

# **A Style-Based Generator Architecture for Generative Adversarial Networks**

---

CVPR 2019

Tero Karras, Samuli Laine, Timo Aila

# Contents

---

**I. Introduction**

**II. Style-based Generator**

**III. Properties of the style-based generator**

**IV. Disentanglement Studies**

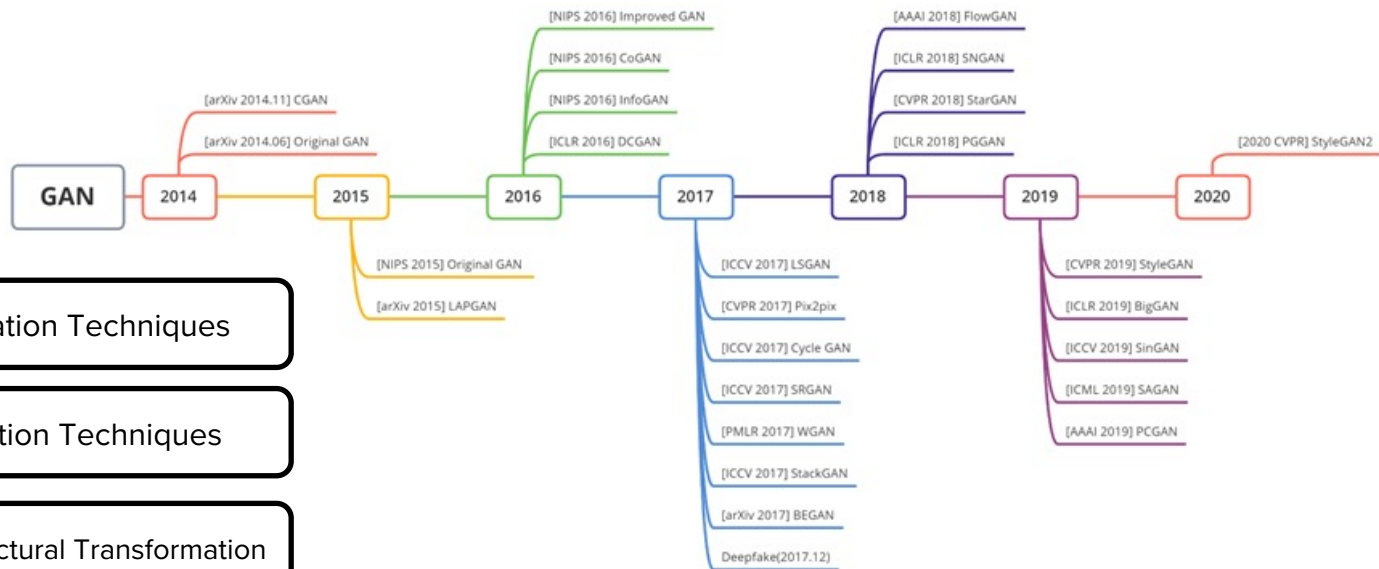
**V. Results**

# Introduction

---

# Introduction

## History of GAN



+ Normalization Techniques

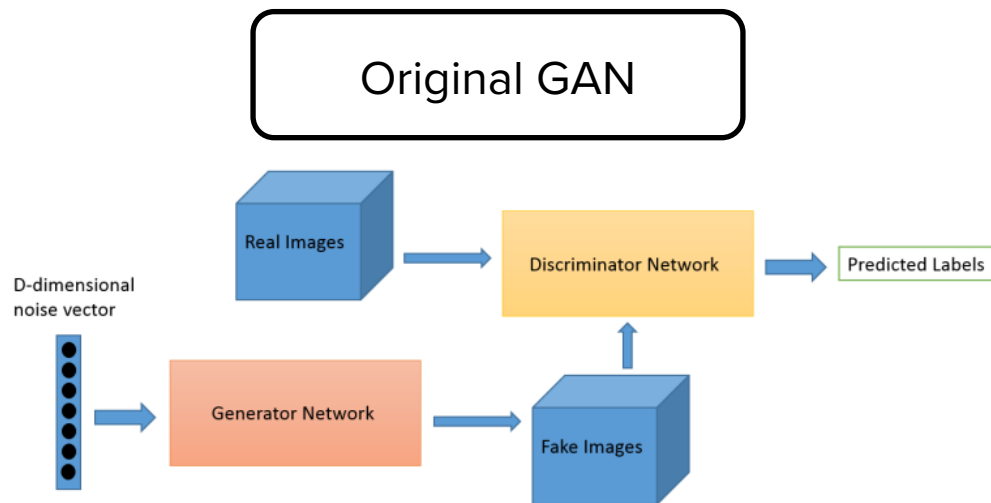
+ Optimization Techniques

+ Auxiliary Structural Transformation

# Introduction

---

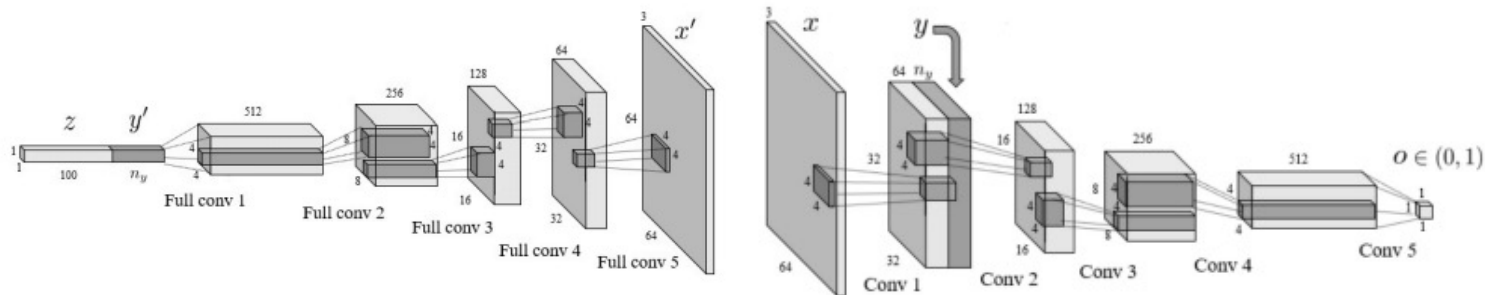
## History of GAN



# Introduction

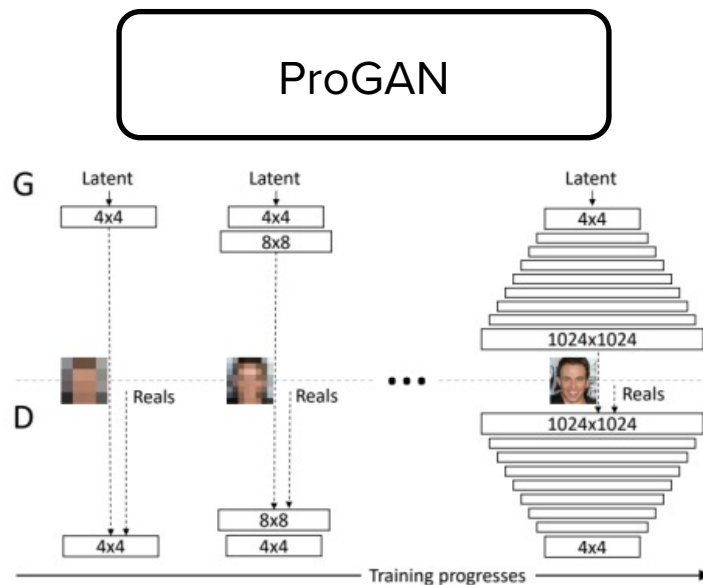
## History of GAN

cGAN



# Introduction

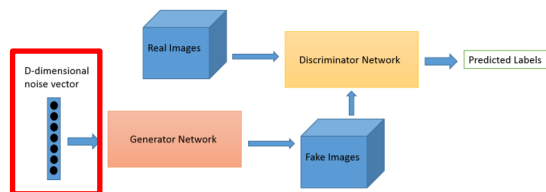
## History of GAN



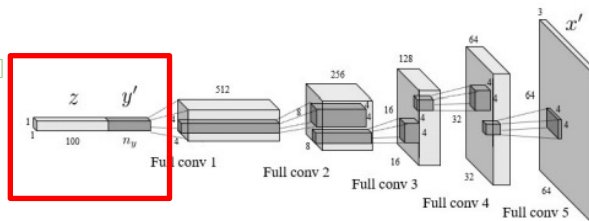
# Introduction

## ❑ Problem of previous models

Original GAN



cGAN



ProGAN

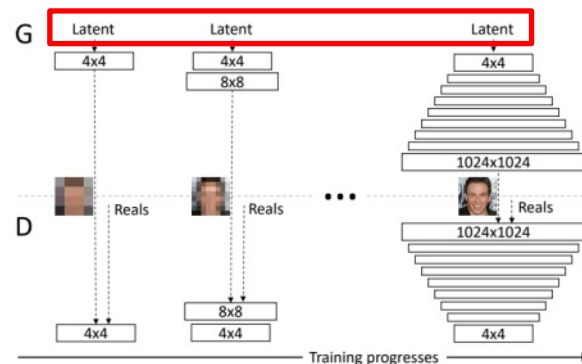


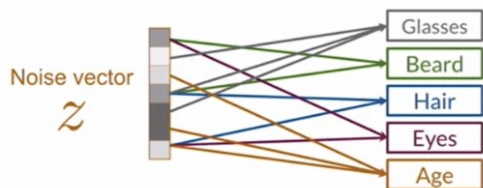
Image Generation based on **single latent space**



