# A U-Net Based Discriminator for Generative Adversarial Networks
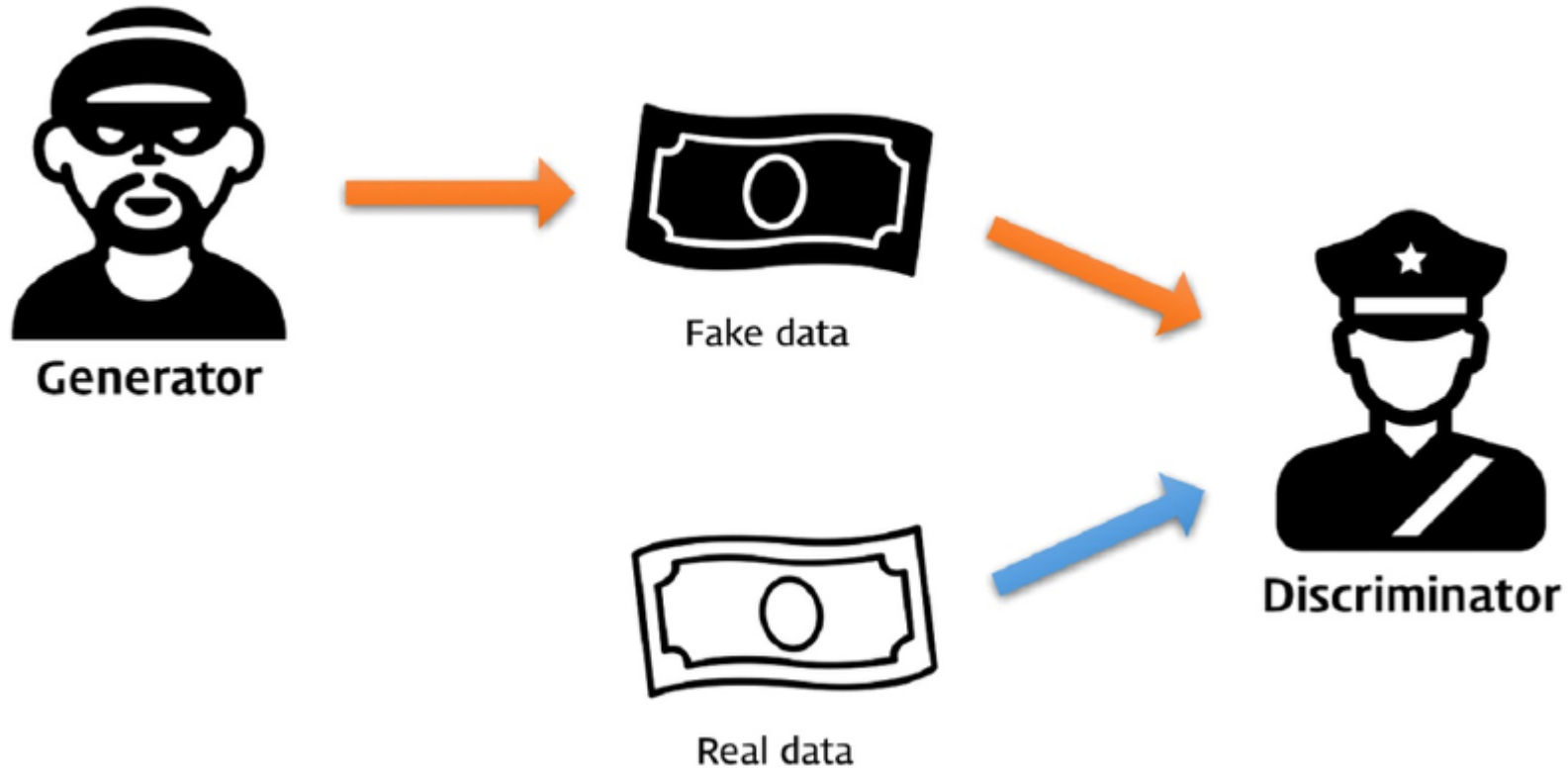
CVPR 2020

**Haneol Lee**

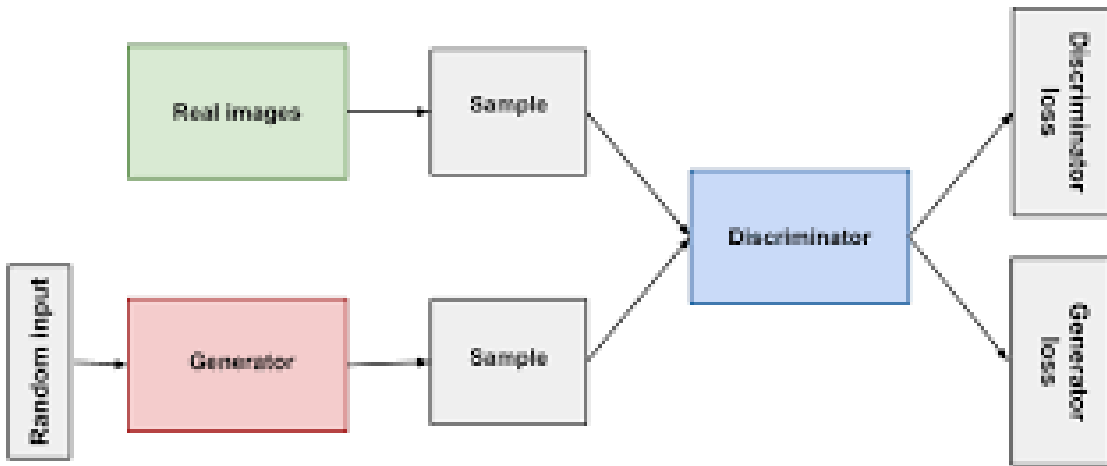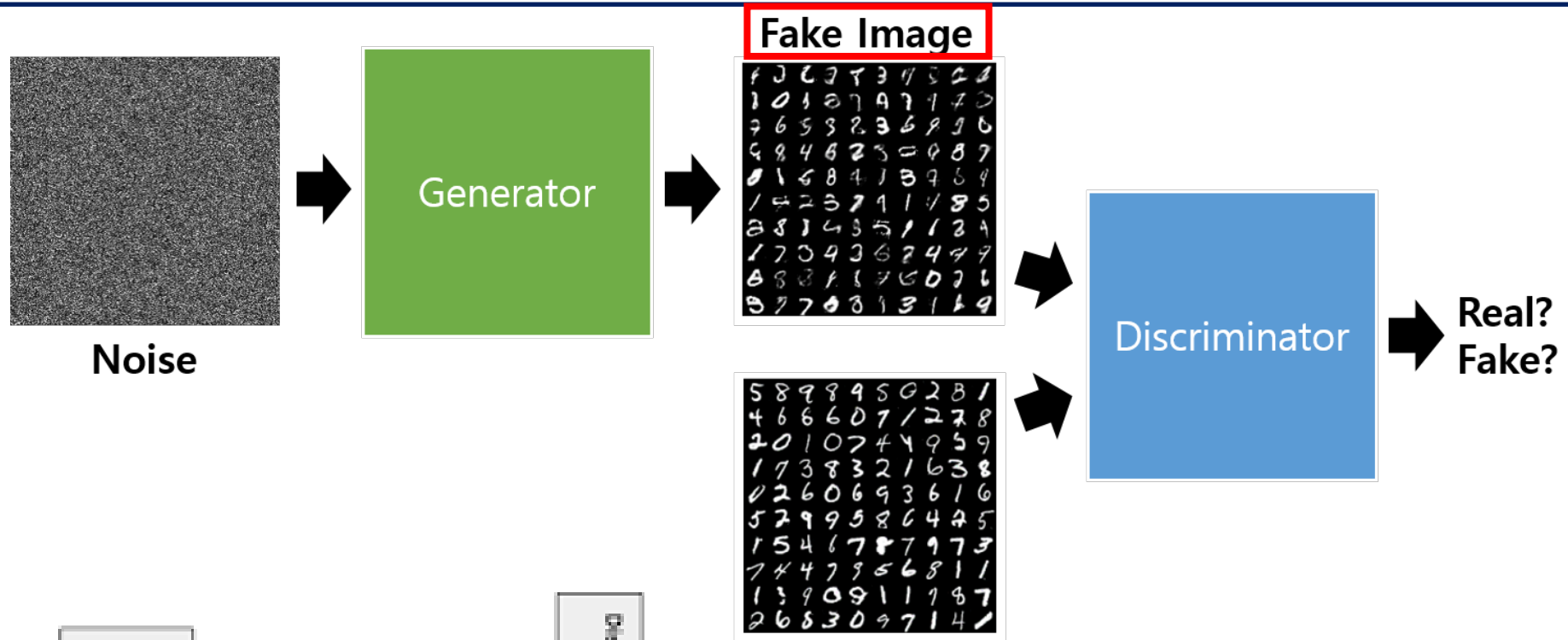2th May, 2022

