

# Deep Equilibrium Diffusion Restoration with Parallel Sampling

(Supplementary Materials)

Jiezhang Cao<sup>1</sup>, Yue Shi<sup>1,2</sup>, Kai Zhang<sup>3</sup>, Yulun Zhang<sup>2</sup>, Radu Timofte<sup>1,4</sup>, Luc Van Gool<sup>1,5</sup>

<sup>1</sup>ETH Zürich, <sup>2</sup>Shanghai Jiao Tong University, <sup>3</sup>Nanjing University, <sup>4</sup>University of Würzburg, <sup>5</sup>KU Leuven

**Organization.** In this paper, we organize our supplementary materials as follows. In Section A, we provide detailed proofs of our proposed Proposition 1. In Section C, we provide more results of our method. In Section B, we provide more visual comparisons. In Section D, we provide more details and visual results of initialization optimization. In Section E, the limitations and future work of our proposed method are discussed.

## A. Proof of Proposition 1

**Proof** Based on existing diffusion model-based IR (e.g., [68]), the estimated  $\mathbf{x}_0$  at time-step  $t$  can be formulated as:

$$\mathbf{x}_{0|t} = \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)). \quad (14)$$

We fix the range-space as  $\mathbf{A}^\dagger \mathbf{y}$  and leave the null-space unchanged, then  $\mathbf{x}_0$  at time-step  $t$  can be further estimated by

$$\begin{aligned} \hat{\mathbf{x}}_{0|t} &= (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \mathbf{x}_{0|t} + \mathbf{A}^\dagger \mathbf{y} \\ &= \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)) + \mathbf{A}^\dagger \mathbf{y} \\ &= \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \mathbf{x}_t - \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) + \mathbf{A}^\dagger \mathbf{y}, \end{aligned} \quad (15)$$

where the second line is based on Eqn. (14). Based on DDIM [61],  $\mathbf{x}_{t-1}$  can be updated by

$$\begin{aligned} \mathbf{x}_{t-1} &= \sqrt{\bar{\alpha}_{t-1}} \hat{\mathbf{x}}_{0|t} + \gamma_t \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_t + \gamma_t \sqrt{1 - \bar{\alpha}_t} \sqrt{1 - \eta^2} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \\ &:= \sqrt{\bar{\alpha}_{t-1}} \hat{\mathbf{x}}_{0|t} + \gamma_t c_t^1 \boldsymbol{\epsilon}_t + \gamma_t c_t^2 \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \\ &= \sqrt{\bar{\alpha}_{t-1}} \left[ \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \mathbf{x}_t - \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) + \mathbf{A}^\dagger \mathbf{y} \right] + \gamma_t c_t^1 \boldsymbol{\epsilon}_t + \gamma_t c_t^2 \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \\ &= \frac{\sqrt{\bar{\alpha}_{t-1}}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \mathbf{x}_t + \left( \gamma_t c_t^2 - \frac{\sqrt{\bar{\alpha}_{t-1}}(1 - \bar{\alpha}_t)}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \right) \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) + \sqrt{\bar{\alpha}_{t-1}} \mathbf{A}^\dagger \mathbf{y} + \gamma_t c_t^1 \boldsymbol{\epsilon}_t \\ &:= \boldsymbol{\Psi}_t^1 \mathbf{x}_t + \boldsymbol{\Psi}_t^2 \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) + \boldsymbol{\Psi}_t^3 \mathbf{y} + \boldsymbol{\Psi}_t^4 \boldsymbol{\epsilon}_t, \end{aligned} \quad (16)$$

where  $\boldsymbol{\epsilon}_t$  is sampled from standard Gaussian distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ , and  $\gamma_t$  is a parameter which can be set as 1. The third line is based on Eqn. (15). Let  $c_t^1 = \sqrt{1 - \bar{\alpha}_t} \eta$  and  $c_t^2 = \sqrt{1 - \bar{\alpha}_t} \sqrt{1 - \eta^2}$ , and we define  $\boldsymbol{\Psi}_t^i, i = 1, \dots, 4$  as

$$\boldsymbol{\Psi}_t^1 = \frac{\sqrt{\bar{\alpha}_{t-1}}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}), \quad \boldsymbol{\Psi}_t^2 = \gamma_t c_t^2 - \frac{\sqrt{\bar{\alpha}_{t-1}}(1 - \bar{\alpha}_t)}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}), \quad \boldsymbol{\Psi}_t^3 = \sqrt{\bar{\alpha}_{t-1}} \mathbf{A}^\dagger, \quad \boldsymbol{\Psi}_t^4 = \gamma_t c_t^1.$$

Based on the formulation of  $\mathbf{x}_{t-1}$ , we can write  $\mathbf{x}_{t-2}$  as

$$\begin{aligned} \mathbf{x}_{t-2} &= \boldsymbol{\Psi}_{t-1}^1 \mathbf{x}_{t-1} + \boldsymbol{\Psi}_{t-1}^2 \boldsymbol{\epsilon}_\theta(\mathbf{x}_{t-1}, t-1) + \boldsymbol{\Psi}_{t-1}^3 \mathbf{y} + \boldsymbol{\Psi}_{t-1}^4 \boldsymbol{\epsilon}_{t-1} \\ &= \boldsymbol{\Psi}_{t-1}^1 (\boldsymbol{\Psi}_t^1 \mathbf{x}_t + \boldsymbol{\Psi}_t^2 \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) + \boldsymbol{\Psi}_t^3 \mathbf{y} + \boldsymbol{\Psi}_t^4 \boldsymbol{\epsilon}_t) + \boldsymbol{\Psi}_{t-1}^2 \boldsymbol{\epsilon}_\theta(\mathbf{x}_{t-1}, t-1) + \boldsymbol{\Psi}_{t-1}^3 \mathbf{y} + \boldsymbol{\Psi}_{t-1}^4 \boldsymbol{\epsilon}_{t-1} \\ &= \boldsymbol{\Psi}_{t-1}^1 \boldsymbol{\Psi}_t^1 \mathbf{x}_t + \boldsymbol{\Psi}_{t-1}^1 \boldsymbol{\Psi}_t^2 \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) + \boldsymbol{\Psi}_{t-1}^1 \boldsymbol{\Psi}_t^3 \mathbf{y} + \boldsymbol{\Psi}_{t-1}^1 \boldsymbol{\Psi}_t^4 \boldsymbol{\epsilon}_t + \boldsymbol{\Psi}_{t-1}^2 \boldsymbol{\epsilon}_\theta(\mathbf{x}_{t-1}, t-1) + \boldsymbol{\Psi}_{t-1}^3 \mathbf{y} + \boldsymbol{\Psi}_{t-1}^4 \boldsymbol{\epsilon}_{t-1} \\ &= \boldsymbol{\Psi}_{t-1}^1 \boldsymbol{\Psi}_t^1 \mathbf{x}_t + (\boldsymbol{\Psi}_{t-1}^1 \boldsymbol{\Psi}_t^2 \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) + \boldsymbol{\Psi}_{t-1}^2 \boldsymbol{\epsilon}_\theta(\mathbf{x}_{t-1}, t-1)) + (\boldsymbol{\Psi}_{t-1}^1 \boldsymbol{\Psi}_t^3 \mathbf{y} + \boldsymbol{\Psi}_{t-1}^3 \mathbf{y}) + (\boldsymbol{\Psi}_{t-1}^1 \boldsymbol{\Psi}_t^4 \boldsymbol{\epsilon}_t + \boldsymbol{\Psi}_{t-1}^4 \boldsymbol{\epsilon}_{t-1}) \\ &= \frac{\sqrt{\bar{\alpha}_{t-2}}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \mathbf{x}_t + \sum_{s=t-2}^{t-1} \frac{\sqrt{\bar{\alpha}_{t-2}}}{\sqrt{\bar{\alpha}_s}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) (\boldsymbol{\Psi}_{s+1}^2 \boldsymbol{\epsilon}_\theta(\mathbf{x}_{s+1}, s+1) + \boldsymbol{\Psi}_{s+1}^3 \mathbf{y} + \boldsymbol{\Psi}_{s+1}^4 \boldsymbol{\epsilon}_{s+1}) \\ &\quad + \mathbf{A}^\dagger \mathbf{A} (\boldsymbol{\Psi}_{t-1}^2 \boldsymbol{\epsilon}_\theta(\mathbf{x}_{t-1}, t-1) + \boldsymbol{\Psi}_{t-1}^3 \mathbf{y} + \boldsymbol{\Psi}_{t-1}^4 \boldsymbol{\epsilon}_{t-1}) \\ &:= \frac{\sqrt{\bar{\alpha}_{t-2}}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \mathbf{x}_t + \sum_{s=t-2}^{t-1} \frac{\sqrt{\bar{\alpha}_{t-2}}}{\sqrt{\bar{\alpha}_s}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \mathbf{z}_{s+1} + \mathbf{A}^\dagger \mathbf{A} \mathbf{z}_{t-1}, \end{aligned} \quad (17)$$

