

Image Generation Using Disentanglement

Ron Mokady

What is Disentanglement?

No formal definition which is widely accepted.

Our definition: a change in a single underlying factor of variation in the sample x should lead to a change in a single factor in the learned representation $r(x)$.

Disentanglement for Generative model

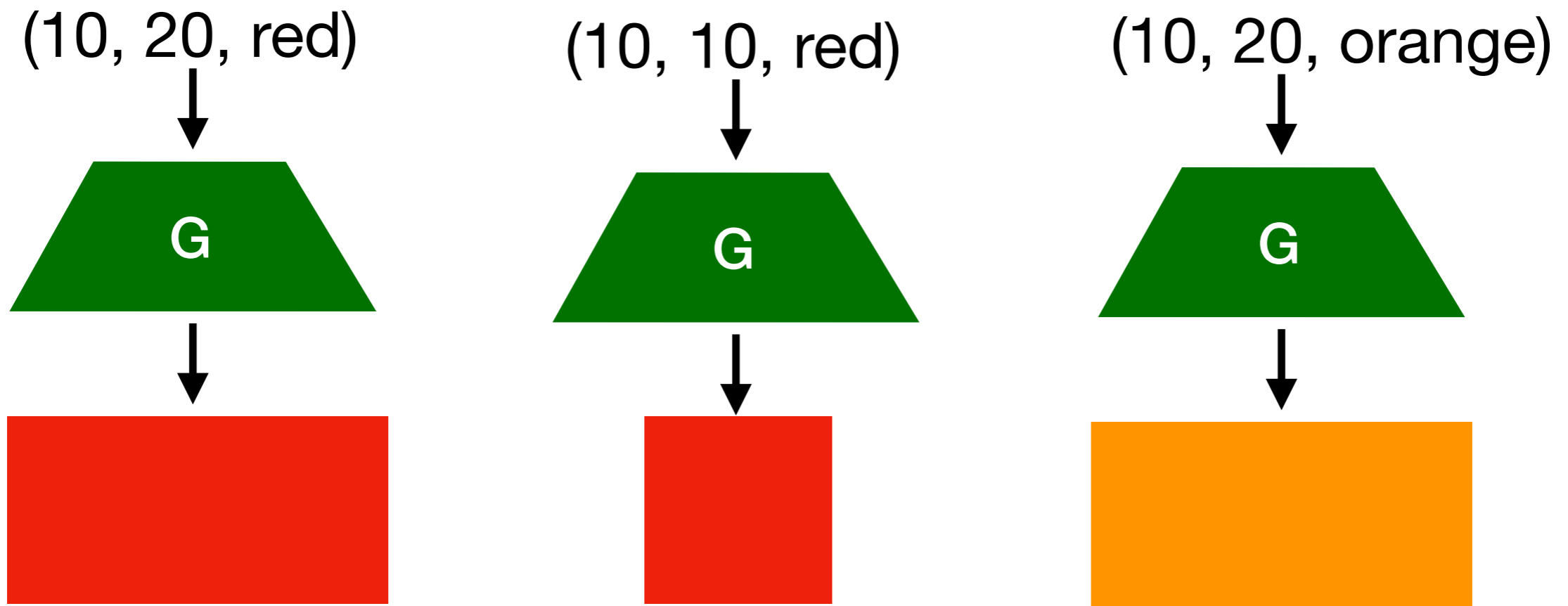
A change in a single underlying factor of variation in the generated image $G(z)$ should lead to a change in a single factor in the learned representation z .

Factors could be style, content, rotation, size, etc.



