prior. Once the generator is tuned, the final HR image is obtained by $\hat{\mathbf{x}} = \mathsf{G}_s(\mathbf{w}^*, \boldsymbol{\theta}^*, \boldsymbol{\eta}^*)$.

## 4. Results

**Experimental setup.** We used a StyleGAN2 generator [13] pre-trained on the FFHQ dataset [12] that includes 70,000 high-quality face images of resolution $1024 \times 1024$. For evaluation, we used the first 2000 samples from the CelebA-HQ test set [11] and simulated degraded faces from the HR images using bicubic downsampling.

The regularization parameters $\lambda_w$, $\lambda_c$ and $\lambda_g$ were set to 0.0002, 0.05, and 0.0004, respectively. The normalizing flow model we use is the MAF, as it tends to work better than RealNVP for density estimation tasks. Five flow blocks are used in our model and all hidden dimensions are set to 1024. For RLS, we used an Adam optimizer over 200 iterations with a learning rate of 0.5 and initialized the search by the mean of 100,000 randomly generated latent vectors. Then, for RLS$^+$, we further run for only 50 iterations with a learning rate of 0.0001. To enforce the $\ell_1$-norm ball constraint, we use Projected Gradient Descent and set the radius of the ball as $\sqrt{d}$. Also, only the first nine input layers are optimized, and the rest are fixed.

We compared our algorithm with state-of-the-art face restoration methods, including PULSE [18], BRGM [17], GPEN [27], GFPGAN [23] and DDRM [14]. We used the original codes and weights from the official paper repositories for all experiments, and replicated the same parameter settings reported in the original papers.

**Quality of the Gaussianization process.** We generated 5000 random samples $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and their associated style vector $\mathbf{w} = \mathsf{G}_m(\mathbf{z})$, then gaussianized distribution of $\mathcal{W}$ by PULSE and our method. We then computed the squared norm for all of these samples (see Figure 3). As expected, the squared norm of the standard gaussian distribution $\mathcal{Z}$ approximately follows $||\mathbf{z}_n||_2^2 \sim \chi_d^2$ and thus forms a narrow distribution around $d = 512$ while the squared norm of the untransformed distribution $\mathcal{W}$ does not, which is inconsistent with the prior assumption held by

| Method | FID↓ | KID$^{(\times 10^3)}$↓ | NIQE↓ | ID↑ | LPIPS↓ | PSNR↑ | MSSIM↑ |
|---|---|---|---|---|---|---|---|
| w/o anchor | 32.4085 | 18.9718 | 4.4316 | 0.7817 | 0.4309 | 22.7645 | 0.7009 |
| w/o noise | 31.1461 | 16.7528 | 4.3895 | 0.7887 | 0.4132 | 22.6009 | 0.6934 |
| w/o w | 36.0475 | 20.5331 | 4.2209 | 0.7790 | 0.4150 | 22.4965 | 0.6932 |
| w/o g | 30.2329 | 15.5456 | 4.2967 | 0.7896 | 0.4121 | 22.4442 | 0.6867 |
| w/o $\ell_1$-ball | 26.8366 | 13.1830 | 4.4294 | 0.8040 | 0.3998 | 23.7940 | 0.7288 |
| RLS$^+$ | 28.3786 | 13.7663 | 4.2878 | 0.7981 | 0.3972 | 23.5242 | 0.7195 |

Table 2. Ablation study on 16x super-resolution. (The best and the second-best are emphasized by blue and red respectively.)
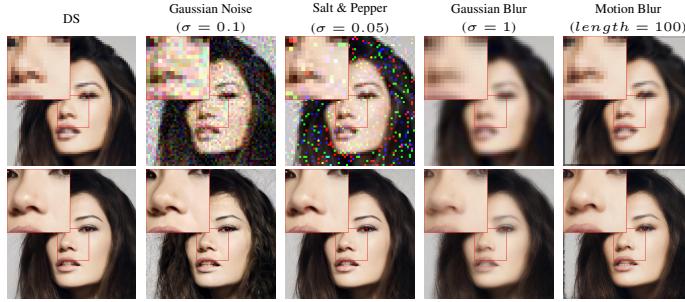


Figure 6. Robustness evaluation: degradation includes downscaling followed by corresponding operations, except for the last column, which consists of motion blur followed by downscaling (16x).