where the second line follows Eqn. (16), the third line holds by the definition of  $\Psi_t^i$ , and in the last line, we define  $z_s$  as:

$$\begin{aligned} z_s &= \Psi_s^2 \epsilon_\theta(\mathbf{x}_s, s) + \Psi_s^3 \mathbf{y} + \Psi_s^4 \epsilon_s \\ &= \left( \gamma_s c_s^2 - \frac{\sqrt{1 - \bar{\alpha}_s}}{\sqrt{\alpha_s}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \right) \epsilon_\theta(\mathbf{x}_s, s) + \sqrt{\bar{\alpha}_{s-1}} \mathbf{A}^\dagger \mathbf{y} + \gamma_s c_s^1 \epsilon_s, \end{aligned} \quad (18)$$

where the second line is based on the definitions of  $c_t^2$  and  $\Psi_t^i, i = 2, 3, 4$ . Similarly, we further write  $\mathbf{x}_{t-3}$  as

$$\begin{aligned} \mathbf{x}_{t-3} &= \Psi_{t-2}^1 \mathbf{x}_{t-2} + \Psi_{t-2}^2 \epsilon_\theta(\mathbf{x}_{t-2}, t-2) + \Psi_{t-2}^3 \mathbf{y} + \Psi_{t-2}^4 \epsilon_{t-2} \\ &= \Psi_{t-2}^1 \Psi_{t-1}^1 \Psi_t^1 \mathbf{x}_t + \left( \Psi_{t-2}^1 \Psi_{t-1}^1 \Psi_t^2 \epsilon_\theta(\mathbf{x}_t, t) + \Psi_{t-2}^1 \Psi_{t-1}^2 \epsilon_\theta(\mathbf{x}_{t-1}, t-1) + \Psi_{t-2}^2 \epsilon_\theta(\mathbf{x}_{t-2}, t-2) \right) \\ &\quad + \left( \Psi_{t-2}^1 \Psi_{t-1}^1 \Psi_t^3 \mathbf{y} + \Psi_{t-2}^1 \Psi_{t-1}^2 \Psi_t^3 \mathbf{y} + \Psi_{t-2}^3 \mathbf{y} \right) + \left( \Psi_{t-2}^1 \Psi_{t-1}^1 \Psi_t^4 \epsilon_t + \Psi_{t-2}^1 \Psi_{t-1}^2 \Psi_t^4 \epsilon_{t-1} + \Psi_{t-2}^4 \epsilon_{t-2} \right) \\ &= \frac{\sqrt{\bar{\alpha}_{t-3}}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \mathbf{x}_t + \sum_{s=t-3}^{t-1} \frac{\sqrt{\bar{\alpha}_{t-3}}}{\sqrt{\bar{\alpha}_s}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \left( \Psi_{s+1}^2 \epsilon_\theta(\mathbf{x}_{s+1}, s+1) + \Psi_{s+1}^3 \mathbf{y} + \Psi_{s+1}^4 \epsilon_{s+1} \right) \\ &\quad + \mathbf{A}^\dagger \mathbf{A} \left( \Psi_{t-2}^2 \epsilon_\theta(\mathbf{x}_{t-2}, t-2) + \Psi_{t-2}^3 \mathbf{y} + \Psi_{t-2}^4 \epsilon_{t-2} \right) \end{aligned} \quad (19)$$

$$:= \frac{\sqrt{\bar{\alpha}_{t-3}}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \mathbf{x}_t + \sum_{s=t-3}^{t-1} \frac{\sqrt{\bar{\alpha}_{t-3}}}{\sqrt{\bar{\alpha}_s}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) z_{s+1} + \mathbf{A}^\dagger \mathbf{A} z_{t-2}, \quad (20)$$

where the second line is according to the formulation of  $\mathbf{x}_{t-2}$ , and the third line is based on the definition of  $\Psi_t^i$ . In the line of Eqn. (19), we can derive the following

$$\begin{aligned} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \Psi_{t+1}^3 \mathbf{y} &= \sqrt{\bar{\alpha}_t} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \mathbf{A}^\dagger \mathbf{y} = \mathbf{0}, \\ \mathbf{A}^\dagger \mathbf{A} \Psi_{T-k+1}^3 \mathbf{y} &= \sqrt{\bar{\alpha}_{T-k}} \mathbf{A}^\dagger \mathbf{A} \mathbf{A}^\dagger \mathbf{y} = \sqrt{\bar{\alpha}_{T-k}} \mathbf{A}^\dagger \mathbf{y}, \\ (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \Psi_{t+1}^2 &= (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \left( \gamma_t c_t^2 - \frac{\sqrt{\bar{\alpha}_{t-1}(1 - \bar{\alpha}_t)}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \right) \\ &= \gamma_t c_t^2 (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) - \frac{\sqrt{\bar{\alpha}_{t-1}(1 - \bar{\alpha}_t)}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \\ &= \left( \gamma_t c_t^2 - \frac{\sqrt{\bar{\alpha}_{t-1}(1 - \bar{\alpha}_t)}}{\sqrt{\bar{\alpha}_t}} \right) (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}), \\ \mathbf{A}^\dagger \mathbf{A} \Psi_{T-k+1}^2 &= \mathbf{A}^\dagger \mathbf{A} \left( \gamma_{T-k+1} c_{T-k+1}^2 - \frac{\sqrt{\bar{\alpha}_{T-k}(1 - \bar{\alpha}_{T-k+1})}}{\sqrt{\bar{\alpha}_{T-k+1}}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \right) \\ &= \gamma_{T-k+1} c_{T-k+1}^2 \mathbf{A}^\dagger \mathbf{A}. \end{aligned}$$

Based on the induction, we can write  $\mathbf{x}_{T-k}$  as:

$$\mathbf{x}_{T-k} = \frac{\sqrt{\bar{\alpha}_{T-k}}}{\sqrt{\bar{\alpha}_T}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \mathbf{x}_T + \sum_{s=T-k}^{T-1} \frac{\sqrt{\bar{\alpha}_{T-k}}}{\sqrt{\bar{\alpha}_s}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) z_{s+1} + \mathbf{A}^\dagger \mathbf{A} z_{T-k+1}. \quad (21)$$

We complete the proof. Actually, the above equation can also be written as

$$\mathbf{x}_{T-k} = \frac{\sqrt{\bar{\alpha}_{T-k}}}{\sqrt{\bar{\alpha}_T}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A}) \mathbf{x}_T + \sum_{s=T-k}^{T-1} \frac{\sqrt{\bar{\alpha}_{T-k}}}{\sqrt{\bar{\alpha}_s}} (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A})^{T-k-s} z_{s+1}. \quad (22)$$

□

We extend Proposition 1 to the noisy inverse problem (i.e.,  $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{n}_\sigma$ ) as follows.

**Proposition 2 (Parallel sampling for noisy inverse problem)** *Given a degradation matrix  $\mathbf{A}$ , a degraded image  $\mathbf{y}$  and a Gaussian noise image  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and assume that the scale factor  $\Sigma_t$  for range-space is a constant diagonal matrix, i.e.,  $\Sigma = \Sigma_t$ , for every state  $t$ , for  $k \in [1, \dots, T]$ , the state  $\mathbf{x}_{T-k}$  can be predicted by previous states  $\{\mathbf{x}_{T-k+1}, \dots, \mathbf{x}_T\}$ , i.e.,*

$$\mathbf{x}_{T-k} = \frac{\sqrt{\bar{\alpha}_{T-k}}}{\sqrt{\bar{\alpha}_T}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A})^k \mathbf{x}_T + \sum_{s=T-k}^{T-1} \frac{\sqrt{\bar{\alpha}_{T-k}}}{\sqrt{\bar{\alpha}_s}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A})^{T-k-s} \mathbf{z}_{s+1}, \quad (23)$$

where  $\mathbf{z}_s = c_s^0 \epsilon_\theta(\mathbf{x}_s, s) + \sqrt{\bar{\alpha}_{s-1}} \Sigma \mathbf{A}^\dagger \mathbf{y} + c_s^1 \epsilon_s$ , the coefficients are defined as  $c_s^0 := c_s^2 - \sqrt{(1 - \bar{\alpha}_s)/\alpha_s} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A})$ ,  $c_s^1 := \sqrt{1 - \bar{\alpha}_s} \eta$  and  $c_s^2 := \sqrt{1 - \bar{\alpha}_s} \sqrt{1 - \eta^2}$ ,  $0 \leq \eta < 1$ .

