# Image Super-Resolution via Iterative Refinement

Chitwan Saharia[†] Jonathan Ho, William Chan, Tim Salimans, David J. Fleet, Mohammad Norouzi

`{sahariac, jonathanho,williamchan,salimans,davidfleet,mnorouzi}@google.com`
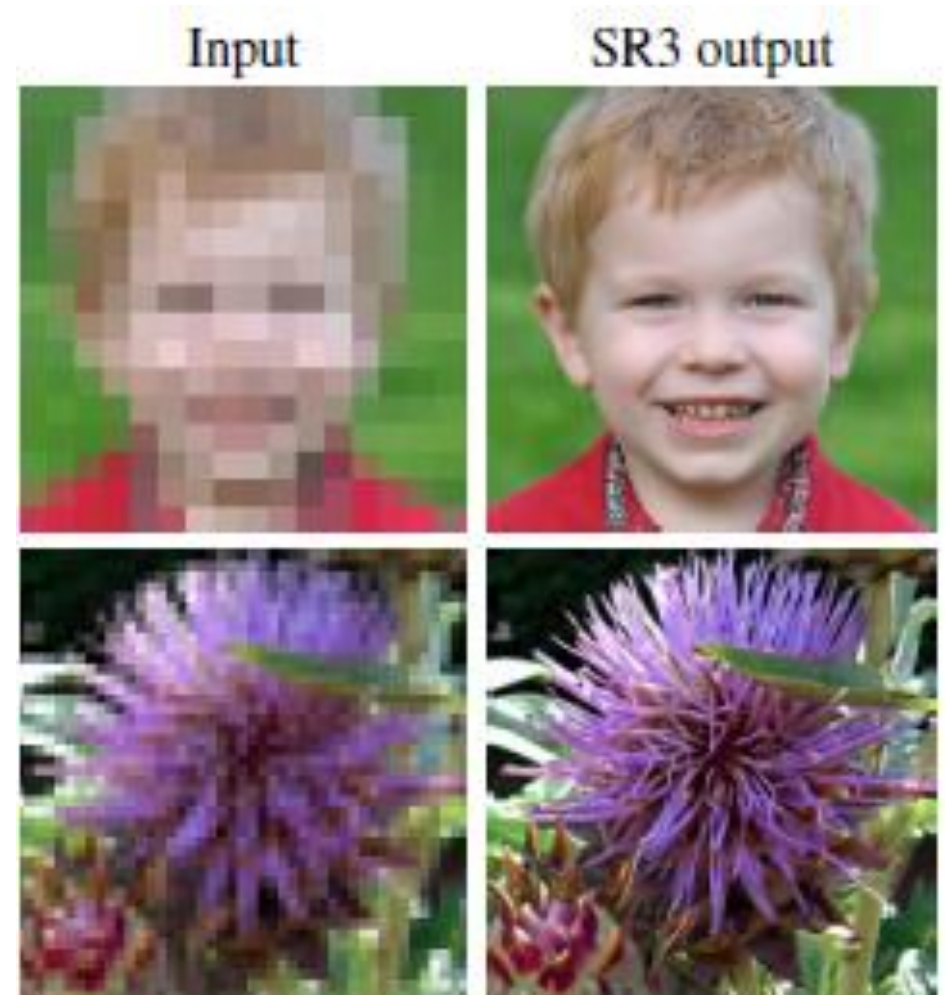
Google Research, Brain Team

2022.06.27(월)

남궁영수

# Contents

# Introduction

- Single-image super-resolution

- It is a process of generating a high-resolution image that is consistent with an input low-resolution image.

- It is inverse problem

- It is challenging because multiple output images may be consistent with a single input image.



Input    SR3 output

# Introduction

- Generative Model

- Autoregressive model :

- VAEs :

- Flows :

- GAN :

# Introduction

- Generative Model

- Autoregressive model : expensive for high-resolution image generation

- VAEs : often yield sub-optimal sample quality

- Flows : often yield sub-optimal sample quality

- GAN : require carefully designed regularization and optimization trick to tame optimization instability and mode collapse

# Introduction

- SR3

- Inspired by Denoising Diffusion Probabilistic Models(DDPM), and denoising score matching.

- Works by learning to transform a standard normal distribution into an empirical data distribution through a sequence of refinement steps.
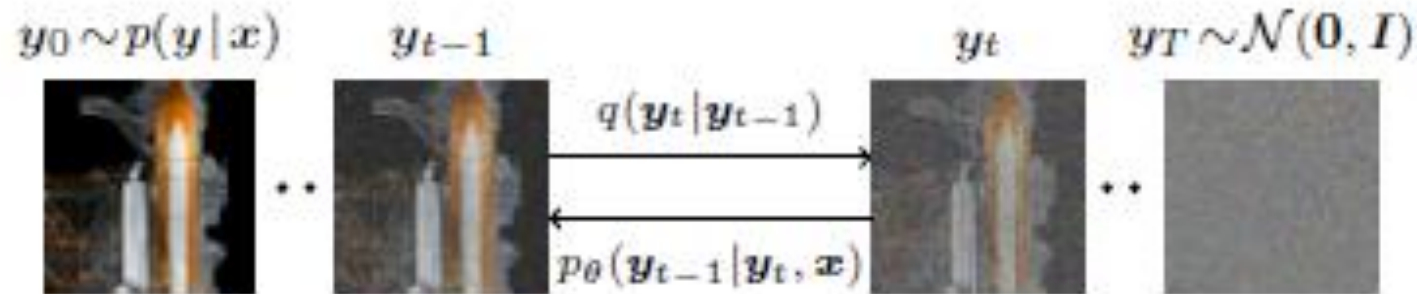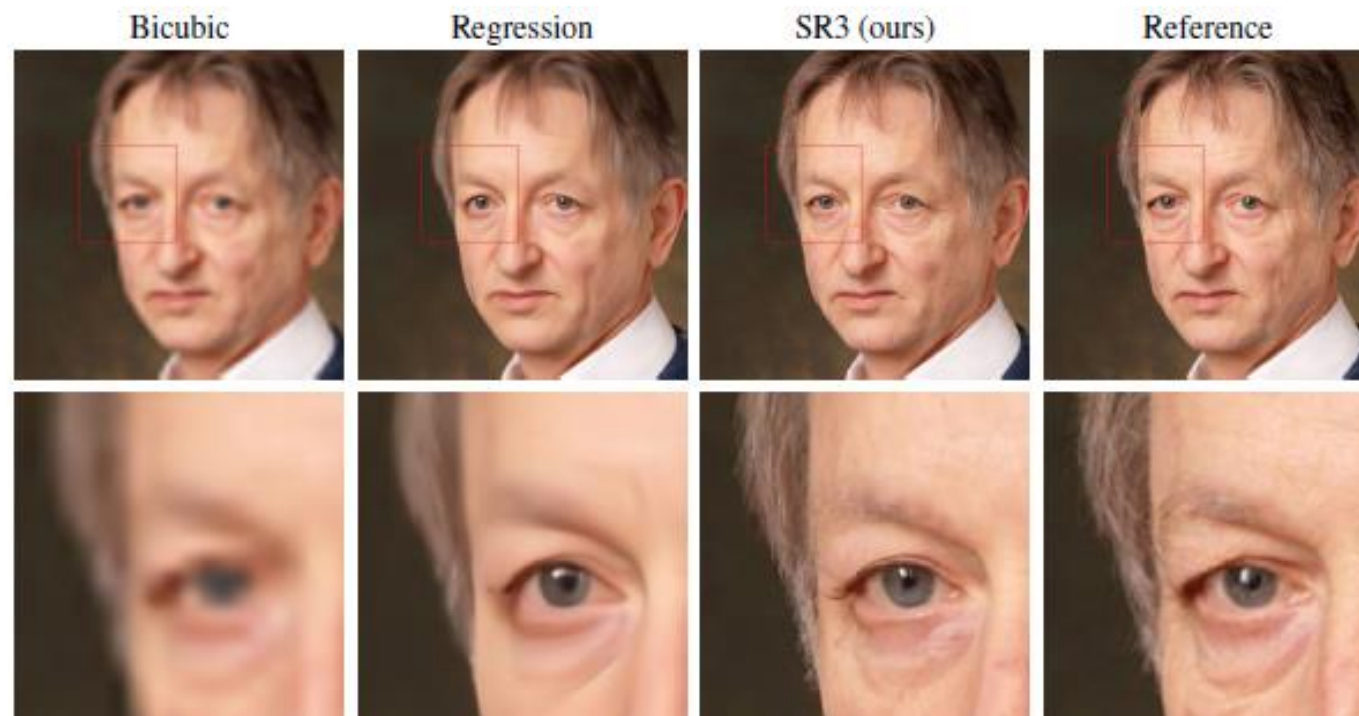


**Figure 2:** The forward diffusion process $q$ (left to right) gradually adds Gaussian noise to the target image. The reverse inference process $p$ (right to left) iteratively denoises the target image conditioned on a source image $x$. Source image $x$ is not shown here.

# Introduction

- Fool rate

- PSNR and SSIM don't reflect human preference well when the input resolution is low and the magnification ratio is large

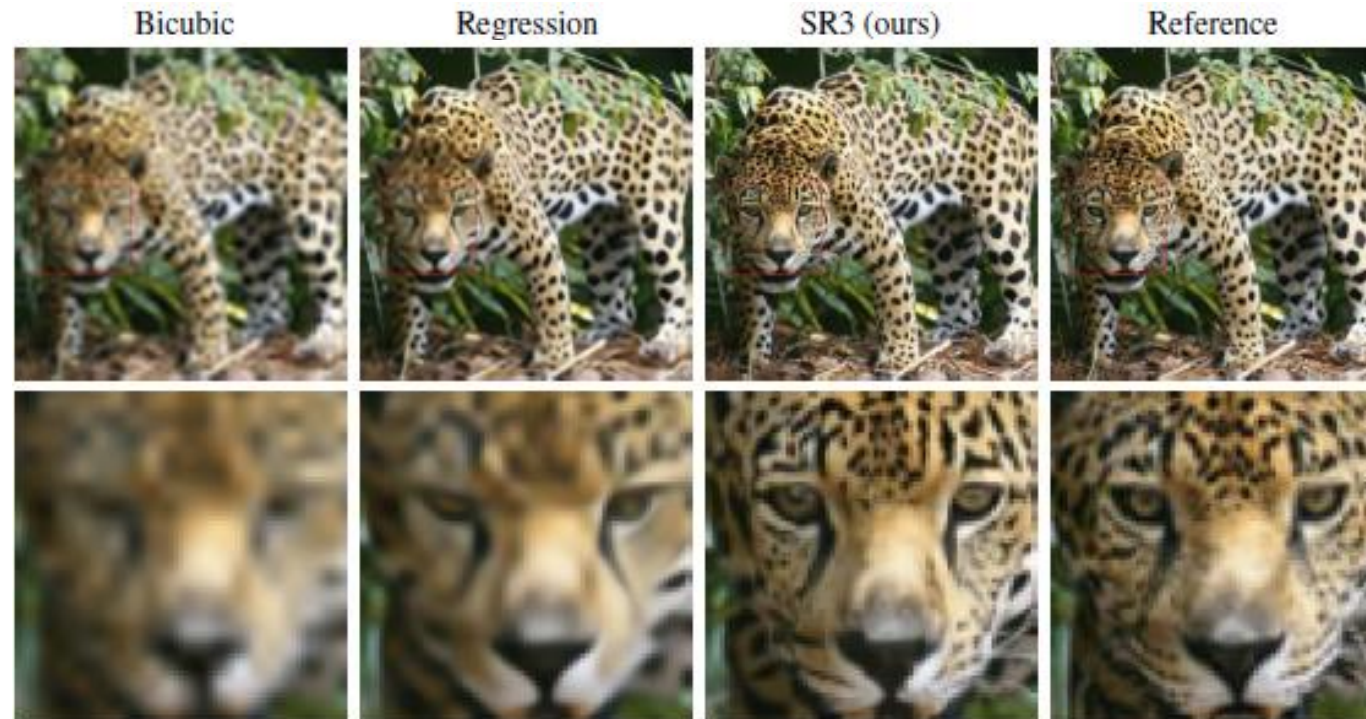- PSNR and SSIM often penalize synthetic high-frequency details.

