

## Paper summary

AI VISION Lab

1. 공부한 논문의 제목, 게재된 학회 혹은 저널 등 논문 기본 정보를 적으세요.
  - A. 이름: **Palette: Image-to-Image Diffusion Models**
  - B. 저널: **SIGGRAPH**
  - C. 도메인: **i2i, Diffusion**
  - D. 출판연도: **2022**
  - E. 저자: **Google Brain Team**
2. 논문에서 제안한 알고리즘 및 프레임워크에 대해 본인이 이해한대로 다이어그램을 그려보세요. **논문 Figure를 그대로 따라 그리면 안됩니다.**
  - A. 구조가 특별하기 보다는 하나의 모델의 구조에서 다양한 Task(Image Inpainting, Uncropping, Colorization, JPEG Restoration)에 모두 적용할 수 있는 구조를 사용 (ADM의 U-Net Architecture를 가져옴, 256\*256)
  - B. 기존 모델들은 Unconditional한 Diffusion  $p(x)$ 를 학습하는 것을 목표로 했으나, 일반적인 경우 Unconditional Diffusion Model이 존재하는 모든 데이터셋을 학습하여 다양한 Task의 결과를 얻어내는 과정은 상당히 어렵기 때문에,  $p(x)$ 가 아닌,  $p(y|x)$ 를 목표로 함.
  - C. 이때, Condition을 주입하는 방법은 Concatenation.
  - D. 결과적으로 다양한 Task의 condition을 줬을 때, condition을 주입하지 않고 학습한 기존 모델들과 달리 다양한 Task(Task마다 다른 Condition주입)에서 높은 성능을 보임.
3. 본인이 생각하는 이 논문의 장점이 무엇이라고 생각하나요? **논문 Contribution bullet를 그대로 따라 적으면 안됩니다.**
  - A. L1 vs L2에 대한 실험
  - B. 다양한 option에 따른 모델 결과 (Unconditional vs Conditional)
  - C. 양적, 질적 관점에서 모두 고려한 결과물

- D. 다양한 Low-level Task측면에서 여전히 고려되고 있는 아키텍처 구조
4. 이 논문을 읽으면서 느낀 점, 혹은 배운 점이 있으면 적어보세요.
- A. 다양한 실험 – metrics 기반으로 진행되는 과정에서 어떤 방식으로 Cloud Removal 실험을 진행하면 좋을지 생각할 수 있는 논문
- B. 또한, 다양한 i2i (i2sb, DBIM)모델들이 해당 논문의 구조(ADM -> Palette로 발전되는 모델 구조)를 기반으로 모델을 작성하고 있기 때문에 모델 설계를 하는 측면에서도 도움이 되었다.
5. 이 논문의 한계점이 있다면 무엇이라고 생각하나요?
- A. Sampling 속도
6. 본인의 연구에 접목시켜볼 점이 있을지 생각하고 적어보세요.
- A. L1 vs L2에 대한 직접적인 성능비교 (L1의 경우 – 다양성 측면에서 낮은 점수를 보였으나 강건한 sample을 보여줌): 실제 Target을 가진 상태에서 Restoration을 목표로 하는 Cloud Removal을 하는 경우 L1을 사용하는 것이 Uncertainty를 상승시키는데 더 효과적이라 판단됨
- B. U-Net의 Attention Layer 성능 비교 실험을 통해 확인. (Global Attention의 결과로 blur한 형태의 output이 나온다고 판단. 그러나 여전히 Global Attention이 Replaced Layer(3x3, 5x5, 7x7 – Receptive Field 확장, Local Attention – head 개수:4개) 보다 성능이 좋게 나오고 있다는 점. 지금 사용하는 모델에서 – Attention을 대체하기 보다는 계산적으로 light한 layer를 사용하는 방향성 고려.
7. 본 Summary를 작성하는 과정에서 생성형AI를 사용했나요?
- A. 아니요

날짜: 2025-08-21

이름: 신준원

## palette: Image-to-Image Diffusion models

### Introduction

Image-to-Image translation  $\Rightarrow$  ex) inpainting, JPEG restoration 등...  
single input & multiple output  $\Rightarrow$  여러 task  
 $\hookrightarrow$  input 조건, output 분포 (Conditioned) 를 학습하는 형태  
 $\hookrightarrow$  Output의 true 분포 학습 = 목적.  
 $\hookrightarrow$  tractable하게 Condition 주입.