**Proof** For the noisy inverse problem, [68] estimates  $\hat{\mathbf{x}}_{0|t}$  as

$$\begin{aligned} \hat{\mathbf{x}}_{0|t} &= \mathbf{x}_{0|t} - \Sigma_t \mathbf{A}^\dagger (\mathbf{A} \mathbf{x}_{0|t} - \mathbf{y}) \\ &= (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A}) \mathbf{x}_{0|t} + \Sigma \mathbf{A}^\dagger \mathbf{y} \end{aligned} \quad (24)$$

$$= \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A}) (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon_\theta(\mathbf{x}_t, t)) + \Sigma \mathbf{A}^\dagger \mathbf{y} \quad (25)$$

$$= \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A}) \mathbf{x}_t - \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A}) \cdot \epsilon_\theta(\mathbf{x}_t, t) + \Sigma \mathbf{A}^\dagger \mathbf{y}. \quad (26)$$

Let  $c_t^0 := c_t^2 - \sqrt{(1 - \bar{\alpha}_t)/\alpha_t} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A})$ ,  $c_t^1 := \sqrt{1 - \bar{\alpha}_t} \eta$  and  $c_t^2 := \sqrt{1 - \bar{\alpha}_t} \sqrt{1 - \eta^2}$ , we define  $\mathbf{z}_t$  with  $\gamma_t = 1$  as

$$\mathbf{z}_t = \left( c_t^2 - \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A}) \right) \epsilon_\theta(\mathbf{x}_t, t) + \sqrt{\bar{\alpha}_{t-1}} \Sigma \mathbf{A}^\dagger \mathbf{y} + c_t^1 \epsilon_t \quad (27)$$

$$:= c_t^0 \epsilon_\theta(\mathbf{x}_t, t) + \sqrt{\bar{\alpha}_{t-1}} \Sigma \mathbf{A}^\dagger \mathbf{y} + c_t^1 \epsilon_t. \quad (28)$$

Similar to Proposition 1, based on DDIM [61],  $\mathbf{x}_{t-1}$ ,  $\mathbf{x}_{t-2}$  and  $\mathbf{x}_{t-3}$  can be calculated by

$$\mathbf{x}_{t-1} = \frac{\sqrt{\bar{\alpha}_{t-1}}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A}) \mathbf{x}_t + \mathbf{z}_t \quad (29)$$

$$= \frac{\sqrt{\bar{\alpha}_{t-1}}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A})^1 \mathbf{x}_t + \sum_{s=t-1}^{t-1} \frac{\sqrt{\bar{\alpha}_{t-1}}}{\sqrt{\bar{\alpha}_s}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A})^{t-1-s} \mathbf{z}_{s+1} \quad (30)$$

$$\mathbf{x}_{t-2} = \frac{\sqrt{\bar{\alpha}_{t-2}}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A})^2 \mathbf{x}_t + \sum_{s=t-2}^{t-1} \frac{\sqrt{\bar{\alpha}_{t-2}}}{\sqrt{\bar{\alpha}_s}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A})^{t-2-s} \mathbf{z}_{s+1} \quad (31)$$

$$\mathbf{x}_{t-3} = \frac{\sqrt{\bar{\alpha}_{t-3}}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A})^3 \mathbf{x}_t + \sum_{s=t-3}^{t-1} \frac{\sqrt{\bar{\alpha}_{t-3}}}{\sqrt{\bar{\alpha}_s}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A})^{t-3-s} \mathbf{z}_{s+1}. \quad (32)$$

According to the induction, we can write

$$\mathbf{x}_{t-k} = \frac{\sqrt{\bar{\alpha}_{t-k}}}{\sqrt{\bar{\alpha}_t}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A})^k \mathbf{x}_t + \sum_{s=t-k}^{t-1} \frac{\sqrt{\bar{\alpha}_{t-k}}}{\sqrt{\bar{\alpha}_s}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A})^{t-k-s} \mathbf{z}_{s+1}. \quad (33)$$

Let  $t$  be  $T$ , the  $\mathbf{x}_{T-k}$  can be estimated by

$$\mathbf{x}_{T-k} = \frac{\sqrt{\bar{\alpha}_{T-k}}}{\sqrt{\bar{\alpha}_T}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A})^k \mathbf{x}_T + \sum_{s=T-k}^{T-1} \frac{\sqrt{\bar{\alpha}_{T-k}}}{\sqrt{\bar{\alpha}_s}} (\mathbf{I} - \Sigma \mathbf{A}^\dagger \mathbf{A})^{T-k-s} \mathbf{z}_{s+1}. \quad (34)$$

We complete the proof of Proposition 2. □

## B. More Implementation Details and Quantitative Results

**More implementation details.** Our method is a zero-shot diffusion model-based IR method and thus does not need training DEQ function and diffusion models. Algorithm 1 is a standard implementation of Anderson acceleration. In Algorithm 1, we sample a noise  $x_T$  from Gaussian distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ , and copy it into multiple copies. The noises are the same and they are the initializations in the fixed point solver, denoted by  $x_{0:T-1}^0$ , as shown in Figure 2. After  $K$  iterations, the sampling chain will converge to  $x_{0:T-1}^*$  where  $x_0^*$  is our desired result. In Algorithm 2, we provide an initialization optimization via DEQ inversion. The initialization is the same as Algorithm 1. For every initialization, we first use  $\text{RootSolve}(\cdot)$  to obtain  $x_0^*$  and its loss. Then we use the DEQ inversion Eqn. (13) to update  $x_T$ . This update process terminates once the gradient norm falls below a default threshold within a sufficiently large step  $S$ . To further improve the performance, we use [18] for SR on ImageNet, and use [84, 86] for colorization, to provide the prior information in the intermediate state.

**Super-Resolution.** In the main paper, we present the SR results for  $\times 2$  and  $\times 4$  scales. In this experiment, we extend to other scales. We compare our method with DPS [20], DiffPIR [87], DDRM [41] and DDNM [68] on ImageNet. Additionally, we use bicubic upscaling as a baseline for SR. The quantitative results are shown in Table B1. For larger scales  $\times 8$  and  $\times 16$ , our method demonstrates significant superiority over most methods across various metrics. Specifically, when compared with the competitive IR method DDNM, our method surpasses it by an LPIPS margin of up to 0.024 and a PSNR margin of up to 1.55 dB. Furthermore, our method requires only 15 sampling steps, compared to other methods.

Methods	$\times 8$ SR			$\times 16$ SR		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
bicubic	21.45	0.512	0.504	19.31	0.431	0.648
DDRM [41]	22.83	0.578	0.444	20.05	0.468	0.577
DPS [20]	18.30	0.382	0.520	16.40	0.315	0.572
DiffPIR [87]	20.40	0.440	0.476	17.72	0.347	0.553
DDNM [68]	22.36	0.558	0.414	20.02	0.459	0.563
<b>DeqIR (Ours)</b>	<b>23.91</b>	<b>0.630</b>	<b>0.390</b>	<b>21.04</b>	<b>0.510</b>	<b>0.546</b>

Table B1. Comparisons of zero-shot **SR** methods on ImageNet.

**Deblurring.** In the main paper, we have provided deblurring results for Gaussian and anisotropic kernels. In this experiment, we further consider a uniform kernel to evaluate the performance of all models. Specifically, we evaluate the zero-shot IR methods, including DDRM [41], DPS [20], DiffPIR [87], and DDNM [68], and employ  $A^\dagger y$  as a baseline. The quantitative results in Table B2 demonstrate that our method outperforms others on all datasets. In comparison with DDNM [68], our method exhibits a PSNR improvement of up to 0.65 dB.

Methods	ImageNet			CelebA-HQ		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
$A^\dagger y$	18.76	0.455	0.574	20.30	0.652	0.446
DDRM [41]	36.97	0.953	0.080	40.72	0.972	0.060
DPS [20]	20.68	0.497	0.414	25.44	0.709	0.260
DiffPIR [87]	22.71	0.494	0.499	28.11	0.777	0.264
DDNM [68]	38.66	0.968	0.050	43.13	0.981	0.044
<b>DeqIR (Ours)</b>	<b>39.31</b>	<b>0.967</b>	<b>0.047</b>	<b>43.32</b>	<b>0.982</b>	<b>0.042</b>

Table B2. Comparisons of zero-shot **deblurring (uniform)** methods.