# Introduction

- Fool rate

- PSNR and SSIM don't reflect human preference well when the input resolution is low and the magnification ratio is large

- PSNR and SSIM often penalize synthetic high-frequency details.

# Introduction

- Fool rate

- <span style="color:red">So we resort to human evaluation.</span>

- Human subjects are shown a low-resolution input and are required to select between a model output and a ground truth image.

Low-resolution input

# Introduction

- Fool rate

- So we resort to human evaluation.
- Human subjects are shown a low-resolution input and are required to select between a model output and a ground truth image.



Select between Model output and ground truth

# Key Contribution

- adapt denoising diffusion models to conditional image generation.

- prove effective on face and natural image super-resolution at different magnification factors. On standard 8x face super-resolution task, SR3 achieves a human fool rate close to 50%, outperforming PSRGAN and PULSE.

- demonstrate unconditional and class-conditional generation by cascading a 64x64 image synthesis model with SR3 models.

# Conditional Denoising Diffusion Model

dataset

Unknown conditional distribution

$$\mathcal{D} = \{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^{N}$$

Sample from

$$p(\boldsymbol{y} \mid \boldsymbol{x})$$

# Conditional Denoising Diffusion Model

dataset

Unknown conditional distribution

$$\mathcal{D} \stackrel{-}{=} \{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^N$$

Sample from

$$p(\boldsymbol{y} \mid \boldsymbol{x})$$

$\boldsymbol{x}$



$\boldsymbol{y}_0 \sim p(\boldsymbol{y} \mid \boldsymbol{x})$

# Conditional Denoising Diffusion Model

dataset

Unknown conditional distribution

$$\mathcal{D} \;\dot{=}\; \{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^{N} \quad \xleftarrow{\text{Sample from}} \quad p(\boldsymbol{y} \mid \boldsymbol{x})$$

$\boldsymbol{x}$

$\boldsymbol{y}_0 \sim p(\boldsymbol{y} \mid \boldsymbol{x})$  $\boldsymbol{y}_{t-1}$  $\boldsymbol{y}_t$  $\boldsymbol{y}_T \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$

$q(\boldsymbol{y}_t \mid \boldsymbol{y}_{t-1})$

$p_\theta(\boldsymbol{y}_{t-1} \mid \boldsymbol{y}_t, \boldsymbol{x})$
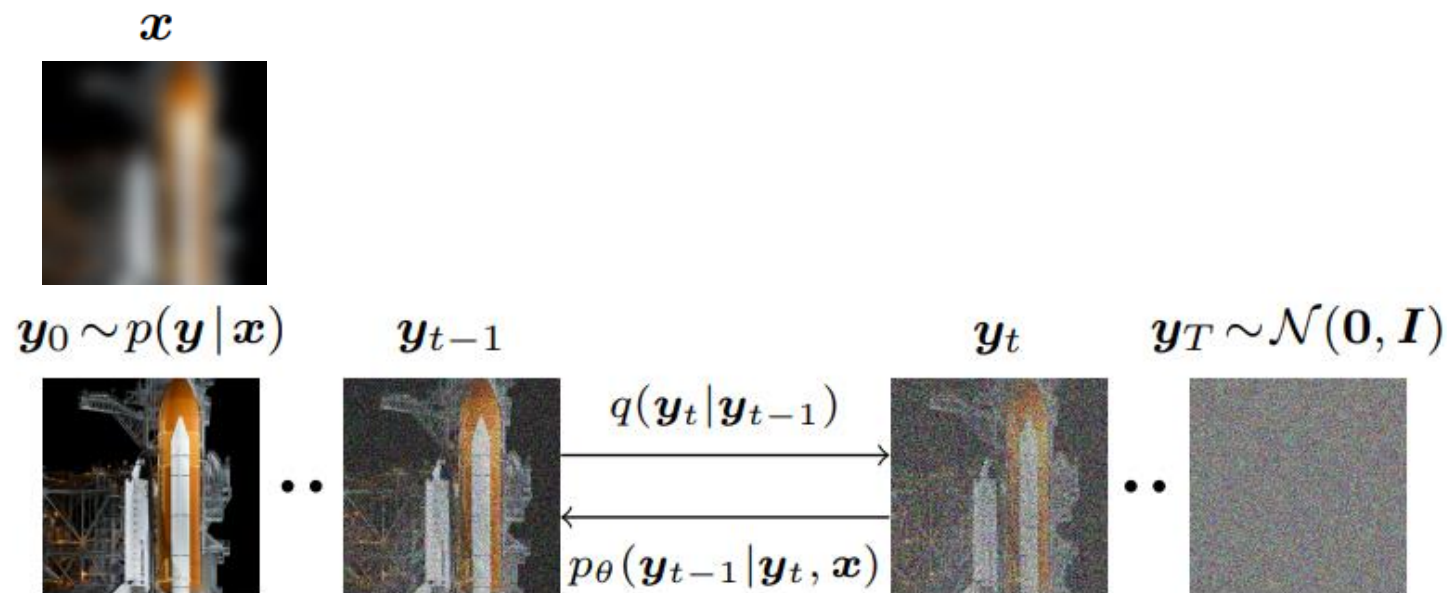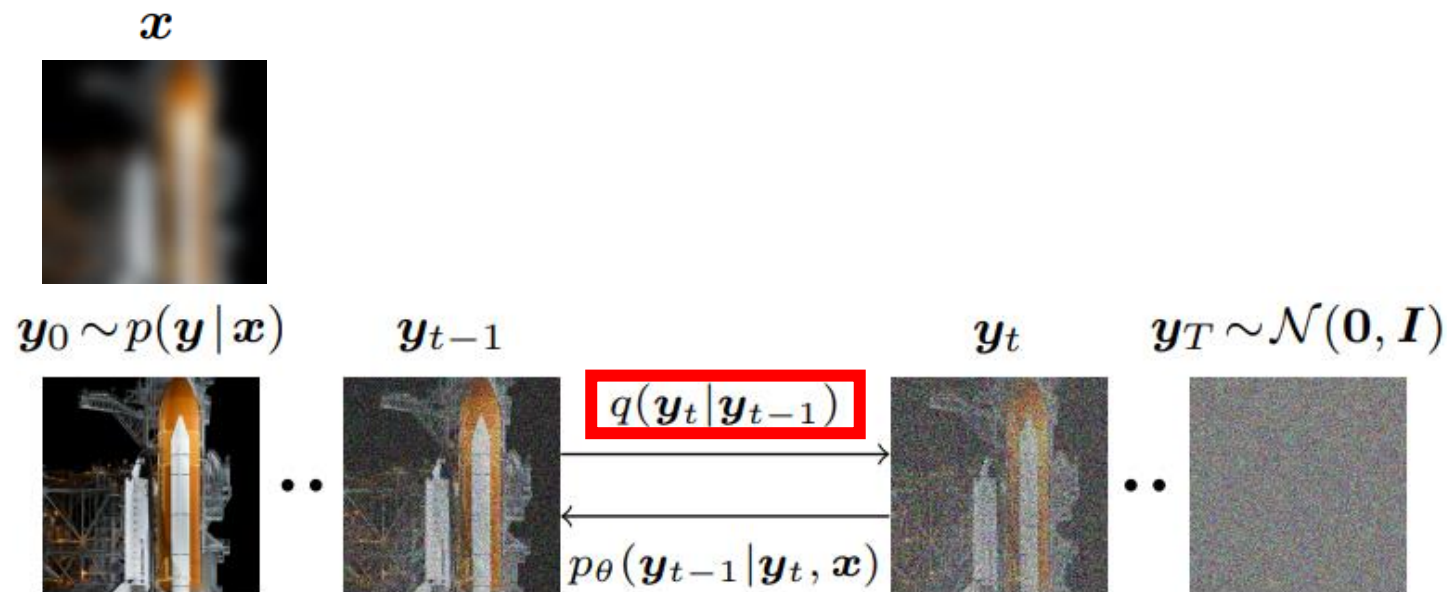
**Figure 2:** The forward diffusion process $q$ (left to right) gradually adds Gaussian noise to the target image. The reverse inference process $p$ (right to left) iteratively denoises the target image conditioned on a source image $\boldsymbol{x}$. Source image $\boldsymbol{x}$ is not shown here.

# Conditional Denoising Diffusion Model

dataset

Unknown conditional distribution

$$\mathcal{D} \doteq \{x_i, y_i\}_{i=1}^{N}$$

Sample from $\longleftarrow$ $p(y \mid x)$

$x$



$y_0 \sim p(y \mid x)$ $\qquad$ $y_{t-1}$ $\qquad$ $y_t$ $\qquad$ $y_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

$q(y_t \mid y_{t-1})$

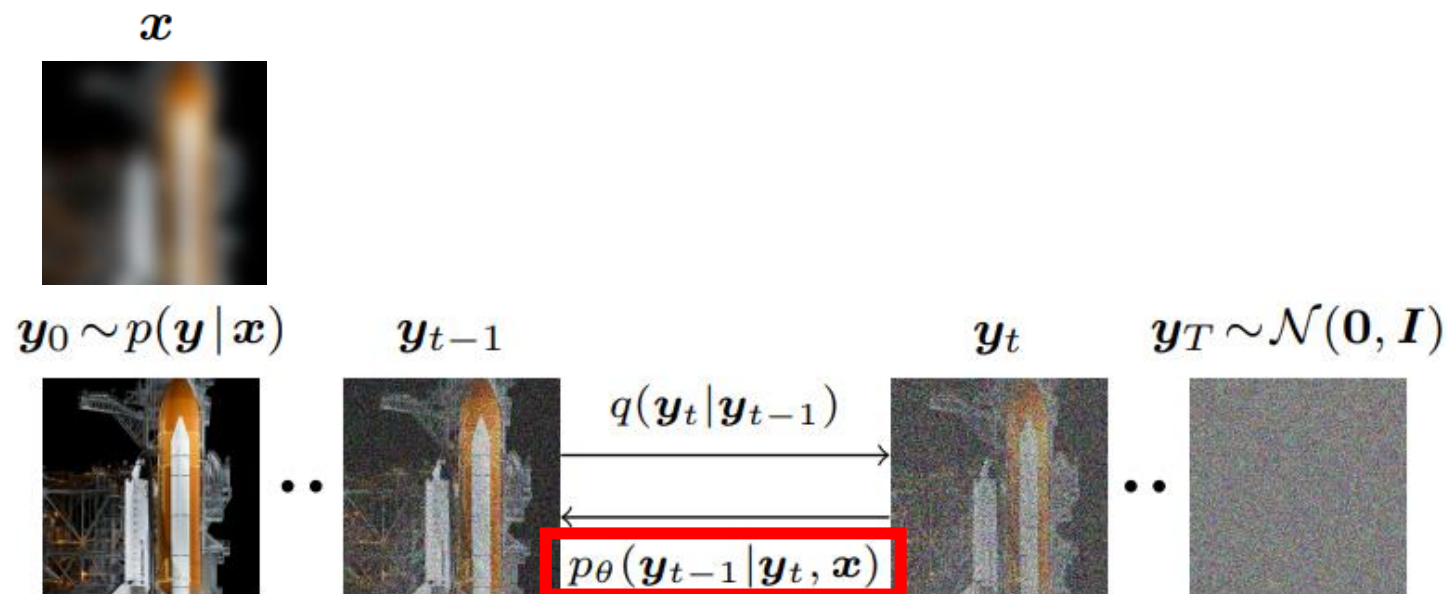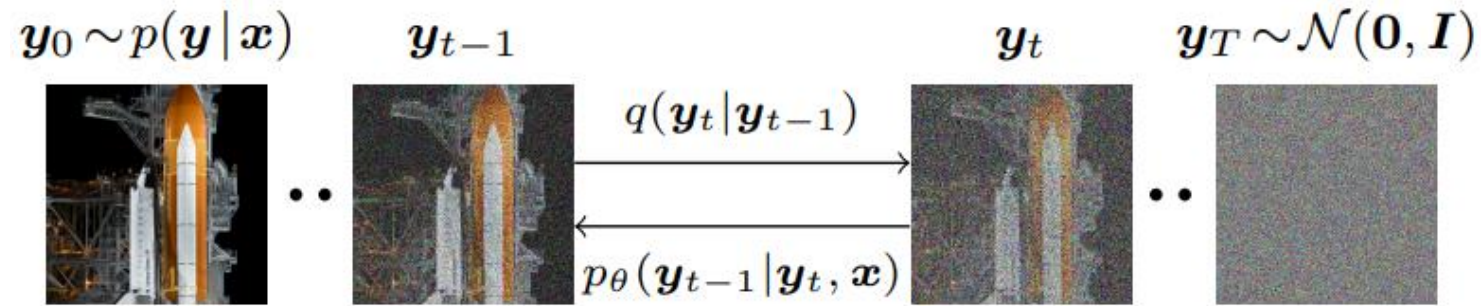$p_\theta(y_{t-1} \mid y_t, x)$