BRGM [17]. Furthermore, PULSE [18] appears to produce a wider squared norm distribution while ours approaches the squared norm distribution of $\mathcal{Z}$ more closely.

**Preservation of domain integrity.** At higher magnification factors (e.g. 64x), it is essential to increase the weight of the image domain prior (e.g., faces) over maintaining LR consistency to avoid generating images that do not resemble a face. Whereas, for lower magnification factors (e.g. 8x), where low-resolution images already contain fine grain details, the focus should be on preserving those details and ensure the LR consistency.

The qualitative comparison in Figure 4 shows that at 64x magnification, competing methods fail to generate reasonable facial details. Optimization-based methods such as PULSE and BRGM tend to generate images that are accurately downscaled to the LR image but at the cost of generating distorted faces, thus moving away from the face image manifold. This could be attributed to a lack of proper regularization, leading to over-fitting to the input LR image.

Moreover, GPEN and GFPGAN extract multi-resolution features from the LR input image and use them to modulate the intermediate features of the pre-trained StyleGAN model. However, at higher magnification factors, the input image contains limited spatial information, leading to insufficient features to generate fine facial details. In comparison, both RLS and RLS$^+$ succeed to produce plausible and realistic faces by maintaining the overall structure of the face and generating visually accurate details such as

eyes, eyebrows, teeth, mouth, and hair, among others. It is worth noting that while GPEN and GFPGAN yield similar outcomes at low magnifications (8x), they may generate unrealistic facial features when dealing with higher magnification factors. Additional examples can be found in Appendix D.

To quantitatively assess this gain in performance, we use Frechet Inception Distance (FID) [8] and Kernel Inception Distance (KID) [2] to measure the discrepancy between the real HR face images and the reconstructed one. We also employ Natural Image Quality Evaluator (NIQE) [19] to evaluate the naturalness of reconstructed images. As expected, the scores of FID, KID, and NIQE listed in Tab. 1 show that RLS$^+$ improves realism at both large and small magnifications, with improvements especially noticeable at higher magnification factors.

While even on the 8x magnification factor, the baselines still suffer from the realism-fidelity trade-off, RLS$^+$ performs significantly better both at reproducing the distribution of real HR face images and at producing images that have better perceptual quality. Thus, RLS$^+$ does not compromise realism despite modulating the generator. This is in accordance with our strategy that imposes the image to belong to a given domain.

Note that, with 8x magnification, even though RLS achieves the best NIQE scores (showing that the outputs look natural), which is also consistent with the visual results, it is almost the worst in terms of FID and KID. This is explained by the fact that FID and KID measure the dis-

tance between the distributions of real HR images and the reconstructed images, whereas RLS aims to produce a plausible image from the original image domain the StyleGAN was trained on.

Although the super-resolution task is ill-posed and has many plausible solutions, we evaluated the reconstruction quality using perceptual LPIPS [28], PSNR, and MS−SSIM [25]. The scores in Tab. 1 indicate that RLS$^+$ may not perform the best in terms of these metrics at low magnification factors, but it shows significant improvements in reconstruction quality at high magnification factors.

It is worth noting that the naive bicubic interpolator achieves high PSNR and MSSIM scores, but it fails to restore facial details, demonstrating that PSNR and SSIM are inadequate metrics for measuring super-resolution tasks. Moreover, RLS$^+$ achieves low LPIPS scores on both scales, indicating that the generated images are perceptually close to the ground-truth. Using a pre-trained face recognition model, CurricularFace [9], we measure the identity-similarity between the ground-truth and reconstructed images. With the 64x magnification factor, RLS reconstructs realistic images resembling the ground-truth, which is improved slightly by RLS$^+$. Note that, for high magnification factors we do not expect the output to perfectly match the ground-truth image as there are many plausible outputs.

