

Accelerating **diffusion models** for **inverse problems** through **stochastic contraction**

Jong Chul Ye

In collaboration with
Hyungjin Chung

Professor
Graduate School of Artificial Intelligence
KAIST, Korea



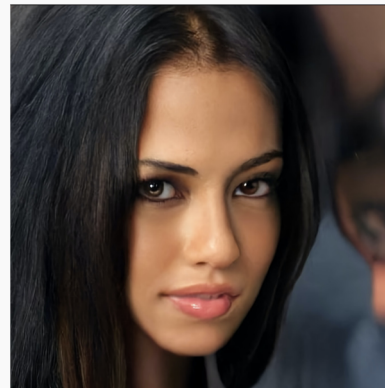
Diffusion-based Generative Models



StyleGAN2-ADA
(Karras et al., 2020)

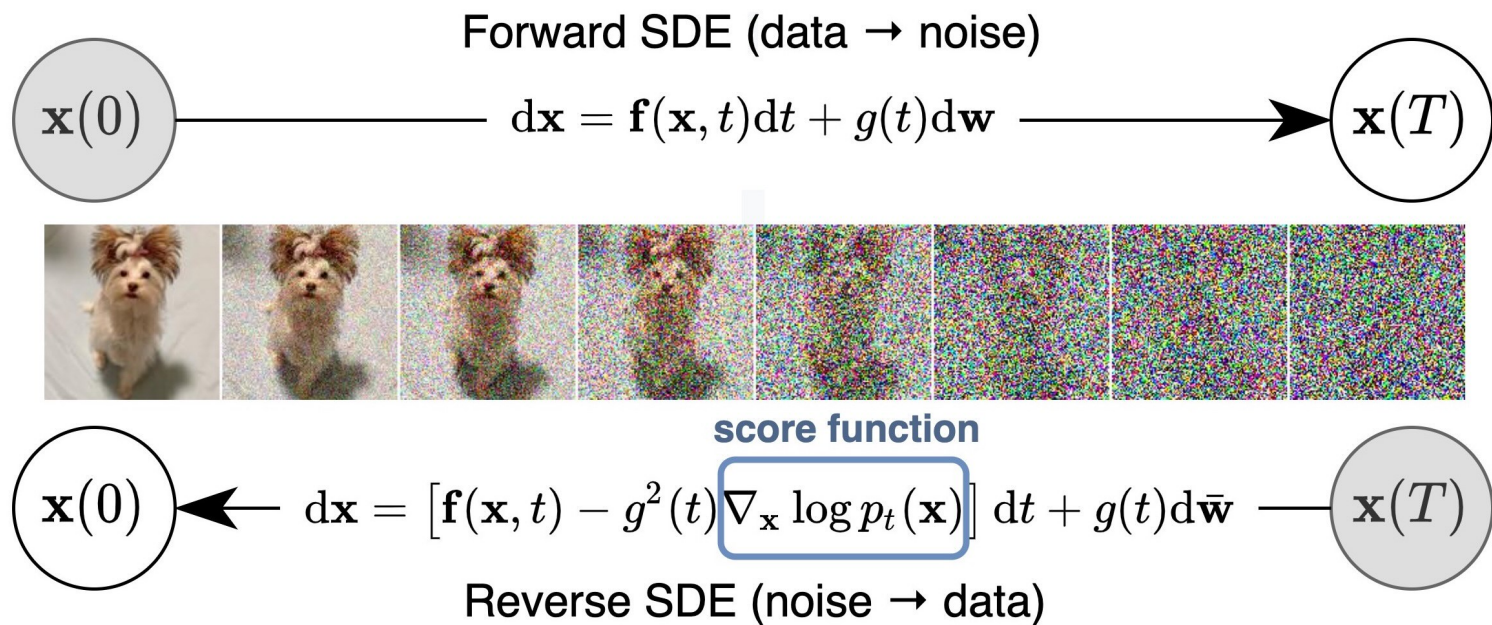


DDPM
(Ho et al., 2020)

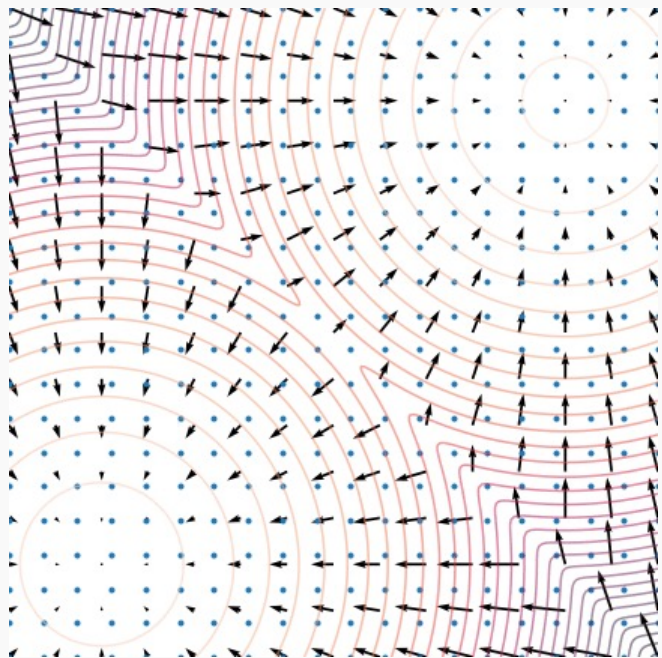


Reverse SDE
(Song et al., 2020)

Score-based Generative Models through SDE



Score-based Generative Models through SDE



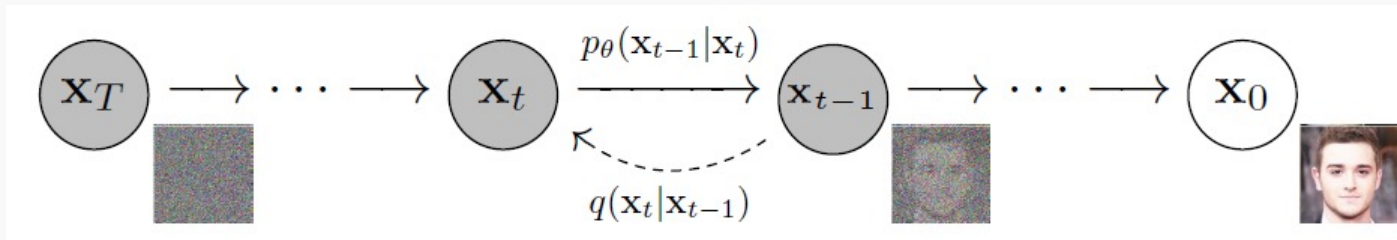
- Once the score model is trained to optimality,
 - i.e. $s_\theta(\mathbf{x}) \simeq \nabla_{\mathbf{x}} p(\mathbf{x})$
- Use **Langevin dynamics** to draw samples

$$\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + \epsilon \nabla_{\mathbf{x}} \log p(\mathbf{x}) + \sqrt{2\epsilon} \mathbf{z}_i$$

$$i = 0, 1, \dots, K$$

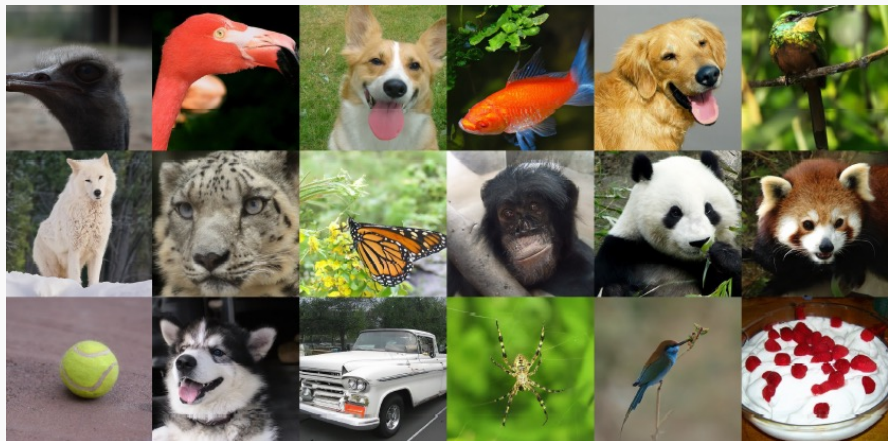
Diffusion Denoising Probabilistic Models (DDPMs)

Ho et al. *NeurIPS*, 2020



- **Train** with **variational lower bound**
- Follow the **reverse markov chain** at **inference**

$$L_{\text{simple}}(\theta) := \mathbb{E}_{t, x_0, \epsilon} \left[\left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t) \right\|^2 \right]$$