# Introduction: GAN (Generative Adversarial Networks)

# Introduction: Generator Loss VS Discriminator Loss



$$\mathcal{L}_D = -\mathbb{E}_x[\log D(x)] - \mathbb{E}_z[\log(1 - D(G(z)))],$$

$$\mathcal{L}_G = -\mathbb{E}_z[\log D(G(z))]^1$$
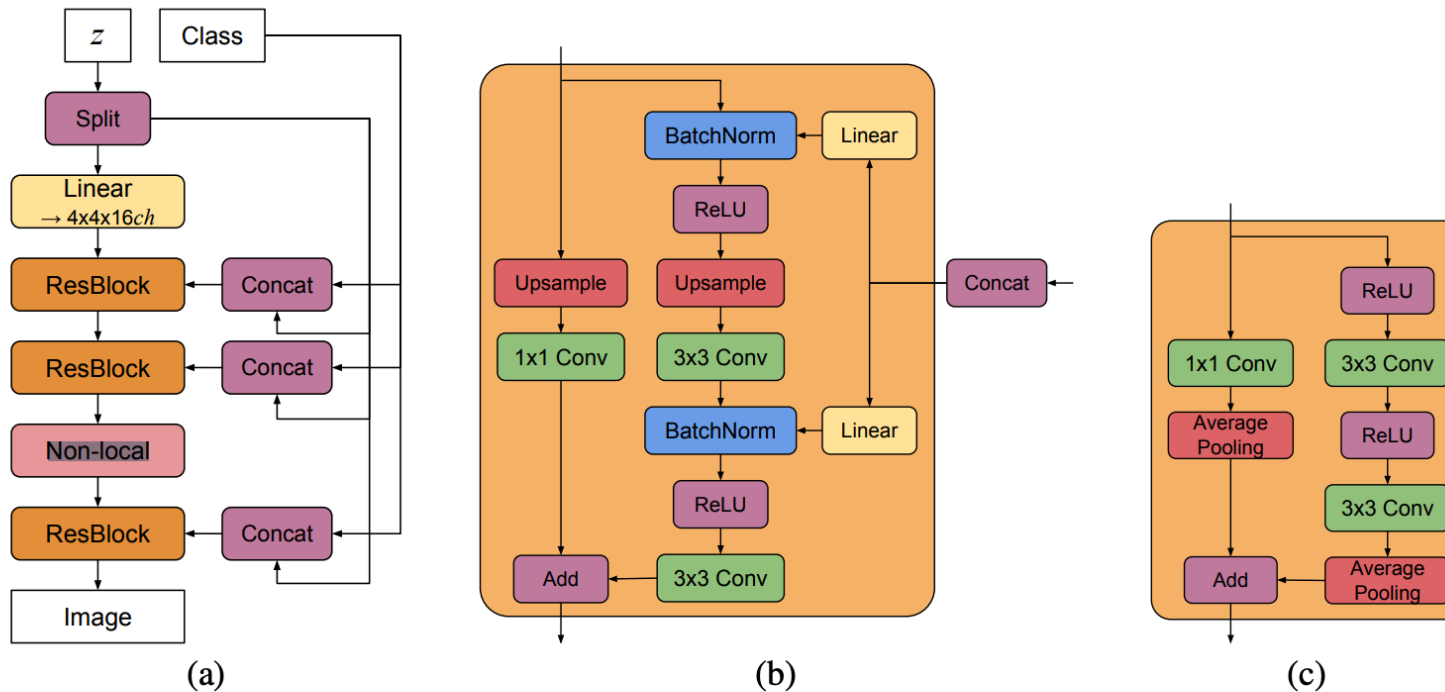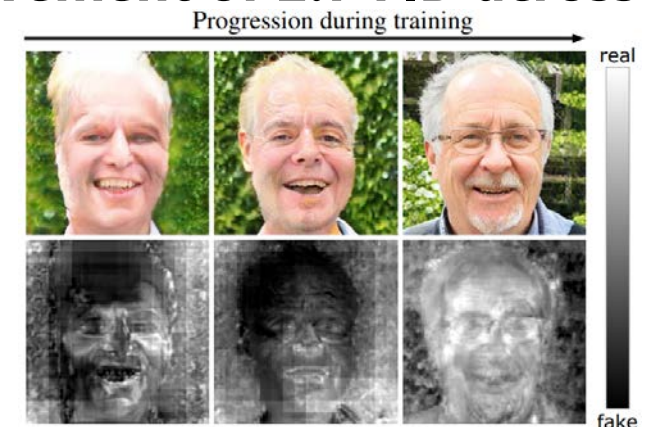
# Introduction: BigGAN



Figure 15: (a) A typical architectural layout for BigGAN's **G**; details are in the following tables. (b) A Residual Block (*ResBlock up*) in BigGAN's **G**. (c) A Residual Block (*ResBlock down*) in BigGAN's **D**.

Class-conditional samples generated by our model.

Link: *https://arxiv.org/pdf/1809.11096.pdf*

# Abstract

- **Borrowing <u>U-Net segmentation network</u> for the <u>discriminator</u> to target the issue of the <span style="color:red">capacity to synthesize globally</span> and <span style="color:red"><u>locally coherent images</u></span> with object shapes and textures indistinguishable from real images.**

- **U-Net based discriminator architecture allows to provide detailed <span style="color:red">per-pixel feedback to the generator</span> while <span style="color:red">maintaining the global coherence</span> of synthesized images.**

- **Propose a <span style="color:red">per-pixel consistency regularization</span> technique based on the <span style="color:red">CutMix</span> data augmentation <span style="color:red">encouraging discriminator to focus</span> more <span style="color:red">on semantic and structural changes</span> <u>between real and fake images.</u>**

- **Finally, the <u>U-Net based discriminator</u> <span style="color:red"><u>enables</u> the <u>generator</u></span> <u>to enhance the quality of generated samples</u>, and <span style="color:red">to synthesize images</span> with <span style="color:red">varying structure, appearance and levels of detail</span>.**

- **Compared to the BigGAN baseline, we achieve an average improvement of 2.7 FID across FFHQ, CelebA, and the newly introduced COCO-Animals dataset.**

Progression during training

# Introduction

**Challenging Limitation issues**:

1. **Global semantic coherence**

   ➔ Detailed and spatially coherent response to the generator

2. **Long-range structure**

   ➔ U-Net based architecture

3. **Exactness of detail**

   ➔Better qualitative performance
   of details

# Introduction

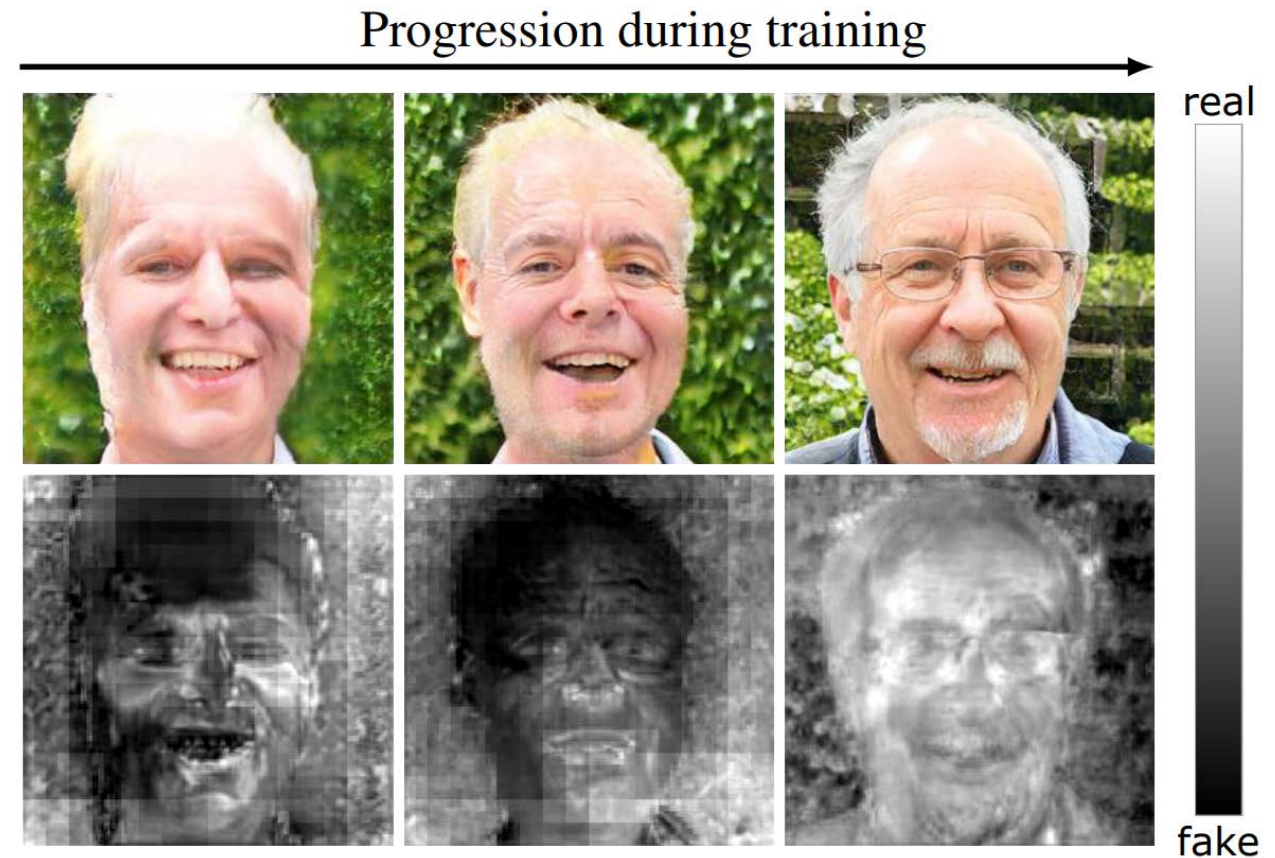**Challenging Limitation issues**:

1. **Global semantic coherence**

   ➔ Detailed and spatially coherent response to the generator

2. **Long-range structure**

   ➔ U-Net based architecture
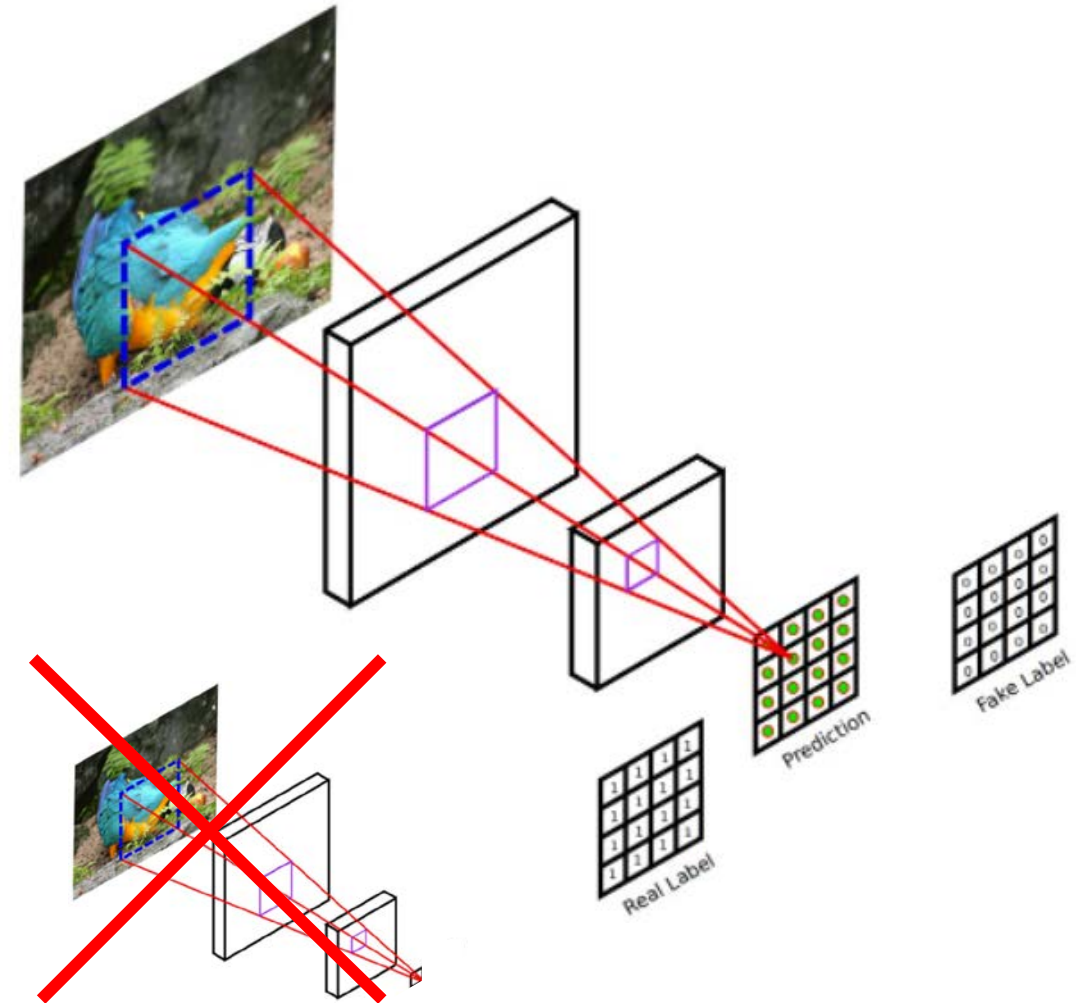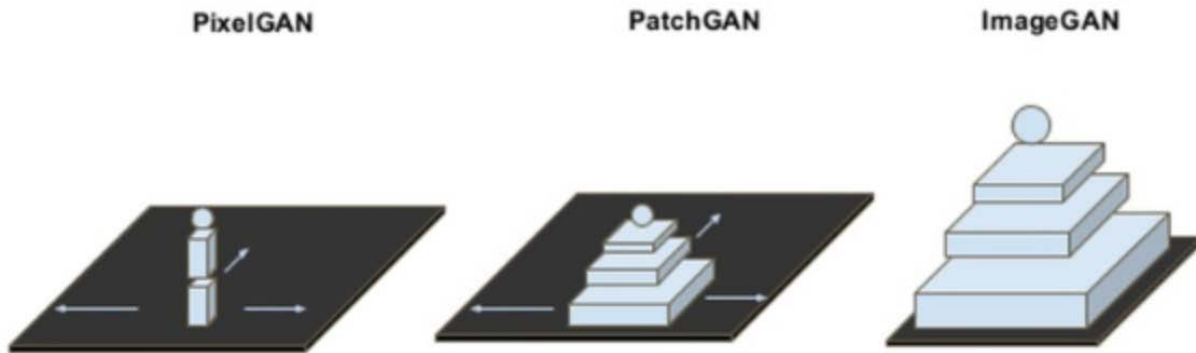
3. **Exactness of detail**

   ➔Better qualitative performance
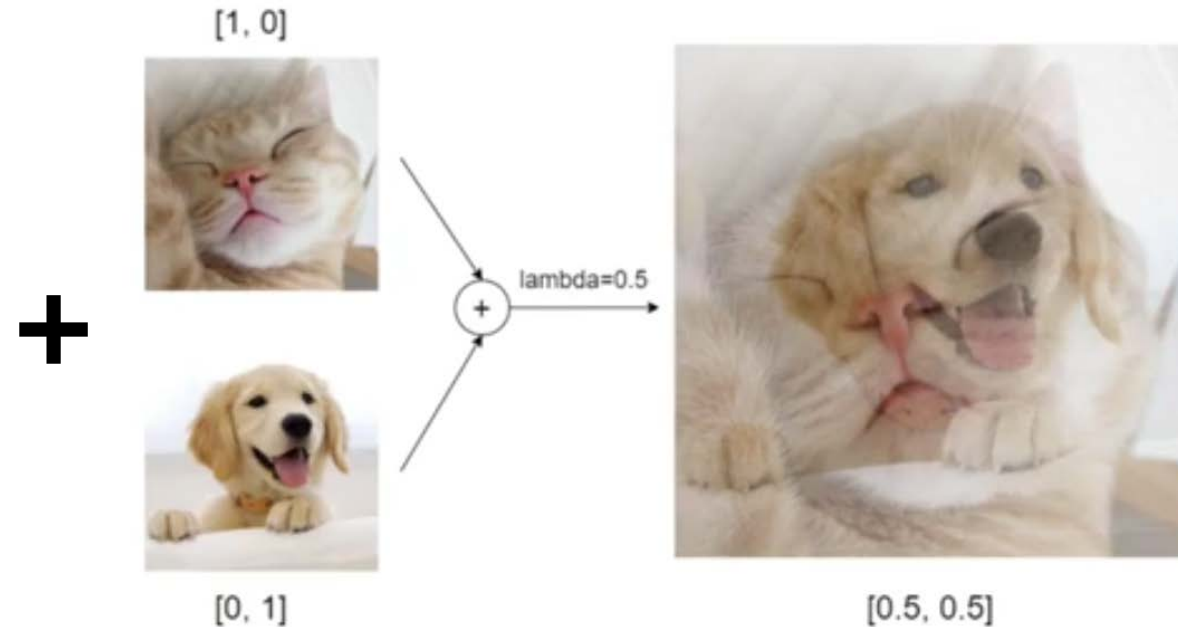   of details

# Introduction: PatchGAN

1. **Lighter parameters**

2. **Proper Receptive Field Size**

# Introduction: CutMix
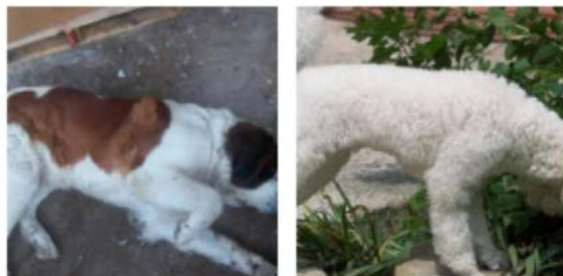
**CutMix**  =  **CutOut**  +  **MixUp**



**+**

[1, 0]

lambda=0.5

[0, 1]

[0.5, 0.5]

Mix-up works by blending 2 images with alpha % from image_1 and (1-alpha) % from image_2

Link1(MixUp): https://arxiv.org/pdf/1710.09412.pdf
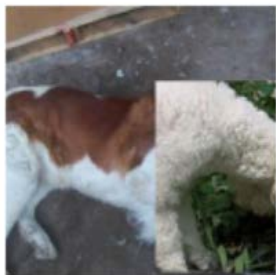Link2(CutOut): https://arxiv.org/pdf/1708.04552.pdf

# Introduction: CutMix

1. Regularization method to **increase generalization** and **localization performance** by **allowing models** to view and **learn less distinctive parts** and the **overall region** of the image, <u>rather than focusing on where they can identify differences</u> in objects
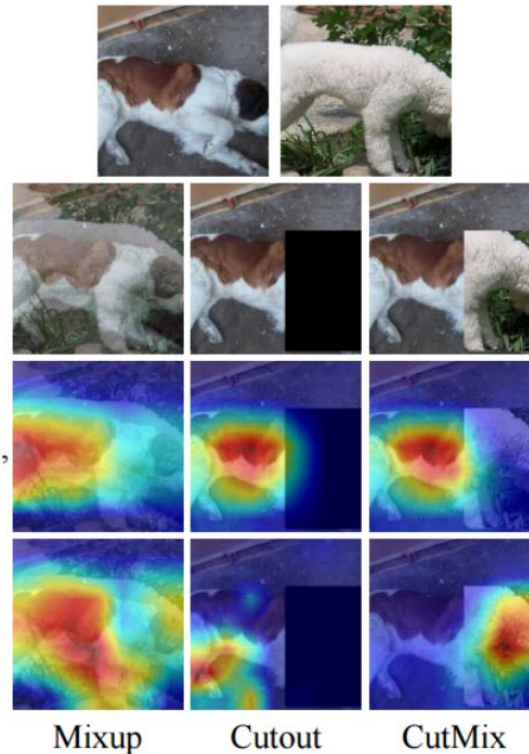


original image

cutmix

| | Original Samples |
| Input Image |
| CAM for 'St. Bernard' |
| CAM for 'Poodle' |

