

# SAM-DiffSR : Structure-Modulated Diffusion Model for Image Super-Resolution

## ソース

- <https://arxiv.org/pdf/2402.17133>
- <https://github.com/lose4578/SAM-DiffSR>

## 概要

- SAM(Segment Anythin Model)をdiffusionに取り入れてSuper Resolution
- 低画質画像のsegmentation maskを学習時(拡散過程)のノイズに加える

## 提案手法

- $\mathbf{x}_0$ をSAMに入力して得られるsegmentation maskを $\mathbf{M}_{\text{SAM}}$ とする
- $\mathbf{M}_{\text{SAM}}$ に対してencoder, embeddingを加えたものを $\mathbf{E}_{\text{SAM}}$ とする

## diffusion modelの変更

拡散過程を以下のように定義する

$$q(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \mathbf{E}_{\text{SAM}}) = \mathcal{N}\left(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1} + \sqrt{\beta_t}\mathbf{E}_{\text{SAM}}, \beta_t \mathbf{I}\right) \quad (1)$$

このとき、 $q(\mathbf{x}_t \mid \mathbf{x}_0, \mathbf{E}_{\text{SAM}})$ を計算すると

$$q(\mathbf{x}_t \mid \mathbf{x}_0, \mathbf{E}_{\text{SAM}}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \phi_t \mathbf{E}_{\text{SAM}}, (1 - \bar{\alpha}_t) \mathbf{I}) \quad (2)$$

$$\alpha_t = \beta_t \quad (3)$$

$$\bar{\alpha}_t = \prod_{i=1}^t \alpha_i \quad (4)$$

$$\phi_t = \sum_{i=1}^t \sqrt{\bar{\alpha}_t \frac{\beta_i}{\bar{\alpha}_i}} \quad (5)$$

さらに、 $p(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0, \mathbf{E}_{\text{SAM}})$ については以下の式になる

$$\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0, \mathbf{E}_{\text{SAM}}) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \left( \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\beta_t}} \mathbf{E}_{\text{SAM}} + \epsilon \right) \right) \quad (6)$$

$$\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t \quad (7)$$

$$p(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0, \mathbf{E}_{\text{SAM}}) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0, \mathbf{E}_{\text{SAM}}), \tilde{\beta}_t \mathbf{I}) \quad (8)$$

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (9)$$

そこで、損失関数を以下で定義する

$$\mathcal{L}(\theta) = \mathbb{E}_{t, \mathbf{x}_0, \epsilon} \left[ \left\| \frac{\sqrt{1 - \bar{\alpha}}}{\sqrt{\beta_t}} \mathbf{E}_{\text{SAM}} + \epsilon - \epsilon_\theta(\mathbf{x}_t, t) \right\|_2^2 \right]$$

## segmentation maskに対するencoding, embedding

画像 $\mathbf{x} (C \times H \times W)$ に対してposition embedding(RoPE)により $\mathbf{x}_{\text{RoPE}} \in \mathbb{R}^{1 \times H \times W}$ を得る。

一方で $\mathbf{x} (C \times H \times W)$ に対してSAMにより $K$ 個のsegmentation mask  $M_{\text{SAM}, i} \in \{0, 1\}^{1 \times H \times W}$ を得る。 $(i = 1, 2, \dots, K)$

これらに対して以下の計算式でStructural position encodingを行い $\mathbf{E}_{\text{SAM}}$ を得る

$$\mathbf{E}_{\text{SAM}} = \sum_i M_{\text{SAM}, i} \cdot \text{mean}(\mathbf{x}_{\text{RoPE}, i}) \quad (10)$$