Tab. 3 also presents the results of our approach on a 16$x$ super-resolution task and compares it with two baselines that achieved similar results on an 8x task, as shown in Tab. 1. The table demonstrates that while our approach produces fidelity values (such as LPIPS and PSNR) that are comparable to the baselines, it significantly improves realism (as indicated by FID and KID).

Additionally, we provide further examples that compare the output of RLS$^+$ with the baseline. These examples include out-of-domain images with challenging features such as extreme poses, heavy makeup, or occluded faces. Figure 8 shows that our approach can faithfully reconstruct images with challenging features, whereas the baselines struggle with such images. It produces realistic facial details and accurately preserves the individual's facial features.

| Method | FID↓ | KID$^{(\times 10^3)}$↓ | LPIPS↓ | PSNR↑ |
|---|---|---|---|---|
| GPEN [27] | 36.0024 | 29.5146 | 0.3945 | 25.0845 |
| GFPGAN [23] | 29.7925 | 18.3028 | 0.4221 | 24.2659 |
| RLS$^+$ | 26.3786 | 13.7664 | 0.3972 | 23.5242 |

Table 3. Quantitative comparison on CelebA for 16x super-resolution.

**Ablation** In this section, we first investigate the impact of parameters in the refinement step. To do so we compare

| Method | NIQE↓ | ID↑ | LPIPS↓ | MSSIM↑ |
|---|---|---|---|---|
| RLS$^+$ | 4.2975 | 0.7976 | 0.4004 | 0.7223 |
| Gaussian Noise | 4.0545 | 0.7056 | 0.4390 | 0.6675 |
| Salt and Pepper | 4.2427 | 0.8045 | 0.4036 | 0.7215 |
| Gaussian Blur | 4.4677 | 0.7901 | 0.4414 | 0.6819 |
| Motion Blur | 4.3160 | 0.8077 | 0.4237 | 0.6821 |

Table 4. Quantitative evaluation of robustness on 1000 images of CelebA (16x)

four variants of RLS$^+$ and show the results both quantitatively (Tab. 2) and qualitatively (Figure 5).

Firstly, "w/o anchor" is an optimization performed over $\mathbf{w}$, $\boldsymbol{\theta}$, and $\boldsymbol{\eta}$ that uses a mean value of random samples as initialization rather than the anchor point. The results produce faces with artifacts which suggests that it is crucial to first find the anchor point in the first step and then optimize from it.

Secondly, "w/o noise" and "w/o g" are variants where noise inputs and respectively generator's weights are removed from the optimization. Quantitative results show that training without noise inputs or generator weights achieves comparable performance. However, without optimizing these parameters, both reconstruction quality, and realism drop, suggesting they are necessary to synthesize facial details. This is further illustrated by qualitative results.

Thirdly, in the variant "w/o w", the latent code is fixed on the anchor point. One can see that w/o optimizing $\mathbf{w}$ the refinement can generate realistic face images by improving reconstruction loss; however, the identity of the face looks rather different from the ground-truth.

Finally, we evaluate the $\ell_1$-norm ball constraint during the fine-tuning. One can see that this variant harm realism, since without this constraint, the refinement process leads to over-fitting to the LR image query, generating unrealistic facial details. The latter also aligns with quantitative results where w/o $\ell_1$-norm ball outperforms RLS$^+$ on all metrics except NIQE and LPIPS. Overall, RLS$^+$ achieves better quantitative measures than its variants, showing that our choice for refinement better balances realism and fidelity.n

Figure 9 demonstrates another ablation experiments highlighting the function of other components of RLS image prior. First, "w/o Regu." searches the latent space without any regularization for the image that, once downscaled matches the LR image. The second variant is denoted "w/o $\mathcal{P}_{\text{cross}}$" refer to the suppression of $\mathcal{P}_{cross}$. The variant "w/o $\mathcal{W}^+$" refers to an optimization in the $\mathcal{W}$ space rather than in the $\mathcal{W}^+$.