**Figure 2:** The forward diffusion process $q$ (left to right) gradually adds Gaussian noise to the target image. The reverse inference process $p$ (right to left) iteratively denoises the target image conditioned on a source image $x$. Source image $x$ is not shown here.

# Conditional Denoising Diffusion Model

dataset

Unknown conditional distribution

$$\mathcal{D} \doteq \{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^{N}$$

Sample from

$$p(\boldsymbol{y} \mid \boldsymbol{x})$$

$\boldsymbol{x}$

$\boldsymbol{y}_0 \sim p(\boldsymbol{y} \mid \boldsymbol{x})$     $\boldsymbol{y}_{t-1}$     $\boldsymbol{y}_t$     $\boldsymbol{y}_T \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$

$q(\boldsymbol{y}_t \mid \boldsymbol{y}_{t-1})$

$p_\theta(\boldsymbol{y}_{t-1} \mid \boldsymbol{y}_t, \boldsymbol{x})$
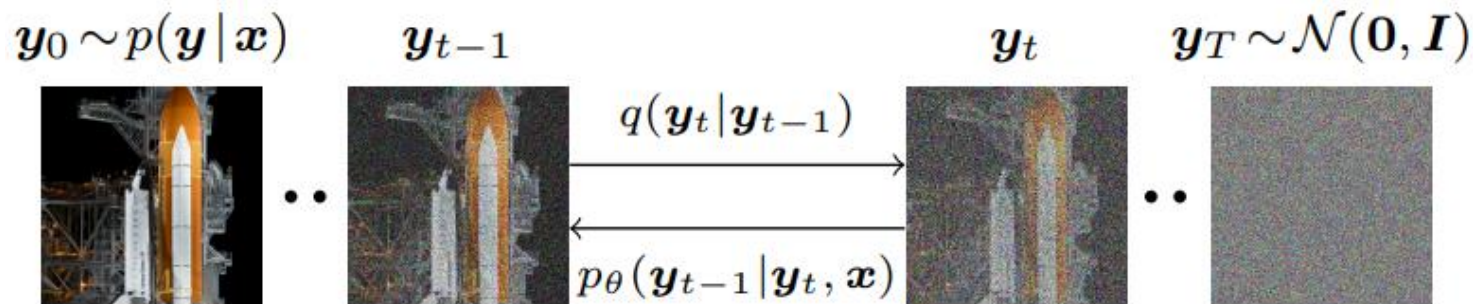
**Figure 2:** The forward diffusion process $q$ (left to right) gradually adds Gaussian noise to the target image. The reverse inference process $p$ (right to left) iteratively denoises the target image conditioned on a source image $\boldsymbol{x}$. Source image $\boldsymbol{x}$ is not shown here.

# Conditional Denoising Diffusion Model



$y_0 \sim p(y|x)$  $\quad y_{t-1}$  $\quad y_t$  $\quad y_T \sim \mathcal{N}(0, I)$

$q(y_t | y_{t-1})$

$p_\theta(y_{t-1} | y_t, x)$

# Conditional Denoising Diffusion Model

$$y_0 \sim p(y|x) \qquad y_{t-1} \qquad\qquad\qquad y_t \qquad y_T \sim \mathcal{N}(0, I)$$



$$q(y_t|y_{t-1})$$

$$p_\theta(y_{t-1}|y_t, x)$$

## 2.1. Gaussian Diffusion Process

$$q(\boldsymbol{y}_{1:T} \mid \boldsymbol{y}_0) = \prod_{t=1}^{T} q(\boldsymbol{y}_t \mid \boldsymbol{y}_{t-1}), \qquad (1)$$

$$q(\boldsymbol{y}_t \mid \boldsymbol{y}_{t-1}) = \mathcal{N}(\boldsymbol{y}_t \mid \sqrt{\alpha_t}\, \boldsymbol{y}_{t-1}, (1 - \alpha_t)\boldsymbol{I}), \qquad (2)$$

# Conditional Denoising Diffusion Model



$$\boldsymbol{y}_0 \sim p(\boldsymbol{y}|\boldsymbol{x}) \qquad \boldsymbol{y}_{t-1} \qquad\qquad\qquad \boldsymbol{y}_t \qquad \boldsymbol{y}_T \sim \mathcal{N}(\boldsymbol{0},\boldsymbol{I})$$

$$q(\boldsymbol{y}_t|\boldsymbol{y}_{t-1})$$

$$p_\theta(\boldsymbol{y}_{t-1}|\boldsymbol{y}_t,\boldsymbol{x})$$

## 2.1. Gaussian Diffusion Process

$$q(\boldsymbol{y}_{1:T} \mid \boldsymbol{y}_0) = \prod_{t=1}^{T} q(\boldsymbol{y}_t \mid \boldsymbol{y}_{t-1}), \qquad (1)$$

$$q(\boldsymbol{y}_t \mid \boldsymbol{y}_{t-1}) = \mathcal{N}(\boldsymbol{y}_t \mid \sqrt{\alpha_t}\, \boldsymbol{y}_{t-1}, (1-\alpha_t)\boldsymbol{I}), \quad (2)$$

$$q(\boldsymbol{y}_t \mid \boldsymbol{y}_0) = \mathcal{N}(\boldsymbol{y}_t \mid \sqrt{\gamma_t}\, \boldsymbol{y}_0, (1-\gamma_t)\boldsymbol{I}), \qquad (3) \quad \text{where } \gamma_t = \prod_{i=1}^{t} \alpha_i$$

# Conditional Denoising Diffusion Model

$$y_0 \sim p(y|x) \qquad y_{t-1} \qquad\qquad y_t \qquad y_T \sim \mathcal{N}(0, I)$$



$$q(y_t|y_{t-1})$$

$$p_\theta(y_{t-1}|y_t, x)$$

## 2.1. Gaussian Diffusion Process

$$q(y_{1:T} \mid y_0) \;=\; \prod_{t=1}^{T} q(y_t \mid y_{t-1}), \qquad\qquad (1)$$

$$q(y_t \mid y_{t-1}) \;=\; \mathcal{N}(y_t \mid \sqrt{\alpha_t}\, y_{t-1}, (1 - \alpha_t)I), \quad (2)$$

$$q(y_t \mid y_0) = \mathcal{N}(y_t \mid \sqrt{\gamma_t}\, y_0, (1 - \gamma_t)I), \qquad (3) \quad \text{where } \gamma_t = \prod_{i=1}^{t} \alpha_i$$

With some algebraic manipulation, one can derive the posterior distribution of $y_{t-1}$ given $(y_0, y_t)$

# Conditional Denoising Diffusion Model

$y_0 \sim p(y|x)$   $y_{t-1}$   $y_t$   $y_T \sim \mathcal{N}(0, I)$



$q(y_t | y_{t-1})$

$p_\theta(y_{t-1} | y_t, x)$

## 2.1. Gaussian Diffusion Process

$$q(y_{1:T} | y_0) = \prod_{t=1}^{T} q(y_t | y_{t-1}), \qquad (1)$$

$$q(y_t | y_{t-1}) = \mathcal{N}(y_t | \sqrt{\alpha_t}\, y_{t-1}, (1-\alpha_t)I), \quad (2)$$

$$q(y_t | y_0) = \mathcal{N}(y_t | \sqrt{\gamma_t}\, y_0, (1-\gamma_t)I), \qquad (3) \quad \text{where } \gamma_t = \prod_{i=1}^{t} \alpha_i$$

With some algebraic manipulation, one can derive the posterior distribution of $y_{t-1}$ given $(y_0, y_t)$

$$q(y_{t-1} | y_0, y_t) = \mathcal{N}(y_{t-1} | \mu, \sigma^2 I)$$

$$\mu = \frac{\sqrt{\gamma_{t-1}}(1-\alpha_t)}{1-\gamma_t} y_0 + \frac{\sqrt{\alpha_t}(1-\gamma_{t-1})}{1-\gamma_t} y_t \qquad (4)$$

$$\sigma^2 = \frac{(1-\gamma_{t-1})(1-\alpha_t)}{1-\gamma_t}.$$

# Conditional Denoising Diffusion Model

**2.2. Optimizing the Denoising Model**

$$\mathbb{E}_{(x,y)}\mathbb{E}_{\epsilon,\gamma}\left\|f_\theta(x,\underbrace{\sqrt{\overline{\gamma}}\,y_0 + \sqrt{1-\gamma}\,\epsilon}_{\tilde{y}},\gamma) - \epsilon\right\|_p^p, \quad (6)$$

$$p \in \{1,2\}$$

# Conditional Denoising Diffusion Model

## 2.2. Optimizing the Denoising Model

$$\mathbb{E}_{(x,y)}\mathbb{E}_{\epsilon,\gamma}\left\|f_\theta(x, \underbrace{\sqrt{\gamma}\, y_0 + \sqrt{1-\gamma}\, \epsilon}_{\tilde{y}}, \gamma) - \epsilon\right\|_p^p, \quad (6) \qquad\qquad p \in \{1, 2\}$$

Input :   $x$

$$\tilde{y} = \sqrt{\gamma}\, y_0 + \sqrt{1-\gamma}\, \epsilon, \qquad \epsilon \sim \mathcal{N}(0, I), \qquad (5)$$

( compatible with   $q(y_t \mid y_0) = \mathcal{N}(y_t \mid \sqrt{\gamma_t}\, y_0, (1-\gamma_t)I)$,   (3)  )

$\gamma \sim p(\gamma)$   ( sampling method :  sample t~{0,...,T} and sample γ $\sim U(\gamma_{t-1}, \gamma_t)$  )
( scaling factor for noise )

Output : estimation of  $\epsilon \sim \mathcal{N}(0, I)$

# Conditional Denoising Diffusion Model

## 2.3. Inference via Iterative Refinement

$$\tilde{y} = \sqrt{\gamma}\, y_0 + \sqrt{1-\gamma}\, \epsilon\,, \qquad \epsilon \sim \mathcal{N}(0, I)\,, \qquad (5)$$

