StyleCLIP: Text-Driven Manipulation of StyleGAN Imagery

Input:

Text prompt: "Mohawk hairstyle"   "Without makeup"   "Cute cat"   "Lion"
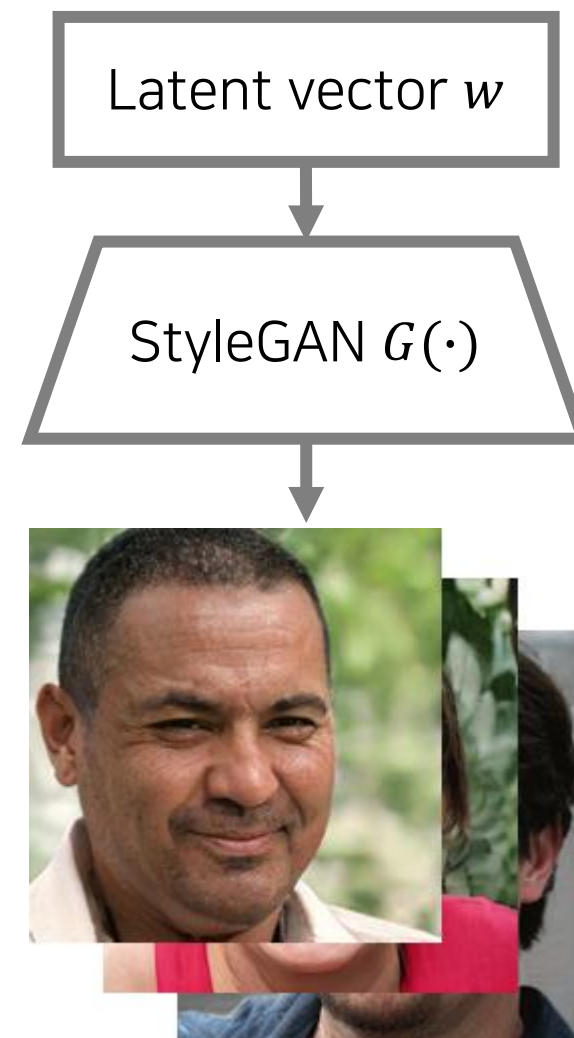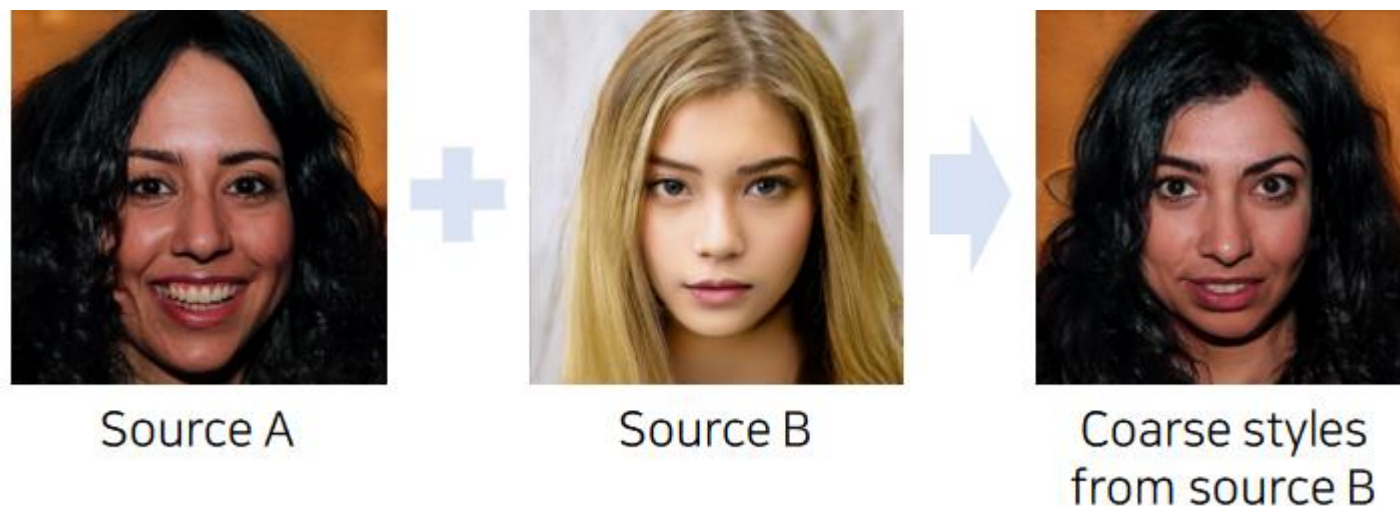
Result:

StyleCLIP: Text-Driven Manipulation of StyleGAN Imagery (2021)

# Background: StyleGAN (CVPR 2019)

- Propose an efficient architecture to generate high-quality images.

- Present a 1024 X 1024 high-quality face dataset (**FFHQ**).

- Improve the **disentanglement** of semantic features.



Source A    Source B    Coarse styles from source B



Latent vector $w$

StyleGAN $G(\cdot)$

Image B

Image A

Coarse styles

Middle styles

Fine styles

512

4

4

10

18 X 512 Latent Vector $w^+$

Background: Face Manipulation Using StyleGAN

① Encoding step    ② Manipulation step

Input image $x$    $G(w_{encoded})$    $G(w_{manipulated})$
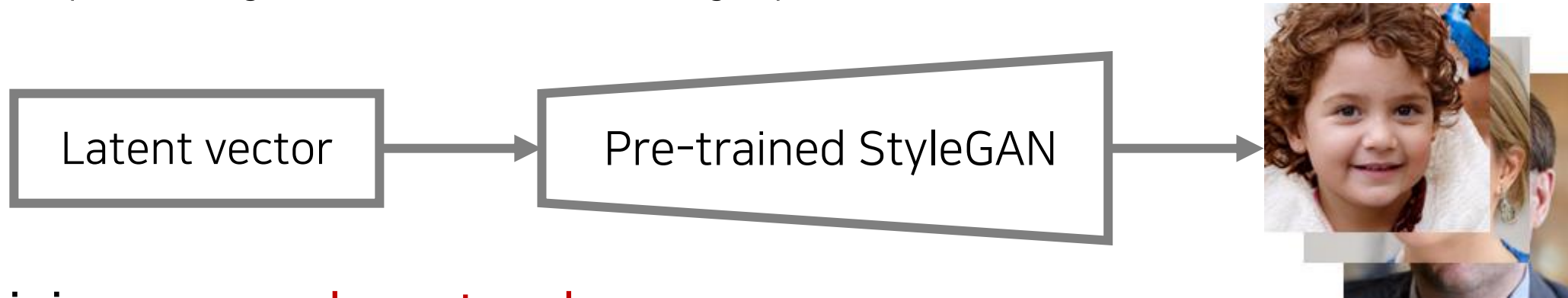
Latent vector $w$

StyleGAN

Generated $G(w)$

Image $x$

VGG16 feature extractor

VGG16 feature extractor

MSE loss

Gradient descent and update latent

Image2StyleGAN: How to Embed Images Into the StyleGAN Latent Space? (ICCV 2019)

## 1) Latent vector dataset generation

- Prepare a large number of (latent, image) pairs.



## 2) Training an encoder network

- We can interpolate between encoded latent $a$ and encoded latent $b$.

1) Learns a boundary of an attribute (such as gender, age).



Random Sample → Image Space $\mathcal{X}$ → Scoring Function → Semantic Space $\mathcal{S}$

Latent Space $\mathcal{Z}$

2) Update a latent vector across the boundary.



$-n$          $n$

- CLIP jointly trains an image encoder and a text encoder using a large dataset.

- CLIP jointly trains an image encoder and a text encoder using a large dataset.