To evaluate the three variants, we use the same set of parameters used in Sec. 4. One can see that searching the latent space without any regularization produces images that do not necessarily belong to the face domain and therefore

Figure 7. Generating multiple solutions for a given LR image. Apart from the first column that displays Low Resolution and Ground Truth images, top row is PULSE and bottom row is RLS (64x).
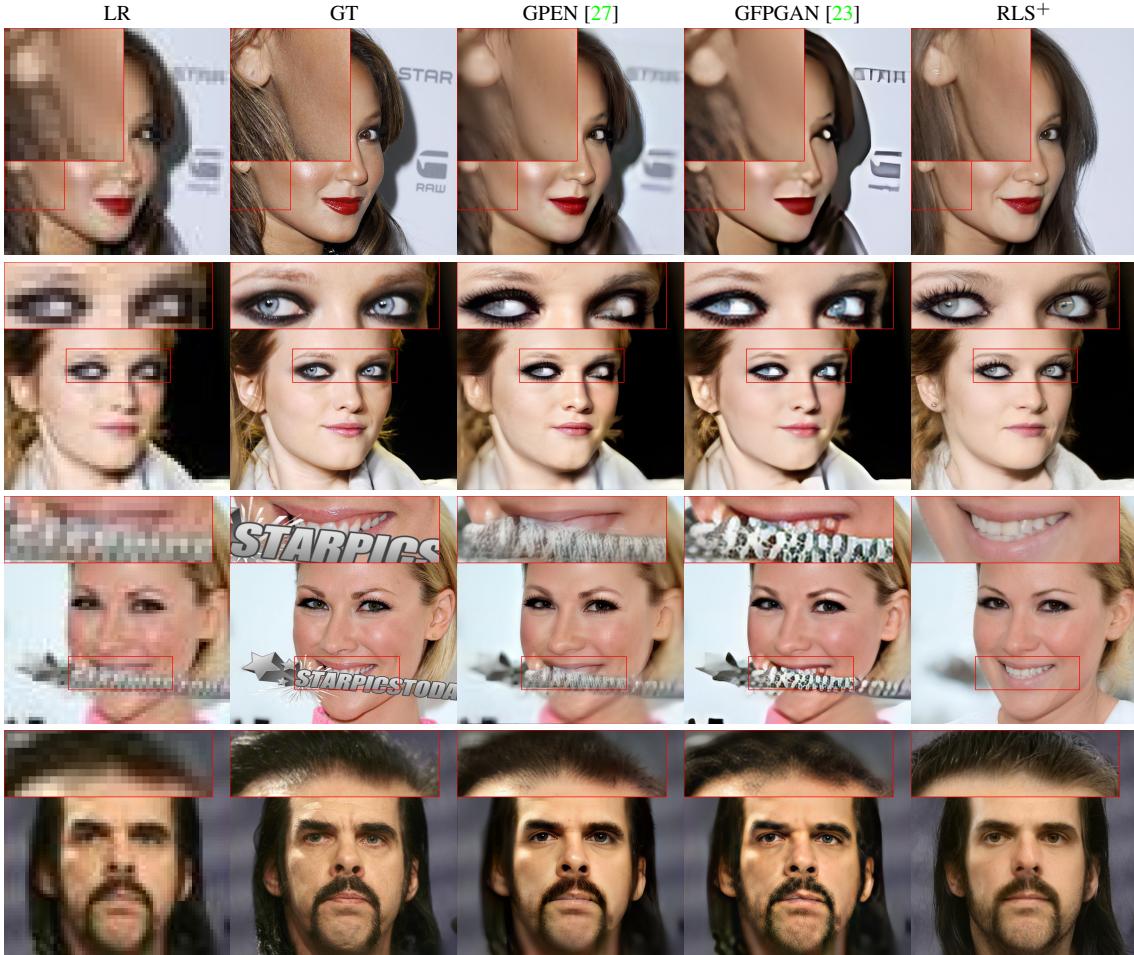