# Conditional Denoising Diffusion Model

**2..3. Inference via Iterative Refinement**

$$\tilde{y} = \sqrt{\gamma}\, y_0 + \sqrt{1 - \gamma}\, \epsilon\,, \qquad \epsilon \sim \mathcal{N}(0, I)\,, \qquad (5)$$

$$\hat{y}_0 = \frac{1}{\sqrt{\gamma_t}} \left( y_t - \sqrt{1 - \gamma_t}\, f_\theta(x, y_t, \gamma_t) \right)\,. \qquad (10)$$

# Conditional Denoising Diffusion Model

## 2..3. Inference via Iterative Refinement

$$\tilde{y} = \sqrt{\bar{\gamma}}\, y_0 + \sqrt{1 - \bar{\gamma}}\, \epsilon\,, \qquad \epsilon \sim \mathcal{N}(0, I)\,, \qquad (5)$$

$$\hat{y}_0 = \frac{1}{\sqrt{\gamma_t}} \left( y_t - \sqrt{1 - \gamma_t}\, f_\theta(x, y_t, \gamma_t) \right)\,. \qquad (10)$$

$$q(y_{t-1} \mid y_0, y_t) = \mathcal{N}(y_{t-1} \mid \mu, \sigma^2 I)$$

$$\mu = \frac{\sqrt{\gamma_{t-1}}\,(1 - \alpha_t)}{1 - \gamma_t}\, y_0 + \frac{\sqrt{\alpha_t}\,(1 - \gamma_{t-1})}{1 - \gamma_t}\, y_t \qquad (4)$$

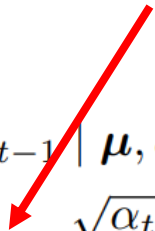$$\sigma^2 = \frac{(1 - \gamma_{t-1})(1 - \alpha_t)}{1 - \gamma_t}\,.$$

# Conditional Denoising Diffusion Model

## 2.3. Inference via Iterative Refinement

$$\tilde{y} = \sqrt{\bar{\gamma}}\, y_0 + \sqrt{1 - \gamma}\, \epsilon \,, \qquad \epsilon \sim \mathcal{N}(0, I) \,, \qquad (5)$$

$$\hat{y}_0 = \frac{1}{\sqrt{\gamma_t}} \left( y_t - \sqrt{1 - \gamma_t}\, f_\theta(x, y_t, \gamma_t) \right) . \qquad (10)$$

$$q(y_{t-1} \mid y_0, y_t) = \mathcal{N}(y_{t-1} \mid \mu, \sigma^2 I)$$

$$\mu = \frac{\sqrt{\gamma_{t-1}}\,(1 - \alpha_t)}{1 - \gamma_t} y_0 + \frac{\sqrt{\alpha_t}\,(1 - \gamma_{t-1})}{1 - \gamma_t} y_t \qquad (4) \qquad \mu_\theta(x, y_t, \gamma_t) = \frac{1}{\sqrt{\alpha_t}} \left( y_t - \frac{1 - \alpha_t}{\sqrt{1 - \gamma_t}} f_\theta(x, y_t, \gamma_t) \right) ,$$

$$\sigma^2 = \frac{(1 - \gamma_{t-1})(1 - \alpha_t)}{1 - \gamma_t} . \qquad (11)$$

# Conditional Denoising Diffusion Model

## 2..3. Inference via Iterative Refinement

$$\tilde{y} = \sqrt{\gamma}\,y_0 + \sqrt{1-\gamma}\,\epsilon\,, \qquad \epsilon \sim \mathcal{N}(0, I)\,, \qquad (5)$$

$$\hat{y}_0 = \frac{1}{\sqrt{\gamma_t}}\left(y_t - \sqrt{1-\gamma_t}\,f_\theta(x, y_t, \gamma_t)\right)\,. \qquad (10)$$

$$q(y_{t-1} \mid y_0, y_t) = \mathcal{N}(y_{t-1} \mid \mu, \sigma^2 I)$$

$$\mu = \frac{\sqrt{\gamma_{t-1}}\,(1-\alpha_t)}{1-\gamma_t}\,y_0 + \frac{\sqrt{\alpha_t}\,(1-\gamma_{t-1})}{1-\gamma_t}\,y_t \qquad (4)$$

$$\sigma^2 = \frac{(1-\gamma_{t-1})(1-\alpha_t)}{1-\gamma_t}\,.$$

$$\mu_\theta(x, y_t, \gamma_t) = \frac{1}{\sqrt{\alpha_t}}\left(y_t - \frac{1-\alpha_t}{\sqrt{1-\gamma_t}}\,f_\theta(x, y_t, \gamma_t)\right)\,, \qquad (11)$$

and we set the variance of $p_\theta(y_{t-1} \mid y_t, x)$ to $(1-\alpha_t)$, a default given by the variance of the forward process [17].

# Conditional Denoising Diffusion Model

## 2..3. Inference via Iterative Refinement

$$\tilde{y} = \sqrt{\gamma}\, y_0 + \sqrt{1 - \gamma}\, \epsilon\, , \qquad \epsilon \sim \mathcal{N}(0, I)\, , \qquad (5)$$

$$\hat{y}_0 = \frac{1}{\sqrt{\gamma_t}} \left( y_t - \sqrt{1 - \gamma_t}\, f_\theta(x, y_t, \gamma_t) \right) . \qquad (10)$$

$$q(y_{t-1} \mid y_0, y_t) = \mathcal{N}(y_{t-1} \mid \mu, \sigma^2 I)$$

$$\mu = \frac{\sqrt{\gamma_{t-1}}\,(1 - \alpha_t)}{1 - \gamma_t}\, y_0 + \frac{\sqrt{\alpha_t}\,(1 - \gamma_{t-1})}{1 - \gamma_t}\, y_t \qquad (4)$$

$$\sigma^2 = \frac{(1 - \gamma_{t-1})(1 - \alpha_t)}{1 - \gamma_t}\, .$$

$$\mu_\theta(x, y_t, \gamma_t) = \frac{1}{\sqrt{\alpha_t}} \left( y_t - \frac{1 - \alpha_t}{\sqrt{1 - \gamma_t}}\, f_\theta(x, y_t, \gamma_t) \right) , \qquad (11)$$

and we set the variance of $p_\theta(y_{t-1} \mid y_t, x)$ to $(1 - \alpha_t)$, a default given by the variance of the forward process [17].

$$y_{t-1} \leftarrow \frac{1}{\sqrt{\alpha_t}} \left( y_t - \frac{1 - \alpha_t}{\sqrt{1 - \gamma_t}}\, f_\theta(x, y_t, \gamma_t) \right) + \sqrt{1 - \alpha_t}\, \epsilon_t$$

# Conditional Denoising Diffusion Model

## 2.3. Inference via Iterative Refinement

$$\tilde{y} = \sqrt{\gamma}\, y_0 + \sqrt{1-\gamma}\, \epsilon\,, \qquad \epsilon \sim \mathcal{N}(0, I)\,, \qquad (5)$$

$$\hat{y}_0 = \frac{1}{\sqrt{\gamma_t}}\left(y_t - \sqrt{1-\gamma_t}\, f_\theta(x, y_t, \gamma_t)\right)\,. \qquad (10)$$

$$q(y_{t-1} \mid y_0, y_t) = \mathcal{N}(y_{t-1} \mid \mu, \sigma^2 I)$$

$$\mu = \frac{\sqrt{\gamma_{t-1}}\,(1-\alpha_t)}{1-\gamma_t}\, y_0 + \frac{\sqrt{\alpha_t}\,(1-\gamma_{t-1})}{1-\gamma_t}\, y_t \qquad (4)$$

$$\sigma^2 = \frac{(1-\gamma_{t-1})(1-\alpha_t)}{1-\gamma_t}\,.$$

$$\mu_\theta(x, y_t, \gamma_t) = \frac{1}{\sqrt{\alpha_t}}\left(y_t - \frac{1-\alpha_t}{\sqrt{1-\gamma_t}}\, f_\theta(x, y_t, \gamma_t)\right)\,,$$

$$(11)$$

and we set the variance of $p_\theta(y_{t-1}|y_t, x)$ to $(1 - \alpha_t)$, a default given by the variance of the forward process [17].

$$y_{t-1} \leftarrow \frac{1}{\sqrt{\alpha_t}}\left(y_t - \frac{1-\alpha_t}{\sqrt{1-\gamma_t}}\, f_\theta(x, y_t, \gamma_t)\right) + \sqrt{1-\alpha_t}\, \epsilon_t$$
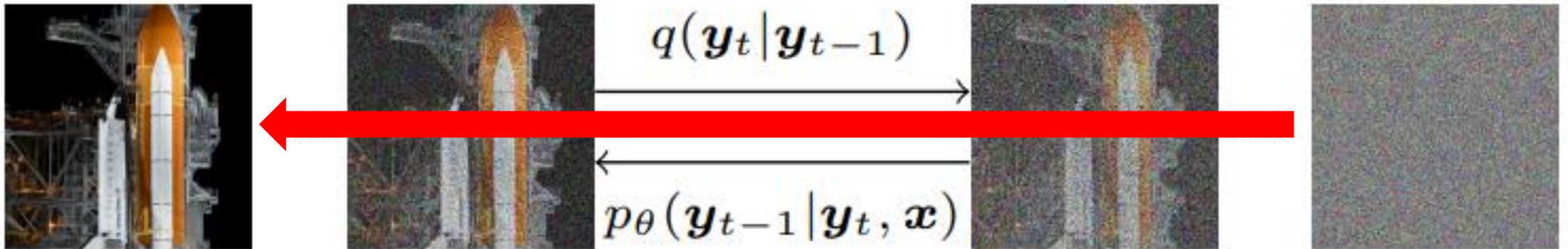
# Conditional Denoising Diffusion Model

## 2.3. Inference via Iterative Refinement

$$\tilde{y} = \sqrt{\gamma}\, y_0 + \sqrt{1-\gamma}\, \epsilon\,, \qquad \epsilon \sim \mathcal{N}(0, I)\,, \qquad (5)$$

$$\hat{y}_0 = \frac{1}{\sqrt{\gamma_t}} \left( y_t - \sqrt{1-\gamma_t}\, f_\theta(x, y_t, \gamma_t) \right)\,. \qquad (10)$$

$$y_0 \sim p(y \mid x) \qquad y_{t-1} \qquad\qquad y_t \qquad y_T \sim \mathcal{N}(0, I)$$



$$q(y_t \mid y_{t-1})$$

$$p_\theta(y_{t-1} \mid y_t, x)$$

$$y_{t-1} \leftarrow \frac{1}{\sqrt{\alpha_t}} \left( y_t - \frac{1-\alpha_t}{\sqrt{1-\gamma_t}} f_\theta(x, y_t, \gamma_t) \right) + \sqrt{1-\alpha_t}\, \epsilon_t$$

