# Denoising Diffusion Probabilistic Models

Горохов Антон, БПМИ212

# Table of contents

1. Intro

2. Diffusion Process

3. Background

4. Model

5. Experiments and results

6. Image interpolation

7. Conclusion

# 1. Intro – Image generation

"Dataset of faces"

# 1. Intro – Known solutions



Generative adversarial networks (GANs)

autoregressive models

flows

variational autoencoders (VAEs)
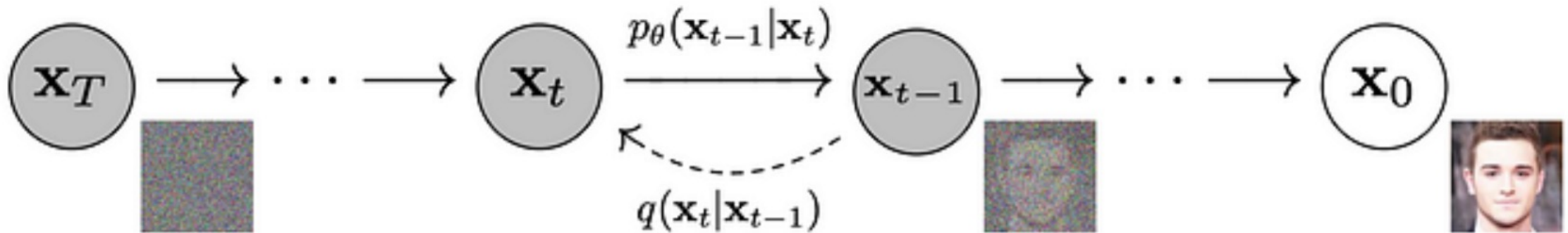
# 2. Diffusion Process as a Markov chain



Image from paper Denoising Diffusion Probabilistic Models, page 2
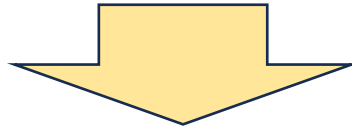
# 2. Diffusion Process forward step



$$q(x_t|x_{t-1}) = \mathcal{N}(x_t, \sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$$

with labels: normal distribution, mean, output, variance

# 3. Background

Training is performed by optimizing the usual variational bound on negative log likelihood:

$$\mathbb{E}\left[-\log p_\theta(\mathbf{x}_0)\right] \leq \mathbb{E}_q\left[-\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\right] = \mathbb{E}_q\left[-\log p(\mathbf{x}_T) - \sum_{t \geq 1}\log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_t|\mathbf{x}_{t-1})}\right] =: L$$

**Difficulties with generating high quality samples**

# 4. Model – loss function improvement

**VLB loss**

$$\mathbb{E}[-\log p_\theta(\mathbf{x}_0)] \le \mathbb{E}_q\left[ L_T + \sum_{t>1} D_{\mathrm{KL}}\left(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \,\|\, p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)\right) + L_0 \right]$$

**DSM loss**

$$\text{constant} * \left\| \epsilon - \epsilon_\theta\left(\sqrt{\bar{\alpha}}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\epsilon, t\right) \right\|^2$$

*This is the final loss function we use to train DDPMs, which is just a "Mean Squared Error" between the noise added in the forward process and the noise predicted by the model. This is the most impactful contribution of the paper Denoising Diffusion Probabilistic Models.*