# Introduction

## ❑ Entanglement in latent space

Z-Space Entanglement



Output Features

Actual Distribution



**Desired**  
(Completely Distinguished)

Correlated Features



Add beard  
Make more masculine



Unexpected Impacts on other features

**Less control** over  
the latent space & the generator

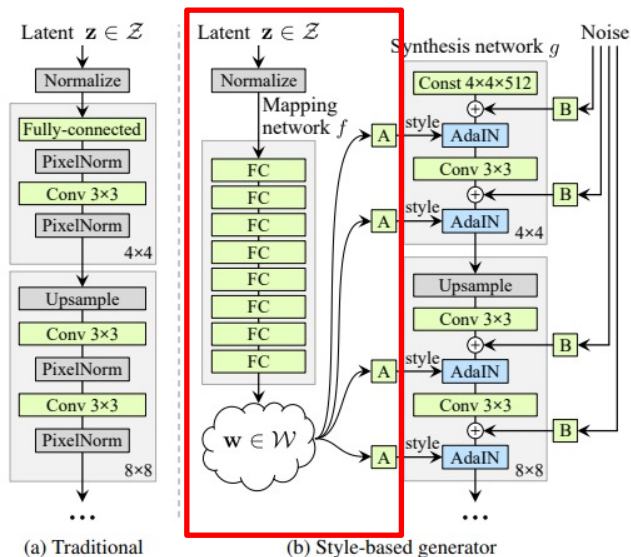
**Motivation for style-based generator**

# Style-based Generator

---

# Style-based Generator

## □ Model Structure



Mapping Network  $f : \mathcal{Z} \rightarrow \mathcal{W}$  (8-layer MLP)

$w \in \mathcal{W}$  : Intermediate Latent Vector

Learnable Affine  
Parameters

$$y = (y_s, y_b)$$

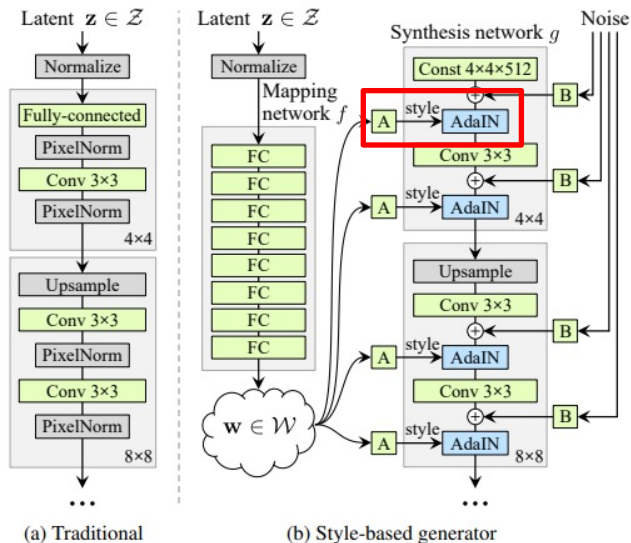
“Style”

$$\text{AdaIN}(x_i, y) = y_{s,i} \frac{x_i - \mu(x_i)}{\sigma(x_i)} + y_{b,i}$$

“Adaptive Instance Normalization”

# Style-based Generator

## □ Model Structure

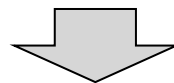


$$\text{AdaIN}(\mathbf{x}_i, \mathbf{y}) = \mathbf{y}_{s,i} \frac{\mathbf{x}_i - \mu(\mathbf{x}_i)}{\sigma(\mathbf{x}_i)} + \mathbf{y}_{b,i}$$

“Adaptive Instance Normalization”