# Experiments

## 4.1. Qualitative Results

Face Images



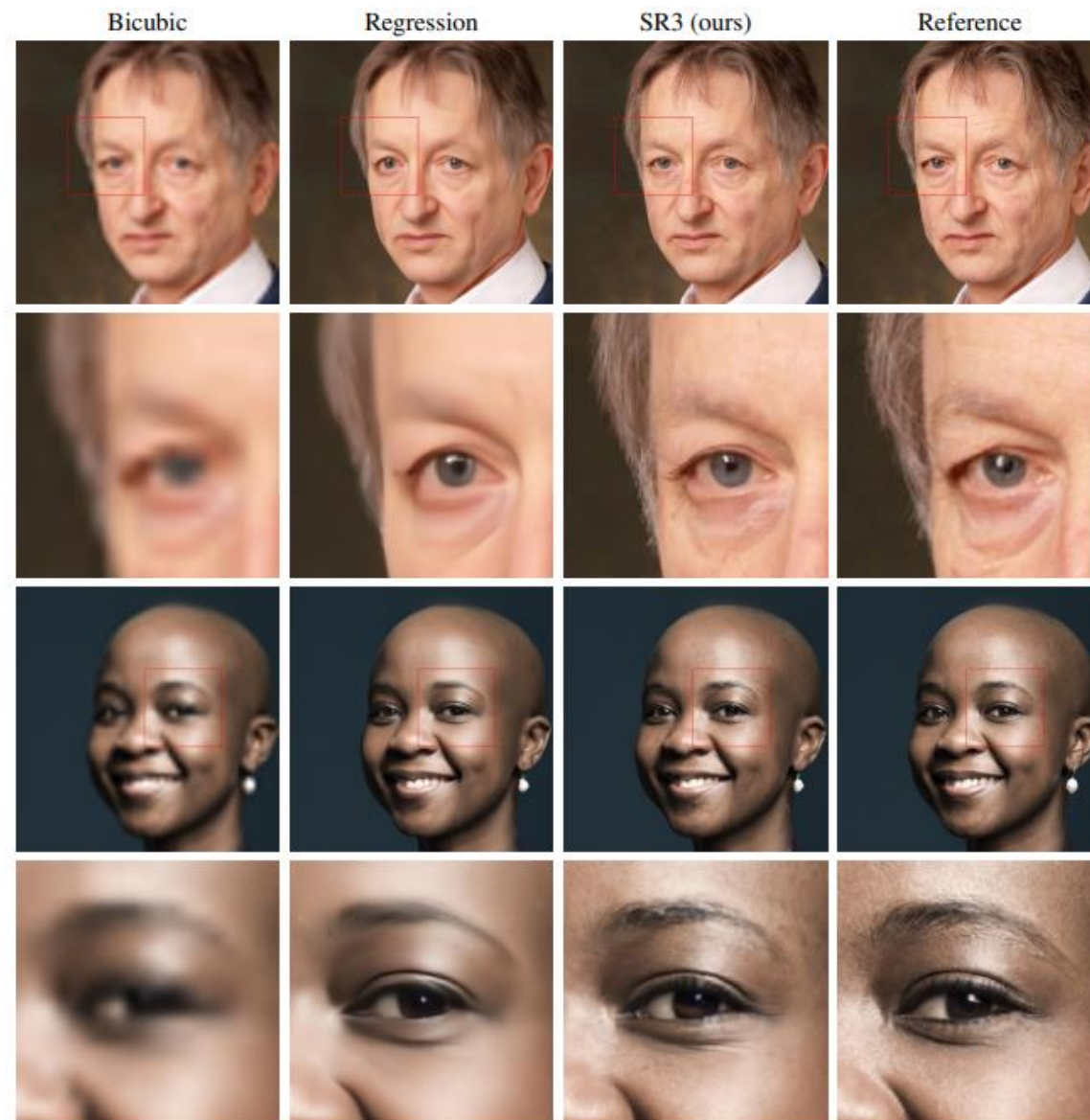**Figure 4:** Results of a SR3 model (64×64 → 512×512), trained on FFHQ, and applied to images outside of the training set, along with enlarged patches to show finer details. Additional results are shown in Appendix C.1 and C.2.

# Experiments

## 4.1. Qualitative Results

Natural Images



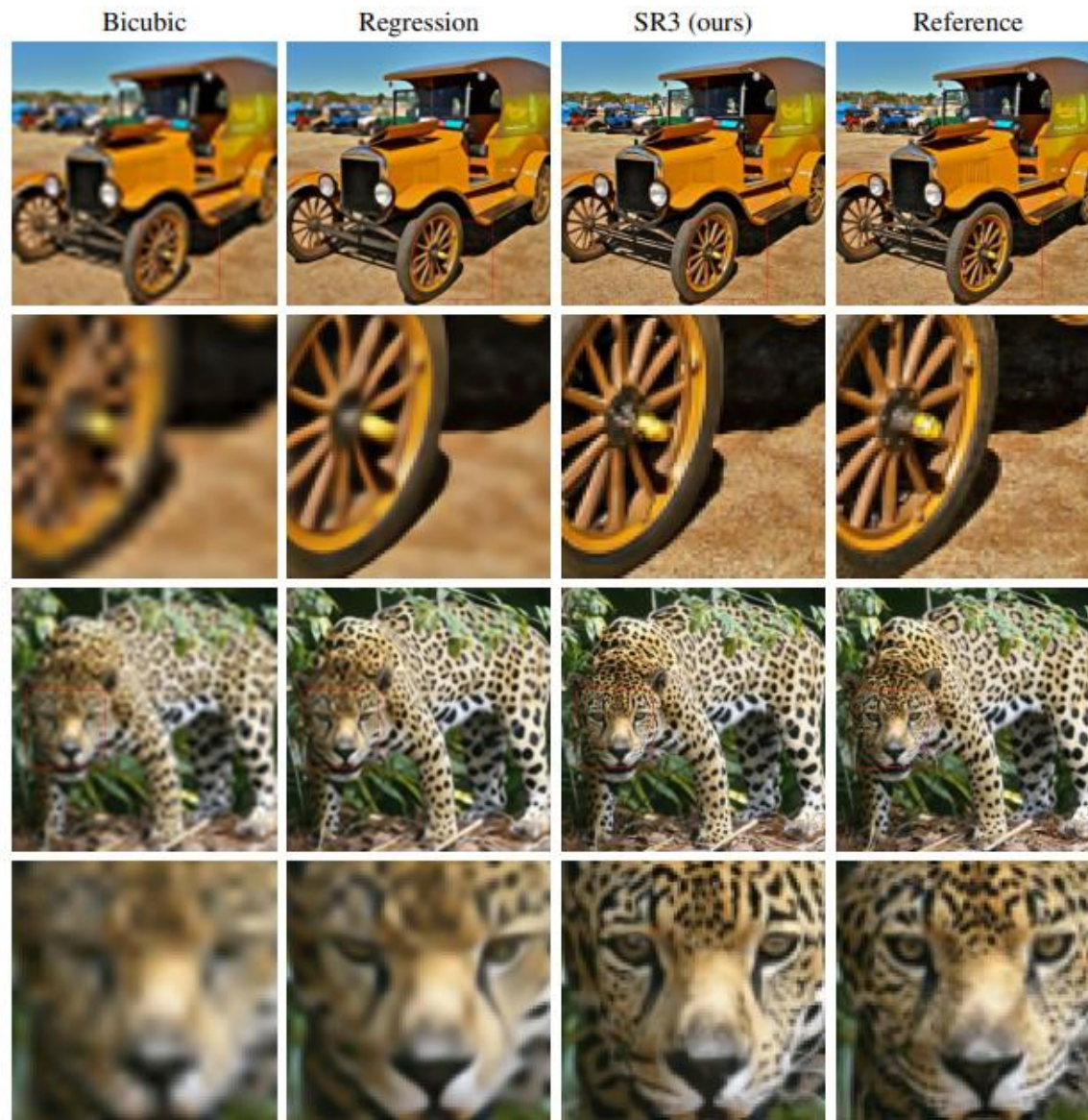**Figure 3:** Results of a SR3 model (64×64 → 256×256), trained on ImageNet and evaluated on two ImageNet test images. For each we also show an enlarged patch in which finer details are more apparent. Additional samples are shown in Appendix C.3 and C.4.

# Experiments

## 4.2. Benchmark Comparison

# Experiments

## 4.2. Benchmark Comparison

### 4.2.1 Automated metrics

| Metric | PULSE [28] | FSRGAN [7] | Regression | SR3 |
|---|---|---|---|---|
| **PSNR** ↑ | 16.88 | 23.01 | **23.96** | 23.04 |
| **SSIM** ↑ | 0.44 | 0.62 | **0.69** | 0.65 |
| **Consistency** ↓ | 161.1 | 33.8 | 2.71 | **2.68** |

| Model | FID ↓ | IS ↑ | PSNR ↑ | SSIM ↑ |
|---|---|---|---|---|
| Reference | 1.9 | 240.8 | - | - |
| Regression | 15.2 | 121.1 | **27.9** | **0.801** |
| SR3 | **5.2** | **180.1** | 26.4 | 0.762 |

**Table 1:** PSNR & SSIM on $16 \times 16 \rightarrow 128 \times 128$ face super-resolution. Consistency measures MSE $(\times 10^{-5})$ between the low-resolution inputs and the down-sampled super-resolution outputs.

**Table 2:** Performance comparison between SR3 and Regression baseline on natural image super-resolution using standard metrics computed on the ImageNet validation set.
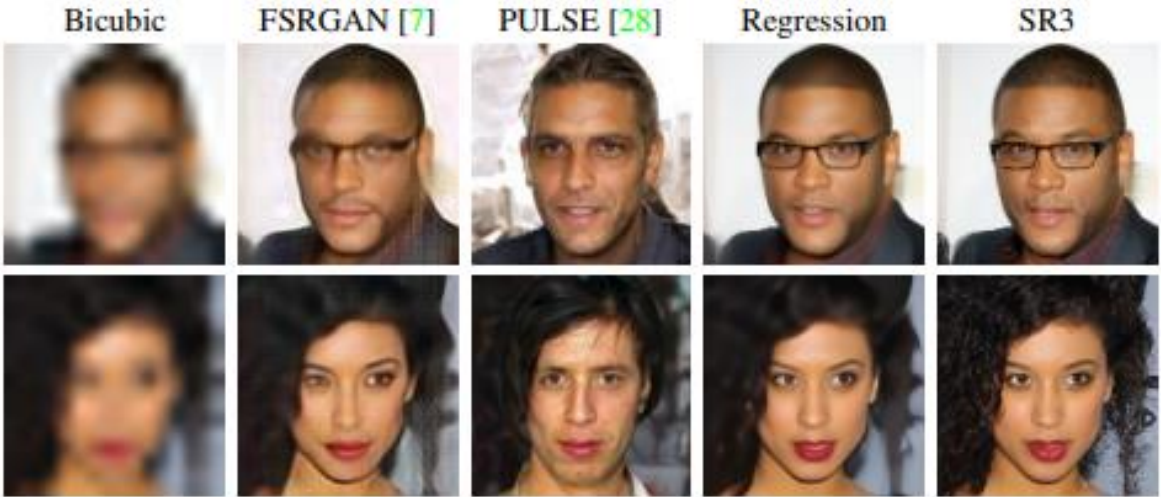


**Figure 5:** Comparison of different methods on the $16 \times 16 \rightarrow 128 \times 128$ face super-resolution task. Reference image has not been included because of privacy concerns. Additional results in Appendix C.5.