# 4. Model

**Algorithm 1** Training

1: **repeat**
2:     $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
3:     $t \sim \text{Uniform}(\{1, \ldots, T\})$
4:     $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
5:     Take gradient descent step on
$$\nabla_\theta \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}, t) \right\|^2$$
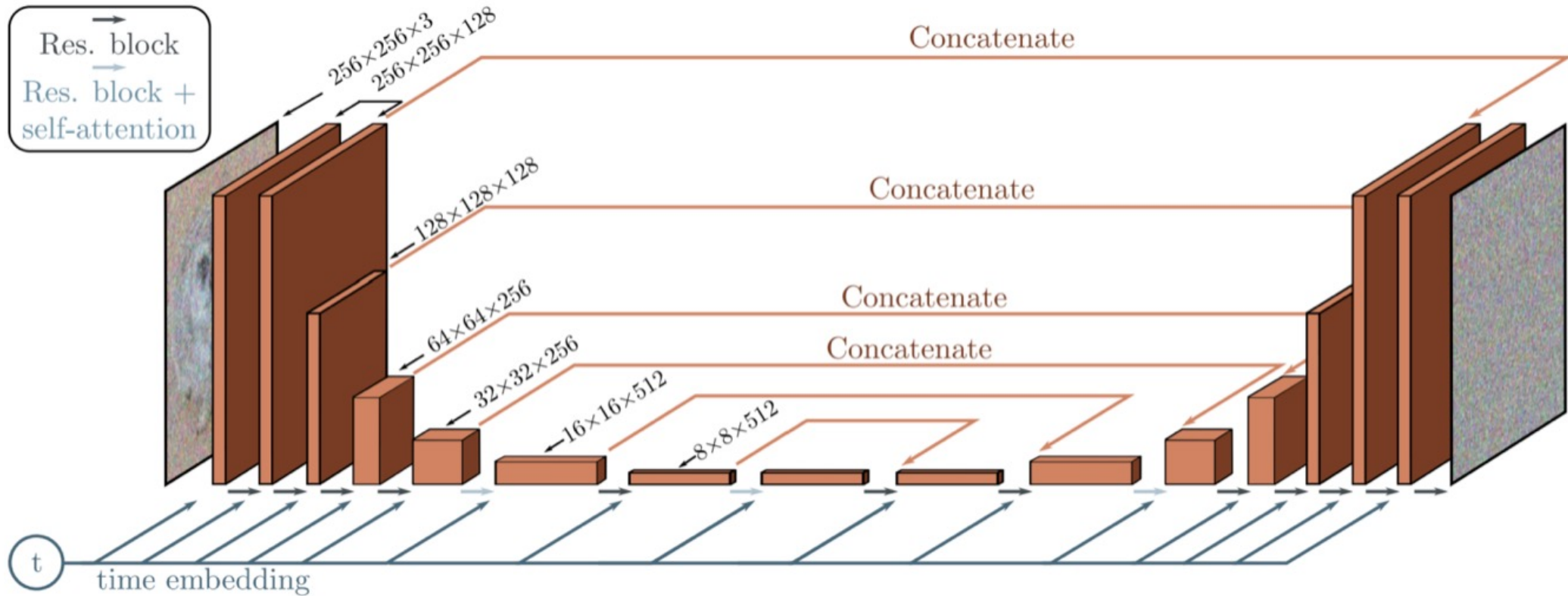6: **until** converged

# 4. Model

**Algorithm 2** Sampling

1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
2: **for** $t = T, \ldots, 1$ **do**
3:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$
4:  $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
5: **end for**
6: **return** $\mathbf{x}_0$

# 4. Model architecture



*The U-Net architecture used in DDPMs*

# 5. Experiments and results

- T = 1000
- Constant variances of Gaussian noise on each step (works better than predictable) from $\beta_1 = 10^{-4}$ to $\beta_T = 2 * 10^{-2}$

# 5. Experiments and results

Table 1: CIFAR10 results. NLL measured in bits/dim.