Figure 8. Comparison of reconstructions of a high resolution face from CelebA (16x). (Zoom-in for best view)

do not appear realistic. It also produces faces with artifacts when $\mathcal{P}_{\text{cross}}$ is discarded. This implies that the cross prior plays an important role in generating realistic facial details. Moreover, we can see that w/o $\mathcal{W}^+$ generates realistic HR face images; however, as it is reported in Tab. 5, RLS acheives higher values in terms of fidelity. This is because the latent space $\mathcal{W}$ is less expressive than $\mathcal{W}^+$, reducing the range of images that can be reconstructed with high fidelity.

We also explored the impact of $P_w$ in Equation (5). Figure 10 shows the impact of $P_w$ on KID, LPIPS, and MSSIM metrics, supporting the claim that the incorporation of $P_w$ improves the overall realism-fidelity trade-off.
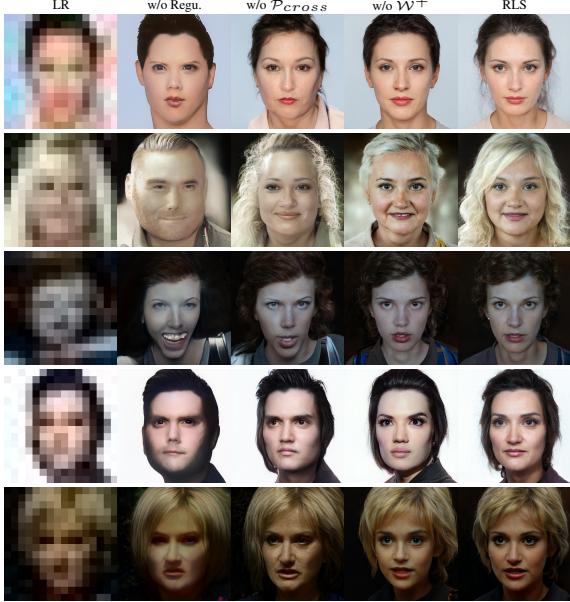
Figure 9. Qualitative comparison of different variants of RLS to evaluate the effectiveness of its image prior (64x). (Zoom-in for best view)

| Method | NIQE↓ | ID↑ | LPIPS↓ | MSSIM↑ |
|---|---|---|---|---|
| w/o Regu. | 4.2970 | 0.7003 | 0.5336 | 0.5508 |
| w/o $\mathcal{P}_{cross}$ | 4.1946 | 0.7078 | 0.4977 | 0.5488 |
| w/o $\mathcal{W}^+$ | 3.7184 | 0.7086 | 0.5106 | 0.5183 |
| RLS | 4.1032 | 0.7210 | 0.5037 | 0.5214 |

Table 5. Quantitative comparison of five variants of RLS (64x). (The best and the second-best are emphasized by blue and red respectively.)
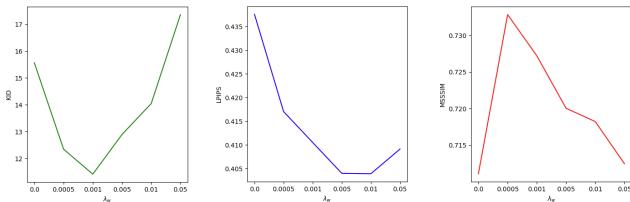


Figure 10. Examining the impact of $P_w$ (16x SR)

**Robustness to artifacts.** In contrast to supervised methods, which are sensitive to the input image domain, this approach is not restricted to a particular degradation operator used during training.