해당 문제  $\Rightarrow$  Diffusion이 Image-to-Image GAN을 학습하는 성능을 보인 것을 바탕으로.

또한, 해당 Task에 맞게  $\Rightarrow$  GAN 방식과 다르게  $\hookrightarrow$  다양성을 보장하는 것.

$\Rightarrow$  ① Colorization / ② un-cropping / ③ inpainting / ④ JPEG restoration 같은  $\Rightarrow$  관련 작업

$\hookrightarrow$  palette model  $\Rightarrow$  특정한 Loss, Condition 하에서 여러 task 성능 구하기!

$\hookrightarrow$  기본 Diffusion model

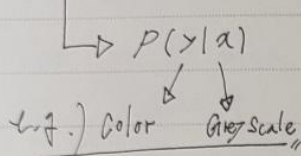
$\hookrightarrow$  train unconditional model. (160/256)

$\Rightarrow$  use in Conditional Task

결과, Conditional 한 train  $\Rightarrow$  for multiple Task

## Palette

Condition diffusion model  $\gamma \in [0, 1]$  (denoising process. Conditional on an input signal)



$\rightarrow 7$  or  $5$  second.

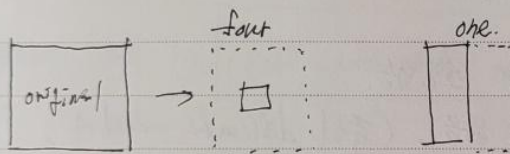
## Appendix) Detail

① human evaluation (Image-to-Image !)

Colorization : grayscale(x), RGB(y) + 6. (random)

inpainting : rectangular mask (10 ~ 40 %) with Lample.

Unwrapping : four direction or one direction



JPEG information : (5, 70)

## Learning process (Denoising)

output :  $y$     noisy :  $\tilde{y}$

$$\|f_{\theta}(x, \sqrt{\gamma}\tilde{y} + \sqrt{1-\gamma}\epsilon, \gamma) - y\|_p^p$$

$\Delta$  for restoration

$\gamma$  : noise indicator

$p=1 \Rightarrow L_1$

$p=2 \Rightarrow L_2$

$\rightarrow L_1$   
hallucination 방지  
diversity 증가

$L_2$   
diversity 증가

## Architecture

based on  $256 \times 256$  Alex

+ ① class condition  $1 \times 1 \times 1$

② Concat을 통한 Conditioning  $1 \times 1$

## Evaluation

12:  $\Rightarrow$  evaluation of  $\forall$  and  $\exists$

- ① FID
- ② human evaluation (3s, 4s)
- ③ Qualitative evaluation
- ④ pGMR, SSIM (reference based score)
- ⑤ Benchmark dataset of 474.

평가 book을  $\Rightarrow$  Imagenet 기반, task (4m) 이 위한 evaluation 이 가능한 protocol 이란  
 $\hookrightarrow$  pixel-level evaluation (예: PSNR, SSIM) 이다.

→ PSNR, SSIM = blur Image에 대한 지표

2.5% human standard at  $\frac{1}{2}$  |  $\frac{1}{2}$ . ( $\frac{1}{2}$  | difficult task A)

→ 4 types of metrics IS, FID, CA, PD  
 ↳ (with test)  
 ↳ (Perceptual distance)

Feature space 상 유클리드 거리

sample diversity  $\frac{1}{2}$  : SSZM & Lpips  $\frac{1}{2}$

human evaluation:  $RAF \rightarrow RSATs$  ( $\Rightarrow$  fool rate)





