

최근 Computer Vision 분야에서 Transformer가 각광받고 있다.

그러면 어려운 task인 GAN에서도 될까?

그래서 완전히 convolution을 사용하지 않은 GAN을 만들어 봤다. (오직 순수한 Transformer만 사용)
TransGAN & memory-friendly transformer를 기반으로 된

생성자이다. 임베딩 차원이 감소하는 동안 feature의 resolution은 점진적으로 증가

패치 수준의 discriminator도 transformer 기반.

Data Augmentation에서 기록된 GAN보다 좋은 이점을 보였다.

+ 생성자에는 multi-task co-training 전략

+ 이미지의 주변영역을 강조하는 self-attention

또, 고화질에 대해서도 훌륭하게 작동.

SOTA에 있는 convolution 기반의 GAN이랑 비교했을 때에도 꽤 좋은 성능을 보였다.

STL-10 dataset에서 SOTA 달성

Introduction

GAN은 성공, 이미지 생성도 성공, 이미지 edit도 성공

그러나 GAN의 훈련 불안정이 문제여서 정규화, 더 좋은 loss function, 더 나은 훈련 방법 연구 진행

그러고 다양한 GAN에 대한 연구...

그러다가 NAS를 GAN에 도입하고 backbone 설계도 GAN을 특가로 개선하는게 중요해 보임

하지만 이런 연구들은 다 convolution 기반의 backbone을 사용해서 한 라는 상식으로 보겠지.

그럼 무진 컨볼루션 아예 안 쓸 강한 GAN을 만들 수 있을까?

- Transformer의 강점

1. 인간 저법의 편향성이 적고 표현력 풍부

2. ad-hoc 백업 복록 제거 가능

- **Model Architecture:** We build the first GAN using purely transformers and no convolution. To avoid overwhelming memory overheads, we create a memory-friendly generator and a patch-level discriminator, both transformer-based without bells and whistles. TransGAN can be effectively scaled up to larger models.
- **Training Technique:** We study a number of techniques to train TransGAN better, ranging from data augmentation, multi-task co-training for generator with self-supervised auxiliary loss, and localized initialization for self-attention. Extensive ablation studies, discussions, and insights are presented. None of them requires any architecture change.
- **Performance:** TransGAN achieves highly competitive performance compared to current state-of-the-art CNN-based GANs. Specifically, it sets new state-of-the-art IS score of 10.10 and FID score of 25.32 on STL-10 and also reaches competitive 8.63 IS score and 11.89 FID score on CIFAR-10, and 12.23 FID score on CelebA 64×64 , respectively. We also summarize the current limitations and future potential of TransGAN.

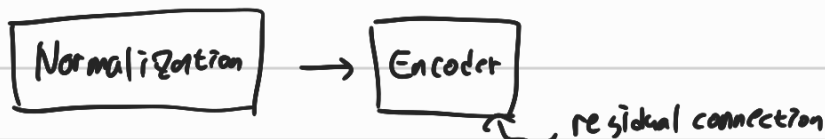
← Contribution

Method

- Basic Block

: Transformer encoder 를 basic block 으로 사용

encoder 구성 { multi-head self-attention module
MLP (feed forward) with GELU



- Memory friendly generator

점진적으로 세분화하는 방법

Input: random noise

→ MLP 통과시킴 → HxWxC vector 만들림 → feature map은 학습 가능한 positional encoding과 결합

비트와 비슷하게 transformer encoder는 토큰도 토큰 명 매칭하고 상관 관계 토큰마다 각각 계산

고 해상도 이미지를 합성하기 위해 각 stage 사이에 pixel shuffle module로 구성됨, upsampling 모듈을 넣는다

pixel shuffle

$$H \times W \times C \rightarrow 2H \times 2W \times C/4 \rightarrow 4H \times 4W \times C/16 \rightarrow \underline{4H \times 4W \times 3}$$

↪ output resolution