To assess this aspect, we applied additional degradation operators such as Gaussian noise, Salt and Pepper, and Gaussian blur to a bicubic downscaled image. We also applied motion blur to the HR image followed by downscaling. Figure 6 demonstrates the robustness evaluation when

the low-resolution image is degraded with various types of noise. RLS$^+$ can still generate realistic images while preserving their identity even in the presence of noise. However, the quality of the output may decrease, particularly when subjected to Gaussian blur, as the generator adapts to produce "noisy" images when the noise becomes more intense. This evaluation justifies the use of the bicubic downscaling operator during training instead of more complicated specific degradations.

It is worth noting that, although the quantitative results presented in Tab. 4 indicate a reduction in identity-similarity when exposed to Gaussian noise, the qualitative outcomes suggest that the additional noise only affects some low-level attributes, such as skin tone and hairstyle, but does not impact the individual's identity. This is illustrated in Figure 11 through additional examples for further support.
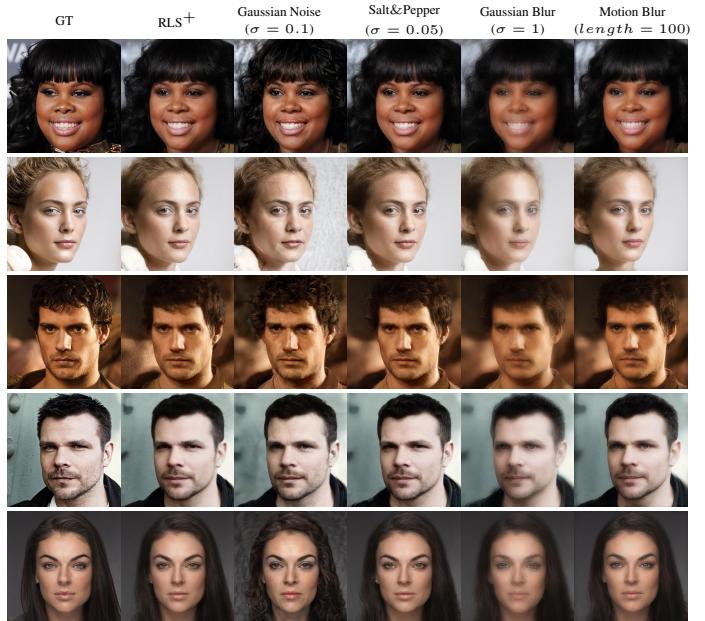


Figure 11. Robustness evaluation on 16x super-resolution.

**Diversity of the output.** At large magnification factors, many SR solutions could match the given LR input. Although these may not precisely match the HR ground-truth image, obtaining a range of diverse and plausible outputs can be advantageous. Different initialization of the latent code can be used to produce diverse SR outputs for the same LR input. In Figure 7, the diversity of outputs generated by RLS and PULSE for a magnification of 64x is compared. It shows that RLS can better produce multiple consistent and realistic HR images from a single LR image. We hypothesize that this could be explained by the fact that PULSE limits the solutions to lie on a hypersphere which does not

cover the whole distribution of images learned during Style-GAN training.

## 5. Conclusion

Super-resolved images reconstructed by existing self-supervised approaches, exploiting a pre-trained style-based generative model, seem to suffer from a lack of realism, especially when the original image is out-of-domain. With this work, we address this issue by first introducing a new regularization of the latent space exploration, which leverages a normalizing flow model, providing a more robust image prior to ensuring that the latent code remains in the original generative model manifold. Then, in order to reconstruct the image with higher fidelity, we slightly fine-tune the generative prior, within a small $\ell_1$-norm ball centered at the latent code obtained at the first step. Doing so enables us to mitigate the fidelity-realness trade-off. We performed extensive experiments demonstrating that our method can generate high-quality face images with clear facial details from severely degraded ones, outperforming prior works.