Mixup    Cutout    CutMix

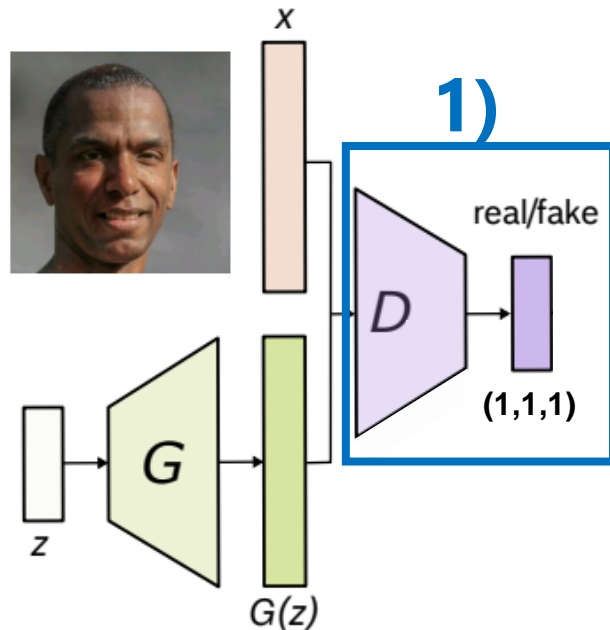| | ResNet-50 | Mixup [48] | Cutout [3] | CutMix |
|---|---|---|---|---|
| Image | | | | |
| Label | Dog 1.0 | Dog 0.5 Cat 0.5 | Dog 1.0 | Dog 0.6 Cat 0.4 |
| ImageNet Cls (%) | 76.3 (+0.0) | 77.4 (+1.1) | 77.1 (+0.8) | **78.6** (+2.3) |
| ImageNet Loc (%) | 46.3 (+0.0) | 45.8 (-0.5) | 46.7 (+0.4) | **47.3** (+1.0) |
| Pascal VOC Det (mAP) | 75.6 (+0.0) | 73.9 (-1.7) | 75.1 (-0.5) | **76.7** (+1.1) |

*Link:* https://arxiv.org/pdf/1905.04899.pdf

# Methods

**Redesign the role of the discriminator as 1) classifier, 2) segmenter**

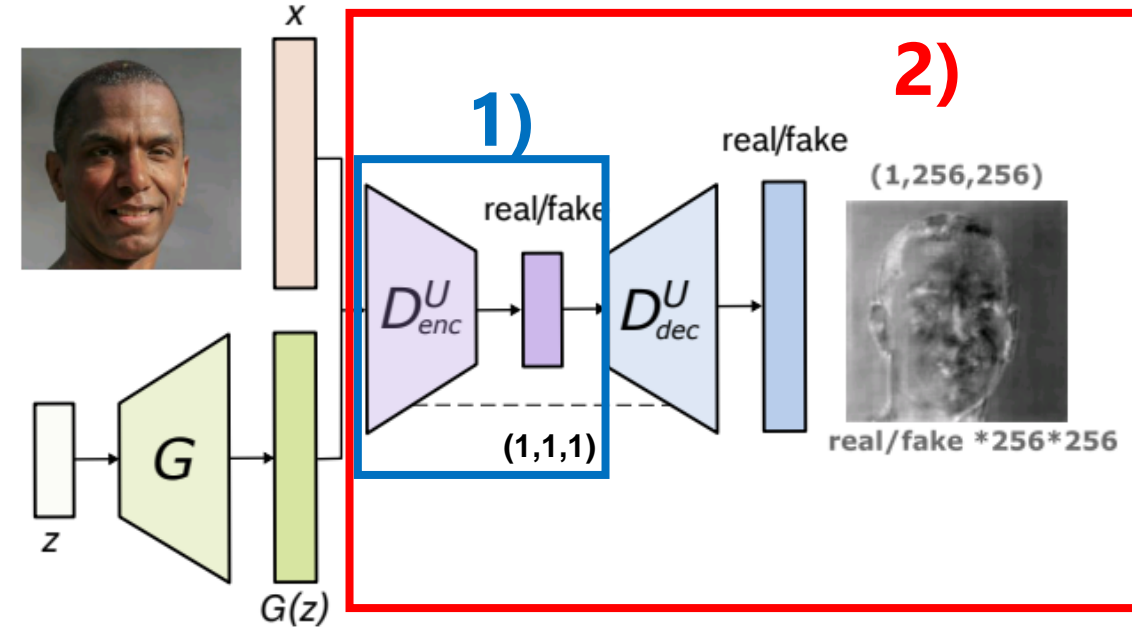**(1) Baseline GAN (Vanilla or BigGAN)  VS.  (2) U-Net GAN**

- Architecture: Same Generator,  Different Discriminator

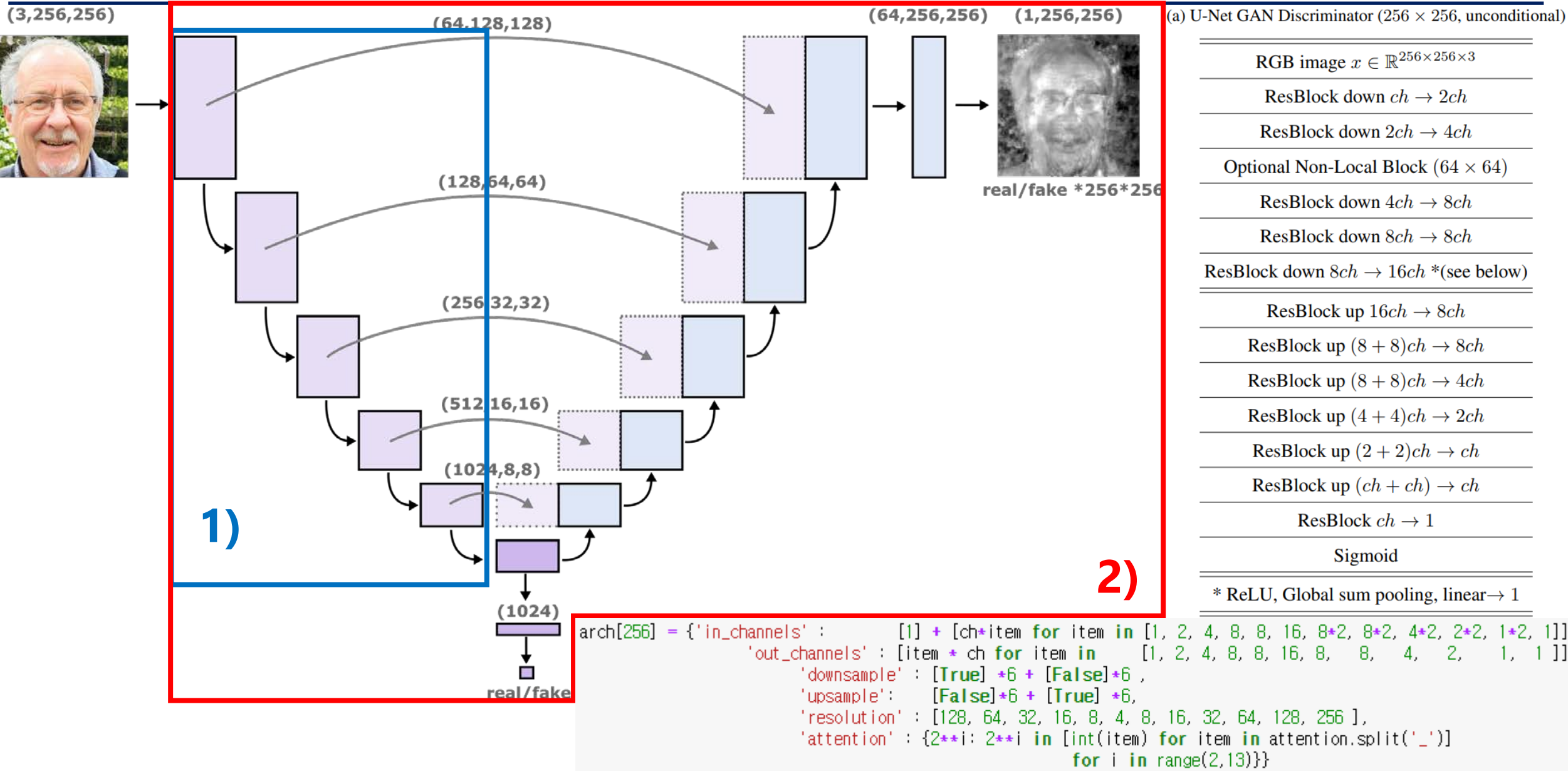**(1)**                                                                                          **(2)**

# Methods: U-Net Based Discriminator Architecture



(3,256,256)     (64,128,128)     (64,256,256)     (1,256,256)

(a) U-Net GAN Discriminator ($256 \times 256$, unconditional)

| |
|---|
| RGB image $x \in \mathbb{R}^{256 \times 256 \times 3}$ |
| ResBlock down $ch \to 2ch$ |
| ResBlock down $2ch \to 4ch$ |
| Optional Non-Local Block ($64 \times 64$) |
| ResBlock down $4ch \to 8ch$ |
| ResBlock down $8ch \to 8ch$ |
| ResBlock down $8ch \to 16ch$ *(see below) |
| ResBlock up $16ch \to 8ch$ |
| ResBlock up $(8 + 8)ch \to 8ch$ |
| ResBlock up $(8 + 8)ch \to 4ch$ |
| ResBlock up $(4 + 4)ch \to 2ch$ |
| ResBlock up $(2 + 2)ch \to ch$ |
| ResBlock up $(ch + ch) \to ch$ |
| ResBlock $ch \to 1$ |
| Sigmoid |
| * ReLU, Global sum pooling, linear$\to 1$ |

real/fake *256*256

(128,64,64)

(256,32,32)

(512,16,16)

(1024,8,8)

(1024)

real/fake

**1)**     **2)**

```
arch[256] = {'in_channels' :       [1] + [ch*item for item in [1, 2, 4, 8, 8, 16, 8*2, 8*2, 4*2, 2*2, 1*2, 1]],
             'out_channels' : [item * ch for item in    [1, 2, 4, 8, 8, 16, 8,   8,   4,   2,   1,  1 ]],
             'downsample' : [True] *6 + [False]*6 ,
             'upsample':    [False]*6 + [True] *6,
             'resolution' : [128, 64, 32, 16, 8, 4, 8, 16, 32, 64, 128, 256 ],
             'attention' : {2**i: 2**i in [int(item) for item in attention.split('_')]
                            for i in range(2,13)}}
```