# Experiments

## 4.2. Benchmark Comparison

### 4.2.1  Automated metrics

| Metric | PULSE [28] | FSRGAN [7] | Regression | SR3 |
|---|---|---|---|---|
| PSNR ↑ | 16.88 | 23.01 | **23.96** | 23.04 |
| SSIM ↑ | 0.44 | 0.62 | **0.69** | 0.65 |
| Consistency ↓ | 161.1 | 33.8 | 2.71 | **2.68** |

Table 1: PSNR & SSIM on $16 \times 16 \rightarrow 128 \times 128$ face super-resolution. Consistency measures MSE $(\times 10^{-5})$ between the low-resolution inputs and the down-sampled super-resolution outputs.

| Model | FID ↓ | IS ↑ | PSNR ↑ | SSIM ↑ |
|---|---|---|---|---|
| Reference | 1.9 | 240.8 | - | - |
| Regression | 15.2 | 121.1 | **27.9** | **0.801** |
| SR3 | **5.2** | **180.1** | 26.4 | 0.762 |

Table 2: Performance comparison between SR3 and Regression baseline on natural image super-resolution using standard metrics computed on the ImageNet validation set.
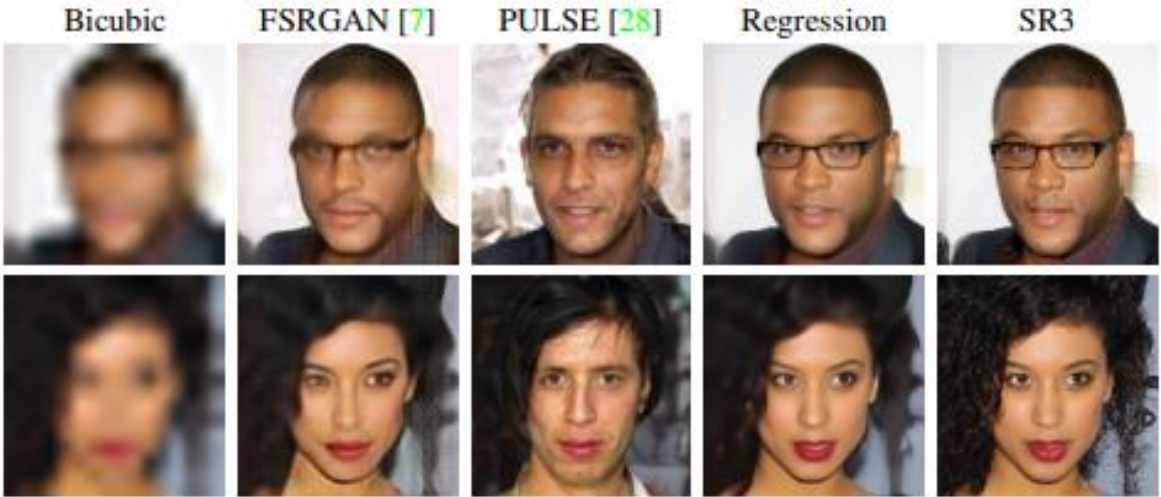


Figure 5: Comparison of different methods on the $16 \times 16 \rightarrow 128 \times 128$ face super-resolution task. Reference image has not been included because of privacy concerns. Additional results in Appendix C.5.

# Experiments

## 4.2. Benchmark Comparison

### 4.2.1 Automated metrics

| Metric | PULSE [28] | FSRGAN [7] | Regression | SR3 |
|---|---|---|---|---|
| PSNR ↑ | 16.88 | 23.01 | **23.96** | 23.04 |
| SSIM ↑ | 0.44 | 0.62 | **0.69** | 0.65 |
| Consistency ↓ | 161.1 | 33.8 | 2.71 | **2.68** |

**Table 1:** PSNR & SSIM on $16 \times 16 \rightarrow 128 \times 128$ face super-resolution. Consistency measures MSE $(\times 10^{-5})$ between the low-resolution inputs and the down-sampled super-resolution outputs.

| Model | FID ↓ | IS ↑ | PSNR ↑ | SSIM ↑ |
|---|---|---|---|---|
| Reference | 1.9 | 240.8 | - | - |
| Regression | 15.2 | 121.1 | **27.9** | **0.801** |
| SR3 | **5.2** | **180.1** | 26.4 | 0.762 |

**Table 2:** Performance comparison between SR3 and Regression baseline on natural image super-resolution using standard metrics computed on the ImageNet validation set.
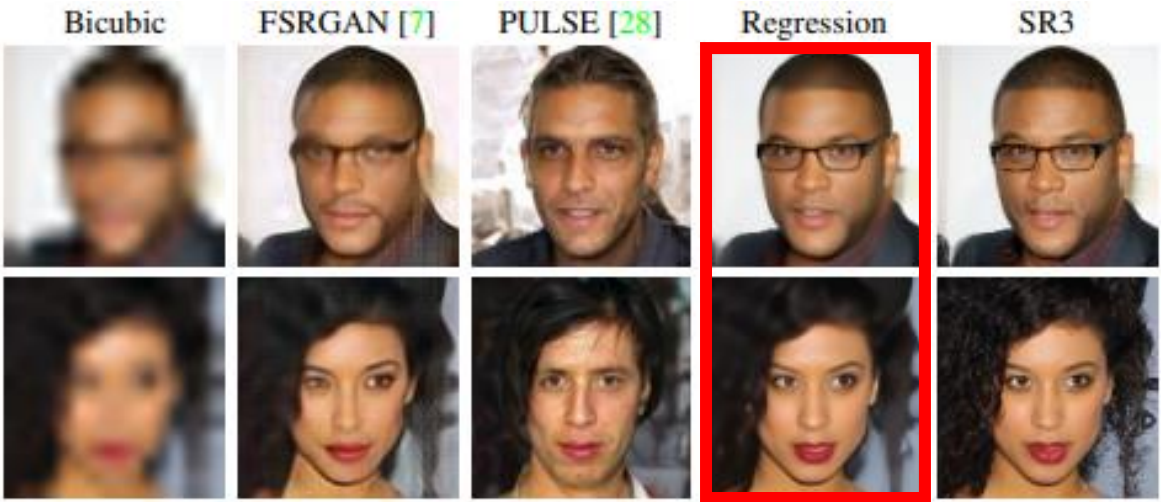


**Figure 5:** Comparison of different methods on the $16 \times 16 \rightarrow 128 \times 128$ face super-resolution task. Reference image has not been included because of privacy concerns. Additional results in Appendix C.5.

# Experiments

## 4.2. Benchmark Comparison

### 4.2.1 Automated metrics

As a measure of the consistency, we compute MSE between the downsampled outputs and the low resolution inputs

| Metric | PULSE [28] | FSRGAN [7] | Regression | SR3 |
|---|---|---|---|---|
| PSNR ↑ | 16.88 | 23.01 | **23.96** | 23.04 |
| SSIM ↑ | 0.44 | 0.62 | **0.69** | 0.65 |
| Consistency ↓ | 161.1 | 33.8 | 2.71 | **2.68** |

**Table 1:** PSNR & SSIM on $16 \times 16 \to 128 \times 128$ face super-resolution. Consistency measures MSE $(\times 10^{-5})$ between the low-resolution inputs and the down-sampled super-resolution outputs.

| Model | FID ↓ | IS ↑ | PSNR ↑ | SSIM ↑ |
|---|---|---|---|---|
| Reference | 1.9 | 240.8 | - | - |
| Regression | 15.2 | 121.1 | **27.9** | **0.801** |
| SR3 | **5.2** | **180.1** | 26.4 | 0.762 |

**Table 2:** Performance comparison between SR3 and Regression baseline on natural image super-resolution using standard metrics computed on the ImageNet validation set.
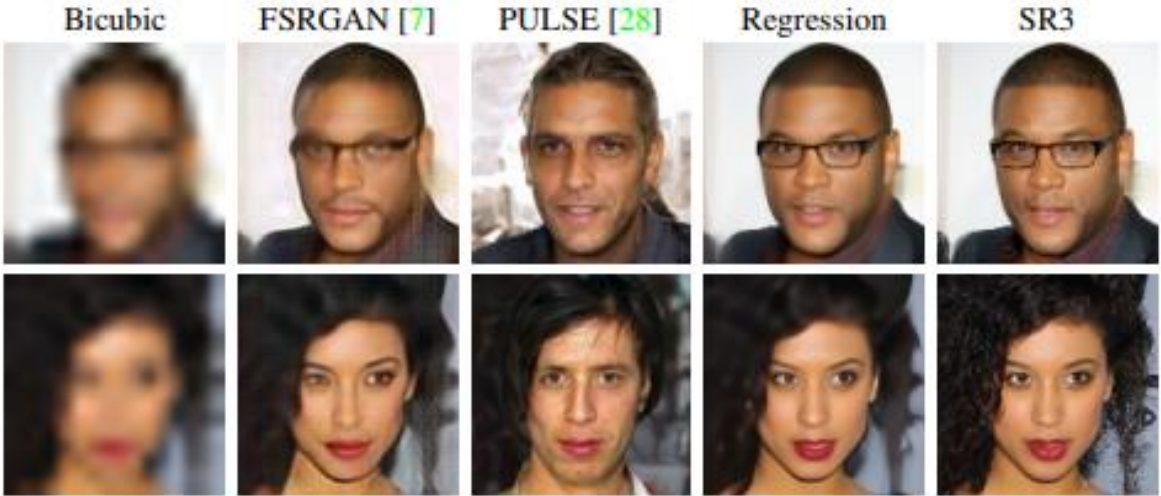
| Bicubic | FSRGAN [7] | PULSE [28] | Regression | SR3 |
|---|---|---|---|---|



**Figure 5:** Comparison of different methods on the $16 \times 16 \to 128 \times 128$ face super-resolution task. Reference image has not been included because of privacy concerns. Additional results in Appendix C.5.

# Experiments

## 4.2. Benchmark Comparison
### 4.2.1 Automated metrics

224x224 ImageNet dataset image

| Method | Top-1 Error | Top-5 Error |
|---|---|---|
| Baseline | 0.252 | 0.080 |
| DRCN [22] | 0.477 | 0.242 |
| FSRCNN [13] | 0.437 | 0.196 |
| PsyCo [35] | 0.454 | 0.224 |
| ENet-E [44] | 0.449 | 0.214 |
| RCAN [64] | 0.393 | 0.167 |
| Regression | 0.383 | 0.173 |
| SR3 | **0.317** | **0.120** |

**Table 3:** Comparison of classification accuracy scores for 4× natural image super-resolution on the first 1K images from the ImageNet Validation set.
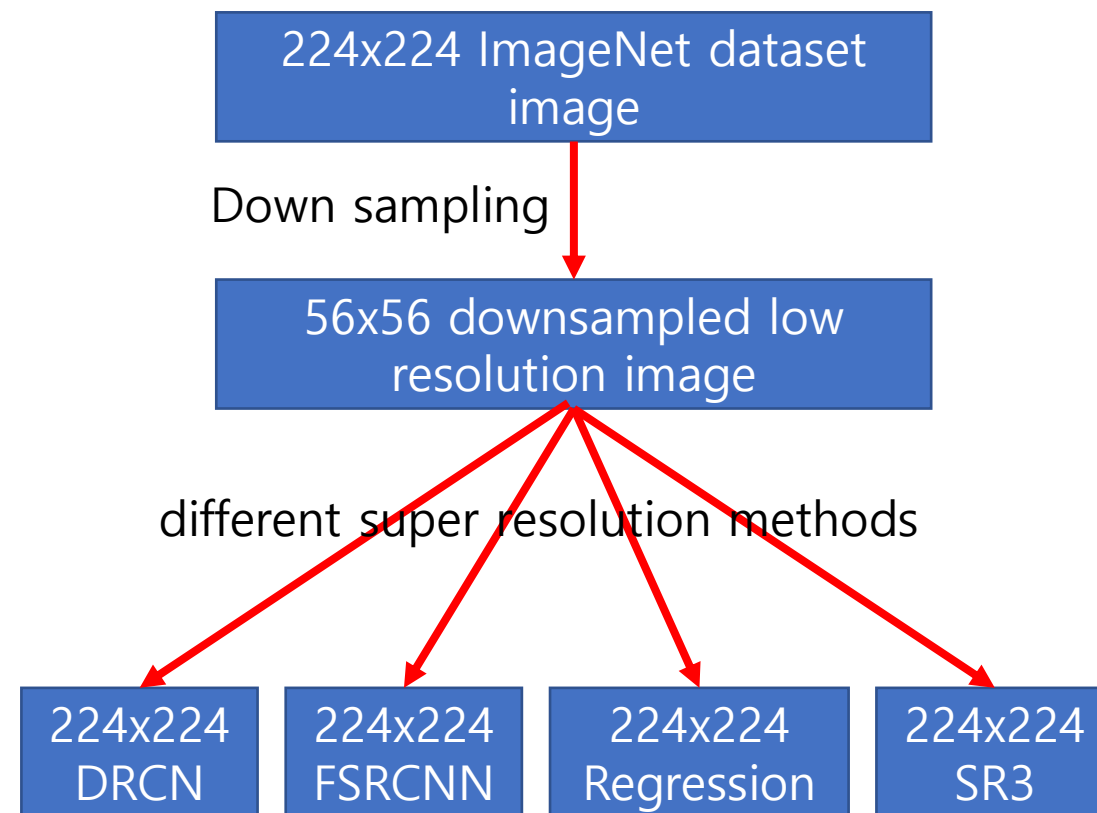
# Experiments