Normalization on each feature map

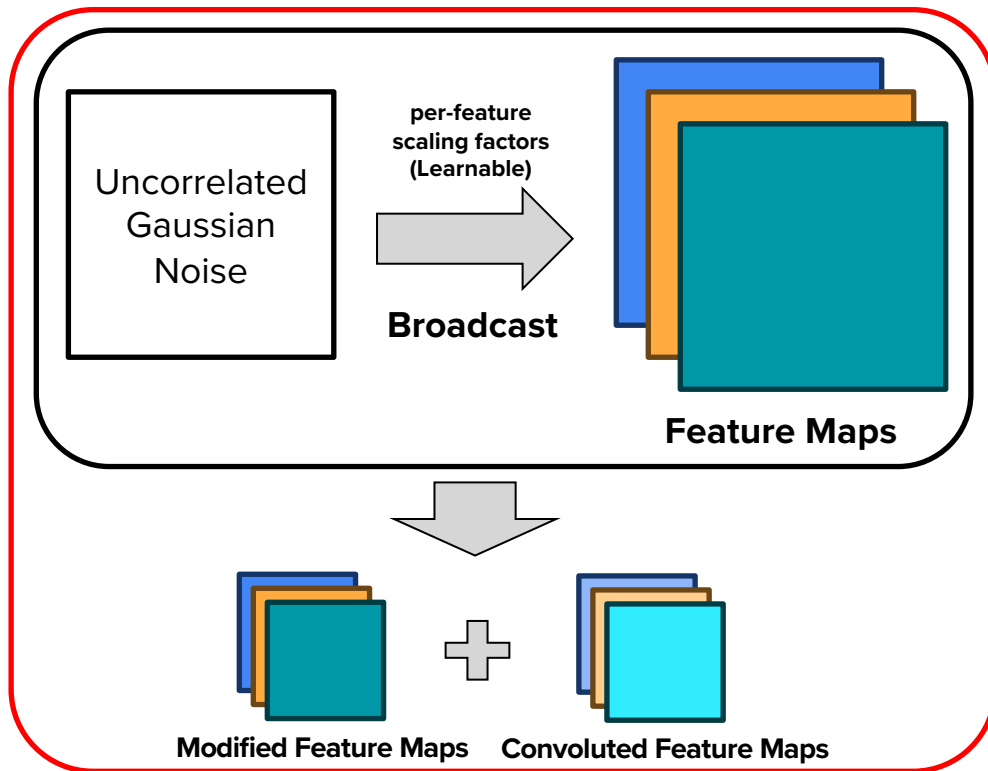
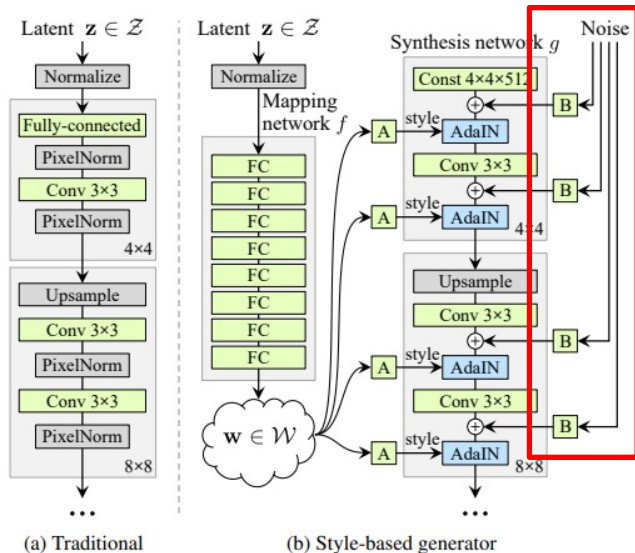
Scaling & Biasing with “Style”



Adding “Style” to the image that is being generated

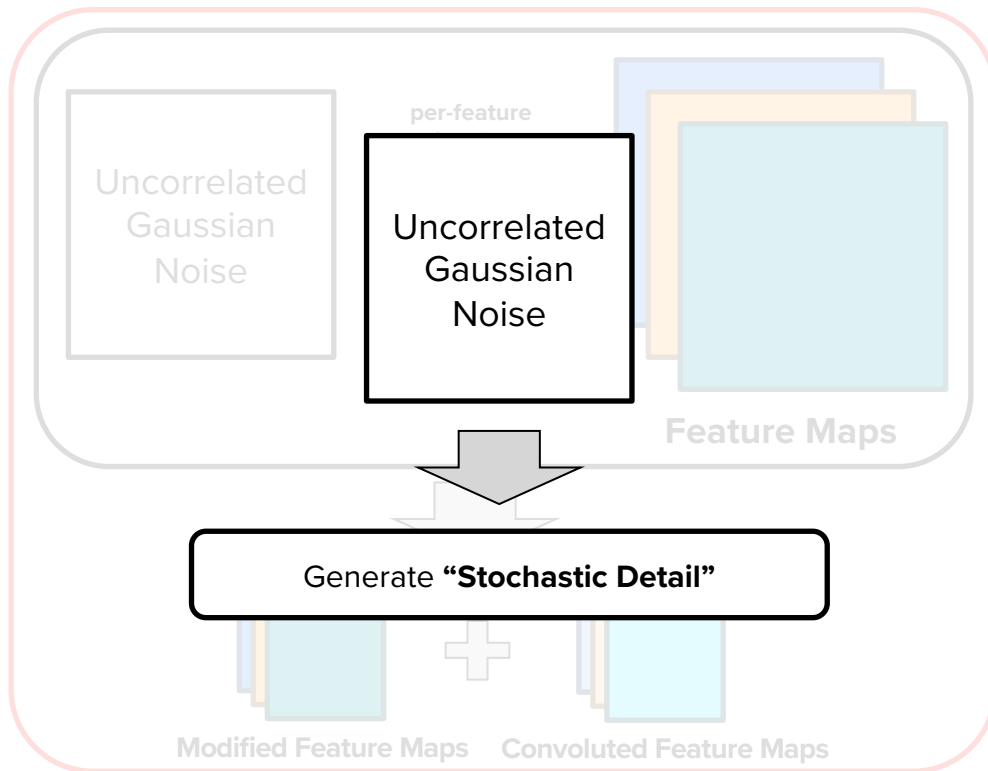
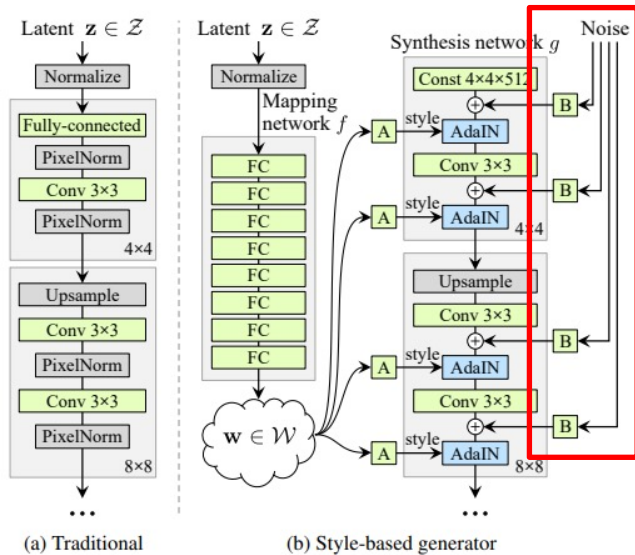
# Style-based Generator