Diffusion Models Beat GANs on Image Synthesis

Dhariwal and Nichol, *NeurIPS*, 2021

Equivalence between the Two Approaches

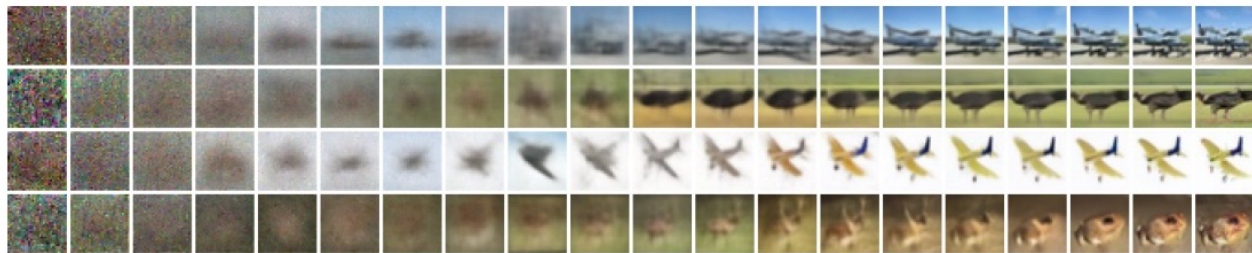
Song et al. *ICLR* 2022

DDPM

- **Training objective:**

$$L_{\text{simple}}(\theta) := \mathbb{E}_{t, \mathbf{x}_0, \epsilon} \left[\left\| \epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) \right\|^2 \right]$$

- **Inference:**



SGM

- **Training objective:**

$$\ell(\theta; \sigma) \triangleq \frac{1}{2} \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{x}, \sigma^2 I)} \left[\left\| \mathbf{s}_{\theta}(\tilde{\mathbf{x}}, \sigma) + \frac{\tilde{\mathbf{x}} - \mathbf{x}}{\sigma^2} \right\|_2^2 \right].$$

- **Inference:**



Reverse Diffusion Through Score-Matching

Noising

$$d\mathbf{x} = \bar{\mathbf{f}}(\mathbf{x}, t) dt + \bar{g}(t)d\mathbf{w}$$

Corresponding reverse SDE

Denoising

$$\begin{aligned}d\mathbf{x} &= [\bar{\mathbf{f}}(\mathbf{x}, t) - \bar{g}(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x})] dt + \bar{g}(t)d\mathbf{w} \\ &\simeq [\bar{\mathbf{f}}(\mathbf{x}, t) - \bar{g}(t)^2 s_{\theta}(\mathbf{x}, t)] dt + \bar{g}(t)d\mathbf{w}\end{aligned}$$

- Solve reverse SDE numerically: Image generation (denoising)

SCORE-BASED DIFFUSION MODELS FOR INVERSE PROBLEMS

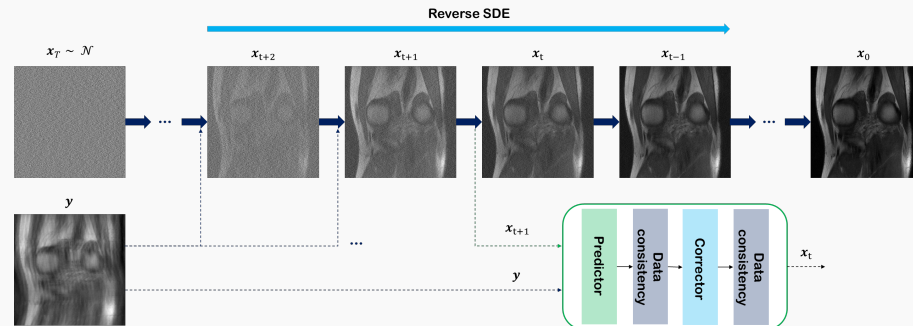
Chung et al, Medical Image Analysis (in revision), 2022

A General Score-based Formula for Inverse Problems

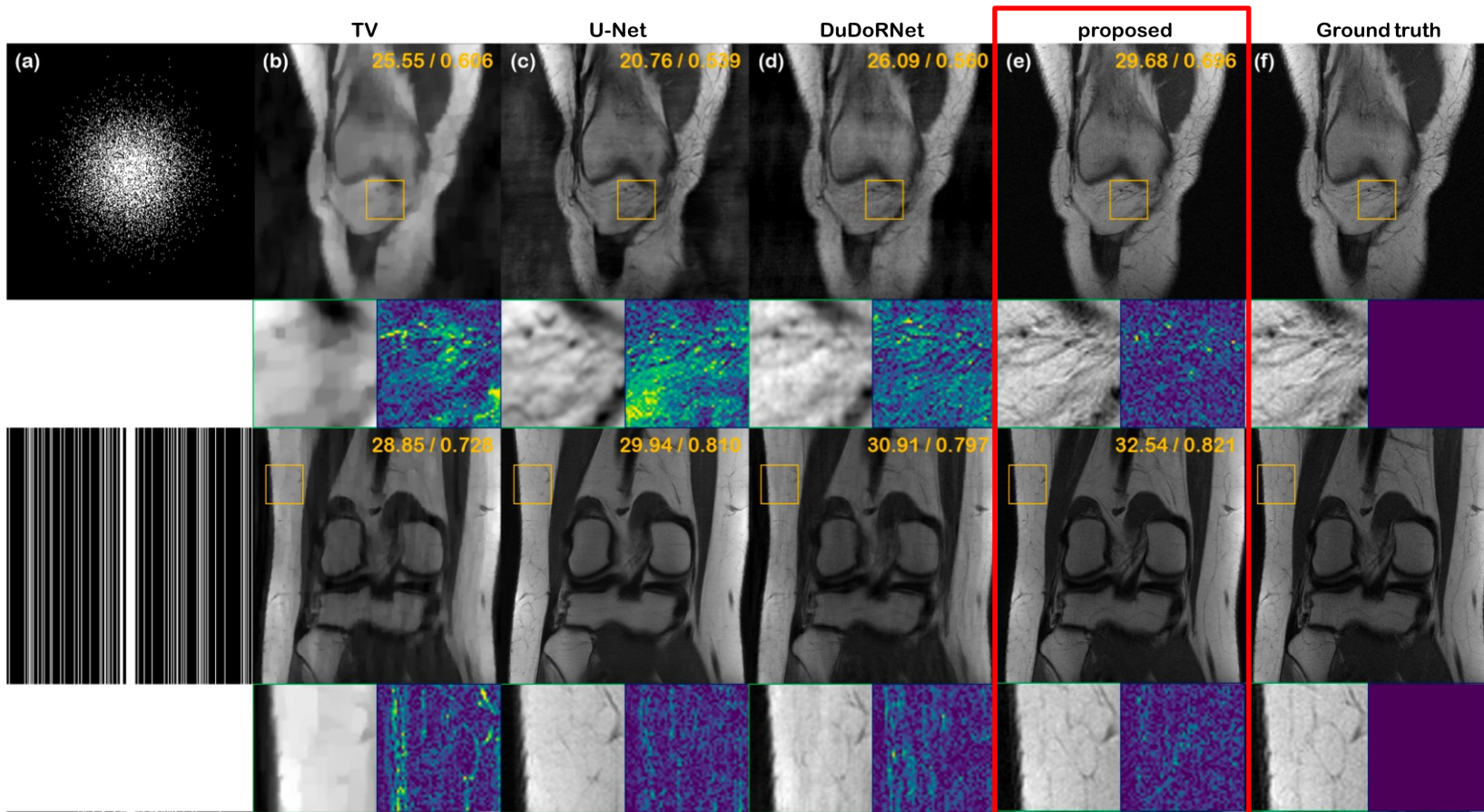
$$\min_{\mathbf{x}} \|y - A\mathbf{x}\|^2$$

$$\mathbf{x}_i \leftarrow \mathbf{x}_{i+1} + \epsilon_i s_{\theta}(\mathbf{x}_{i+1}, \sigma_{i+1}) + \sqrt{2\epsilon_i} \mathbf{z} \quad \text{Denoising step (reverse SDE)}$$