## References

[1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *Int. Conf. Comput. Vis.*, pages 4432–4441, 2019. 3

[2] Mikolaj Binkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. *arXiv preprint arXiv:1801.01401*, 2018. 7

[3] Ashish Bora, Ajil Jalal, Eric Price, and Alexandros G Dimakis. Compressed sensing using generative models. In *Int. Conf. Mach. Learn.*, pages 537–546. PMLR, 2017. 2

[4] Milan Cvitkovic and Gunther Koliander. Minimal achievable sufficient statistic learning. In *Int. Conf. Mach. Learn.*, pages 1465–1474. PMLR, 2019. 4

[5] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. *arXiv preprint arXiv:1605.08803*, 2016. 3

[6] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016. 2

[7] Jingwen He, Wu Shi, Kai Chen, Lean Fu, and Chao Dong. Gcfsr: a generative and controllable face super resolution method without facial and gan priors. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1889–1898, 2022. 1

[8] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Adv. Neural Inform. Process. Syst.*, 30, 2017. 7

[9] Yuge Huang, Yuhan Wang, Ying Tai, Xiaoming Liu, Pengcheng Shen, Shaoxin Li, Jilin Li, and Feiyue Huang. Curricularface: adaptive curriculum learning loss for deep face recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5901–5910, 2020. 8

[10] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1125–1134, 2017. 1

[11] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017. 6

[12] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 4401–4410, 2019. 3, 6

[13] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 8110–8119, 2020. 3, 6

[14] Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. In *Adv. Neural Inform. Process. Syst.*, 2022. 2, 4, 5, 6

[15] Bahjat Kawar, Gregory Vaksman, and Michael Elad. Snips: Solving noisy inverse problems stochastically. *Advances in Neural Information Processing Systems*, 34:21757–21769, 2021. 2

[16] Ivan Kobyzev, Simon JD Prince, and Marcus A Brubaker. Normalizing flows: An introduction and review of current methods. *IEEE Trans. Pattern Anal. Mach. Intell.*, 43(11):3964–3979, 2020. 3

[17] Razvan Marinescu, Daniel Moyer, and Polina Golland. Bayesian image reconstruction using deep generative models. In *NeurIPS Workshop on Deep Generative Models and Downstream Applications*, 2021. 2, 3, 4, 5, 6, 7

[18] Sachit Menon, Alexandru Damian, Shijia Hu, Nikhil Ravi, and Cynthia Rudin. Pulse: Self-supervised photo upsampling via latent space exploration of generative models. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2437–2445, 2020. 2, 4, 5, 6, 7

[19] Anish Mittal, Rajiv Soundararajan, and Alan C Bovik. Making a "completely blind" image quality analyzer. *IEEE Sign. Process. Letters*, 20(3):209–212, 2012. 7

[20] George Papamakarios, Eric T Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *J. Mach. Learn. Res.*, 22(57):1–64, 2021. 3

[21] George Papamakarios, Theo Pavlakou, and Iain Murray. Masked autoregressive flow for density estimation. *NIPS*, 30, 2017. 3

[22] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 8798–8807, 2018. 1

[23] Xintao Wang, Yu Li, Honglun Zhang, and Ying Shan. Towards real-world blind face restoration with generative facial prior. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 9168–9178, 2021. 2, 4, 5, 6, 8, 9

[24] 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 *Eur. Conf. Comput. Vis.*, pages 0–0, 2018. 1

[25] Zhou Wang, Eero P Simoncelli, and Alan C Bovik. Multiscale structural similarity for image quality assessment. In *The Thrity-Seventh Asilomar Conference on Signals, Systems and Computers, 2003*, volume 2, pages 1398–1402. IEEE, 2003. 8

[26] Jonas Wulff and Antonio Torralba. Improving inversion and generation diversity in stylegan using a gaussianized latent space. *arXiv preprint arXiv:2009.06529*, 2020. 3

[27] Tao Yang, Peiran Ren, Xuansong Xie, and Lei Zhang. Gan prior embedded network for blind face restoration in the wild. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 672–681, 2021. 2, 4, 5, 6, 8, 9

[28] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 586–595, 2018. 8