## □ Model Structure



# Style-based Generator

## □ Model Structure

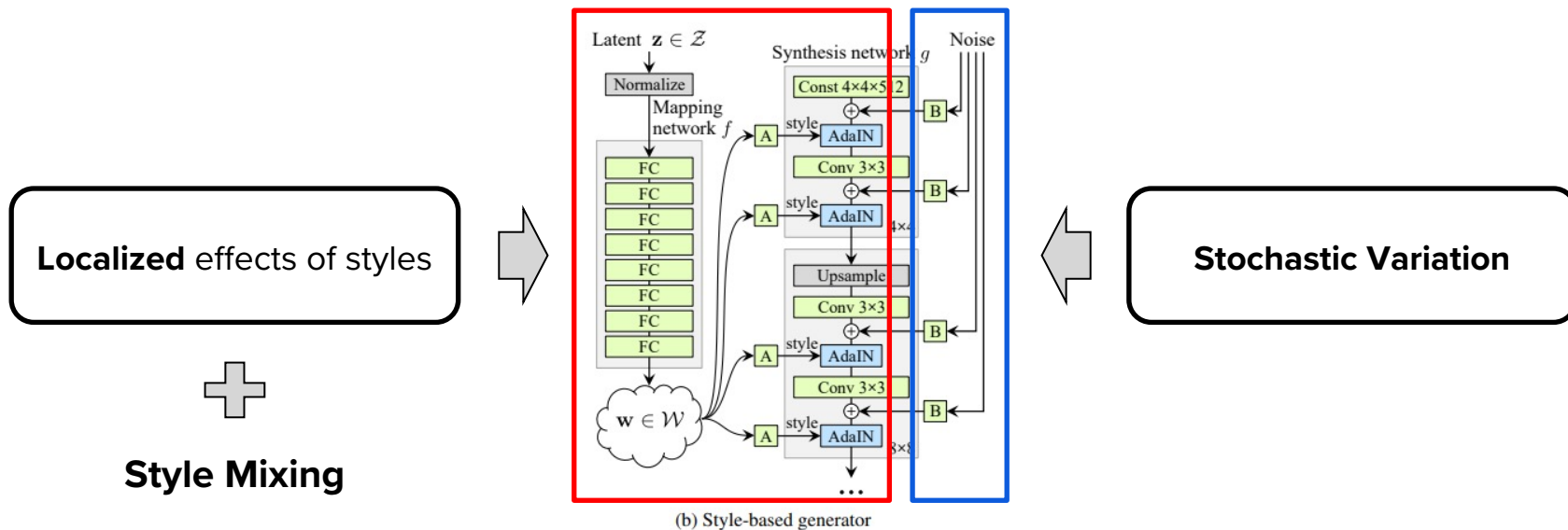


# **Properties of the Style-based Generator**

---

# Properties of the Style-based Generator

## □ Distinct Properties

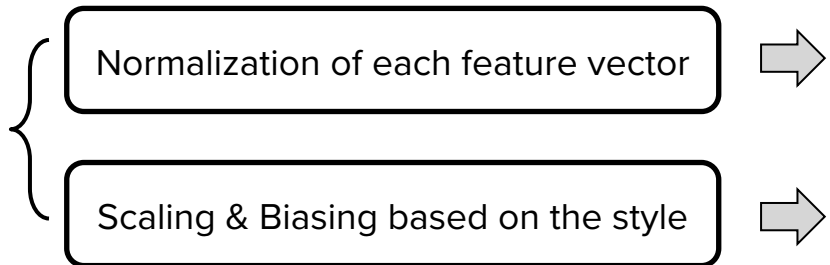




# Properties of the Style-based Generator

## □ Localization of the effects of styles

$$\text{AdaIN}(\mathbf{x}_i, \mathbf{y}) = \mathbf{y}_{s,i} \frac{\mathbf{x}_i - \mu(\mathbf{x}_i)}{\sigma(\mathbf{x}_i)} + \mathbf{y}_{b,i}$$



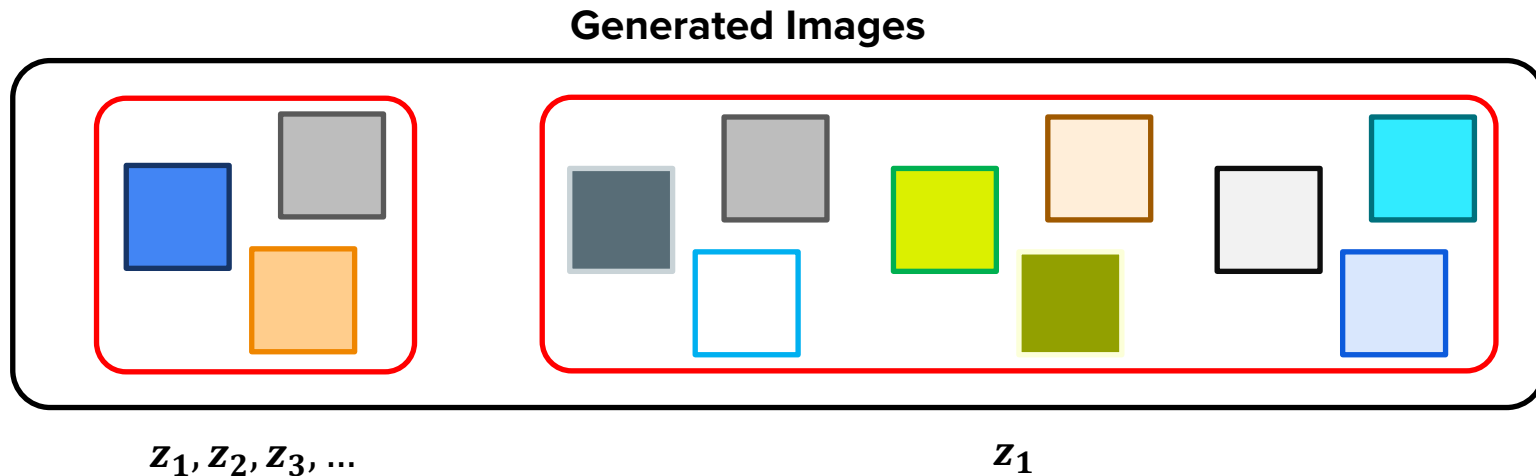
### Achieve Localization

Force the convolution to be performed **only** based on the preceding features

Control the convolution in the desired direction as indicated in the style

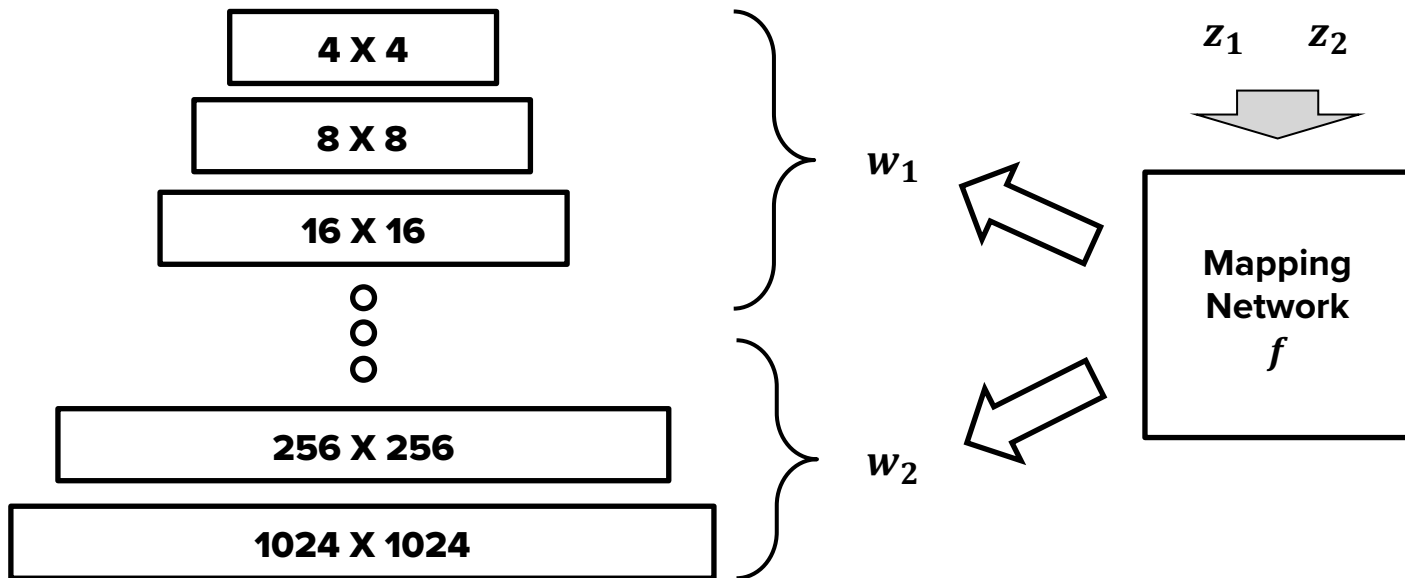
# Properties of the Style-based Generator