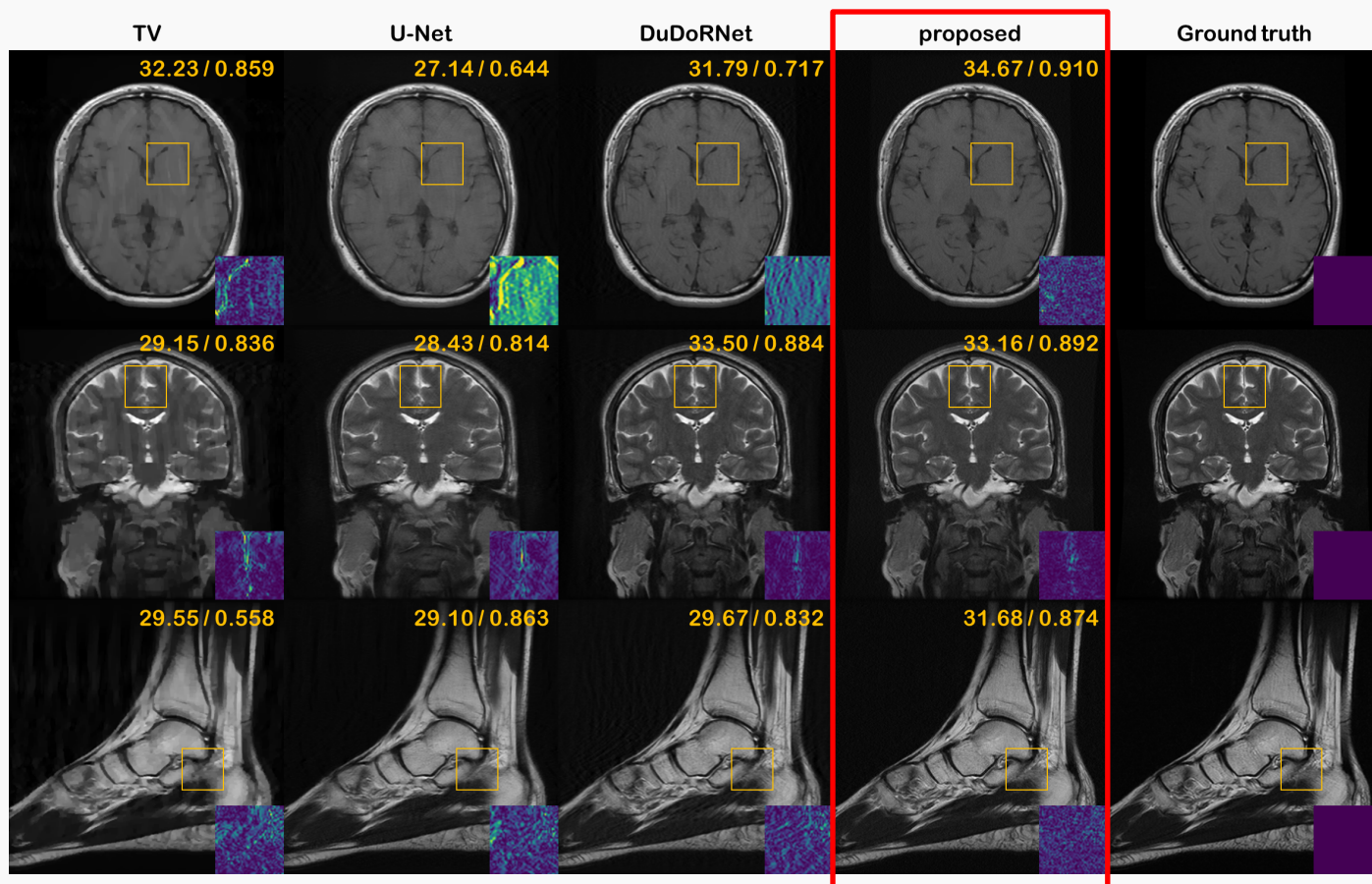
$$\mathbf{x}_i \leftarrow \mathbf{x}_i + \lambda A^*(y - A\mathbf{x}_i), \quad \text{Data consistency step (e.g. GD, POCS)}$$



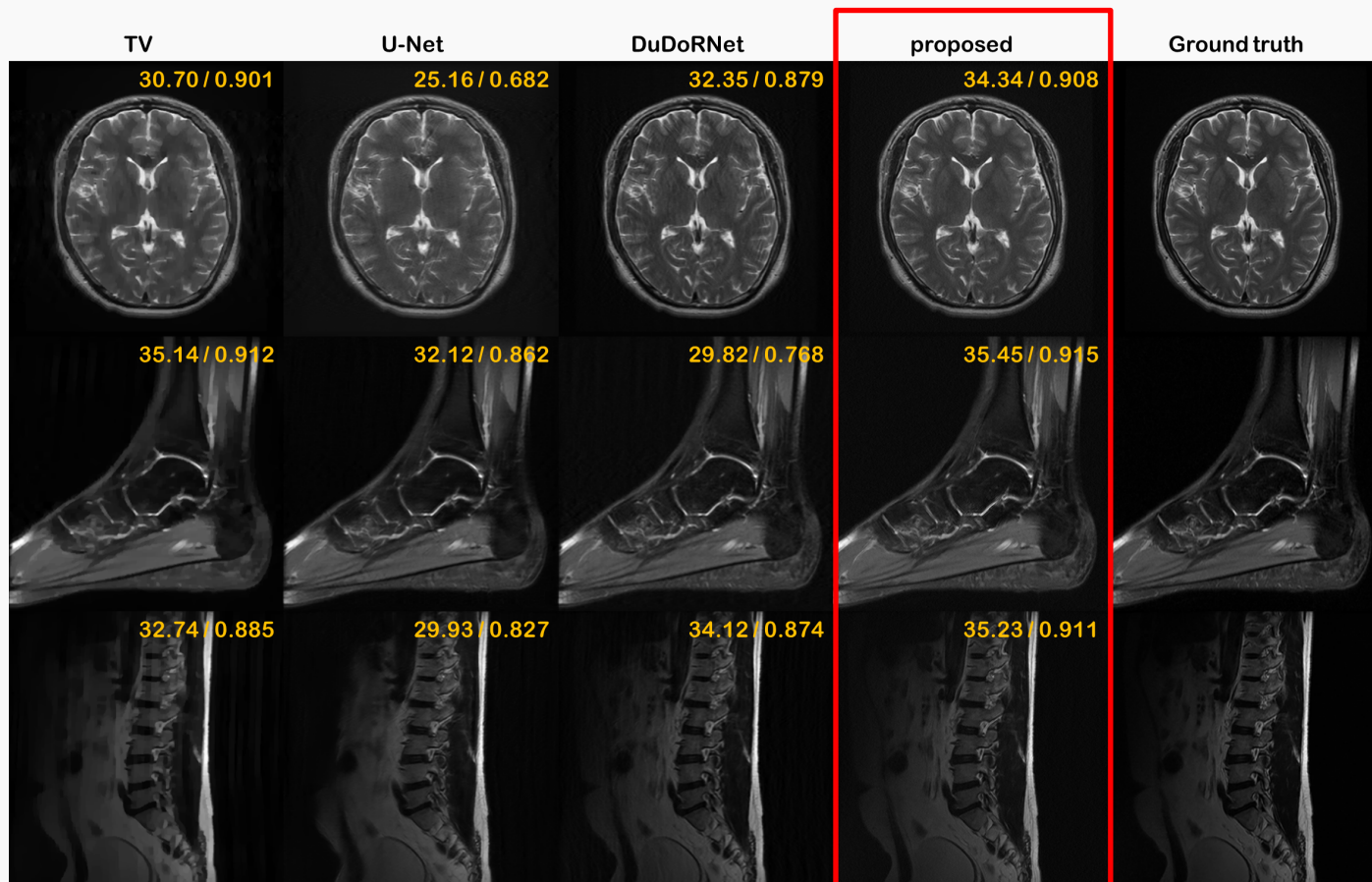
State-of-the-art Performance



Generalization Capability



Generalization Capability



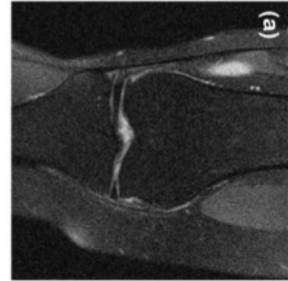
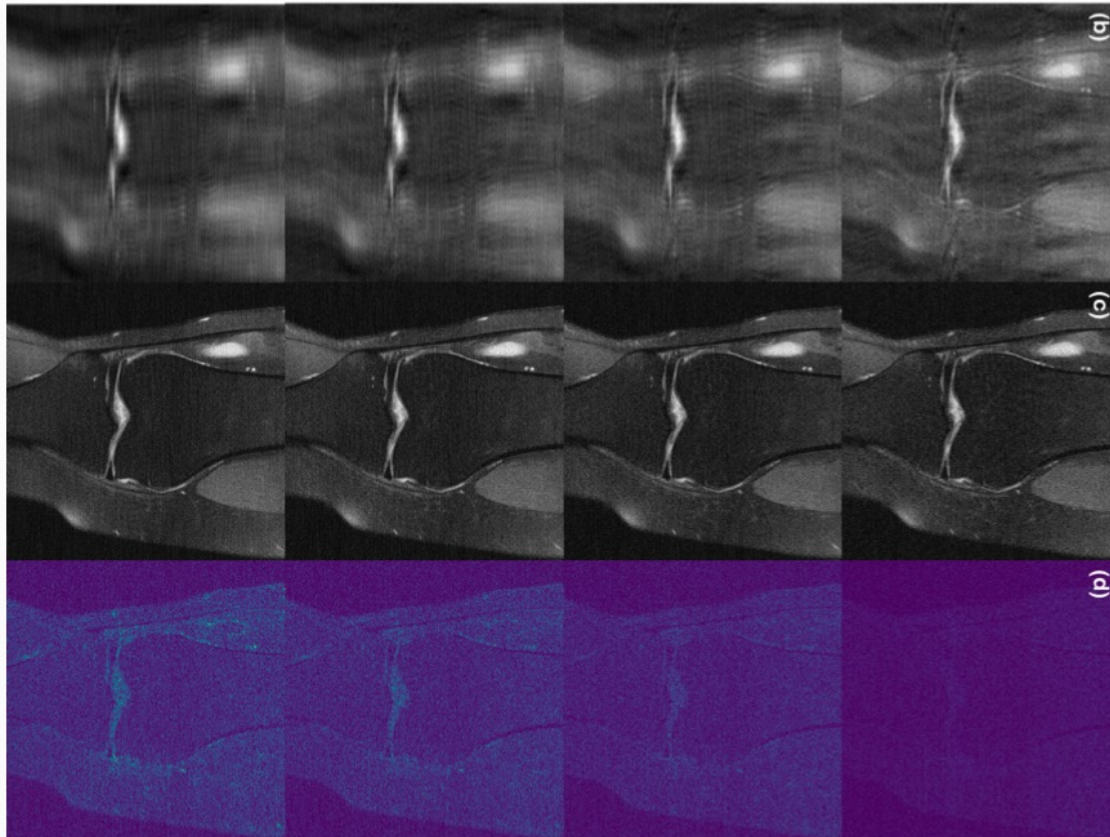
Uncertainty Quantification