# Methods

## (1) Baseline GAN (Vanilla or BigGAN)  VS.  (2) U-Net GAN

- **Architecture: Different Discriminator (Baseline D ➜ U-Net based D)**

**(1)**

(b) BigGAN Discriminator ($128 \times 128$, class-conditional)

| |
| --- |
| RGB image $x \in \mathbb{R}^{128 \times 128 \times 3}$ |
| ResBlock down $ch \rightarrow 2ch$ |
| Non-Local Block ($64 \times 64$) |
| ResBlock down $2ch \rightarrow 4ch$ |
| ResBlock down $4ch \rightarrow 8ch$ |
| ResBlock down $8ch \rightarrow 16ch$ |
| ResBlock down $16ch \rightarrow 16ch$ |
| ReLU, Global sum pooling |
| Embed(y)·$h$ + (linear$\rightarrow$ 1) |

**(2)**

(b) U-Net GAN Discriminator($128 \times 128$, class-conditional)

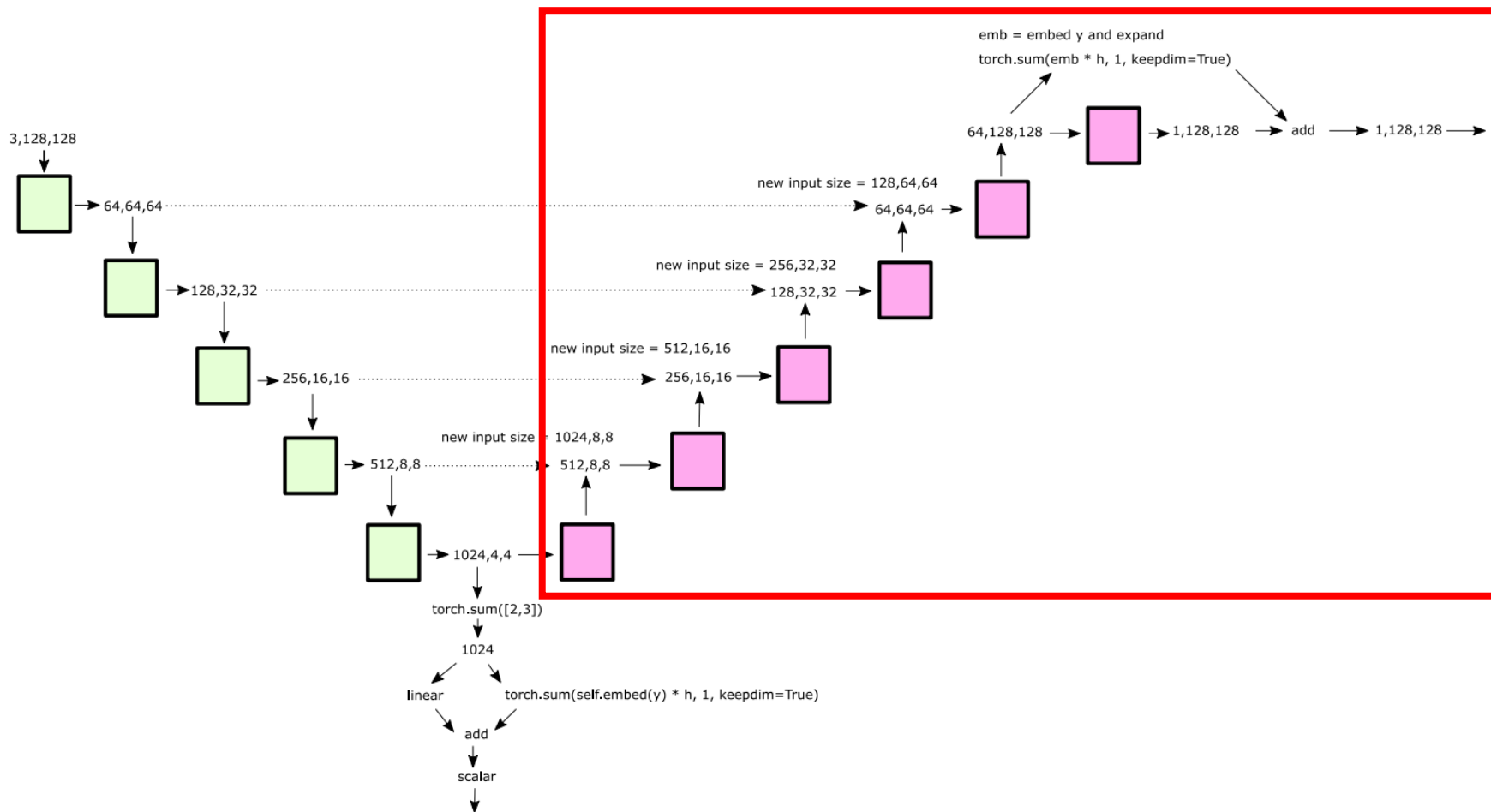| |
| --- |
| RGB image $x \in \mathbb{R}^{128 \times 128 \times 3}$ |
| ResBlock down $ch \rightarrow 2ch$ |
| Optional Non-Local Block ($64 \times 64$) |
| ResBlock down $2ch \rightarrow 4ch$ |
| ResBlock down $8ch \rightarrow 8ch$ |
| ResBlock down $8ch \rightarrow 16ch$ *(see below) |
| ResBlock up $16ch \rightarrow 8ch$ |
| ResBlock up $(8+8)ch \rightarrow 4ch$ |
| ResBlock up $(4+4)ch \rightarrow 2ch$ |
| ResBlock up $(2+2)ch \rightarrow ch$ |
| ResBlock up $(ch+ch) \rightarrow ch$ |
| Embed(y)·$h$ + (Conv $ch \rightarrow 1$) |
| Sigmoid |
| * ReLU, Global sum pooling |
| Embed(y)·$h$ + (linear$\rightarrow$ 1) |

# Methods

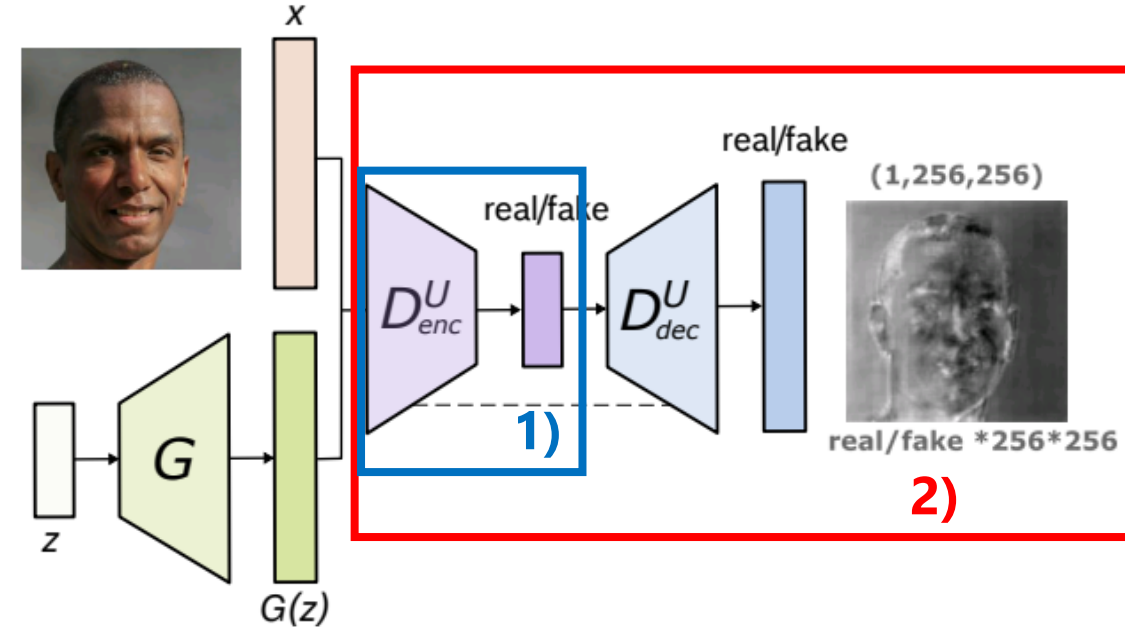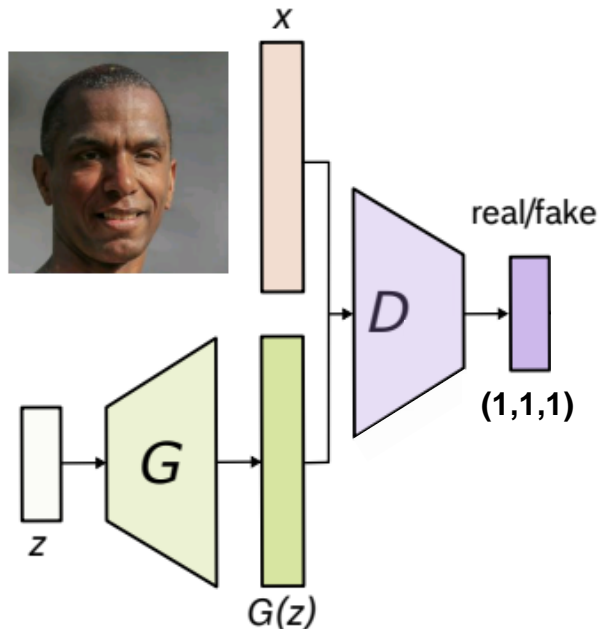## (1) Baseline GAN (Vanilla or BigGAN)  VS.  (2) U-Net GAN