**Inpainting.** In the main paper, we provided quantitative results for the image inpainting task specifically on CelebA-HQ dataset. In this experiment, we extend our evaluation to ImageNet, comparing our method against state-of-the-art (SOTA) inpainting methods, including DDRM [41] and DDNM [68]. Additionally, we establish  $A^\dagger y$  as a baseline for comparison. Furthermore, we explore various inpainting masks, such as text and stripe masks, and present the corresponding results on ImageNet in Table B3. Notably, our method exhibits significant performance improvement over the diffusion model DDRM [41]. Moreover, with a fixed number of timesteps (25), our method surpasses DDNM [68], albeit falls short compared to DDNM with 100 timesteps, as it incorporates more sampling information.

Methods	Text mask			Stripe mask		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
$A^\dagger y$	14.41	0.659	0.461	8.82	0.838	0.197
DDRM [41]	32.22	0.952	0.052	24.07	0.736	0.358
DDNM-25 [68]	32.71	0.958	0.039	25.58	0.788	0.280
DDNM-100 [68]	<b>33.94</b>	<b>0.964</b>	<b>0.031</b>	27.71	<b>0.838</b>	<b>0.197</b>
<b>DeqIR (Ours)</b>	33.03	0.958	0.046	<b>28.53</b>	0.816	0.213

Table B3. Comparisons of zero-shot **inpainting** methods on ImageNet.

**More evaluation metrics (FID).** In addition to PSNR, SSIM and LPIPS, we also consider FID, which is a commonly used evaluation metric for assessing the quality of generated images. Although FID is a useful metric in many real-world scenarios, there are certain circumstances where it may not be appropriate. In particular, I do not use FID to evaluate the image quality for super-resolution and deblurring in the main paper because the ImageNet and CelebA-HQ datasets with 100 classes are relatively small, and FID requires a sufficiently large dataset to accurately estimate the statistics of the data distribution. If the dataset is small, the FID score may not be reliable due to high variance in the estimated statistics. In contrast, we use FID in the colorization task since we can manipulate inputs using various Gaussian noises by altering seeds. Nevertheless, we retain the ability to compute FID scores for super-resolution and deblurring tasks. In general, lower FID indicates higher quality of images. As shown in Table B4, our method has superior FID performance compared to other approaches. Furthermore, we observe a similar trend between FID and LPIPS in our comparisons, as illustrated in Table 1 in the main paper. Hence, we can confidently rely on these FID results for assessing super-resolution and deblurring tasks.

Methods	Bicubic SR		Deblurring	
	2×	4×	Gaussian	anisotropic
Baseline	48.94	134.11	117.87	182.49
DGP [55]	172.72	273.65	231.00	267.71
DPS [20]	142.31	166.36	121.60	119.37
DiffPIR [87]	43.19	96.28	73.50	145.52
DDRM [41]	26.56	94.92	6.43	10.82
DDNM [68]	20.87	86.29	2.57	5.61
<b>DeqIR (Ours)</b>	<b>13.89</b>	<b>59.70</b>	<b>2.41</b>	<b>5.40</b>

Table B4. FID (↓) results of SR and deblurring tasks on ImageNet.

**Application of Proposition 2.** Based on Proposition 2, we apply our method to the task of restoring noisy images. Specifically, our approach addresses noisy super-resolution for scaling factors of 2×, 4×, and 8×, denoted as Ours-σ. In our experiments, we set the noise level to σ = 0.2. As shown in Table B5, our method consistently achieves higher PSNR values compared to DDNM [68] across various scaling factors. Visual comparisons are provided in Figure B1, illustrating our method’s capability to reduce noise and partially restore textures, leveraging the derived sampling formulation in Proposition 2.

Scale	DDNM [68]	Ours-σ
2×SR	24.91	<b>25.23</b>
4×SR	22.13	<b>22.40</b>
8×SR	19.70	<b>19.95</b>

Table B5. PSNR for noisy SR

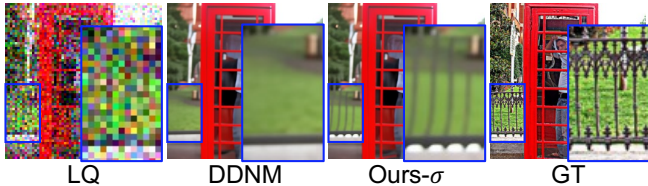


Figure B1. Results for noisy 4×SR.

**GPU VRAM usage.** In addition to comparing running times as presented in the main paper, we also analyze the running time, GPU memory usage, and PSNR values of our method to understand the trade-offs between image quality and computational costs. We conducted experiments with different timesteps for anisotropic deblurring on the ImageNe dataset, as shown in Table B6. As timesteps increase, image quality improves, albeit at the expense of increased inference time and GPU memory usage. Notably, for a timestep of 20, our method achieves the highest PSNR value. However, it also exhibits the longest inference time and requires the largest GPU memory allocation. Given constraints on available memory, one may choose a timestep of 15. This choice allows for the generation of images with a comparable PSNR to that of a timestep of 20, while having more efficient inference times.

Methods	Ours-10	Ours-15	Ours-20
Time (s)	12.11	16.53	21.19
Memory (G)	12.07	16.32	18.06
PSNR (dB)	38.58	39.21	39.47

Table B6. Comparisons of time, memory and PSNR.

## C. More Qualitative Results

### C.1. More Results on Super-Resolution

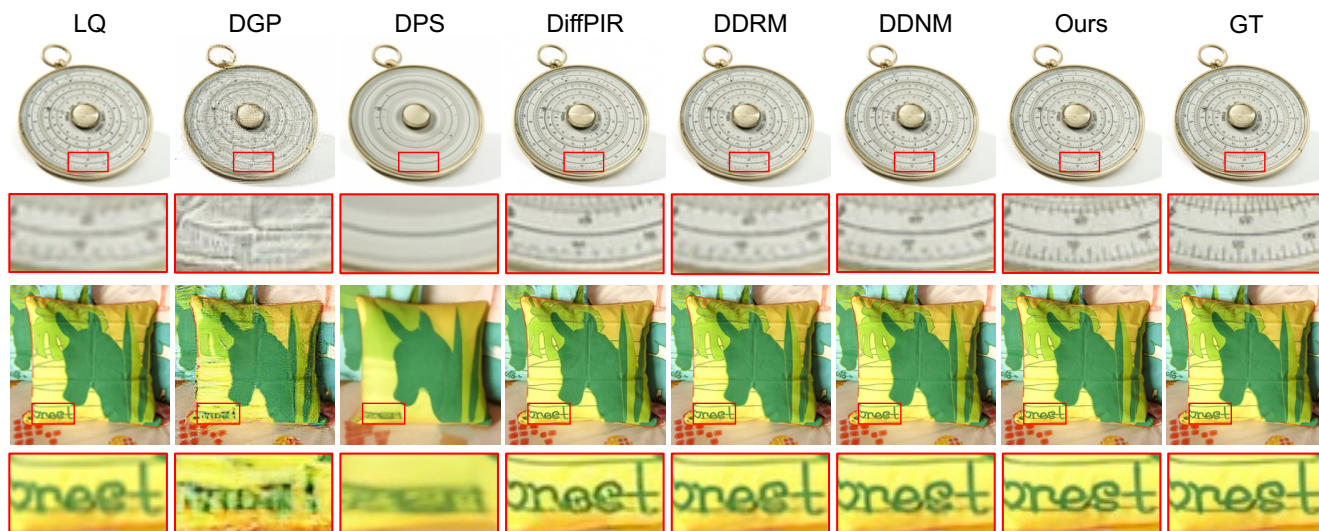


Figure C2. Qualitative results of image super-resolution ( $\times 2$ ) methods on ImageNet.