Uniform x8

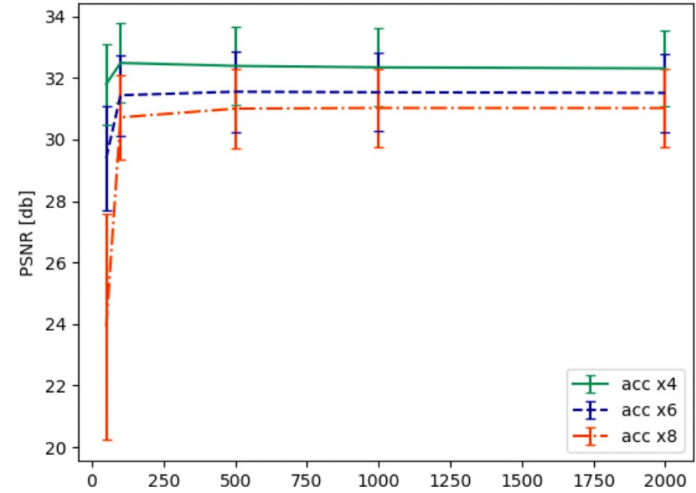
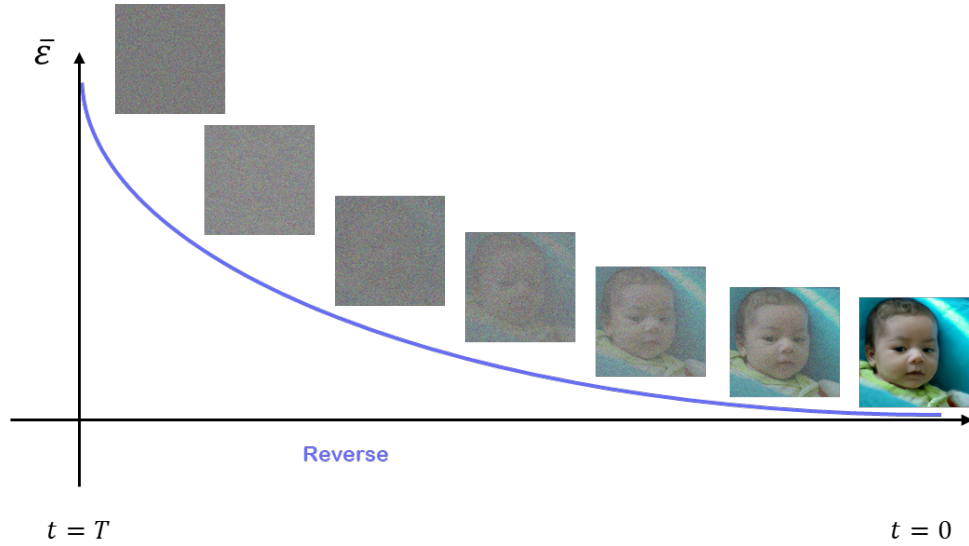
Uniform x6

Uniform x4

Uniform x2



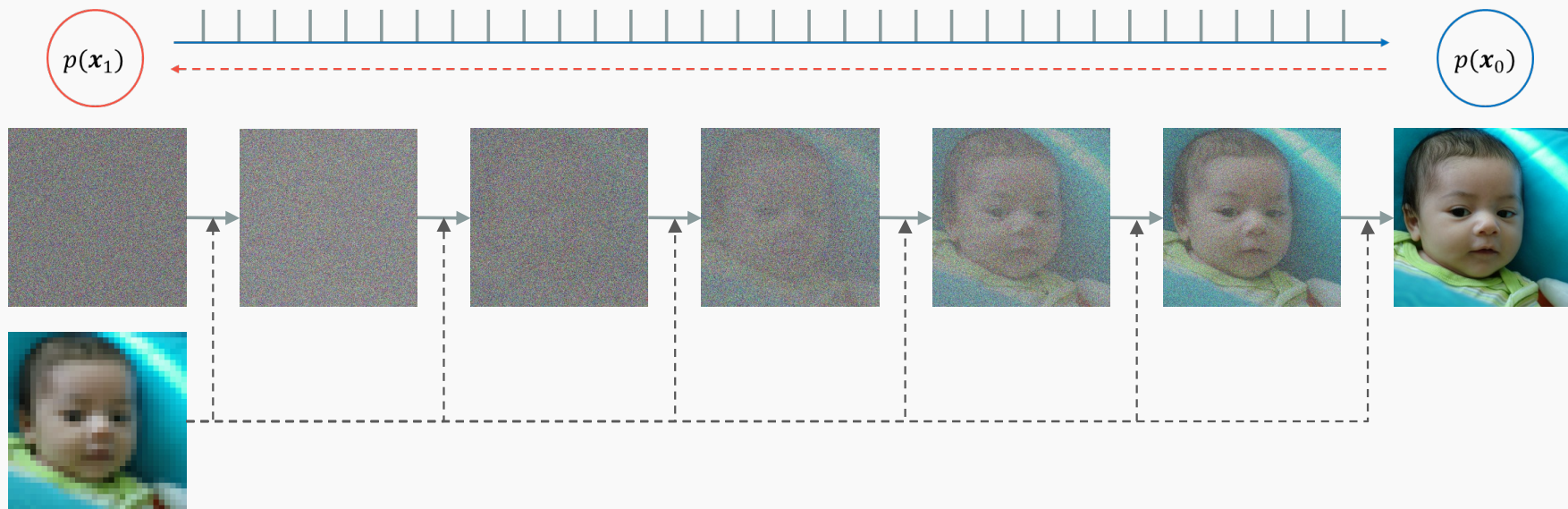
Very Slow Convergence



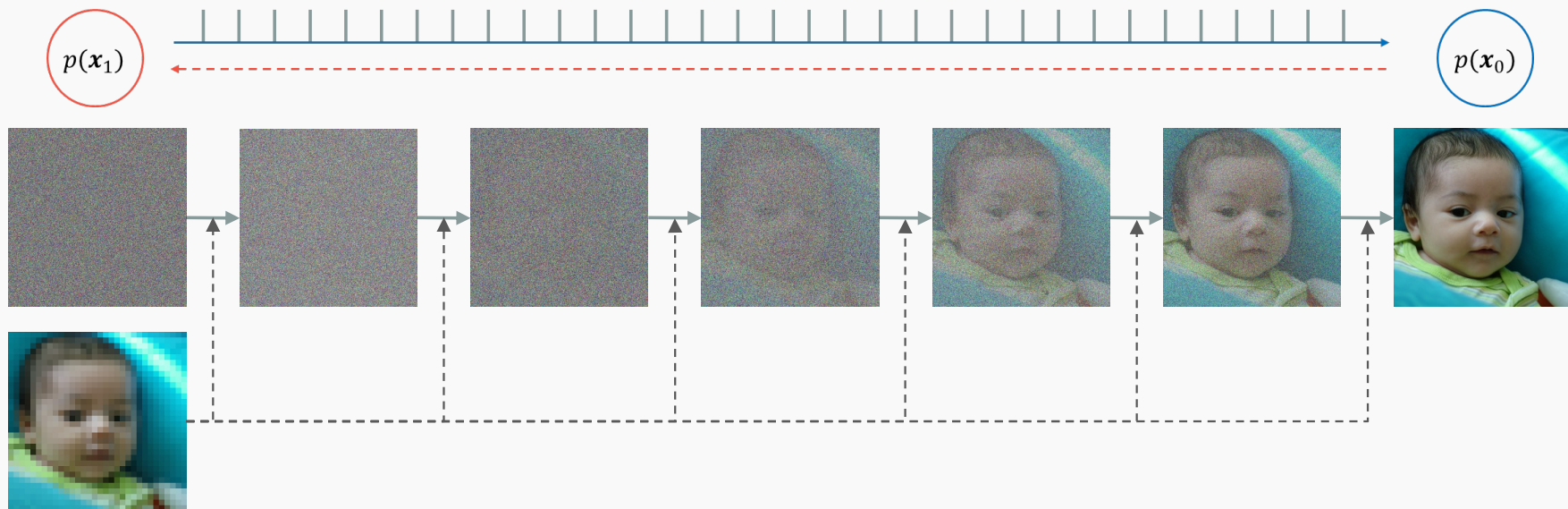
CCDF: COME CLOSER, DIFFUSE FASTER

Chung et al, CVPR, 2022

Intuition: Why use the whole process?



Intuition: Why use the whole process?



Is this part necessary?

Intuition: Why use the whole process?