- **Architecture: Different Discriminator (Baseline D ➡ U-Net based D)**



(b) U-Net GAN Discriminator($128 \times 128$, class-conditional)

| |
|---|
| RGB image $x \in \mathbb{R}^{128 \times 128 \times 3}$ |
| ResBlock down $ch \to 2ch$ |
| Optional Non-Local Block ($64 \times 64$) |
| ResBlock down $2ch \to 4ch$ |
| ResBlock down $8ch \to 8ch$ |
| ResBlock down $8ch \to 16ch$ *(see below) |
| ResBlock up $16ch \to 8ch$ |
| ResBlock up $(8+8)ch \to 4ch$ |
| ResBlock up $(4+4)ch \to 2ch$ |
| ResBlock up $(2+2)ch \to ch$ |
| ResBlock up $(ch+ch) \to ch$ |
| Embed(y)$\cdot h$ + (Conv $ch \to 1$) |
| Sigmoid |
| * ReLU, Global sum pooling |
| Embed(y)$\cdot h$ + (linear$\to 1$) |

# Methods

## (1) Baseline GAN (Vanilla or BigGAN)   VS.   (2) U-Net GAN

- **Architecture:** **Same Generator,**  **Different Discriminator**

**(1)**  $$\mathcal{L}_G = -\mathbb{E}_z[\log D(G(z))]^1$$

**(2)**  $$\mathcal{L}_G = -\mathbb{E}_z\left[\log D^U_{enc}(G(z)) + \sum_{i,j}\log[D^U_{dec}(G(z))]_{i,j}\right]$$

$$\mathcal{L}_D = -\mathbb{E}_x[\log D(x)] - \mathbb{E}_z[\log(1 - D(G(z)))],$$

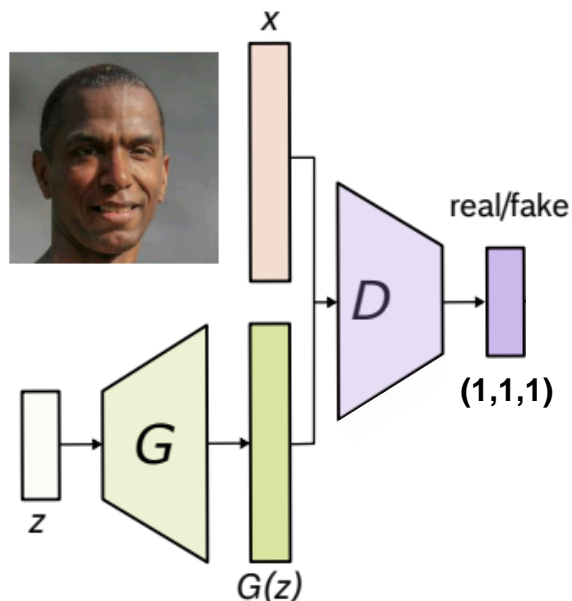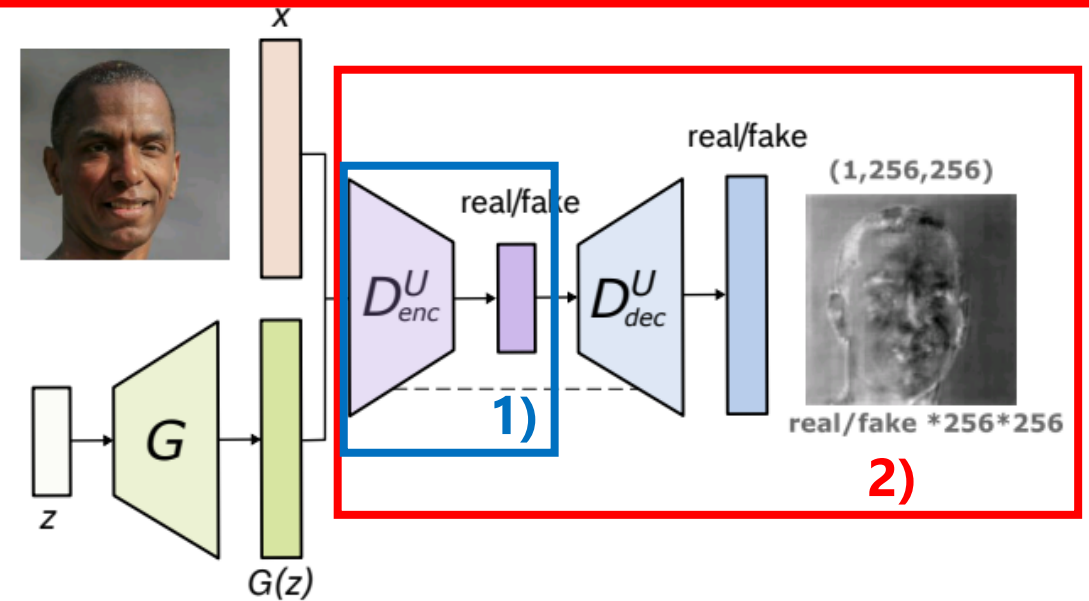$$\mathcal{L}_{D^U} = \mathcal{L}_{D^U_{enc}} + \mathcal{L}_{D^U_{dec}},$$

# Methods

## (1) Baseline GAN (Vanilla or BigGAN)  VS.  (2) U-Net GAN

- **Architecture: Different Discriminator (Baseline D ➔ U-Net based D)**

**(1)**

$$\mathcal{L}_G = -\mathbb{E}_z[\log D(G(z))]^1$$

$$\mathcal{L}_D = -\mathbb{E}_x[\log D(x)] - \mathbb{E}_z[\log(1 - D(G(z)))],$$

**(2)**

$$\mathcal{L}_G = -\mathbb{E}_z\left[\log D_{enc}^U(G(z)) + \sum_{i,j}\log[D_{dec}^U(G(z))]_{i,j}\right]$$
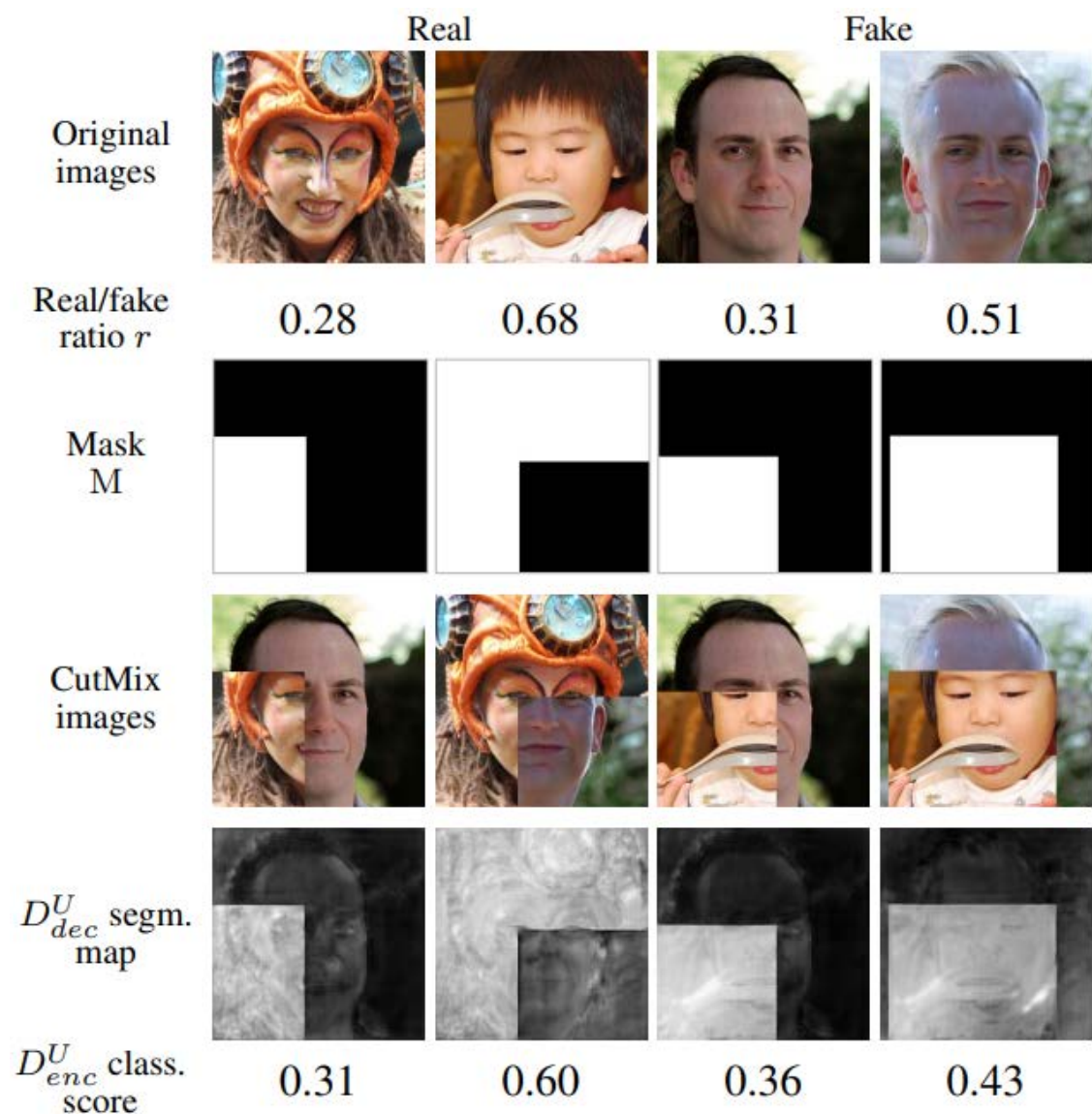
$$\mathcal{L}_{D^U} = \boxed{\mathcal{L}_{D_{enc}^U}} + \boxed{\mathcal{L}_{D_{dec}^U}},$$

$$\mathcal{L}_{D_{enc}^U} = -\mathbb{E}_x[\log D_{enc}^U(x)] - \mathbb{E}_z[\log(1 - D_{enc}^U(G(z)))],$$

$$\mathcal{L}_{D_{dec}^U} = -\mathbb{E}_x\left[\sum_{i,j}\log[D_{dec}^U(x)]_{i,j}\right] - \mathbb{E}_z\left[\sum_{i,j}\log(1 - [D_{dec}^U(G(z))]_{i,j})\right]$$