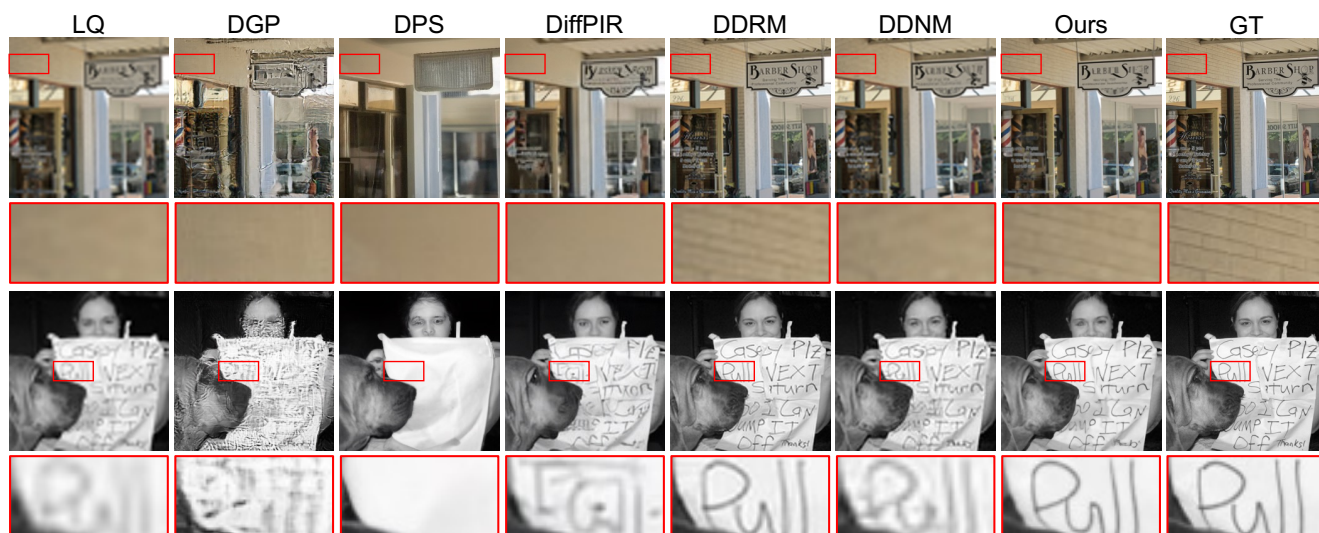


Figure C3. Qualitative results of image super-resolution ( $\times 4$ ) methods on ImageNet.

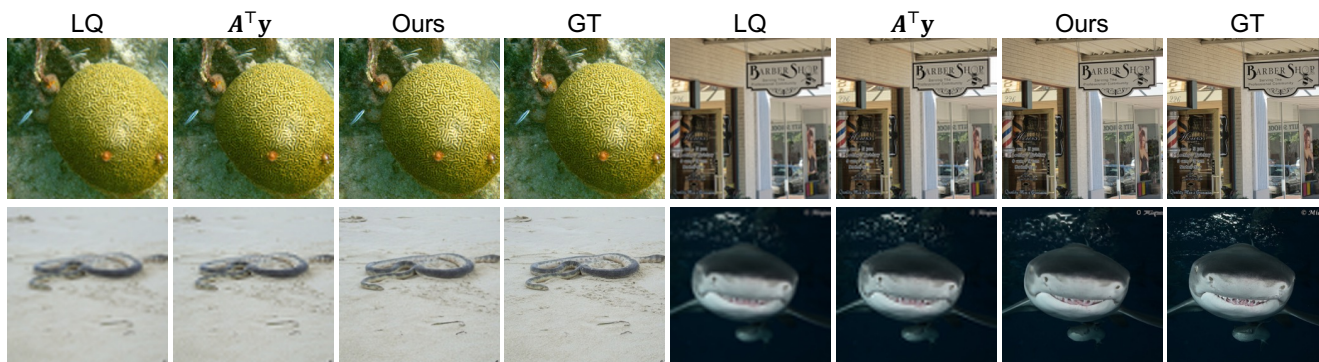


Figure C4. Our qualitative results of image super-resolution ( $\times 2$  (above) and  $\times 4$  (bottom)) on ImageNet.

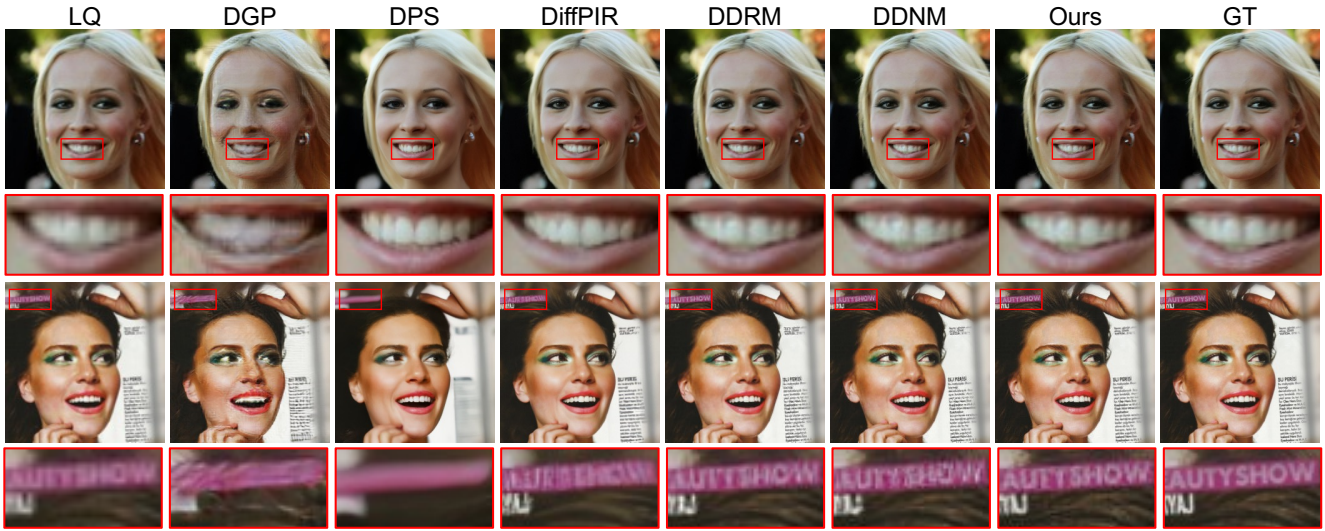


Figure C5. Qualitative results of image super-resolution ( $\times 2$ ) methods on CelebA-HQ.

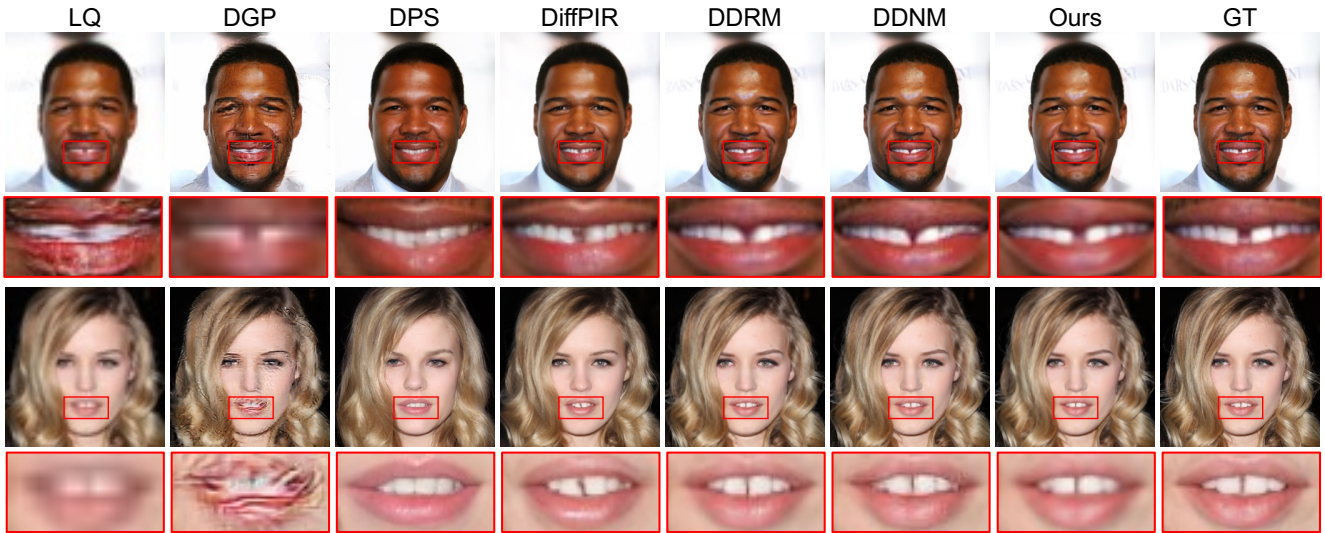


Figure C6. Qualitative results of image super-resolution ( $\times 4$ ) methods on CelebA-HQ.