- $t_0 = 0.3$
- $N' = t_0 N = 300$

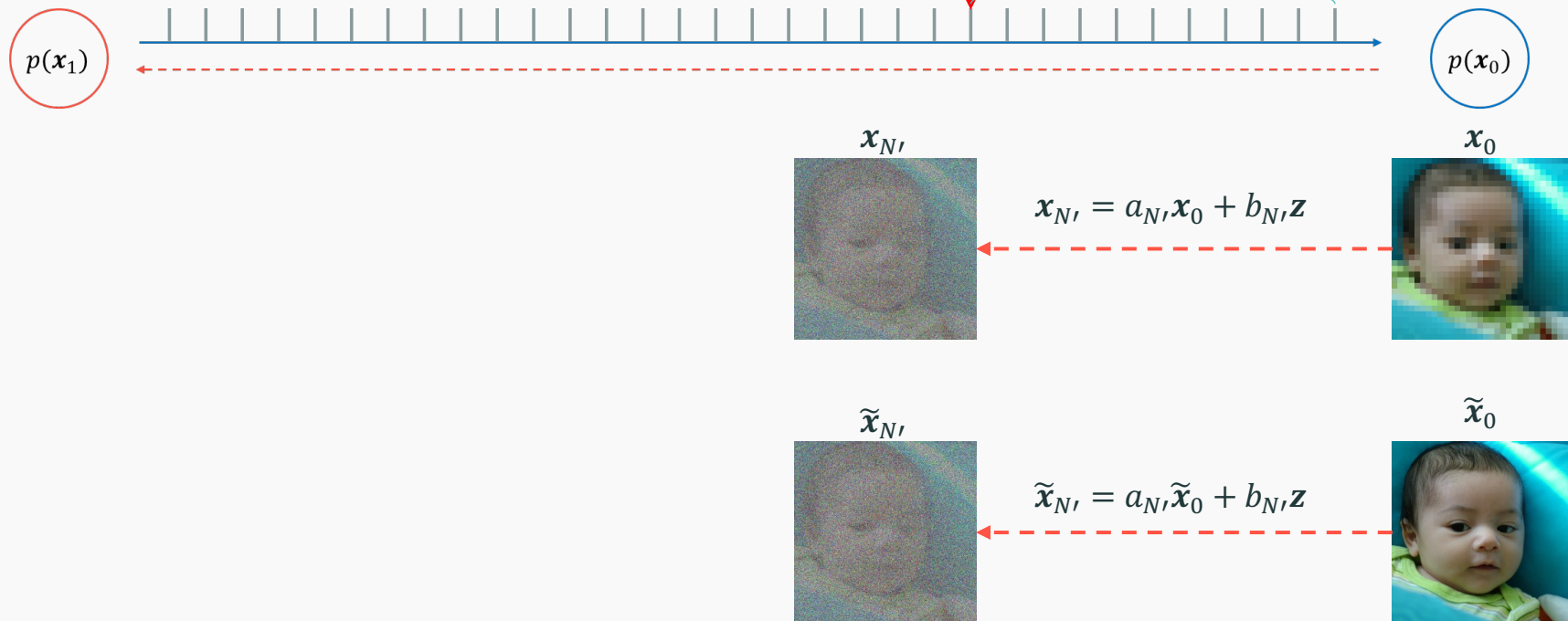
use this much



Intuition: Why use the whole process?

- $t_0 = 0.3$
- $N' = t_0 N = 300$

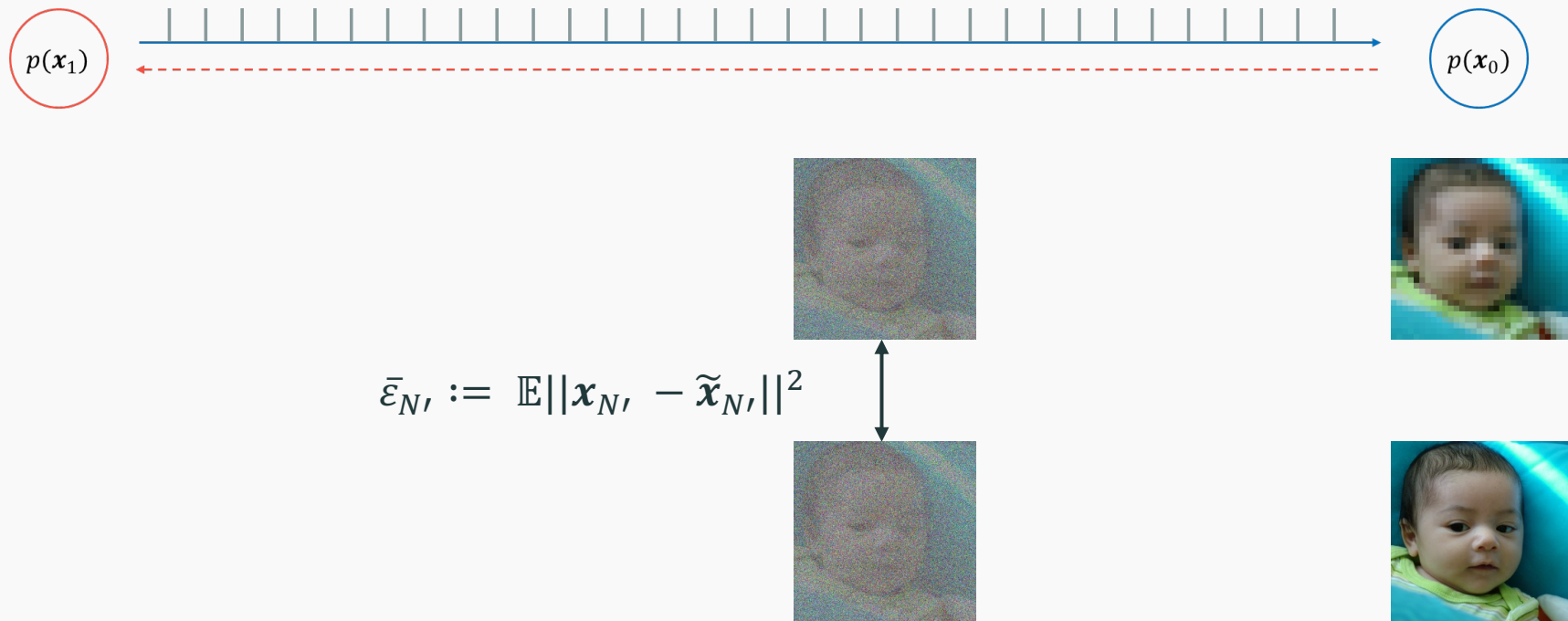
use this much



Intuition of CCDF



Intuition of CCDF



Intuition of CCDF



$$\bar{\epsilon}_{0,r} := \mathbb{E} \|\mathbf{x}_{0,r} - \tilde{\mathbf{x}}_{0,r}\|^2 < \epsilon_0$$



CCDF: The Algorithm



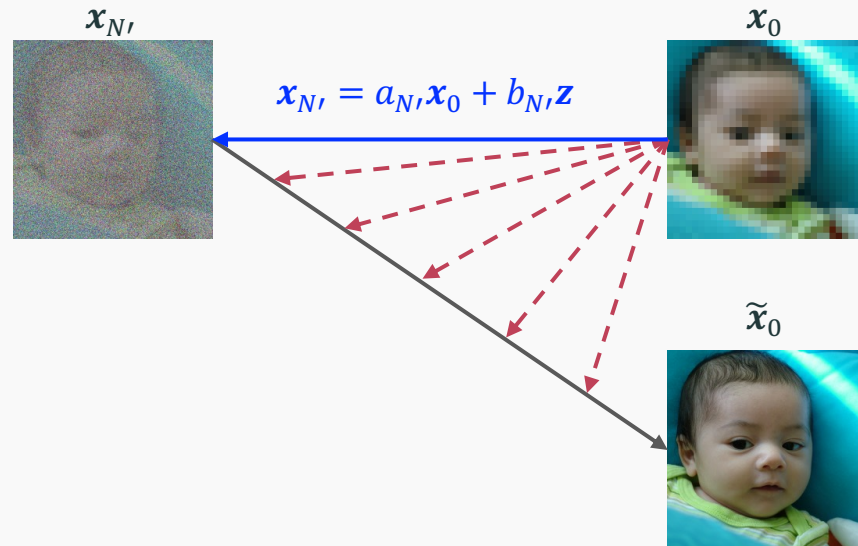
Algorithm 1 Accelerated Super-resolution / inpainting (VP, Markov)

Require: $x_0, \hat{x}_0, N', \{\alpha_i\}_{i=1}^{N'}, \{\sigma_i\}_{i=1}^{N'}, s_\theta$

- 1: $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: $x_{N'} \leftarrow \sqrt{\bar{\alpha}_{N'}}x_0 + \sqrt{1 - \bar{\alpha}_{N'}}z$ \triangleright Forward diffusion
- 3: **for** $i = N'$ to 1 **do** \triangleright Reverse diffusion
- 4: $x'_{i-1} \leftarrow \frac{1}{\sqrt{\alpha_i}}(x_i + (1 - \alpha_i)s_\theta(x_i, i))$
- 5: $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 6: $x_{i-1} \leftarrow x'_{i-1} + \sigma_i z$ \triangleright Unconditional update
- 7: $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 8: $\hat{x}_i \leftarrow \sqrt{\bar{\alpha}_i}\hat{x}_0 + \sqrt{1 - \bar{\alpha}_i}z$
- 9: $x_{i-1} = (\mathbf{I} - \mathbf{P})x_{i-1} + \hat{x}_i$ \triangleright Measurement consistency