## ❑ Mixing Regularization



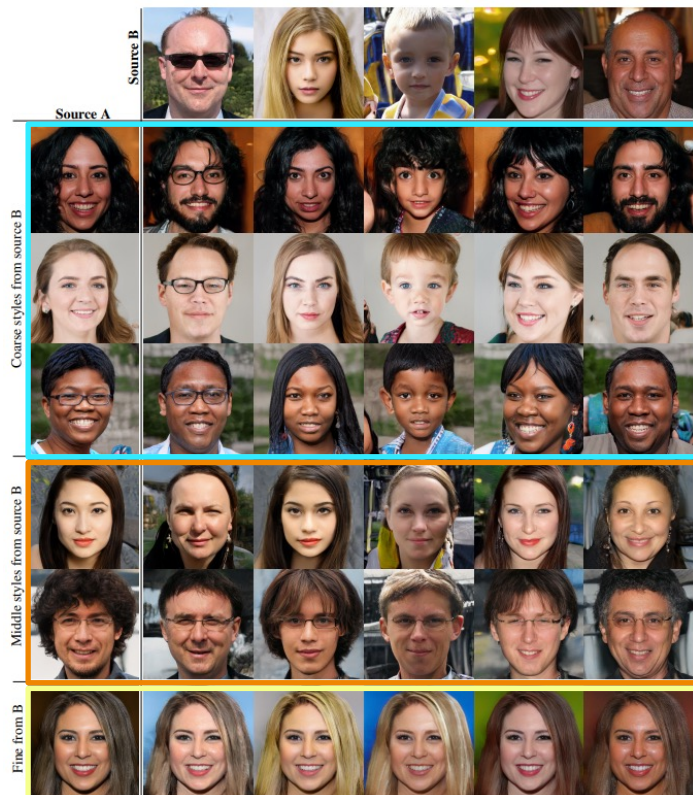
# Properties of the Style-based Generator

## □ Style Mixing



# Properties of the Style-based Generator

## □ Style Mixing



“Coarse Styles” from source B  
“Fine Styles” from source A

“Middle Styles” from source B  
“Middle Styles” from source A

“Fine Styles” from source B  
“Coarse Style” from source A

# Properties of the Style-based Generator

---

## □ Style Mixing

| Mixing<br>regularization | Number of latents during testing |             |             |             |
|--------------------------|----------------------------------|-------------|-------------|-------------|
|                          | 1                                | 2           | 3           | 4           |
| E 0%                     | 4.42                             | 8.22        | 12.88       | 17.41       |
| 50%                      | 4.41                             | 6.10        | 8.71        | 11.61       |
| F 90%                    | <b>4.40</b>                      | <b>5.11</b> | 6.88        | 9.03        |
| 100%                     | 4.83                             | 5.17        | <b>6.63</b> | <b>8.40</b> |

FID scores in FFHQ

\* FFHQ : New dataset of human faces created with style-based generator

# Properties of the Style-based Generator

## ❑ Stochastic Variation



**Stochastic aspects** in the picture of humans



**Can be randomized**  
**without affecting the perception**

# Properties of the Style-based Generator

---

## ❑ Stochastic Variation ➡

Difficult to achieve with traditional generators

- Latent information can be **supplied only through the input layer**
- **Hard to specify when and how many** spatially-varying pseudorandom numbers should be generated



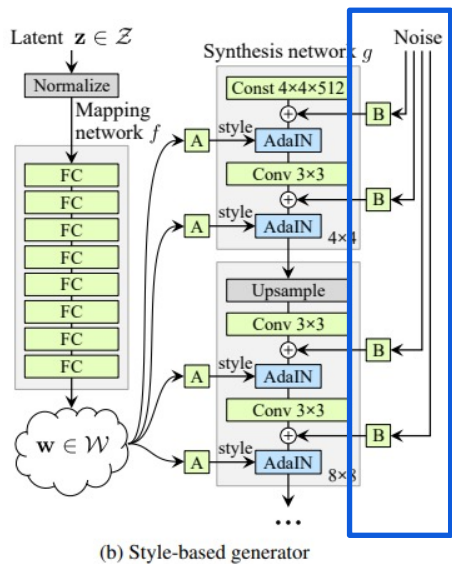
Require more network capacity



Inevitable exposure of the periodicity of generated signal

# Properties of the Style-based Generator

## □ Stochastic Variation



Add **per-pixel noise** after every convolution

With the help of **per-feature scaling factors**

Fresh noise is supplied to every layer (Features)

Only affects stochastic aspects (= Effect is localized)



# Properties of the Style-based Generator

## □ Stochastic Variation



(a) Generated image

(b) Stochastic variation

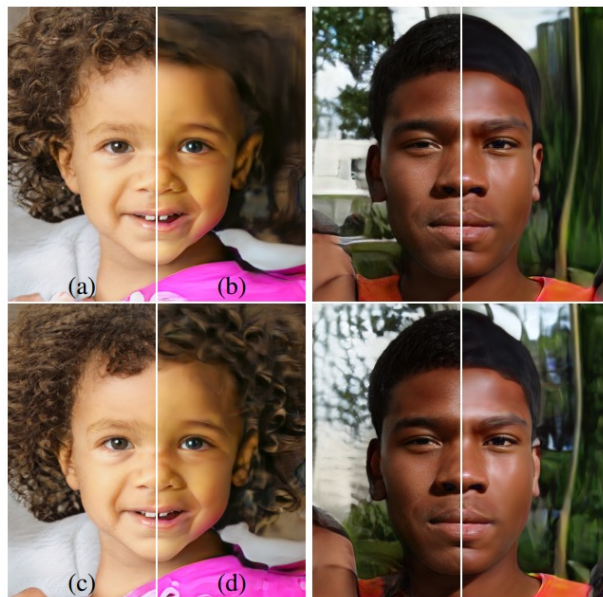
(c) Standard deviation

**White parts**  
→ **High Variations**

**Black parts**  
→ **Low Variations**

# Properties of the Style-based Generator

## □ Stochastic Variation



(a) : All layers  
(b) : No noise  
(c) : Fine Layers  
(d) : Coarse Layers

Different Results according to the stage that noise is injected

# Properties of the Style-based Generator

---

## ❑ Style & Stochasticity

### Style

Affect the entire image → Features are **controlled coherently** with  $y = (y_s, y_b)$

### Stochasticity

- Noise added to all pixels **independently** → **Only handle** stochastic variation