(1) Contrastive pre-training

(2) Create dataset classifier from label text

(3) Use for zero-shot prediction

StyleCLIP (StyleGAN + CLIP): Main Idea

Latent vector w

StyleGAN

Generated $G(w)$

Text Encoder

Image Encoder

Update

"Without makeup"

Similarity

CLIP embedding space

Explained by Dongbin Na

StyleCLIP: Text-Driven Manipulation of StyleGAN Imagery (2021)

- StyleCLIP is a text-based interface for StyleGAN image manipulation.
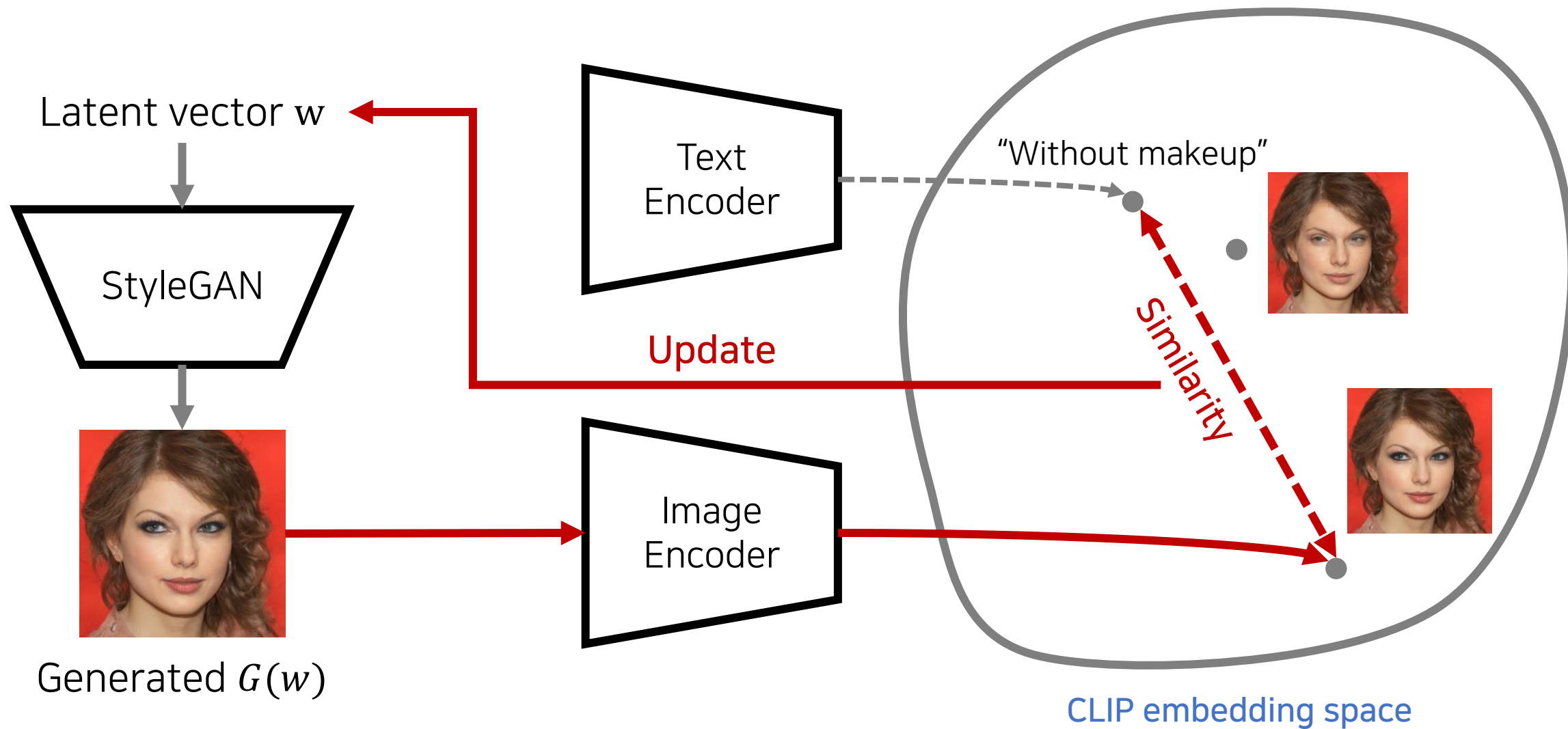
- Propose three methods that do not require such manual effort.