10: **end for**

11: **return** x_0



CCDF: The Algorithm

General form

$$\mathbf{x}_{N'} = a_{N'}\mathbf{x}_0 + b_{N'}\mathbf{z}$$

$$\mathbf{x}'_{i-1} = \mathbf{f}(\mathbf{x}_i, i) + \mathbf{g}(\mathbf{x}_i, i)\mathbf{z}_i$$

$$\mathbf{x}_{i=1} = A\mathbf{x}'_{i-1} + \mathbf{b}$$

: 1-step noising

: Iterative denoising

- Denoising step (reverse SDE)
- Data consistency step (e.g. GD, POCS)

Constraint

$$\|A\mathbf{x} - A\mathbf{x}'\| \leq \|\mathbf{x} - \mathbf{x}'\| \quad \forall \mathbf{x}, \mathbf{x}'$$

• Non-expansive mapping

Key Idea: Stochastic Contraction

Contraction on \mathbb{R}^n

A function $f: \mathbb{R}^n \mapsto \mathbb{R}^n$ **contraction mapping**,
if there exists $0 \leq \lambda < 1$ s.t. $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

$$\|f(\mathbf{x}) - f(\mathbf{y})\| \leq \lambda \|\mathbf{x} - \mathbf{y}\|$$

$$\sigma_{\max} \left(\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right) \leq \lambda < 1$$

Theorem A.1. (Pham et al. 2008)

$$\mathbf{x}_{i+1} = f(\mathbf{x}_i, i) + g(\mathbf{x}_i, i)\mathbf{z}$$

- f is contracting with λ
- $\text{Tr}(g(\mathbf{x}, i)I g(\mathbf{x}, i)) \leq C \quad \forall \mathbf{x}, i$

Then,
$$\mathbb{E}\|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|^2 \leq \frac{2C}{1 - \lambda^2} + \lambda^{2i} \mathbb{E}\|\mathbf{x}_0 - \tilde{\mathbf{x}}_0\|^2$$

Reverse SDE is Contracting!

Proof of Theorem 1. (VE-SDE; SMLD)

Forward

$$\mathbf{x}_i = \mathbf{x}_0 + \sigma_i \mathbf{z}$$

Reverse SDE

$$\mathbf{x}'_{i-1} = \boxed{\mathbf{x}_i + (\sigma_i^2 - \sigma_{i-1}^2) s_\theta(\mathbf{x}_i, i)} + \sqrt{\sigma_i^2 - \sigma_{i-1}^2} \mathbf{z} \quad \text{: Stochastically contracting}$$

f(x_i, i) →

Proof.

$$\frac{\partial \mathbf{f}^T(\mathbf{x}_i, i)}{\partial \mathbf{x}_i} = \mathbf{I} + (\sigma_i^2 - \sigma_{i-1}^2) \frac{\partial s_\theta(\mathbf{x}_i, i)}{\partial \mathbf{x}_i} = \frac{\sigma_{i-1}^2 - \sigma_0^2}{\sigma_i^2 - \sigma_0^2} \mathbf{I}$$

$$\lambda = \max_{i \in [N']} \frac{\sigma_{i-1}^2 - \sigma_0^2}{\sigma_i^2 - \sigma_0^2} < 1$$

$$C = \max_{i \in [N']} \sigma_i^2 - \sigma_{i-1}^2$$

Non-expansiveness is Sufficient!

Corollary 1.

$$\mathbf{x}'_{i+1} = f(\mathbf{x}_i, i) + g(\mathbf{x}_i, i)\mathbf{z}_i$$

$$\mathbf{x}_{i+1} = \mathbf{A}\mathbf{x}'_{i+1} + \mathbf{b}$$

Non-expansive mapping

Proof.

$$\mathbf{x}_{i+1} = \underbrace{\mathbf{A}f(\mathbf{x}_i, i) + \mathbf{b}}_{\tilde{\mathbf{f}}(\mathbf{x}_i, i)} + \sigma(\mathbf{x}_i, i)\mathbf{A}\mathbf{z}_i$$

$$\mathbb{E}\|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|^2 \leq \frac{2C\tau}{1 - \lambda^2} + \lambda^{2i}\mathbb{E}\|\mathbf{x}_0 - \tilde{\mathbf{x}}_0\|^2$$

$$\tau = \frac{\text{Tr}(\mathbf{A}^T \mathbf{A})}{n}$$

$$\sigma_{\max}\left(\frac{\partial \tilde{\mathbf{f}}(\mathbf{x}, i)}{\partial \mathbf{x}}\right) \leq \sigma_{\max}(\mathbf{A})\sigma_{\max}\left(\frac{\partial f(\mathbf{x}, i)}{\partial \mathbf{x}}\right) \leq \lambda$$

$$\text{Tr}(g(\mathbf{x}, i)\mathbf{A}^T \mathbf{A}g(\mathbf{x}, i)) = g(\mathbf{x}, i)^2 \text{Tr}(\mathbf{A}^T \mathbf{A}) = C\tau$$

Theoretical Findings

Theorem 1.

$$\bar{\varepsilon}_{0,r} \leq \frac{2C\tau}{1 - \lambda^2} + \lambda^{2N'} \bar{\varepsilon}_{N'}$$

Error decreases **exponentially**
with reverse diffusion!

$$\lambda = \begin{cases} \max_{i \in [N']} \sqrt{\alpha_i} \left(\frac{1 - \bar{\alpha}_{i-1}}{1 - \bar{\alpha}_i} \right) & (DDPM) \\ \max_{i \in [N']} \frac{\sigma_{i-1}^2 - \sigma_0^2}{\sigma_i^2 - \sigma_0^2} & (SMLD) \\ \max_{i \in [N']} \frac{\sigma_{i-1}}{\sigma_i} & (DDIM) \end{cases}$$

$$C = \begin{cases} n(1 - \alpha_N) & (DDPM) \\ n \max_{i \in [N']} \sigma_i^2 - \sigma_{i-1}^2 & (SMLD) \\ 0 & (DDIM) \end{cases}$$

$$\tau = \frac{\text{Tr}(\mathbf{A}^T \mathbf{A})}{n}$$

Theoretical Findings

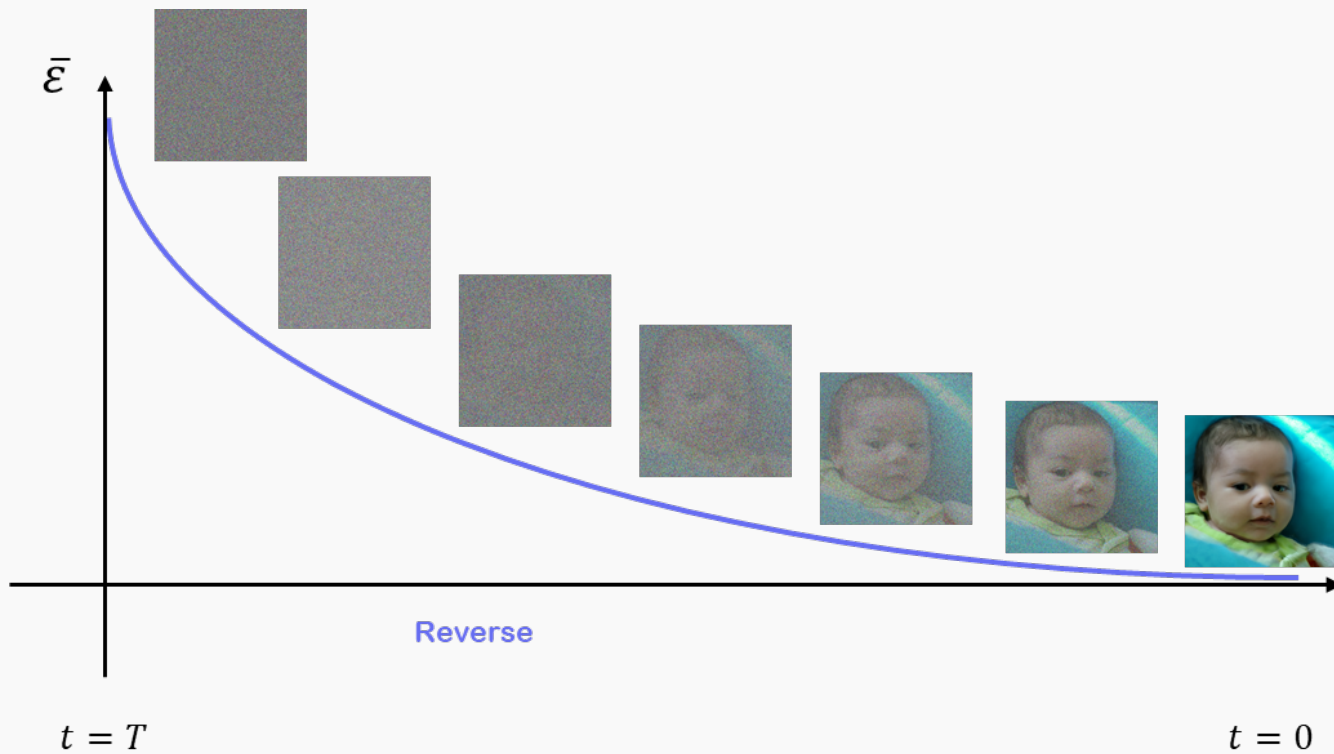
Theorem 2. (shortcut path)

- For any $0 < \mu \leq 1$, there exists a **minimum N'** s.t.