## 4.2. Benchmark Comparison

### 4.2.1 Automated metrics

| Method | Top-1 Error | Top-5 Error |
|--------|-------------|-------------|
| Baseline | 0.252 | 0.080 |
| DRCN [22] | 0.477 | 0.242 |
| FSRCNN [13] | 0.437 | 0.196 |
| PsyCo [35] | 0.454 | 0.224 |
| ENet-E [44] | 0.449 | 0.214 |
| RCAN [64] | 0.393 | 0.167 |
| Regression | 0.383 | 0.173 |
| SR3 | 0.317 | 0.120 |

**Table 3:** Comparison of classification accuracy scores for 4× natural image super-resolution on the first 1K images from the ImageNet Validation set.
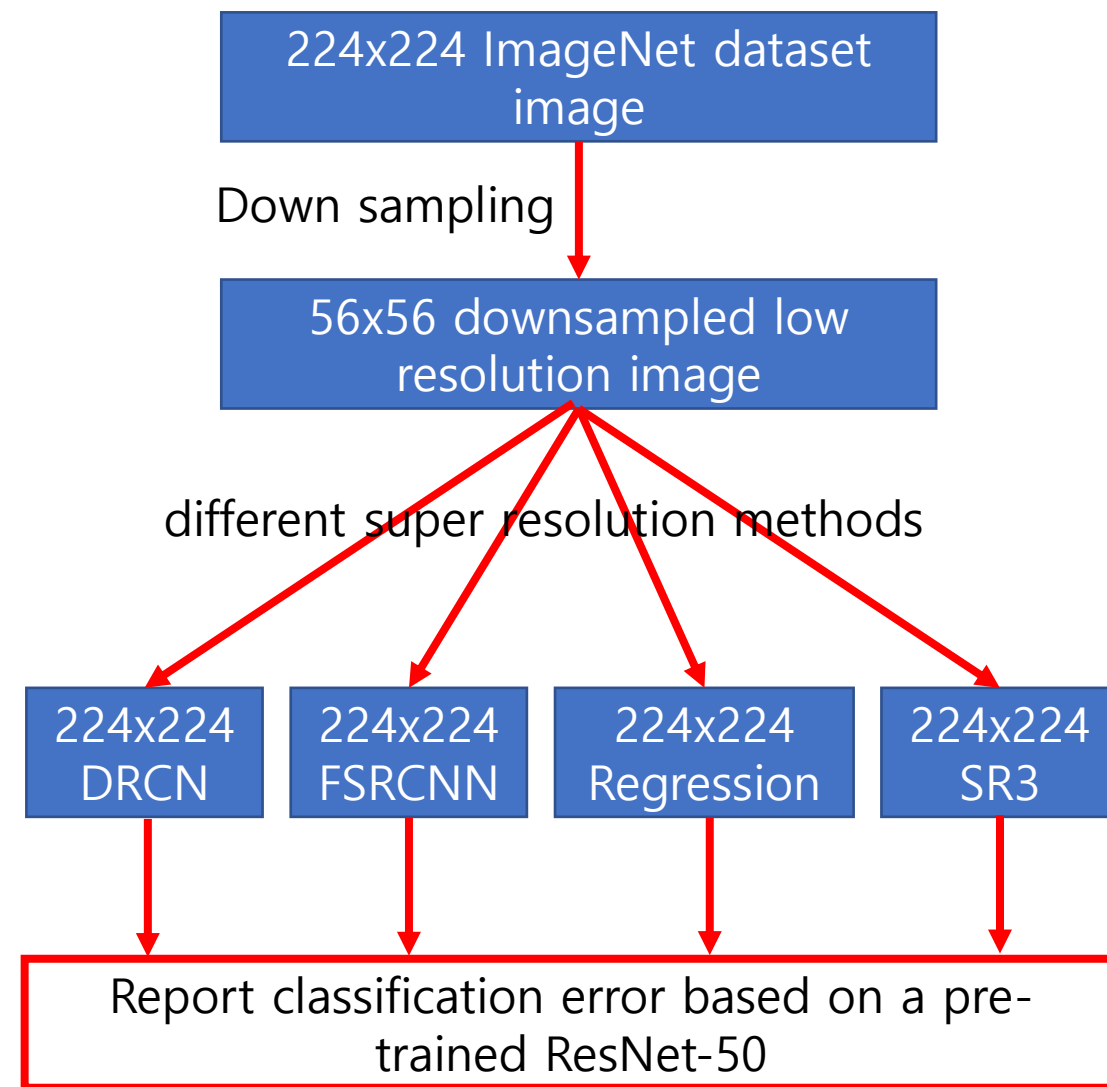
224x224 ImageNet dataset image

Down sampling

56x56 downsampled low resolution image

# Experiments

## 4.2. Benchmark Comparison

### 4.2.1 Automated metrics

| Method | Top-1 Error | Top-5 Error |
|--------|-------------|-------------|
| Baseline | 0.252 | 0.080 |
| DRCN [22] | 0.477 | 0.242 |
| FSRCNN [13] | 0.437 | 0.196 |
| PsyCo [35] | 0.454 | 0.224 |
| ENet-E [44] | 0.449 | 0.214 |
| RCAN [64] | 0.393 | 0.167 |
| Regression | 0.383 | 0.173 |
| SR3 | **0.317** | **0.120** |

**Table 3:** Comparison of classification accuracy scores for 4× natural image super-resolution on the first 1K images from the ImageNet Validation set.

# Experiments

## 4.2. Benchmark Comparison

### 4.2.1   Automated metrics

| Method | Top-1 Error | Top-5 Error |
|---|---|---|
| Baseline | 0.252 | 0.080 |
| DRCN [22] | 0.477 | 0.242 |
| FSRCNN [13] | 0.437 | 0.196 |
| PsyCo [35] | 0.454 | 0.224 |
| ENet-E [44] | 0.449 | 0.214 |
| RCAN [64] | 0.393 | 0.167 |
| Regression | 0.383 | 0.173 |
| SR3 | **0.317** | **0.120** |

**Table 3:** Comparison of classification accuracy scores for 4× natural image super-resolution on the first 1K images from the ImageNet Validation set.

# Experiments

## 4.2. Benchmark Comparison
### 4.2.2 Human Evaluation (2AFC)

Human subjects are shown a low-resolution input and are required to select between a model output and a ground truth image.
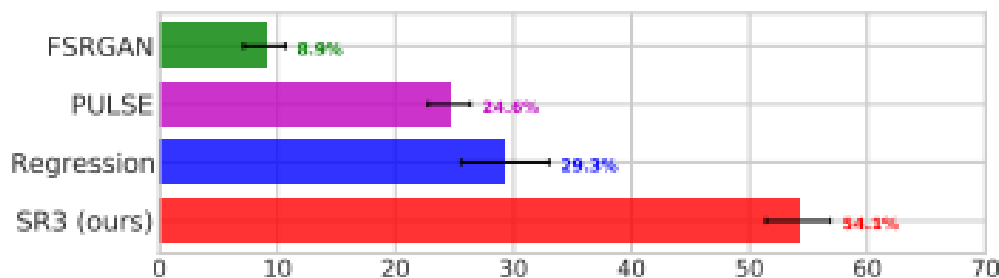
Low-resolution input

# Experiments

## 4.2. Benchmark Comparison
### 4.2.2 Human Evaluation (2AFC)

Human subjects are shown a low-resolution input and are required to select between a model output and a ground truth image.



Select between Model output and ground truth

# Experiments

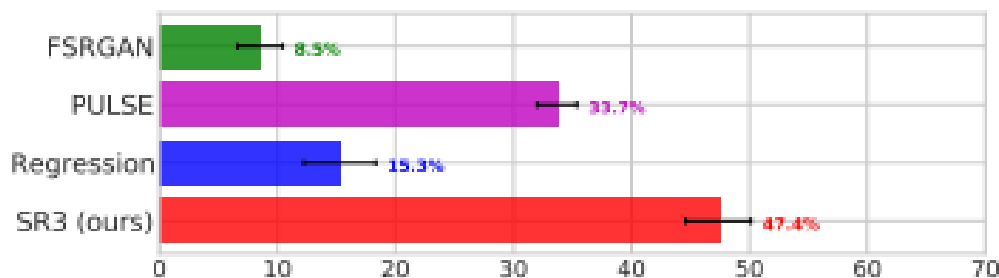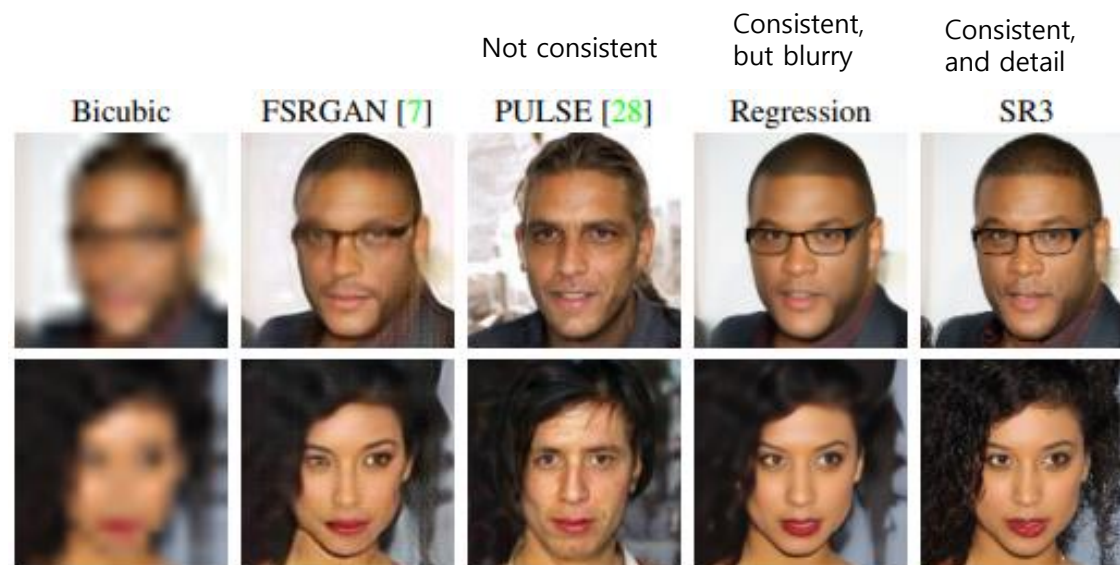## 4.2. Benchmark Comparison
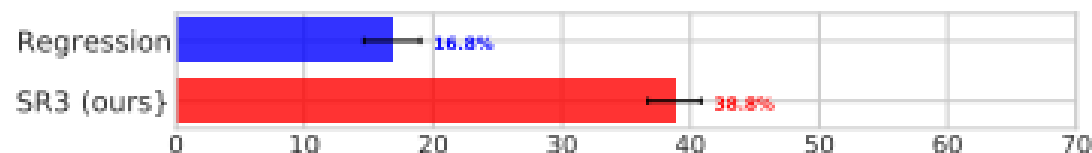
## 4.2.2 Human Evaluation (2AFC)