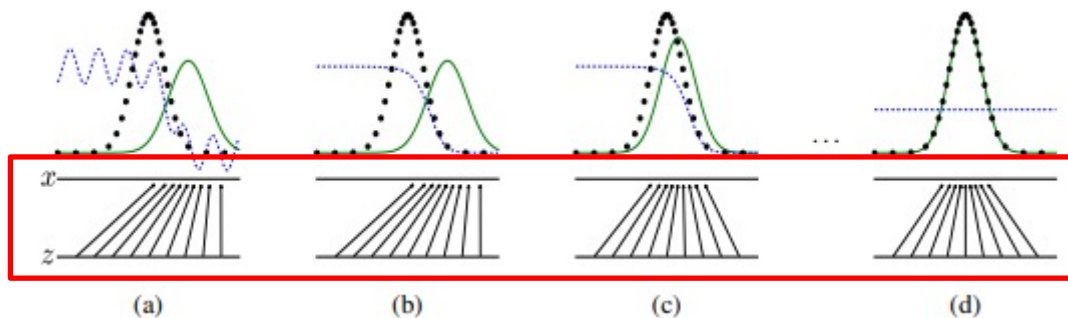
Unsupervised learning with appropriate use of **global and local channels**

# Disentanglement Studies

---

# Disentanglement Studies

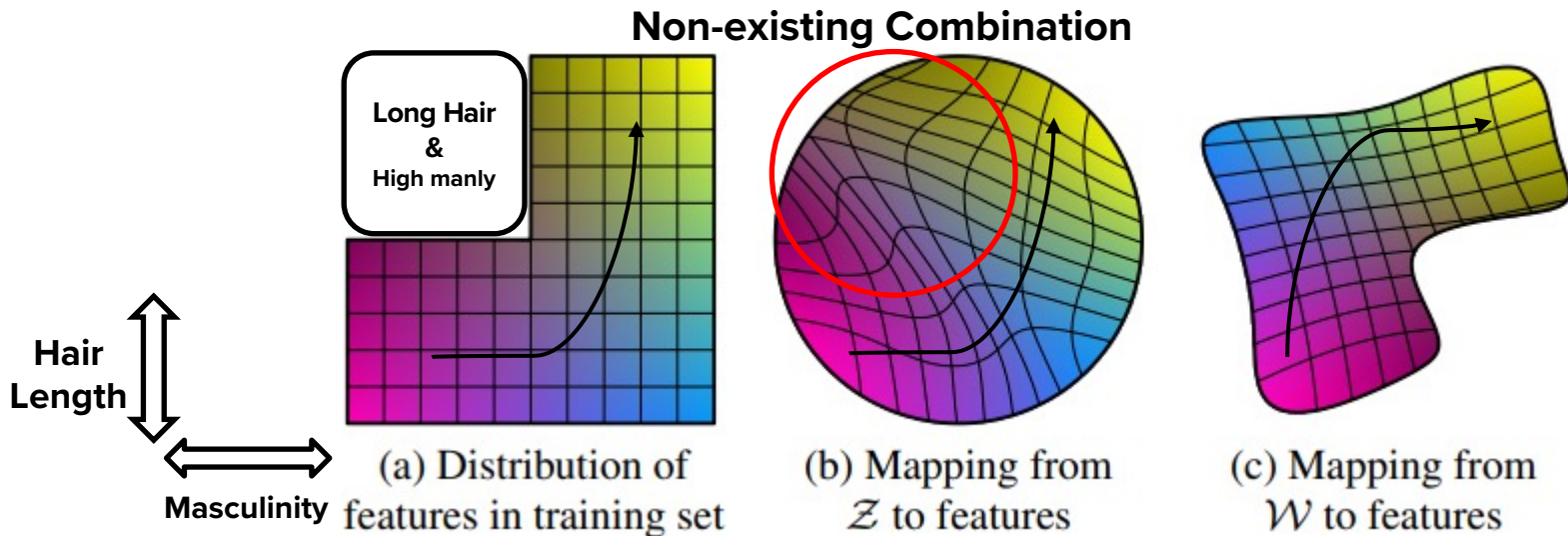
## □ Entanglement in $\mathcal{Z}$



Sampling probability is forced to follow the density function of training data

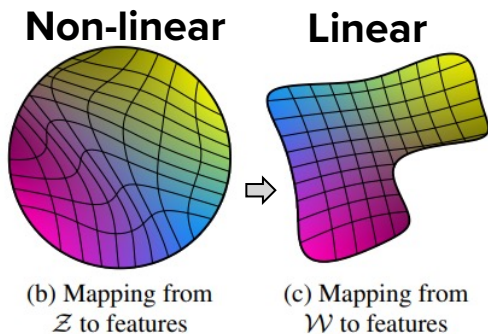
# Disentanglement Studies

## □ Achieve Disentanglement with $\mathcal{W}$



# Disentanglement Studies

## □ Achieve Disentanglement with $\mathcal{W}$



$$\mathbf{z} \rightarrow f(\mathbf{z}) \rightarrow \mathbf{w}$$

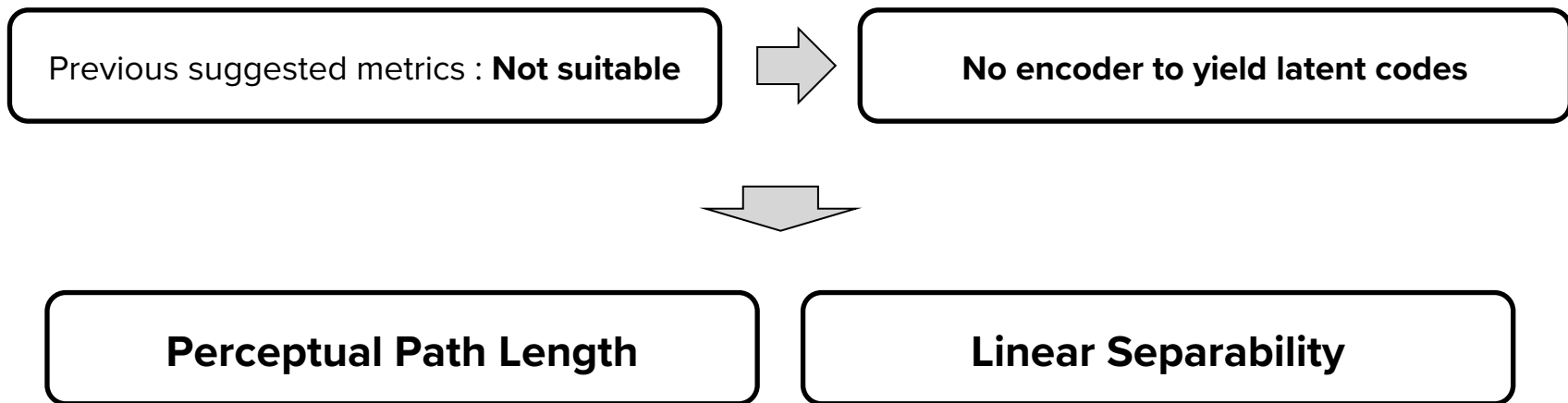
“Unwarp”  $\mathcal{W}$   $\rightarrow$  Factors of variation become **more linear**

**Mapping function**  $\rightarrow$  Help disentanglement through training

# Disentanglement Studies

---

## ❑ New metrics for quantifying disentanglement





# Disentanglement Studies

## ❑ Perceptual Path Length

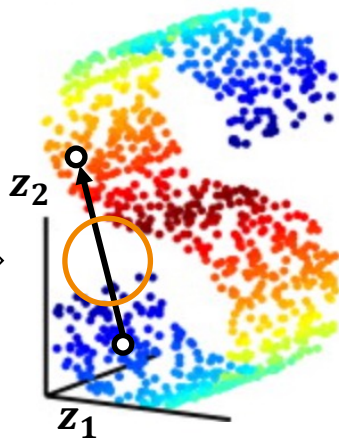
How to quantify disentanglement?