| Model | IS | FID | NLL Test (Train) |
|---|---|---|---|
| **Conditional** | | | |
| EBM [11] | 8.30 | 37.9 | |
| JEM [17] | 8.76 | 38.4 | |
| BigGAN [3] | 9.22 | 14.73 | |
| StyleGAN2 + ADA (v1) [29] | **10.06** | **2.67** | |
| **Unconditional** | | | |
| Diffusion (original) [53] | | | $\leq 5.40$ |
| Gated PixelCNN [59] | 4.60 | 65.93 | 3.03 (2.90) |
| Sparse Transformer [7] | | | **2.80** |
| PixelIQN [43] | 5.29 | 49.46 | |
| EBM [11] | 6.78 | 38.2 | |
| NCSNv2 [56] | | 31.75 | |
| NCSN [55] | 8.87±0.12 | 25.32 | |
| SNGAN [39] | 8.22±0.05 | 21.7 | |
| SNGAN-DDLS [4] | 9.09±0.10 | 15.42 | |
| StyleGAN2 + ADA (v1) [29] | **9.74** $\pm$ 0.05 | 3.26 | |
| Ours ($L$, fixed isotropic $\Sigma$) | 7.67±0.13 | 13.51 | $\leq 3.70$ (3.69) |
| **Ours** ($L_{\mathrm{simple}}$) | 9.46±0.11 | **3.17** | $\leq 3.75$ (3.72) |

# 5. Experiments and results

Table 2: Unconditional CIFAR10 reverse process parameterization and training objective ablation. Blank entries were unstable to train and generated poor samples with out-of-range scores.