  1. Introduce an **optimization method** that utilizes a CLIP-based loss.

  2. Introduce a **latent mapper** that infers a text-guided latent manipulation.

  3. Present a method for mapping text prompts to input-agnostic **global directions**.

| | pre-proc. | train time | infer. time | input image dependent | latent space |
|---|---|---|---|---|---|
| optimizer | – | – | 98 sec | yes | $\mathcal{W}+$ |
| mapper | – | 10 – 12h | 75 ms | yes | $\mathcal{W}+$ |
| global dir. | 4h | – | 72 ms | no | $\mathcal{S}$ |

- Latent optimization: a simple approach for leveraging CLIP to guide image manipulation.

$$\arg\min_{w \in \mathcal{W}+} D_{\text{CLIP}}(G(w), t) + \lambda_{\text{L2}} \|w - w_s\|_2 + \lambda_{\text{ID}} \mathcal{L}_{\text{ID}}(w)$$

For manipulation     For similarity to the input image

$$\mathcal{L}_{\text{ID}}(w) = 1 - \langle R(G(w_s)), R(G(w)) \rangle$$
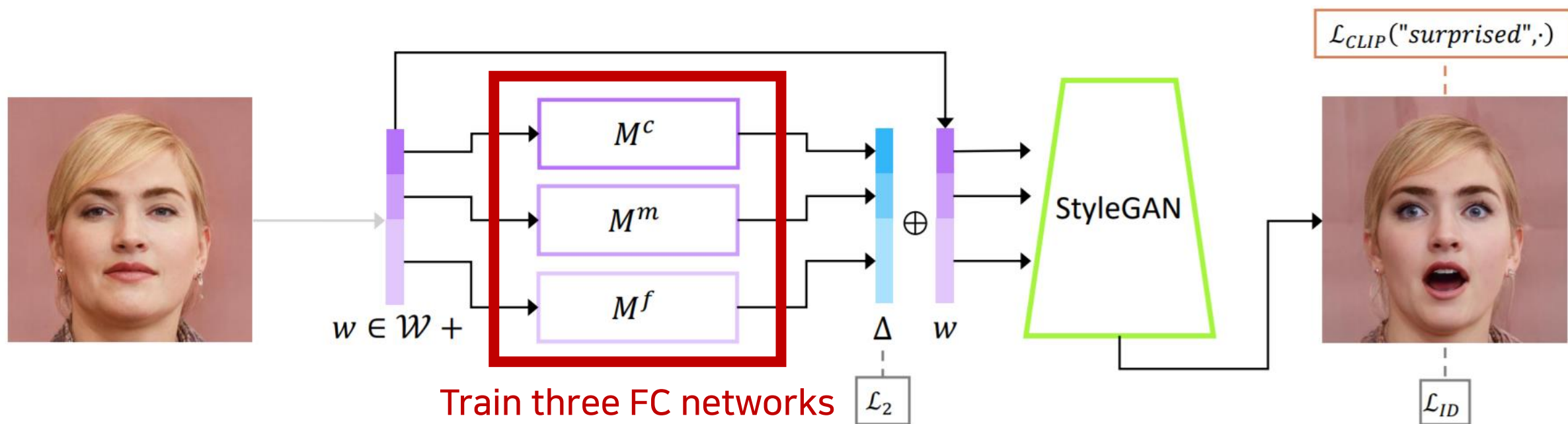
$R$: Pretrained *ArcFace* network