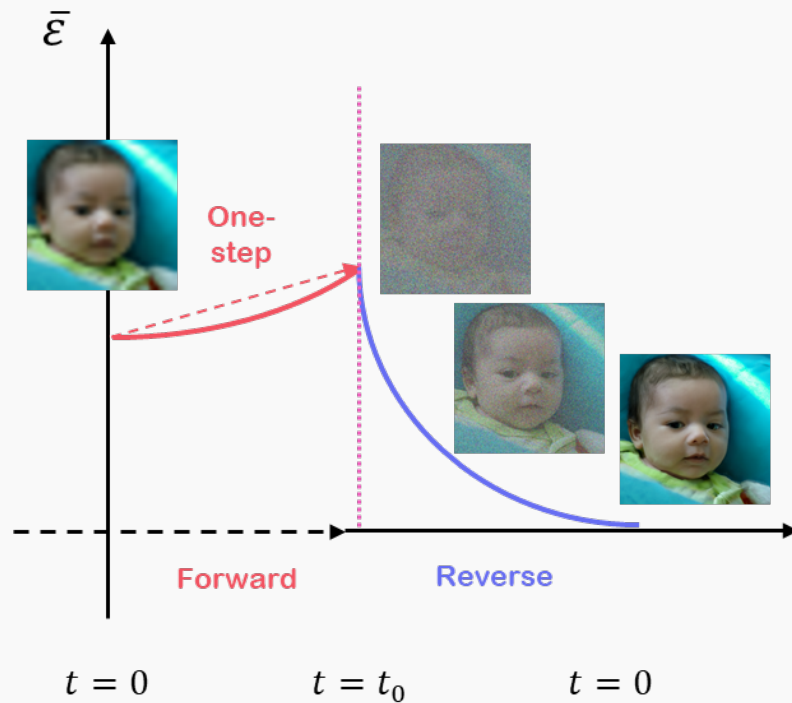
$$\bar{\varepsilon}_{0,r} \leq \mu \varepsilon_0$$

- **Optimal N' decreases** as ε_0 **gets smaller**

Come Closer, Diffuse Faster

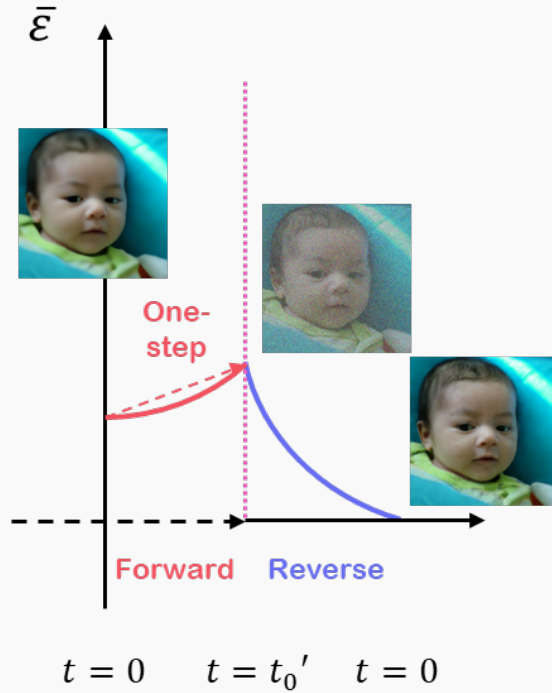


Come Closer, Diffuse Faster

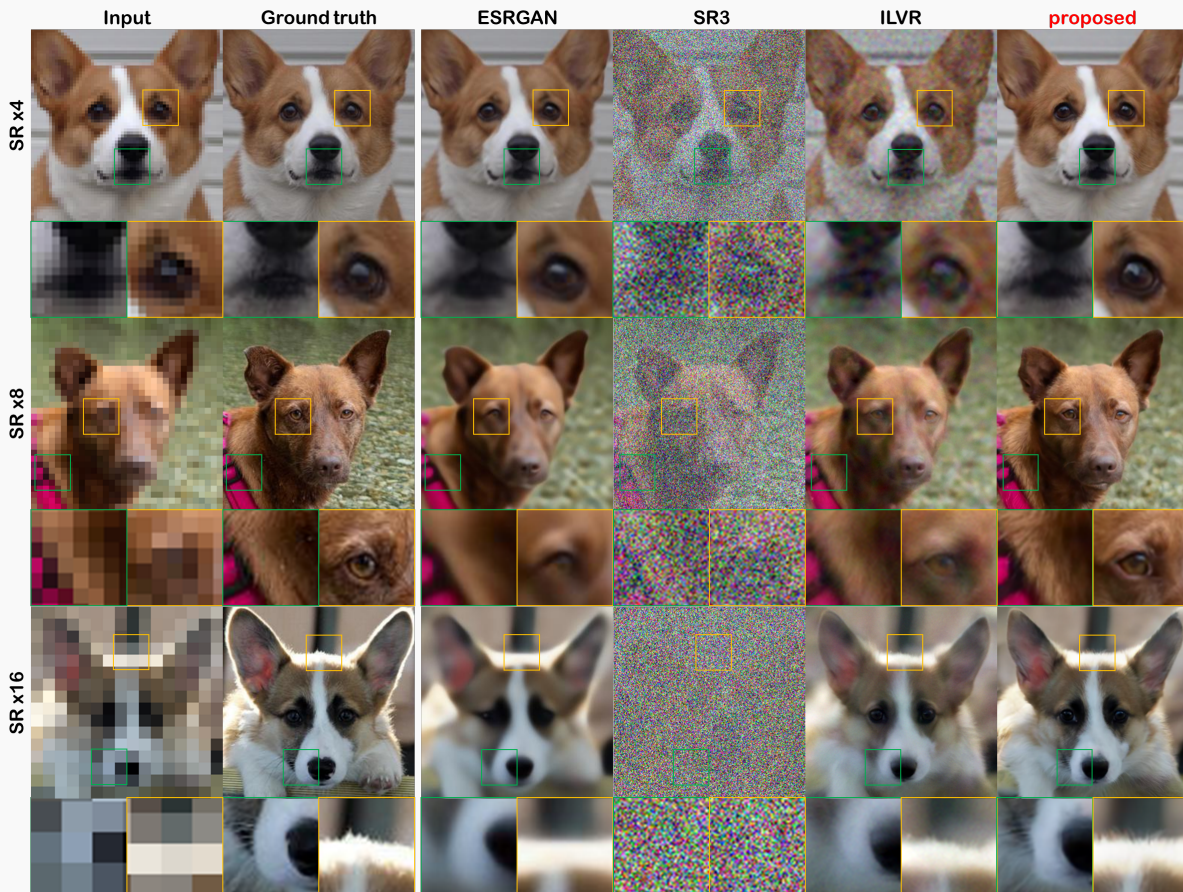


Come Closer, Diffuse Faster

Feed-forward
network correction



Experimental Results: SR



20 step diffusion

- ILVR, SR3

$$N = 20, \quad t_0 = 1.0$$

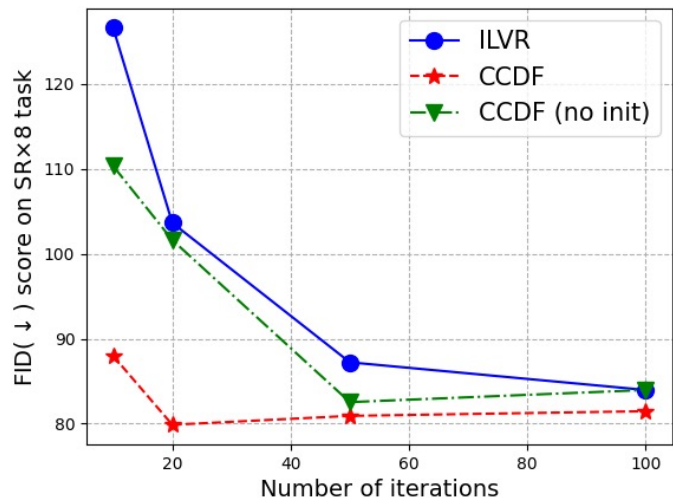
- proposed

$$N = 100, \quad t_0 = 0.2$$

| t_0 | 0.05 | 0.1 | 0.2 | 0.5 | 0.75 | 1.0 [5] |
|----------------|--------|--------------|--------------|--------------|-------|--------------|
| SR $\times 4$ | 63.90 | 60.90 | 60.91 | 64.04 | 64.14 | 63.31 |
| SR $\times 8$ | 85.21 | 78.13 | 75.76 | 79.34 | 79.67 | <u>77.34</u> |
| SR $\times 16$ | 116.37 | 101.79 | 92.59 | 88.09 | 92.12 | <u>88.49</u> |

Table 1. FID(\downarrow) scores on FFHQ test set for SR task with $N = 1000$, and varying t_0 values. $t_0 = 1.0$ is the baseline method without any acceleration used in [5]. Numbers in boldface, and underline indicate the best, and the second best scores.