Interpolation path in latent space

Non-existing new feature in the latent space

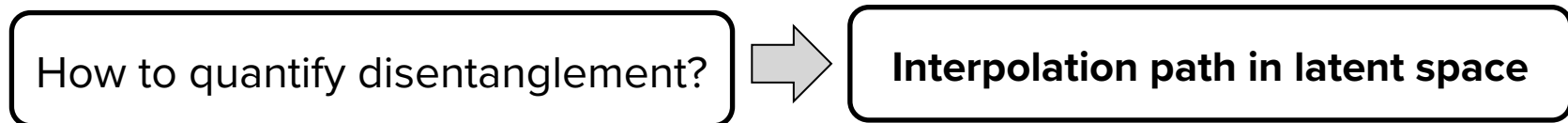
Entangled Latent Space



# Disentanglement Studies

---

## ❑ Perceptual Path Length



### Proposed Metric

The **amount of change** while performing **interpolation in the latent space**

# Disentanglement Studies

---

## ❑ Perceptual Path Length

Perceptually-based pairwise image distance



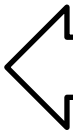
Weighted Difference between 2 VGG16 embeddings

Human-like Similarity Judgement on images

---

< Interpretation to the perceptual path length >

Low



High

High Similarity → Smooth Change

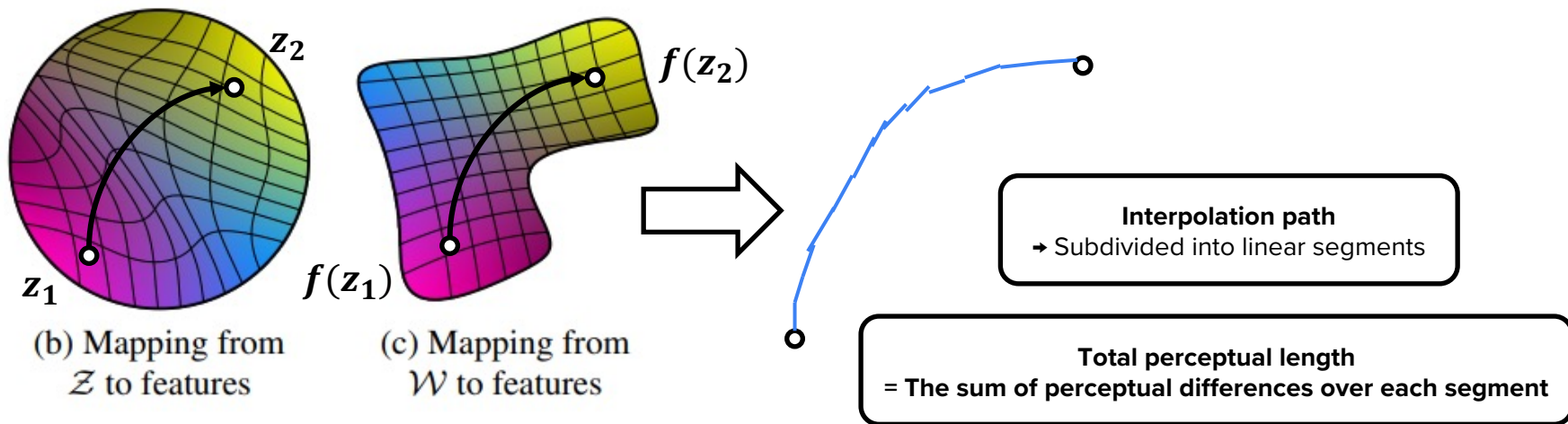
“More disentangled”

Low Similarity → Drastic Change

“Less disentangled”

# Disentanglement Studies

## □ Perceptual Path Length

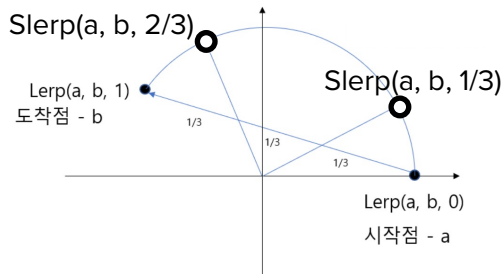


# Disentanglement Studies

## □ Perceptual Path Length

$$l_{\mathcal{Z}} = \mathbb{E} \left[ \frac{1}{\epsilon^2} d(G(\text{slerp}(\mathbf{z}_1, \mathbf{z}_2; t)), G(\text{slerp}(\mathbf{z}_1, \mathbf{z}_2; t + \epsilon))) \right]$$

Perceptual Path Length in  $\mathcal{Z}$

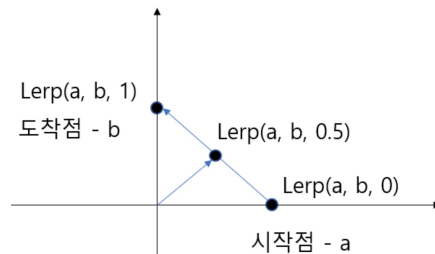


\* **slerp** : Spherical Interpolation

\* **lerp** : Linear Interpolation

$$l_{\mathcal{W}} = \mathbb{E} \left[ \frac{1}{\epsilon^2} d(g(\text{lerp}(f(\mathbf{z}_1), f(\mathbf{z}_2); t)), g(\text{lerp}(f(\mathbf{z}_1), f(\mathbf{z}_2); t + \epsilon))) \right]$$

Perceptual Path Length in  $\mathcal{W}$



# Disentanglement Studies

## □ Perceptual Path Length

| Method                                | Path length  |              | Separability |
|---------------------------------------|--------------|--------------|--------------|
|                                       | full         | end          |              |
| B Traditional generator $\mathcal{Z}$ | 412.0        | 415.3        | 10.78        |
| D Style-based generator $\mathcal{W}$ | 446.2        | 376.6        | 3.61         |
| E + Add noise inputs $\mathcal{W}$    | <b>200.5</b> | <b>160.6</b> | 3.54         |
| + Mixing 50% $\mathcal{W}$            | 231.5        | 182.1        | <b>3.51</b>  |
| F + Mixing 90% $\mathcal{W}$          | 234.0        | 195.9        | 3.79         |

Perceptual Path Length and Separability score  
depending on the architecture

| Method                        | FID         | Path length  |              | Separability |
|-------------------------------|-------------|--------------|--------------|--------------|
|                               |             | full         | end          |              |
| B Traditional 0 $\mathcal{Z}$ | 5.25        | 412.0        | 415.3        | 10.78        |
| Traditional 8 $\mathcal{Z}$   | 4.87        | 896.2        | 902.0        | 170.29       |
| Traditional 8 $\mathcal{W}$   | 4.87        | 324.5        | 212.2        | 6.52         |
| Style-based 0 $\mathcal{Z}$   | 5.06        | 283.5        | 285.5        | 9.88         |
| Style-based 1 $\mathcal{W}$   | 4.60        | 219.9        | 209.4        | 6.81         |
| Style-based 2 $\mathcal{W}$   | 4.43        | <b>217.8</b> | 199.9        | 6.25         |
| F Style-based 8 $\mathcal{W}$ | <b>4.40</b> | 234.0        | <b>195.9</b> | <b>3.79</b>  |

Results depending on  
the architecture and depth of mapping network

\* full :  $t \sim U(0, 1)$

\* end :  $t \in \{0, 1\}$

# Disentanglement Studies

---

## ❑ Linear Separability

If sufficiently disentangled..



