

Paper summary

AI VISION Lab

1. 공부한 논문의 제목, 게재된 학회 혹은 저널 등 논문 기본 정보를 적으세요.
 - A. 이름: **SR3: Image Super-Resolution via Iterative Refinement**
 - B. 저널: **IEEE TPAMI**
 - C. 도메인: **Super Resolution**
 - D. 출판연도: **2021**
 - E. 저자: **Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J. Fleet, Mohammad Norouzi (Google Research, Brain Team)**
2. 논문에서 제안한 알고리즘 및 프레임워크에 대해 본인이 이해한대로 다이어그램을 그려보세요. 논문 **Figure**를 그대로 따라 그리면 안됩니다.
 - A. 일반적인 구조는 Diffusion(DDPM)과 동일함 (UNet사용, 충분한 T Size, step 당 Noise 크기는 최소화)
 - B. 다만 차이점은 다음과 같음. 우선, Low Resolution Image x 에 대해서 Resample를 진행함. (목표는: High Resolution Image와 동일한 크기의 이미지로 맞춰주는 것. Bicubic Interpolation -> 해당 Interpolation 기법의 경우 다른 SR에서도 자주 사용됨.)
 - C. 이때 Resampled LR 과 y_t (이때 y_t 의 경우- Reversed process -> Gaussian Noise인 y_T 에서 Iterative하게 진행 -> Iterative Refinement)을 Concat (Concat말고 다른 방식도 시도해봤으나 Concat대비 성능이 나오지 않았다고 함- FiLM)
 - D. 이후 학습 진행 (Diffusion과 동일, 이때 모델의 목표는 $p(x)$ 가 아닌 $p(y|x)$ 라는 점을 반영하여, x 를 condition으로 주입하는 구조임. x : source | y = target) 마치 x_0 를 condition으로 주는 DDIM의 구조와 상당히 유사함. 다만 Forward 과정을 x_0 conditioning하게 진행하지 않았다는 점 -> 여전히 Markovian하다고 볼 수 있음.)
 - E. 추가적으로 Cascaded에 대해서도 다뤘는데, 즉각 High resolution으로 진행하는 것보다, 64 -> 256 -> 1024(Model 두 번)로 진행하는게 비용, 성능 측

면에서 장점이 있음을 보였음.

3. 본인이 생각하는 이 논문의 장점이 무엇이라고 생각하나요? **논문 Contribution bullet을 그대로 따라 적으면 안됩니다.**
 - A. Condition, Uncondition, Face, Natural 다양한 조건에서 실험을 진행했고 성능을 입증했다는 점이 가장 큰 장점이 아닐까 생각됨.
4. 이 논문을 읽으면서 느낀 점, 혹은 배운 점이 있으면 적어보세요.
 - A. 실제 사람의 판단을 기반으로 평가하는 점: (2AFC)
 - B. SR의 목표 Distribution. 그리고 Regression(SR model)과의 차이점
5. 이 논문의 한계점이 있다면 무엇이라고 생각하나요?
 - A. Step size의 과도한 크기 ($T = 2000$)
6. 본인의 연구에 접목시켜볼 점이 있을지 생각하고 적어보세요.
 - A. Diffusion, SR 선행연구
7. 본 Summary를 작성하는 과정에서 생성형AI를 사용했나요?
 - A. 아니요

날짜: 2025-08-12

이름: 신준원

Image Super-Resolution via Iterative Refinement

용어 정리

x : Source Image

y : target Image (+ noise 제거)

model 목표: $P(Y|x)$ 찾기

과정 (Reverse process): $P(x_{t-1} | x_t, x)$

Conditional DDPM 과정

direction

(x_0, x, y)

SR3 Architecture

Similar to DDPM.

① U-Net Base

② ResBlock (from BiG GAN)

③ Re-scale skip Connection $\Rightarrow \frac{1}{\sqrt{2}}$

④ Block # 증가

⑤ bicubic interpolation

\Rightarrow low resolution image \Rightarrow target과 동일한 resolution까지 interpolation

\rightarrow low resolution (R size는 target과 동일) + target

(channel concat) \Rightarrow ex) $B, C, H, W \rightarrow B, 2C, H, W$

($\frac{1}{\sqrt{2}} \Rightarrow$ Film for Conditioning. 다른 Concat Best)

⑥ $T=2000$

⑦ Variance (Uniform sampling) \rightarrow speech.

* 설명.

Regression based of H.

$$y = Hx + n$$

\downarrow low \downarrow high \downarrow noise

Down-sampling

ex) bicubic interpolation

$$F(x) = x$$

$$\rightarrow F(x) \text{ fit}$$

x or y fit $f \Rightarrow$ loss: MSE

Regression ex) SRCNN.

Generative based? Regression

\rightarrow $\int d p(d|x) d d$ distribution

$$\int d p(d|x) d d$$

high, low

$$\Rightarrow \begin{bmatrix} \lambda_{pm} \\ \text{posterior mean} \end{bmatrix}$$

\rightarrow regression based

IBM Generative

\rightarrow posterior

fit

Experiments

Dataset

① Flickr-Faces-HQ (training) \rightarrow for Face (uncondition)
② CelebA - HQ (evaluation)

③ ImageNet 1k \rightarrow for natural image super-resolution (condition)

Step) Image \rightarrow biCubic interpolation.
(with anti-aliasing)

Step) $64 \times 64 \rightarrow \{256 \times 256, 512 \times 512\}$

$16 \times 16 \rightarrow \{28 \times 28\}$

$256 \times 256 \rightarrow \{1024 \times 1024\}$

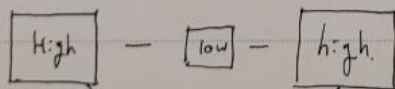
Result

기준 MOS (mean opinion Score)

평가 \rightarrow resolution 반영 X / 풀 크기 상에 한해서 평가

SR과 관련 사항

step1)



1) 어떤 high가 더 나은가?
(low에러 생김)

step2)



2) Camera가 더 나은 이미지는 둘중 어떤 것일까?

(구분 안함)

fail rate

select model's
선택

→ 50 % 42

→ 왜

low 확률

→ 표준 42이 높을지

제약이 후속 → image 2장 인식

한번이 더 쉬움

↓
image 2장 → 40 x

Cascaded generation

상황 → 곧바로 Generate 할 수 있음

예) 64^v (input) → model 1 (96^v) 256^v → model 2 (96^v) 1024^v