# Methods: Consistency Regularization



Real      Fake

Original images

Real/fake ratio $r$    0.28    0.68    0.31    0.51

Mask M

CutMix images

$D_{dec}^{U}$ segm. map

$D_{enc}^{U}$ class. score    0.31    0.60    0.36    0.43

$$\text{M} \in \{0,1\}^{W \times H}$$
$$\text{fake}(\text{M}_{i,j} = 0)$$
$$\text{real}(\text{M}_{i,j} = 1)$$

$$\tilde{x} = \text{mix}(x, G(z), \text{M}),$$
$$\text{mix}(x, G(z), \text{M}) = \text{M} \odot x + (1 - \text{M}) \odot G(z), \quad (6)$$

$$D_{dec}^{U}\big(\text{mix}(x, G(z), \text{M})\big) \approx \text{mix}\big(D_{dec}^{U}(x), D_{dec}^{U}(G(z)), \text{M}\big),$$

$$\mathcal{L}_{D_{dec}^{U}}^{cons} = \left\| D_{dec}^{U}\Big(\text{mix}(x, G(z), \text{M})\Big) \right.$$
$$\left. - \text{mix}\Big(D_{dec}^{U}(x), D_{dec}^{U}(G(z)), \text{M}\Big) \right\|^{2}, \quad (7)$$

# Methods: Consistency Regularization



| | Real | | Fake | |
|---|---|---|---|---|
| Original images | | | | |
| Real/fake ratio $r$ | 0.28 | 0.68 | 0.31 | 0.51 |
| Mask M | | | | |
| CutMix images | | | | |
| $D_{dec}^U$ segm. map | | | | |
| $D_{enc}^U$ class. score | 0.31 | 0.60 | 0.36 | 0.43 |

$$\mathrm{M} \in \{0,1\}^{W \times H}$$
$$\mathrm{fake}(\mathrm{M}_{i,j} = 0)$$
$$\mathrm{real}(\mathrm{M}_{i,j} = 1)$$

$$r = \frac{|\mathrm{M}|}{W * H}$$

$$\tilde{x} = \mathrm{mix}(x, G(z), \mathrm{M}),$$
$$\mathrm{mix}(x, G(z), \mathrm{M}) = \mathrm{M} \odot x + (1 - \mathrm{M}) \odot G(z), \quad (6)$$

$$D_{dec}^U(\mathrm{mix}(x, G(z), \mathrm{M})) \approx \mathrm{mix}(D_{dec}^U(x), D_{dec}^U(G(z)), \mathrm{M}),$$

$$\mathcal{L}_{D_{dec}^U}^{cons} = \left\| D_{dec}^U\big(\mathrm{mix}(x, G(z), \mathrm{M})\big) \right.$$
$$\left. - \mathrm{mix}\big(D_{dec}^U(x), D_{dec}^U(G(z)), \mathrm{M}\big) \right\|^2, \quad (7)$$

# Experiments

# Experiments



Figure 5: Images generated with U-Net GAN trained on COCO-Animals with resolution $128 \times 128$.

# Experiments

| Method | FFHQ Best FID↓ | FFHQ Best IS↑ | FFHQ Median FID↓ | FFHQ Median IS↑ | COCO-Animals Best FID↓ | COCO-Animals Best IS↑ | COCO-Animals Median FID↓ | COCO-Animals Median IS↑ |
|---|---|---|---|---|---|---|---|---|
| BigGAN [5] | 11.48 | 3.97 | 12.42 | 4.02 | 16.37 | 11.77 | 16.55 | 11.78 |
| U-Net GAN | **7.48** | **4.46** | **7.63** | **4.47** | **13.73** | **12.29** | **13.87** | **12.31** |

Table 1: Evaluation results on FFHQ and COCO-Animals. We report the best and median FID score across 5 runs and its corresponding IS, see Section 4.2 for discussion.

| Method | Dataset | FID Best | FID Median | FID Mean | FID Std |
|---|---|---|---|---|---|
| BigGAN | COCO-Animals | 16.37 | 16.55 | 16.62 | 0.24 |
| U-Net GAN | | **13.73** | **13.87** | **13.88** | **0.11** |
| BigGAN | FFHQ | 11.48 | 12.42 | 12.35 | 0.67 |
| U-Net GAN | | **7.48** | **7.63** | **7.73** | **0.56** |
| BigGAN | CelebA | 3.70 | 3.89 | 3.94 | 0.16 |
| U-Net GAN | | **2.03** | **2.07** | **2.08** | **0.04** |

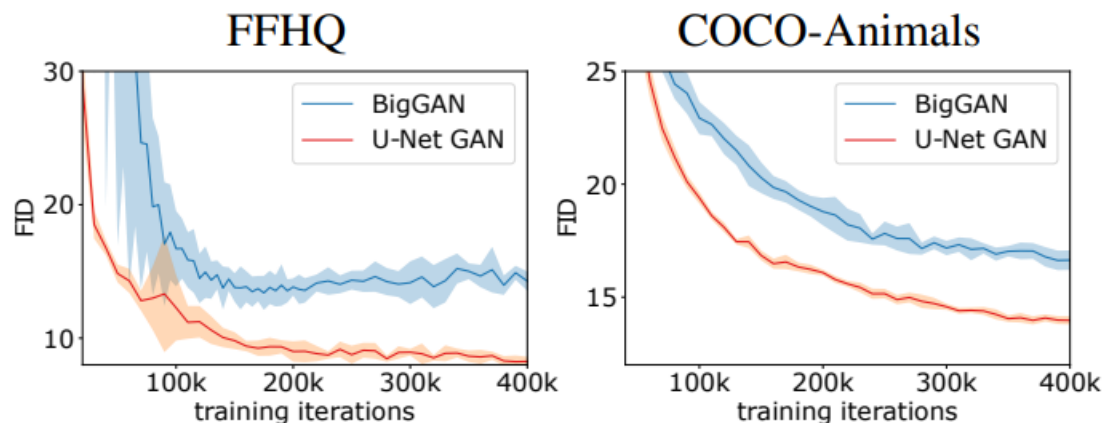Table S2: Best, median, mean and std of FID values across 5 runs.



Figure 6: FID curves over iterations of the BigGAN model (blue) and the proposed U-Net GAN (red). Depicted are the FID mean and standard deviation across 5 runs per setting.

# Experiments

| Method | COCO-Animals | FFHQ |
|---|---|---|
| BigGAN [5] | 16.55 | 12.42 |
| U-Net based discriminator | 15.86 | 10.86 |
| + CutMix augmentation | 14.95 | 10.30 |
| + Consistency regularization | **13.87** | **7.63** |

Table 2: Ablation study of the U-Net GAN model on FFHQ and COCO-Animals. Shown are the median FID scores. The proposed components lead to better performance, on average improving the median FID by 3.7 points over Big-GAN [5]. See Section 4.2 for discussion.

| Method | FID ↓ | IS ↑ |
|---|---|---|
| PG-GAN [19] | 7.30 | – |
| COCO-GAN [27] | 5.74 | – |
| BigGAN [5] | 4.54 | 3.23 |
| U-Net GAN | **2.95** | **3.43** |

Table 3: Comparison with the state-of-the-art models on CelebA ($128 \times 128$). See Section 4.2 for discussion.

| Dataset | Method | PyTorch | | TensorFlow | |
|---|---|---|---|---|---|
| | | FID ↓ | IS ↑ | FID ↓ | IS ↑ |
| FFHQ | BigGAN [5] | 11.48 | 3.97 | 14.92 | 3.96 |
| ($256 \times 256$) | U-Net GAN | **7.48** | **4.46** | **8.88** | **4.50** |
| COCO-Animals | BigGAN [5] | 16.37 | 11.77 | 16.42 | 11.34 |
| ($128 \times 128$) | U-Net GAN | **13.73** | **12.29** | **13.96** | **11.77** |
| | PG-GAN [19] | – | – | 7.30 | – |
| CelebA | COCO-GAN [27] | – | – | 5.74 | – |
| ($128 \times 128$) | BigGAN [5] | 3.70 | 3.08 | 4.54 | 3.23 |
| | U-Net GAN | **2.03** | **3.33** | **2.95** | **3.43** |

Table S1: Evaluation results on FFHQ, COCO-Animals and CelebA with PyTorch and TensorFlow FID/IS scores. The difference lies in the choice of framework in which the inception network is implemented, which is used to extract the inception metrics. See Section A for discussion.

# Experiments



Figure S2: Samples generated by U-Net GAN and the corresponding real-fake predictions of the U-Net decoder. Brighter colors correspond to the discriminator confidence of pixel being real (and darker of being fake).

Figure S7: Generated samples on COCO-Animals and the corresponding U-Net decoder predictions.

# Experiments



Figure S4: Qualitative comparison of uncurated images generated with the unconditional BigGAN model (top) and our U-Net GAN (bottom) on FFHQ with resolution 256 × 256. Note that the images generated by U-Net GAN exhibit finer details and maintain better local realism.

# Experiments



Figure 7: Visualization of the predictions of the encoder $D_{enc}^U$ and decoder $D_{dec}^U$ modules during training, within a batch of 50 generated samples. For visualization purposes, the $D_{dec}^U$ score is averaged over all pixels in the output. Note that quite often decisions of $D_{enc}^U$ and $D_{dec}^U$ are not coherent with each other. As judged by the U-Net discriminator, samples in the upper left consist of locally plausible patterns, while not being globally coherent (example in orange), whereas samples in the lower right look globally coherent but have local inconsistencies (example in purple: giraffe with too many legs and vague background).



Figure 8: Comparison of the generator and discriminator loss behavior over training for U-Net GAN and BigGAN. The generator and discriminator loss of U-Net GAN is additionally split up into its encoder- and decoder components.