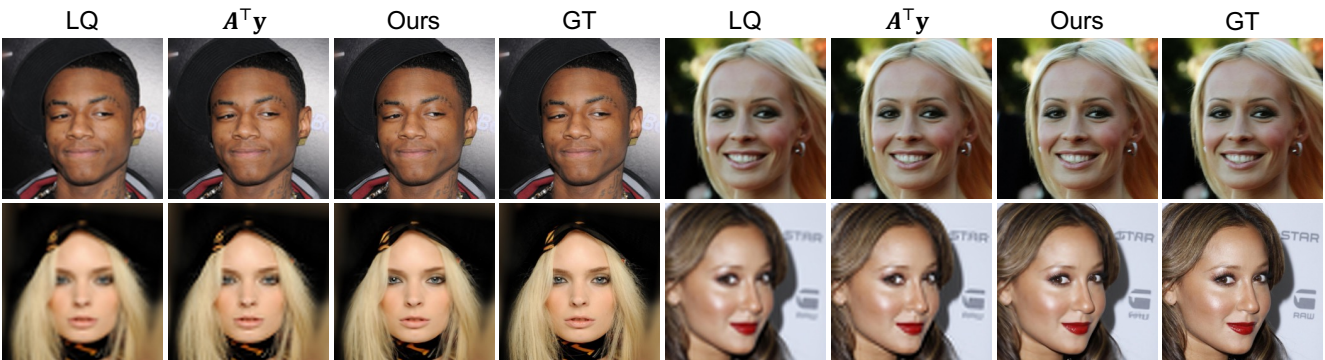


Figure C7. Our qualitative results of image super-resolution ( $\times 2$  (above) and  $\times 4$  (bottom)) on CelebA-HQ.

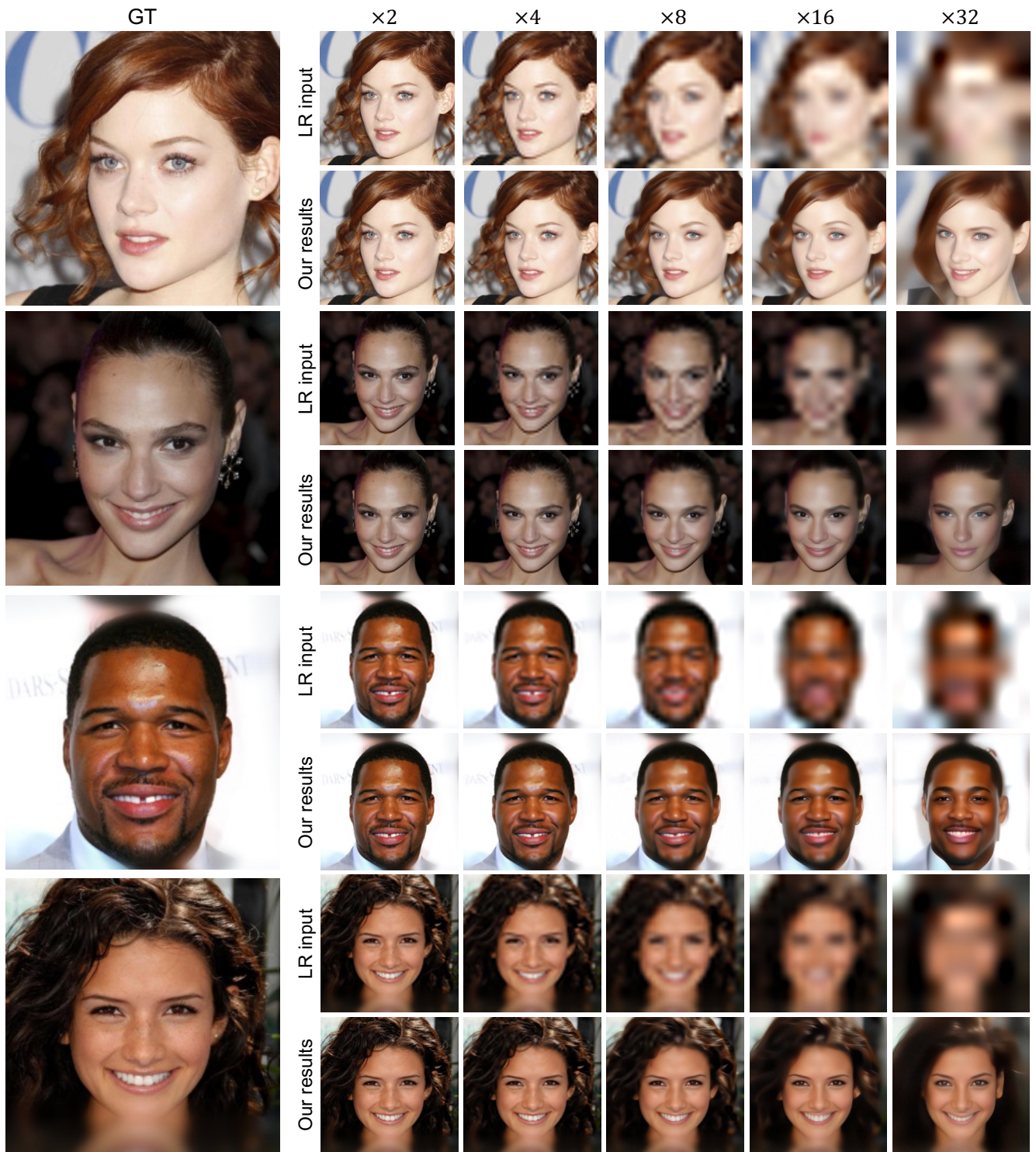


Figure C8. Our qualitative results of image super-resolution for different scales on CelebA-HQ.



## C.2. More Results on Image Deblurring



Figure C9. Our qualitative results of image deblurring (gauss) on ImageNet.

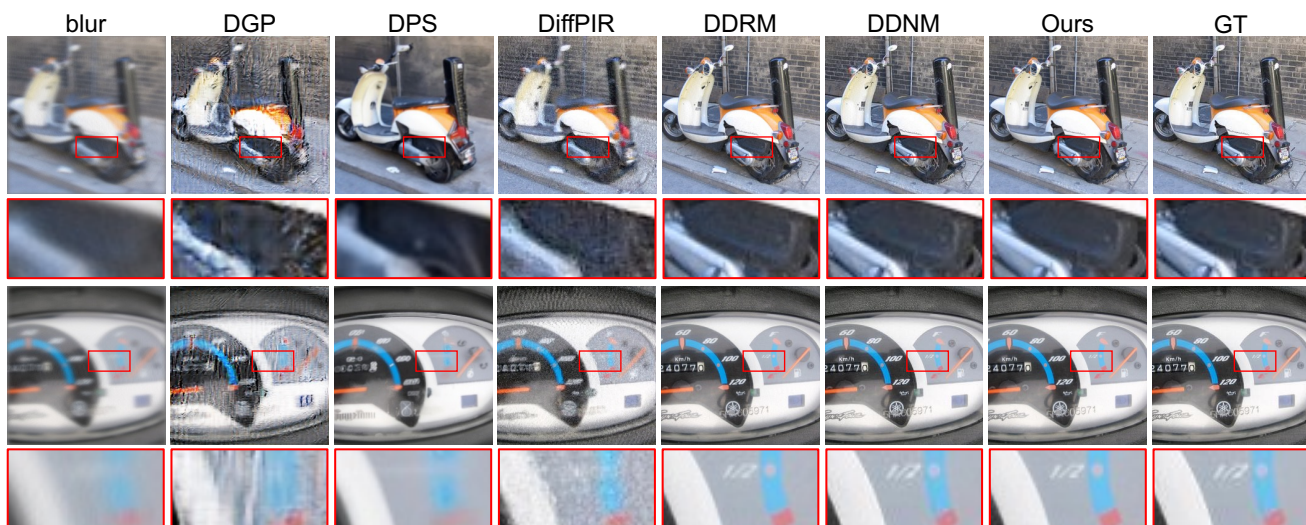


Figure C10. Our qualitative results of image deblurring (anisotropic) on ImageNet.

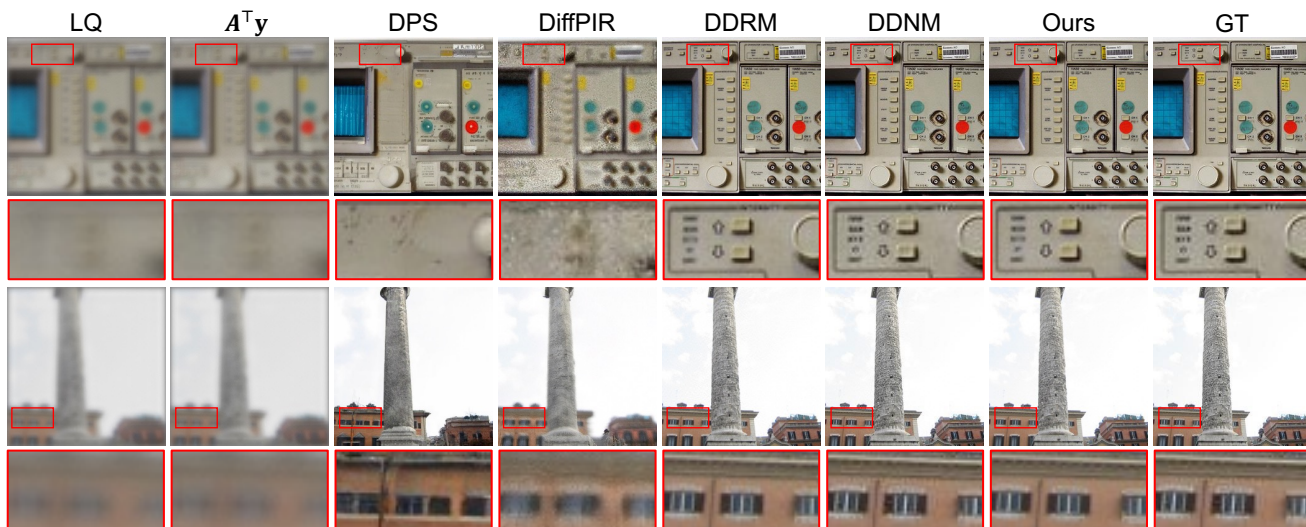


Figure C11. Our qualitative results of image deblurring (uniform) on ImageNet.