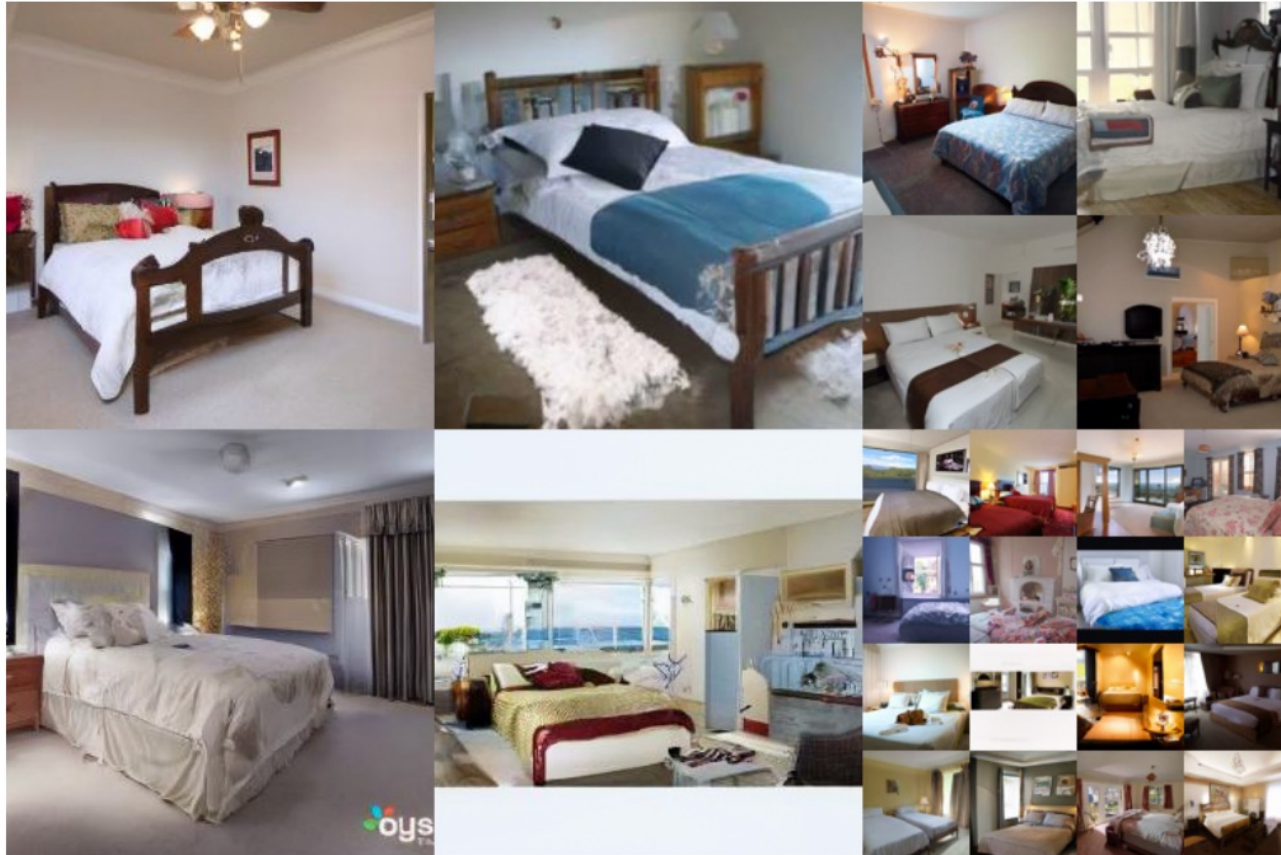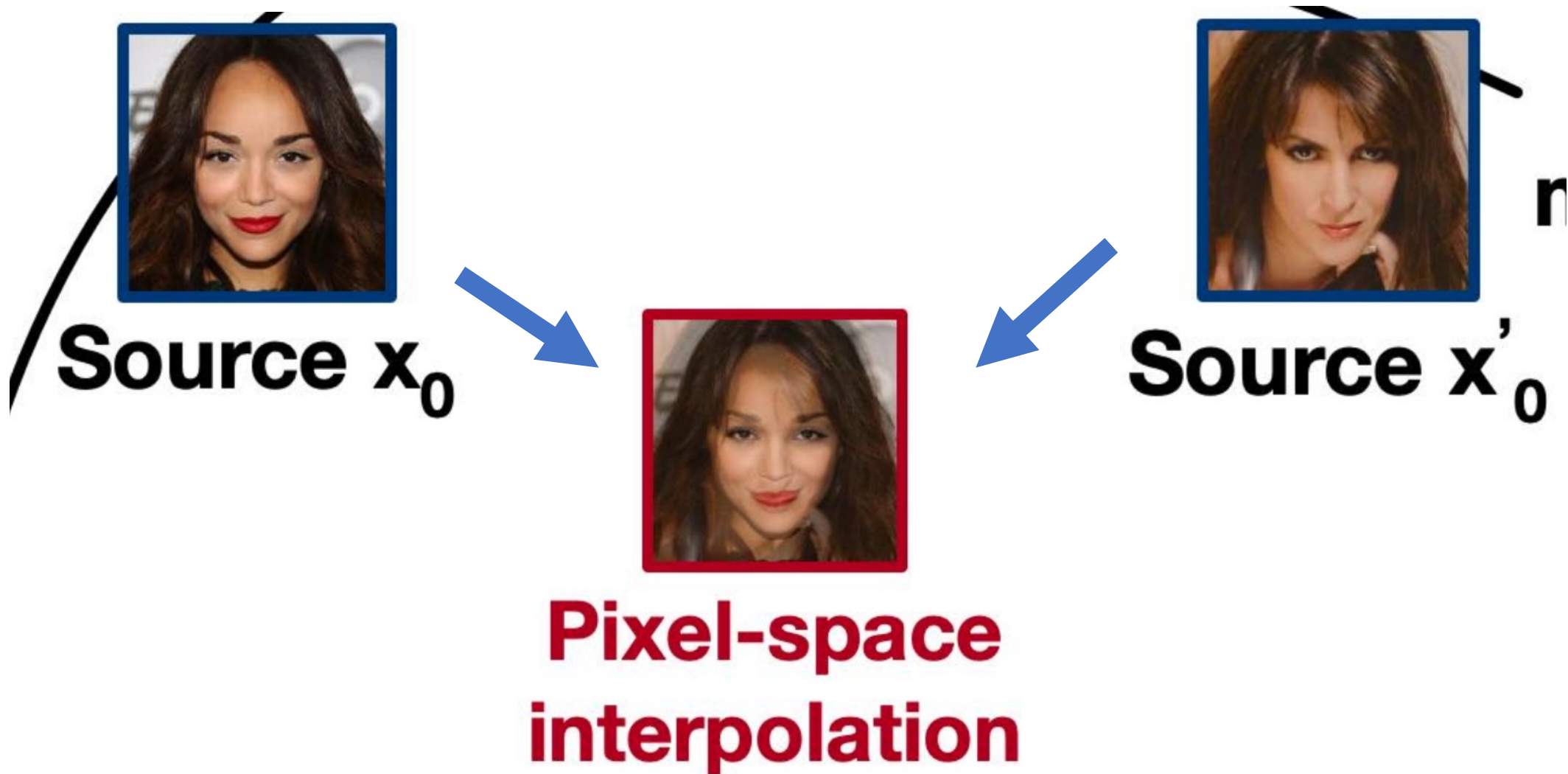| Objective | IS | FID |
|---|---|---|
| $\tilde{\boldsymbol{\mu}}$ **prediction (baseline)** | | |
| $L$, learned diagonal $\boldsymbol{\Sigma}$ | $7.28 \pm 0.10$ | 23.69 |
| $L$, fixed isotropic $\boldsymbol{\Sigma}$ | $8.06 \pm 0.09$ | 13.22 |
| $\|\tilde{\boldsymbol{\mu}} - \tilde{\boldsymbol{\mu}}_\theta\|^2$ | – | – |
| $\epsilon$ **prediction (ours)** | | |
| $L$, learned diagonal $\boldsymbol{\Sigma}$ | – | – |
| $L$, fixed isotropic $\boldsymbol{\Sigma}$ | $7.67 \pm 0.13$ | 13.51 |
| $\|\tilde{\epsilon} - \epsilon_\theta\|^2$ ($L_{\mathrm{simple}}$) | **9.46 $\pm$ 0.11** | **3.17** |