**Figure 6:** Face super-resolution human fool rates (higher is better, photo-realistic samples yield a fool rate of 50%). Outputs of 4 models are compared against ground truth. (top) Subjects are shown low-resolution inputs. (bottom) Inputs are not shown.
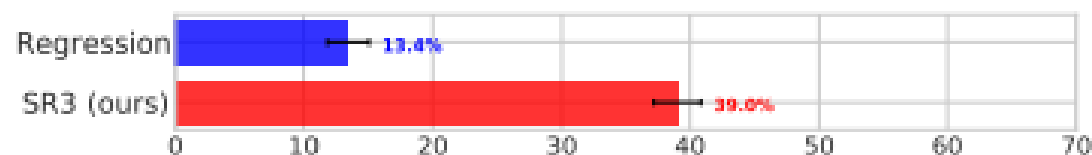


**Figure 7:** ImageNet super-resolution fool rates (higher is better, photo-realistic samples yield a fool rate of 50%). SR3 and Regression outputs are compared against ground truth. (top) Subjects are shown low-resolution inputs. (bottom) Inputs are not shown.
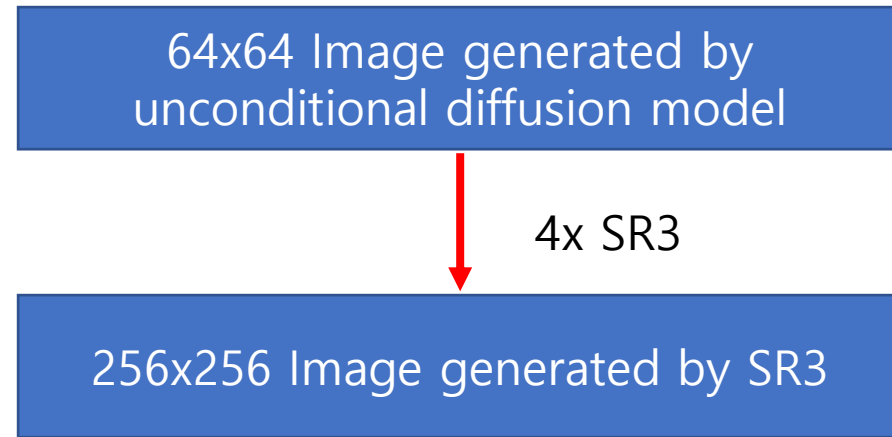
# Experiments

## 4.3. Cascaded High-Resolution Image Synthesis

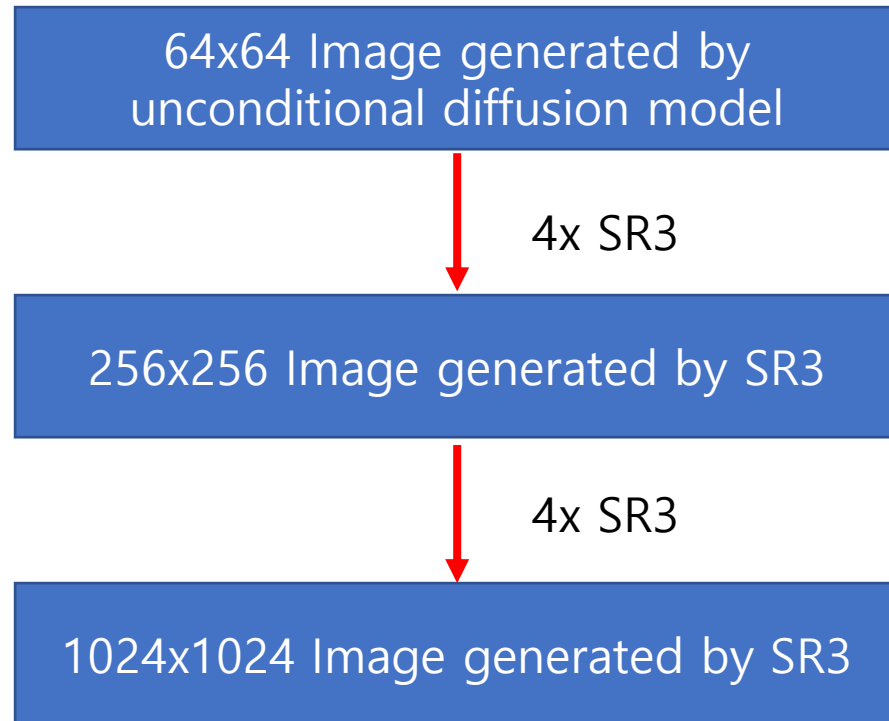64x64 Image generated by unconditional diffusion model

# Experiments

## 4.3. Cascaded High-Resolution Image Synthesis

| 64x64 Image generated by unconditional diffusion model |
|---|

↓ 4x SR3

| 256x256 Image generated by SR3 |
|---|

# Experiments

## 4.3. Cascaded High-Resolution Image Synthesis

| 64x64 Image generated by unconditional diffusion model |
|---|

↓ 4x SR3

| 256x256 Image generated by SR3 |
|---|

↓ 4x SR3

| 1024x1024 Image generated by SR3 |
|---|

# Experiments

## 4.3. Cascaded High-Resolution Image Synthesis

| 64x64 Image generated by unconditional diffusion model |
|---|

↓ 4x SR3

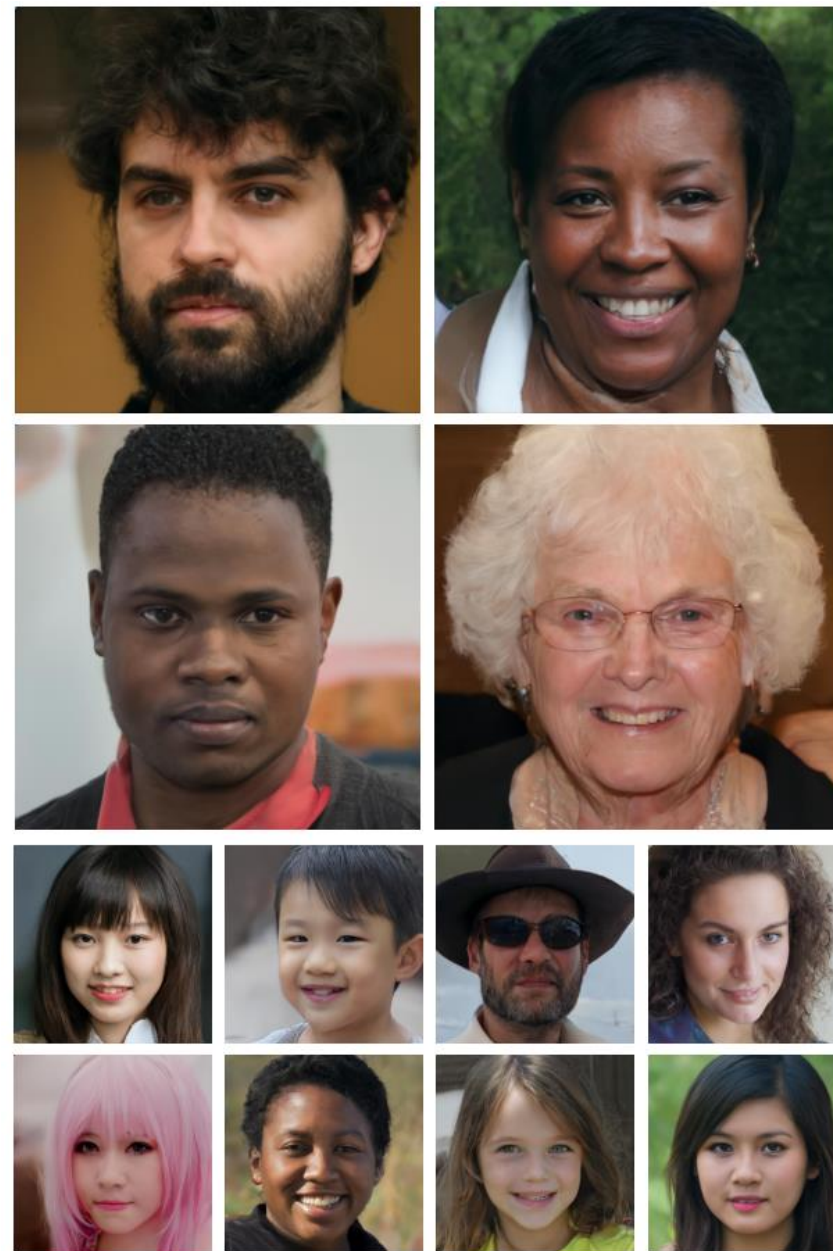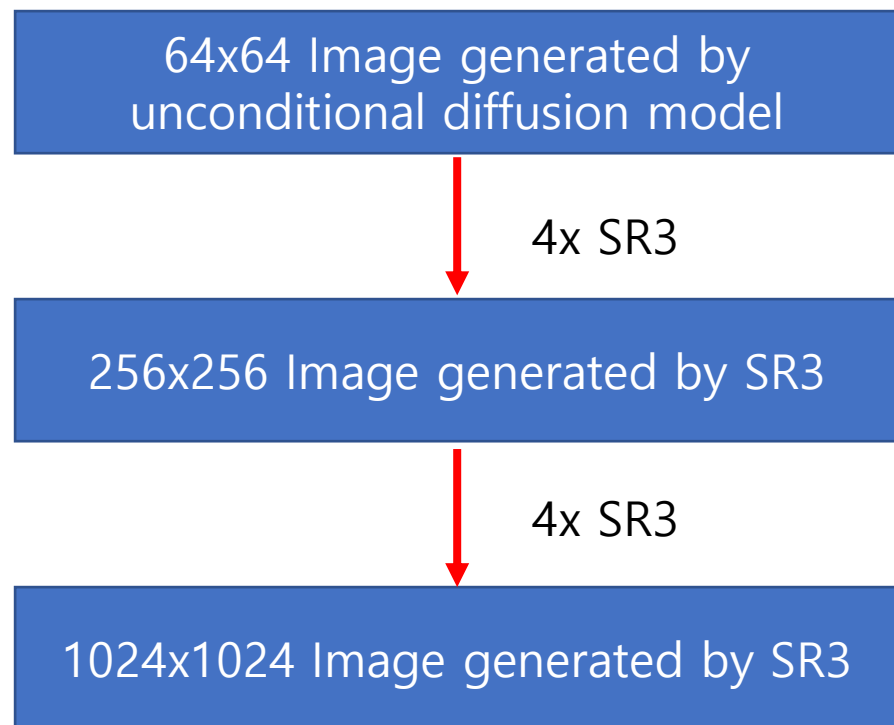| 256x256 Image generated by SR3 |
|---|

↓ 4x SR3

| 1024x1024 Image generated by SR3 |
|---|



**Figure 8:** Synthetic 1024×1024 face images. We first sample from an unconditional 64×64 diffusion model, then pass the samples through two 4× SR3 models, *i.e.*, 64×64 → 256×256 → 1024×1024. Additional samples in Appendix C.7, C.8 and C.9.

# Discussion and Conclusion

Bias is an important problem in all generative models.

SR3 also suffers from bias issues.

We believe it is likely our diffusion-based models drop modes.

We also observed the model to generate very continuous skin texture in face super-resolution, dropping moles, pimples, and pierceings found in the reference.

SR3 should not be used for any real world super-resolution tasks, until these biases are thoroughly understood and mitigated.

# Discussion and Conclusion

In conclusion,

SR3 is an approach to image super-resolution via diffusion model.

SR3 can be used in a cascaded high-resolution image synthesis.

SR3 achieves a human fool rate close to 50%, suggesting photo-realistic outputs.

# Q&A

Thank you