### C.3. More Results on Image Inpainting

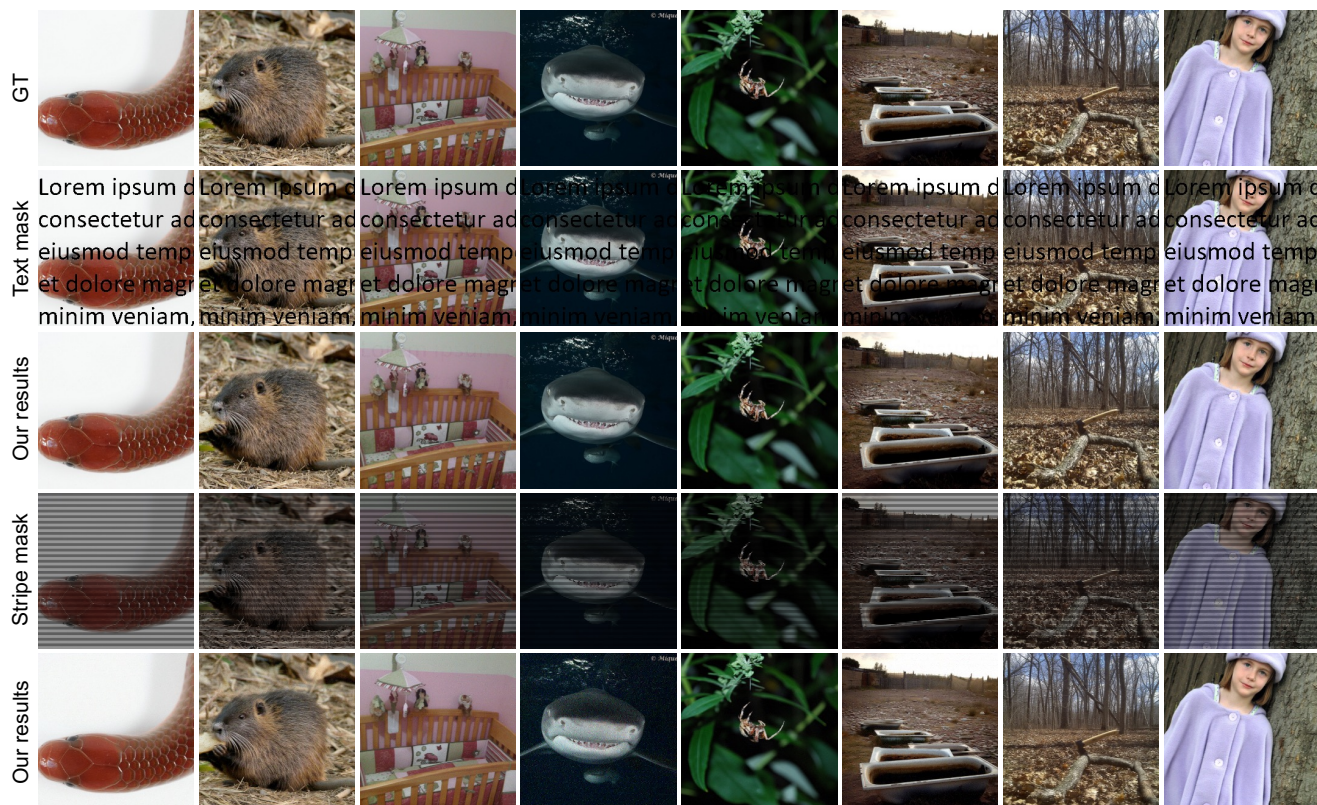


Figure C12. Our qualitative results of image inpainting on ImageNet.



Figure C13. Our qualitative results of image inpainting on CelebA-HQ.

### C.4. More Results on Image Colorization

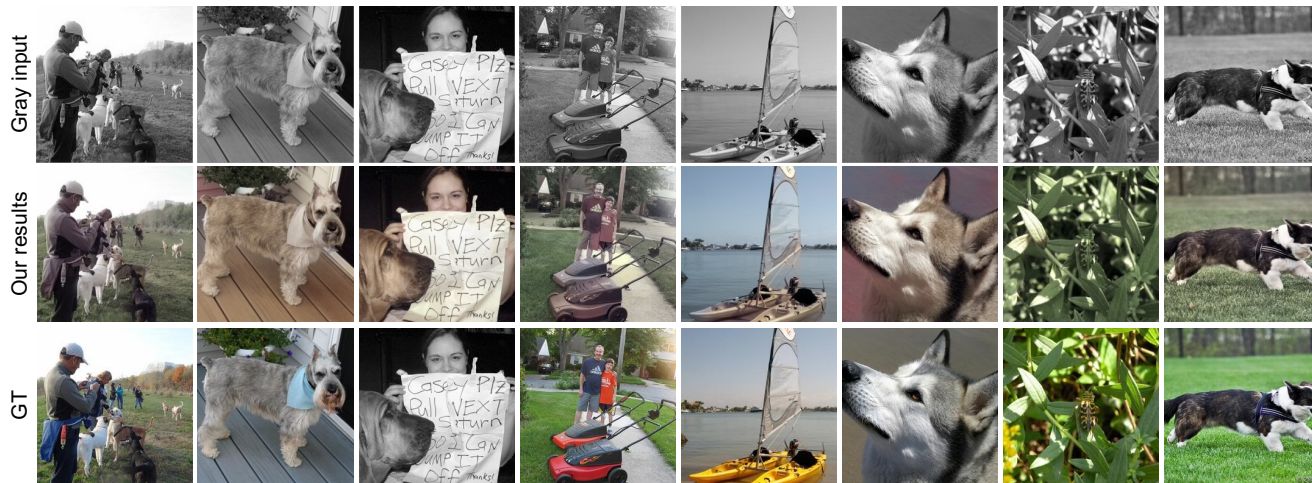


Figure C14. Our qualitative results of image colorization on ImageNet.



Figure C15. Our qualitative results of image colorization on CelebA-HQ.

### C.5. More Results on Old Photo Restoration

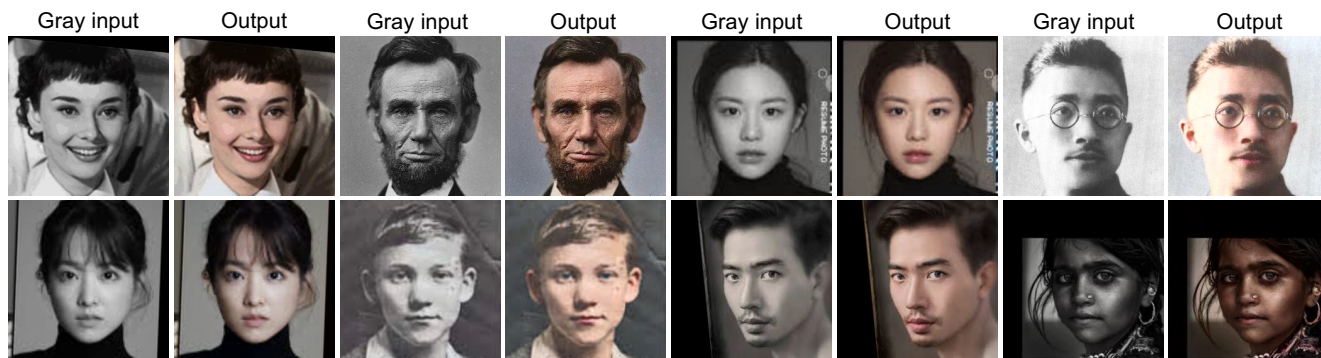


Figure C16. Our qualitative results of image colorization on real-world images.