# 5. Experiments and results



Figure 4: LSUN Bedroom samples. FID=4.90

# 5. Experiments and results



Figure 3: LSUN Church samples. FID=7.89

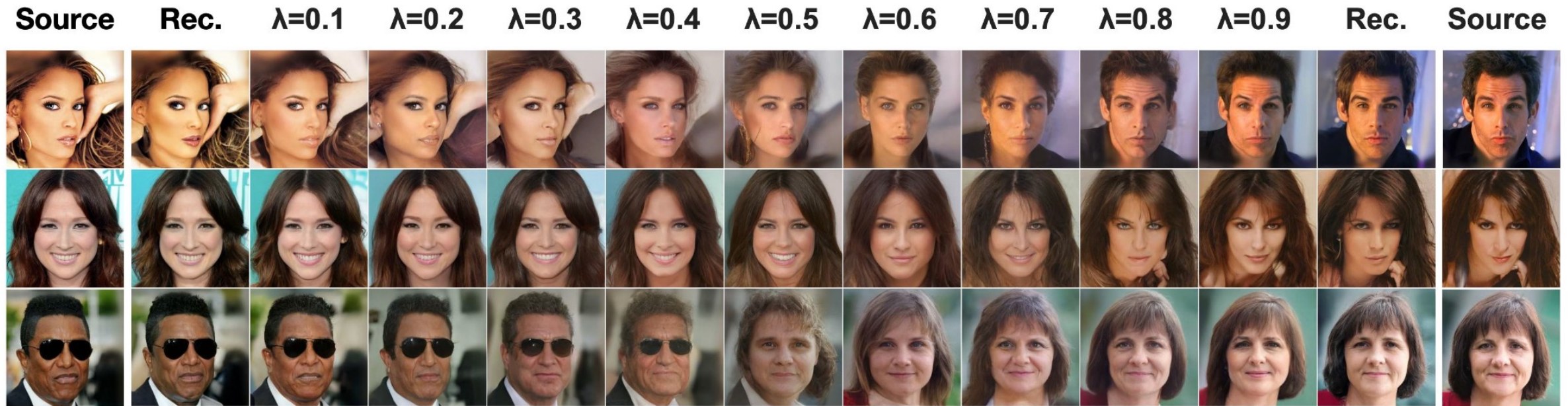# 6. Image interpolation



Source $x_0$

Pixel-space interpolation

Source $x'_0$

# 6. Image interpolation

# 6. Image interpolation

# 7. Conclusion

- Great potential shown
- Appliances in data compression
- Possible risks of malicious usage