Direction vectors consistently pointing to the corresponding factors of variation should exist

### Proposed Metric

How well latent-space points can be separated into distinct 2 sets

= How clearly each set can represent binary attribute

# Disentanglement Studies

---

## ❑ Linear Separability

Train auxiliary **classification** networks

Same as the discriminator in ProGAN



Obtain **binary attributes**

Generate 200000 images to evaluate one of the attributes



**Perform classification**



Leave **the most confident** samples  
(100000 samples)



# Disentanglement Studies

## ❑ Linear Separability

$H(X|Y)$  : Conditional Entropy

$X$  : The class **predicted** by linear SVM based on  $\mathbf{z}$

$Y$  : The **true class** predicted by pre-trained classifier

➡ The **amount of additional information** to determine the true class

➡  $\exp(\sum_i H(X_i|Y_i))$  : Final **Separability Score**

# Disentanglement Studies

## □ Linear Separability

| Method                                | Path length  |              | Separability |
|---------------------------------------|--------------|--------------|--------------|
|                                       | full         | end          |              |
| B Traditional generator $\mathcal{Z}$ | 412.0        | 415.3        | 10.78        |
| D Style-based generator $\mathcal{W}$ | 446.2        | 376.6        | 3.61         |
| E + Add noise inputs $\mathcal{W}$    | <b>200.5</b> | <b>160.6</b> | 3.54         |
| + Mixing 50% $\mathcal{W}$            | 231.5        | 182.1        | <b>3.51</b>  |
| F + Mixing 90% $\mathcal{W}$          | 234.0        | 195.9        | 3.79         |

Perceptual Path Length and Separability score  
depending on the architecture

| Method                        | FID         | Path length  |              | Separability |
|-------------------------------|-------------|--------------|--------------|--------------|
|                               |             | full         | end          |              |
| B Traditional 0 $\mathcal{Z}$ | 5.25        | 412.0        | 415.3        | 10.78        |
| Traditional 8 $\mathcal{Z}$   | 4.87        | 896.2        | 902.0        | 170.29       |
| Traditional 8 $\mathcal{W}$   | 4.87        | 324.5        | 212.2        | 6.52         |
| Style-based 0 $\mathcal{Z}$   | 5.06        | 283.5        | 285.5        | 9.88         |
| Style-based 1 $\mathcal{W}$   | 4.60        | 219.9        | 209.4        | 6.81         |
| Style-based 2 $\mathcal{W}$   | 4.43        | <b>217.8</b> | 199.9        | 6.25         |
| F Style-based 8 $\mathcal{W}$ | <b>4.40</b> | 234.0        | <b>195.9</b> | <b>3.79</b>  |

Results depending on  
the architecture and depth of mapping network

# Results

---

# Results

---

## □ Comparing with other generators

| Method                              | CelebA-HQ   | FFHQ        |
|-------------------------------------|-------------|-------------|
| A Baseline Progressive GAN [30]     | 7.79        | 8.04        |
| B + Tuning (incl. bilinear up/down) | 6.11        | 5.25        |
| C + Add mapping and styles          | 5.34        | 4.85        |
| D + Remove traditional input        | 5.07        | 4.88        |
| E + Add noise inputs                | <b>5.06</b> | 4.42        |
| F + Mixing regularization           | 5.17        | <b>4.40</b> |

FID scores for various generators

# Results

## □ Perceptual Path Length & Separability

| Method                                | Path length  |              | Separa-<br>bility |
|---------------------------------------|--------------|--------------|-------------------|
|                                       | full         | end          |                   |
| B Traditional generator $\mathcal{Z}$ | 412.0        | 415.3        | 10.78             |
| D Style-based generator $\mathcal{W}$ | 446.2        | 376.6        | 3.61              |
| E + Add noise inputs $\mathcal{W}$    | <b>200.5</b> | <b>160.6</b> | 3.54              |
| + Mixing 50% $\mathcal{W}$            | 231.5        | 182.1        | <b>3.51</b>       |
| F + Mixing 90% $\mathcal{W}$          | 234.0        | 195.9        | 3.79              |

| Method                        | FID         | Path length  |              | Separa-<br>bility |
|-------------------------------|-------------|--------------|--------------|-------------------|
|                               |             | full         | end          |                   |
| B Traditional 0 $\mathcal{Z}$ | 5.25        | 412.0        | 415.3        | 10.78             |
| Traditional 8 $\mathcal{Z}$   | 4.87        | 896.2        | 902.0        | 170.29            |
| Traditional 8 $\mathcal{W}$   | 4.87        | 324.5        | 212.2        | 6.52              |
| Style-based 0 $\mathcal{Z}$   | 5.06        | 283.5        | 285.5        | 9.88              |
| Style-based 1 $\mathcal{W}$   | 4.60        | 219.9        | 209.4        | 6.81              |
| Style-based 2 $\mathcal{W}$   | 4.43        | <b>217.8</b> | 199.9        | 6.25              |
| F Style-based 8 $\mathcal{W}$ | <b>4.40</b> | 234.0        | <b>195.9</b> | <b>3.79</b>       |



Even the traditional model has **performed better** when  $\mathcal{W}$  is introduced



Best results are derived from the proposed model →  $\mathcal{W}$  is **less entangled** than  $\mathcal{Z}$

# Results

## ❑ Unresolved problem



Figure 6: Drop-like artefacts in generated images

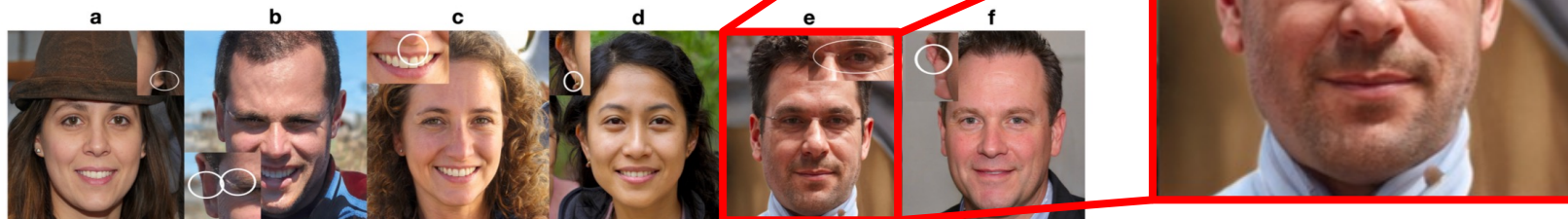


Figure 7: Other visual artefacts in the CelebA dataset. Notice the lack of symmetry and visible marks in ears. There's also some artefacts in the eyes, glasses, and lips

# Thank you

---