# Conclusion

• **Proposing an alternative U-Net based architecture for the discriminator, which allows to provide both global and local feedback to the generator.**

• **U-Net based discriminator architecture allows to provide detailed <span style="color:red">per-pixel feedback to the generator</span> while <span style="color:red">maintaining the global coherence</span> of synthesized images.**

• **Introducing a consistency regularization technique for the U-Net discriminator based on the CutMix data augmentation**

• **Showing that all the proposed changes result in a stronger discriminator, enabling the generator to synthesize images with varying levels of detail, maintaining global and local realism.**

• **We demonstrate the improvement over the state-of-the-art BigGAN model [5] in terms of the FID score on three different datasets.**

pr
di
lo

# Q&A

# Question 1

- 논문의 **3.2 Consistency regularization** 챕터에서 **Class-domain-perturbation**이 무슨 말인가요?
    - **Class가 정확하게 무엇을 지칭하나요?**
        - 해당 **class-domain-perterbation**에서 명시적으로 **calss c={0,1}**라고 적혀 있고, **perturbation**이 될 만한 부분은 **real/fake class**이라고 생각합니다.
        - 논문의 해당 파트에 명시적으로 **Class c = {0,1}**로 적혀 있고, **BigGAN**의 **condition** 이 명확히 언급되지 않은 것으로 보아서도 **e.g. real=1, fake=0** 가 맞는 것 같습니다.
    - **Class-domain-perturbation**이 무슨 말인가요? 뒤에서 설명드리겠습니다.

- **Consistency regularization의 motivation이 무엇이고 구체적인 효과(성능이 왜 높아지는지)에 대해서 자세히 설명해주세요.**

  - **Real/Fake Domain의 bias 관점에서:**

  - 이미지를 **[A] = [cutmix하고 D에 넣은 결과]와 [B] = [D에 넣은 결과를 cutmix]** 한 것을 같도록 제약을 주는 것이 **consistency regulation loss term**인데요. **CutMix output**결과가 **(real)M, (fake)1-M**에 마스크 영역 의해 **real/fake**가 나뉘어지므로 **real image pixel**들은 **real**끼리 **fake image pixel**들은 **fake**끼리 뭉쳐 있는 상황이겠습니다. 만약 **real pixels domain, / fake pixels domain**별로 **domain** 차이에 의한 **pixel data** 분포에 **bias**있다고 가정해봅시다.

  - 이중 **Cutmix**결과 이미지가 **Enc(U-Net D ConvNet Encdoer)** 이미지로 들어갈 때, 최종 **Discriminator output(1,256,256)**의 각 **output pixel**별로 대응되는 **input**의 **Receptive field**를 생각해보면, (앞에서 말한 **real domain, fake domain** 사이의 **bias**에 의해서) **Receptive field**에서 **real pixel**의 비율이 많은 곳은 **pixel**을 **real**이라고 판별할 가능성이 높고, **fake pixel**들의 비율이 많은 곳은 **output**을 **fake**라고 판별할 가능성이 높아질 수 있을 것 같습니다. **(domain**간의 **bias** 정보를 **Discriminator**가 활용할 수 있기 때문이에요.)

- **Class-domain-perturbation이 무슨 말인가요?**

  - **Class-domain-perturbation**이라는 것은 바로 **Output**의 특정 **pixel**에 대응하는 **Receptive field**에서 **fake pixels, real pixels** 수의 비율 차이에 따라 각 개별 **pixel** 정보에 알게 모르게 남아있던 일종의 **domain** 별 **bias**값을 **Discriminator**가 활용하는 문제를 지칭하는 것이 아닐까 생각합니다. 이런 문제는 **Cutmix**를 수행해주는 의도에 맞지 않으니 제약해주는 게 맞지 않을까 합니다.
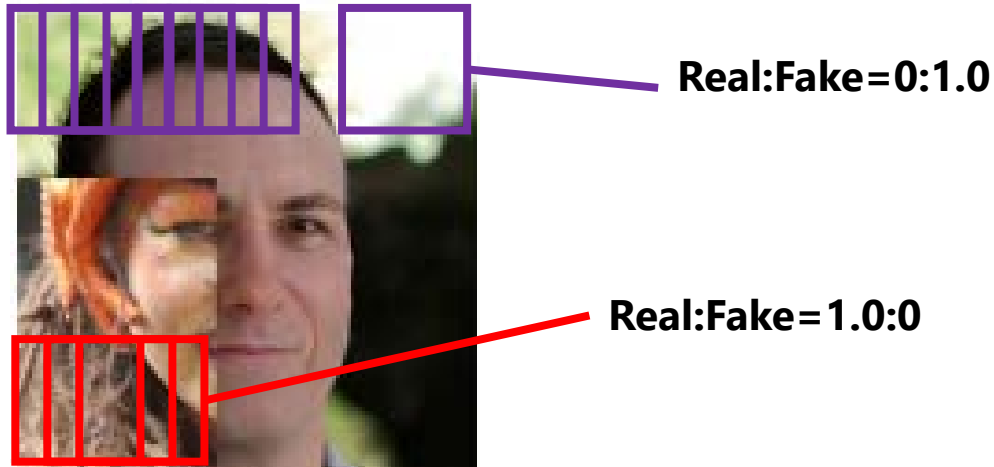
- **Consistency regularization의 motivation**이 무엇이고 구체적인 효과(성능이 왜 높아지는지)에 대해서 자세히 설명해주세요.
  - 반면에 정보를 학습할 때 주변영역에 대한 **dependence**를 학습하지 못하게 하는 효과가 조금이라도 더 생기고 동시에 **pixel**별 독립적인 판단을 하도록 제약을 걸어 주게 되는 것 같습니다.
  - 왜냐하면 **consistency** 제약 **Loss**를 전체 **loss term**에 추가해주었을 때, 구체적으로 (**cutmix**를 하는 순서와 **D**에 통과시키는 순서가 달라지는 것)[A], [B]사이에 어떤 차이가 있느냐 자세히 살펴본다면, ([A]일 때와 [B]일 때 모두) **receptive field**에 **real pixel**만 있거나 **fake pixel**만 있을 때는 동일한 연산과정을 거쳐 동일한 값이 나오게 되므로, **Discriminator**입장에서 **Q2 (3/4)** 그림에서처럼 결과가 동일하게 됩니다.
  - 그러나 **output**의 특정 **pixel** 지점에 대응하는 **Receptive field**에서 **real/fake mask** 주변의 경계선에서 **real pixels**과 **fake pixels**이 섞여 있다면, **Q2 (4/4)** 그림에서처럼, **convNet window filter**가 **sliding**함에 따라 **real/fake** 자체의 비율이 그 때 그 때 달라지므로, 모델 입장에서 **real**과 **fake**의 비율 정보를 **discriminating** 결과에 덜 참고하게 학습이 되는 제약 효과가 주어질 것 같습니다.
  - 결과적으로 **Discriminator**가 **real/fake** 비율이나 그 비율에 따른 **bias** 정보를 컨닝(?)하기보다는 개별 **pixel** 정보들을 **independent**하게 좀 더 집중하게 하는 효과를 주지 않나 싶습니다.
    따라서, **consistency** 제약을 주어 학습을 하면 **Discriminator**가 (**pixel** 단위로 독립적으로) **independently pixel-wise**하게 **pixel** 정보만을 활용해서 **Discriminator**가 집중해서 학습할 수 있도록 유도할 수 있고, 그에 따라 **Generator**도 전체적으로 합성 이미지가 **semantic**하게 결과를 판별하도록 학습 할 수 있을 것 같습니다.
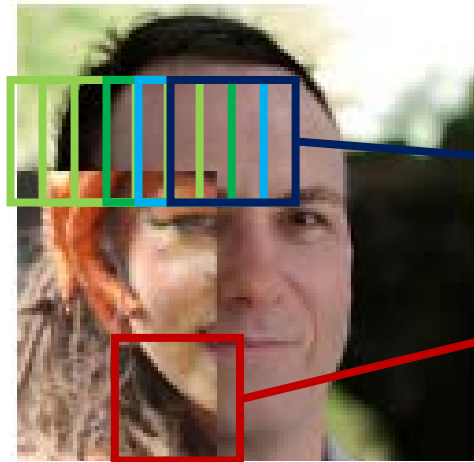
**[A]**

**= Cumix를 먼저해서 D에 통과시키는 경우**



**Real:Fake=0:1.0**

**Real:Fake=1.0:0**

**Receptive field안에 real pixels만, 혹은 fake pixes만 있는 경우,**

**[A], [B]의 연산결과는 동일하다.**

**[B]**

**= real/fake를 각각 D에 통과시키고 나서, 그 다음에 Cutmix하는 경우**



**Real:Fake=0:1.0**

**Real:Fake=1.0:0**

**[A]**

   = **Cumix**를 먼저해서 **D**에 통과시키는 경우

**[B]**

   = **real/fake**를 각각 **D**에 통과시키고 나서, 그 다음에 **Cutmix**하는 경우



**Real:Fake=0.1:0.9**

**Real:Fake=0.8:0.2**



**Real:Fake=0:1.0**

**Real:Fake=0:1.0**

**Real:Fake=1.0:0**

**Receptive field**안에 **real pixels** 일부와 **fake pixes**일부가 섞여 있는 경우, **[A]**, **[B]**의 연산 결과는 다르다.

**ConvLayer Filter**가 input에 대해서 **sliding** 하게 **Conv** 연산을 수행할 때, 한 **receptive field** 안의 **real**, **fake**의 비율도그 때 그 때 달라진다. ➔ 즉, **real/fake**의 비율이 **sliding** 할 때마다 달라지게 되면서, **Discriminator**입장에서는 **real**과 **fake**의 비율의 정보를 활용하는 것이 **loss**를 줄이는 방향과 직결되지 않는다.

# References

[1] "A U-Net Based Discriminator for Generative Adversarial Networks", CVPR 2019

[2] "mixup: BEYOND EMPIRICAL RISK MINIMIZATION"

[3] "Improved Regularization of Convolutional Neural Networks with Cutout"

[4]"Large Scale GAN Training for High Fidelity Natural Image Synthesis" (BigGAN)