$D_{clip}$: Cosine distance between the CLIP embeddings

The optimization method requires **200 – 300 iterations** that spend <span style="color:red">several minutes</span>.
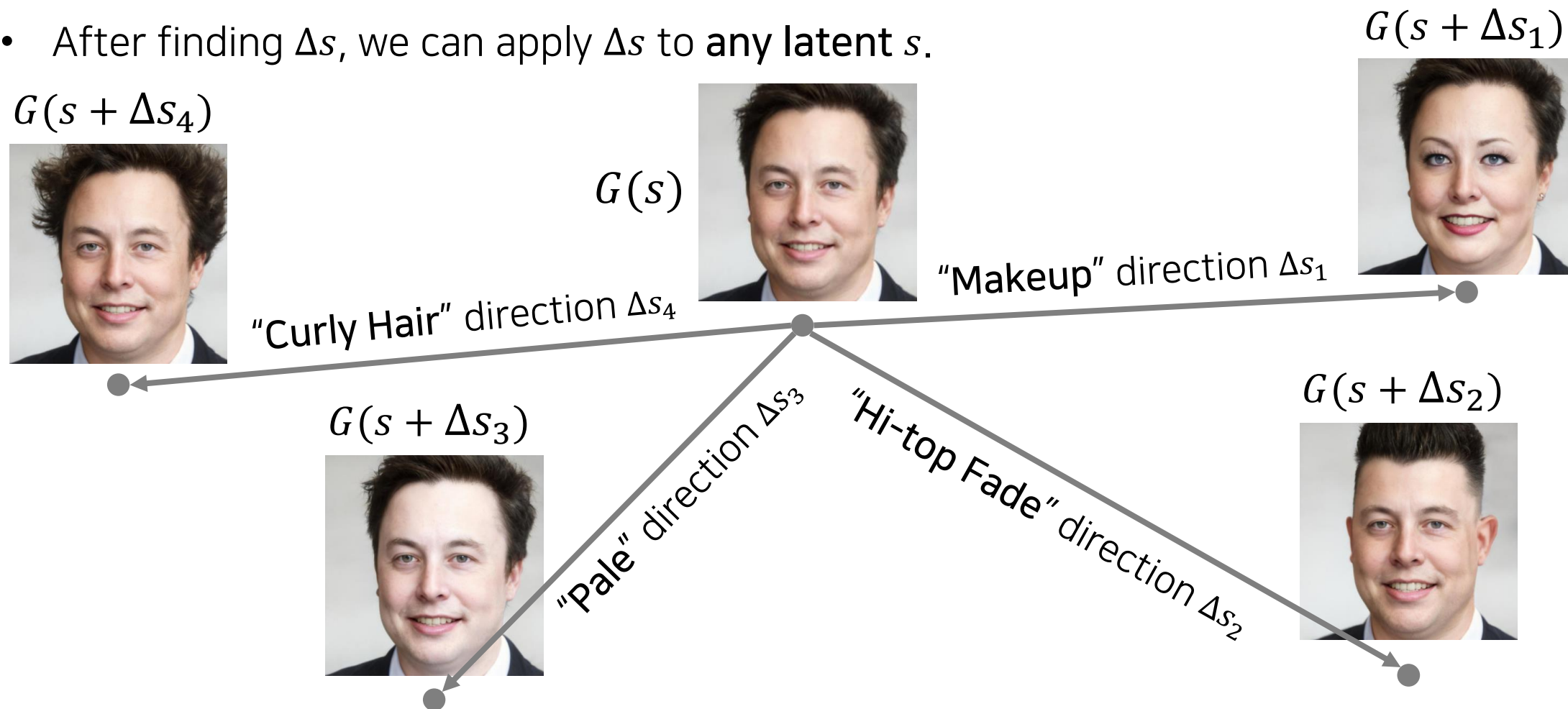
- **Latent mapper** is trained **to manipulate the desired attributes of the image** as indicated by the text prompt $t$, while preserving the other visual attributes of the input image.



Train three FC networks

After trained per text prompt (10 hours), the mapper manipulates attributes in one forward.

- Find a **global direction** $\Delta s \in S$ in a StyleGAN's style space $S$.

  - After finding $\Delta s$, we can apply $\Delta s$ to **any latent** $s$.



$G(s + \Delta s_4)$

$G(s)$

$G(s + \Delta s_1)$

"Makeup" direction $\Delta s_1$

"Curly Hair" direction $\Delta s_4$

$G(s + \Delta s_3)$

$G(s + \Delta s_2)$

"Pale" direction $\Delta s_3$

"Hi-top Fade" direction $\Delta s_2$

# Comparisons

- StyleCLIP is **not limited** to preset manipulation directions.
  - However, StyleCLIP shows **competitive results** on even common attributes.

StyleCLIP: Text-Driven Manipulation of StyleGAN Imagery (2021)