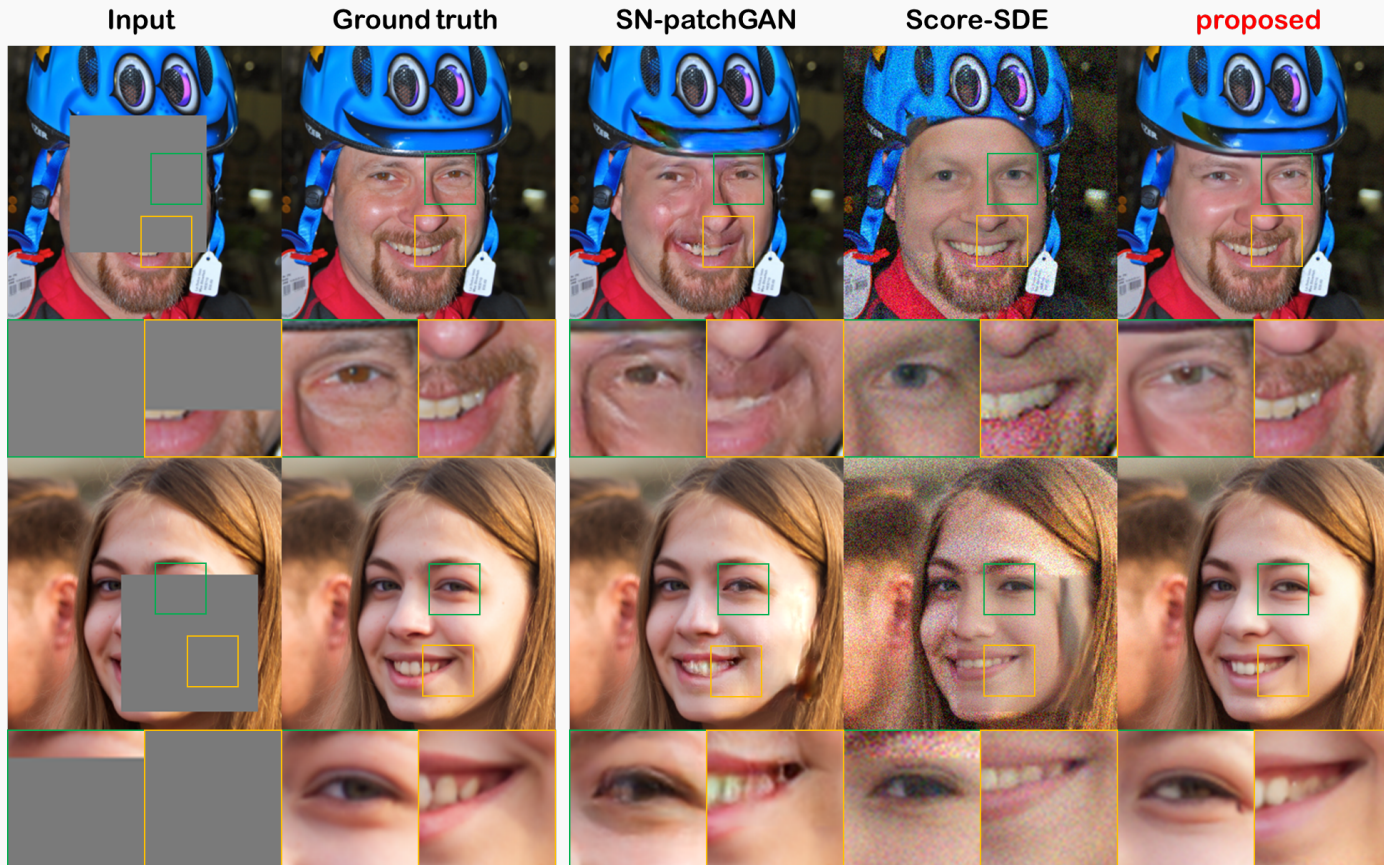
Experimental Results: SR



| | SR factor | ESRGAN [36] | SR3* [25] | ILVR [5] | CCDF (ours) |
|------|-----------|-------------|-----------|--------------|--------------|
| FFHQ | ×4 | 81.14 | 66.79 | 63.14 | 60.90 |
| | ×8 | 108.96 | 80.27 | 81.85 | 75.76 |
| | ×16 | 143.80 | 99.46 | 92.32 | 88.39 |
| AFHQ | ×4 | 24.52 | 20.68 | 18.70 | 15.53 |
| | ×8 | 51.84 | 30.23 | 34.85 | 32.30 |
| | ×16 | 98.22 | 60.76 | 47.28 | 48.77 |

Table 2. Comparison of FID(↓) scores on FFHQ and AFHQ test set. t_0 values used for the proposed method is 0.1, 0.2, 0.3 for ×4, ×8, ×16 SR, respectively. Numbers in boldface represent the best results among the row. (*unofficial re-implementation)

Experimental Results: Inpainting



20 step diffusion

- **Score-SDE**

$$N = 20, \quad t_0 = 1.0$$

- **proposed**

$$N = 100, \quad t_0 = 0.2$$

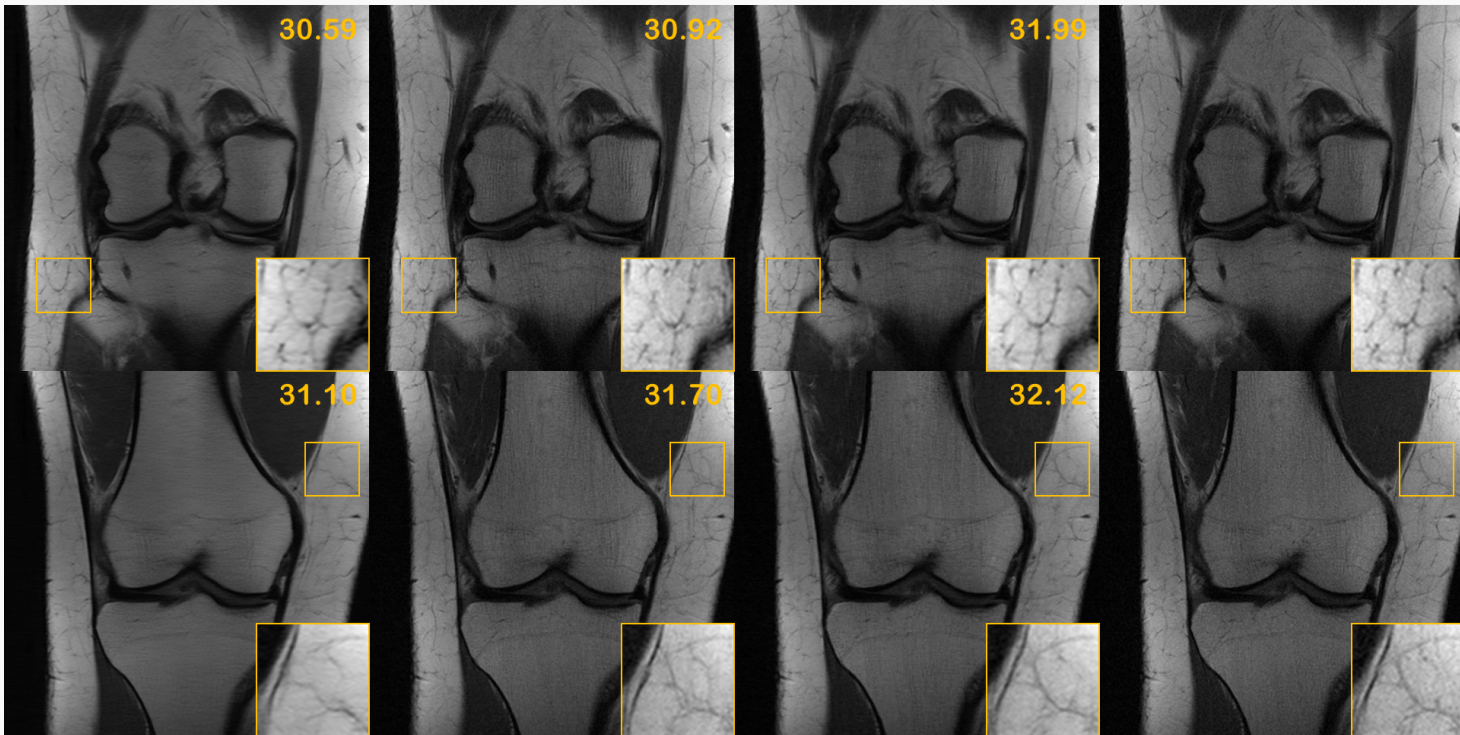
Experimental Results: Fast MRI

U-Net

Chung et al.

proposed

Ground truth



20 step diffusion

- Chung et al.

$N = 1000$, $t_0 = 1.0$

- proposed

$N = 1000$, $t_0 = 0.02$

Summary

- **Diffusion models**: Exciting new path for solving inverse problems
- **Universal solver** without knowledge about the problem a priori
- Great **generalization** capacity
- Acceleration through **stochastic contraction** theory





Questions?