### C.6. Results on Other Image Restoration Tasks

In the main paper, we have conducted some typical image restoration (IR) tasks, including super-resolution, image deblurring, image inpainting and colorization. For other IR tasks, we show that our method can be used in the compressed sensing task. Specifically, we use the Walsh-Hadamard sampling matrix with a 0.5 compression ratio. We show the visual results on ImageNet and CelebA-HQ in Figures C17 and C18. As we can see, with the severely compressed inputs, our method is able to recover the high-frequency details while preserving the inherent identity of the content.



Figure C17. Our qualitative results of compressed sensing on ImageNet.

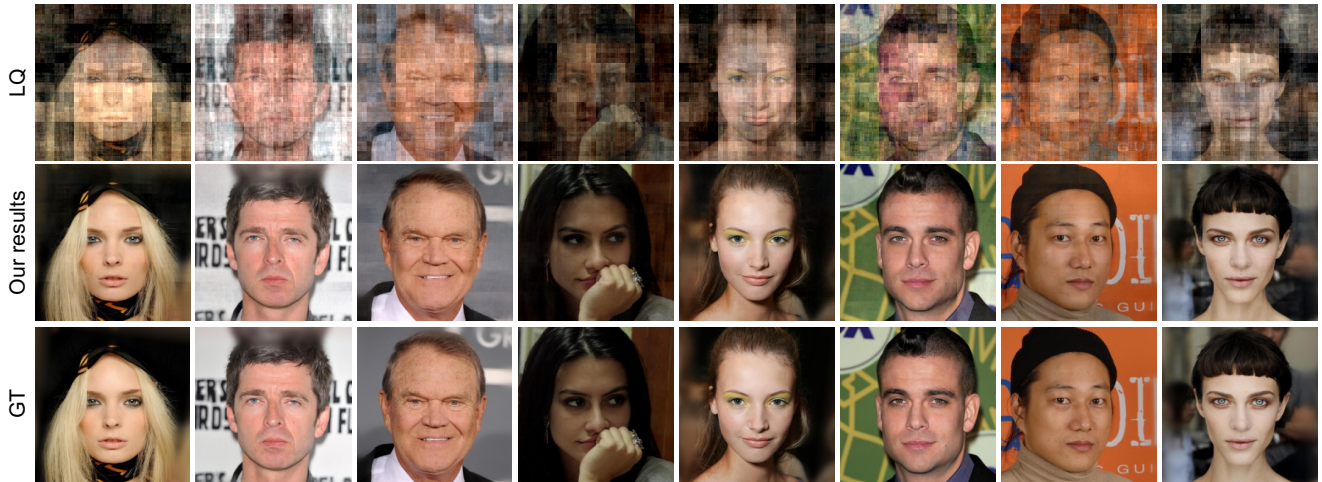


Figure C18. Our qualitative results of compressed sensing on CelebA-HQ.

### C.7. Arbitrary Size Image Restoration

In our primary experiments, the input size is set at  $3 \times 256 \times 256$ . However, for real-world applications, the input sizes may vary. To address this variability, we demonstrate the adaptability of our model to accommodate inputs of arbitrary sizes. As an illustration, we use an input size of  $3 \times 256 \times 1024$ , a methodology that can be extended to accommodate diverse sizes. In line with the approaches in [48, 68], we segment larger images into multiple overlapping patches and conduct individual tests on each. Our method differs from [68] as our overlapped patches are smaller, resulting in reduced computational costs. Finally, we consolidate the generated outputs to form conclusive results.



Figure C19. Our qualitative results of arbitrary-size image restoration on real-world images.

## D. More Results on Initialization Optimization

In the super-resolution, we set  $S$  as 500 for convergence, and we use the  $\ell_2$  loss to guide the generation process. In the loss, we use the classical IR model (e.g., SwinIR [48]) as the supervision information such that the generated images can be close to the given supervision information. In the colorization, we set  $S$  as 2k for convergence, and we use the perceptual loss to guide the generation process. The results are shown in Figures D20 and D21. With the initialization optimization, PSNR can be further improved, and the colorization has guidance to generate according to the reference images.

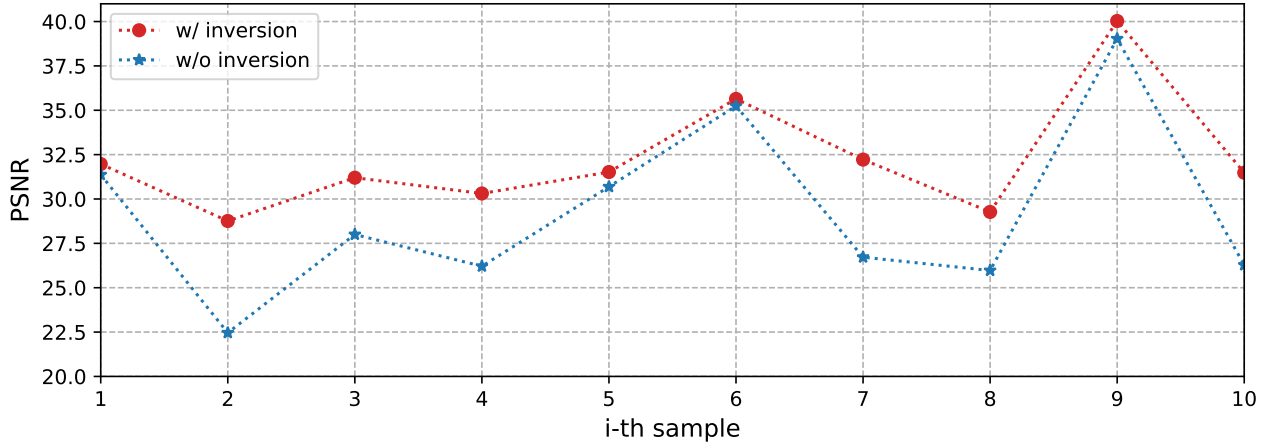


Figure D20. PSNR improvement of using initialization optimization on ImageNet (we show the first 10 samples).

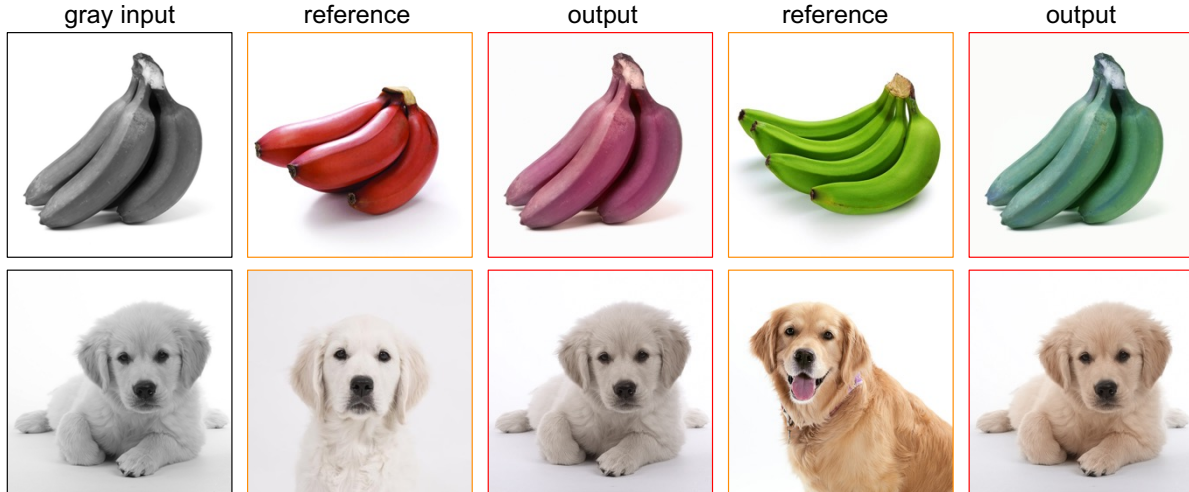


Figure D21. Our qualitative results of colorization with initialization optimization.

## E. Limitations and Future Work

There remain many limitations that can be studied in the future.

- Our method processes multiple timestep images in parallel, leading to large memory in GPU. It is necessary to reduce the memory and boost the inference time.
- Our parallel sampling requires explicit forms of the degradation matrix  $A$  which is linear. For unknown, complex and non-linear, we need to design pseudo-inverse by hand. This way requires multiple attempts and is cumbersome. In addition, one can approximate the degradation matrix  $A$  by training a network on constructed pair data.
- Our initialization optimization via DEQ inversion needs many iterations for convergence. How to reduce the iterations and accelerate remains a future work.
- Our method is a new zero-shot image restoration method. Moreover, the IR performance depends on the pre-trained denoiser. The performance is not better than some supervised learning methods on a specific task. With the help of our model inversion, one can train our DeqIR model to improve the performance of IR.