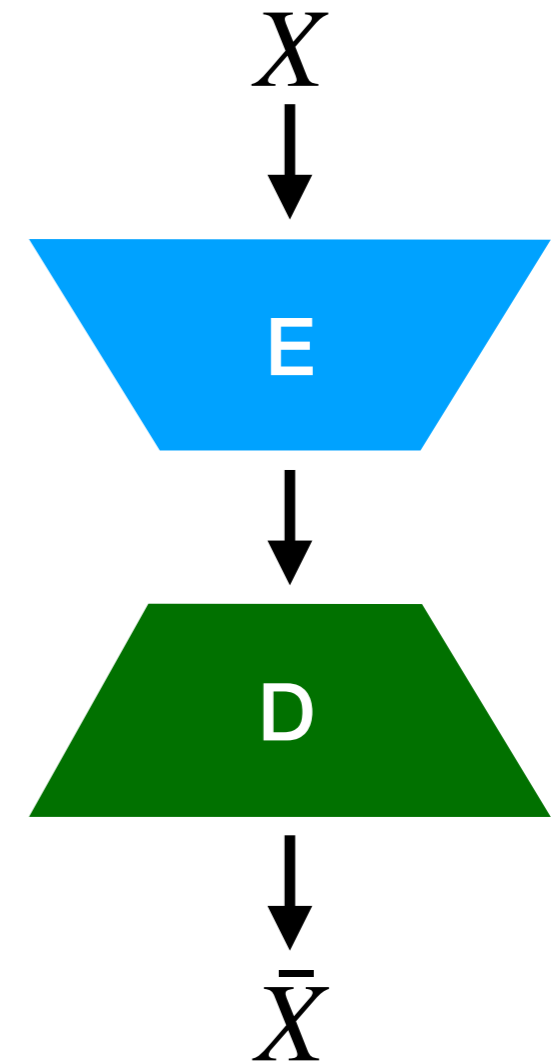
Example - Rectangles

$z = (\text{height}, \text{width}, \text{color})$



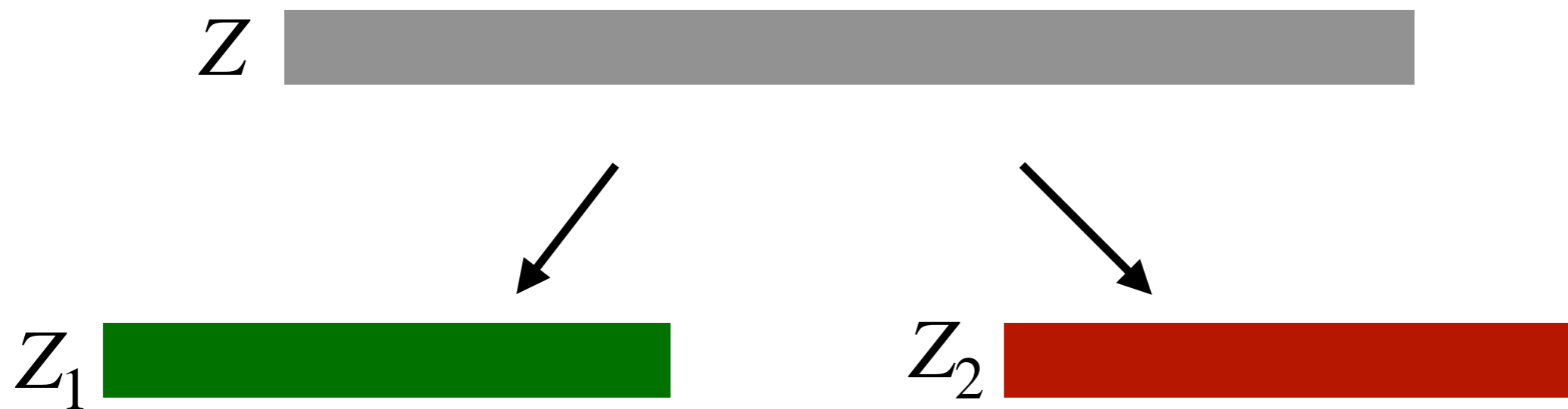
Example - Autoencoder

Using standard Autoencoder, we likely to get entangled latent space.



Separate the Latent Space

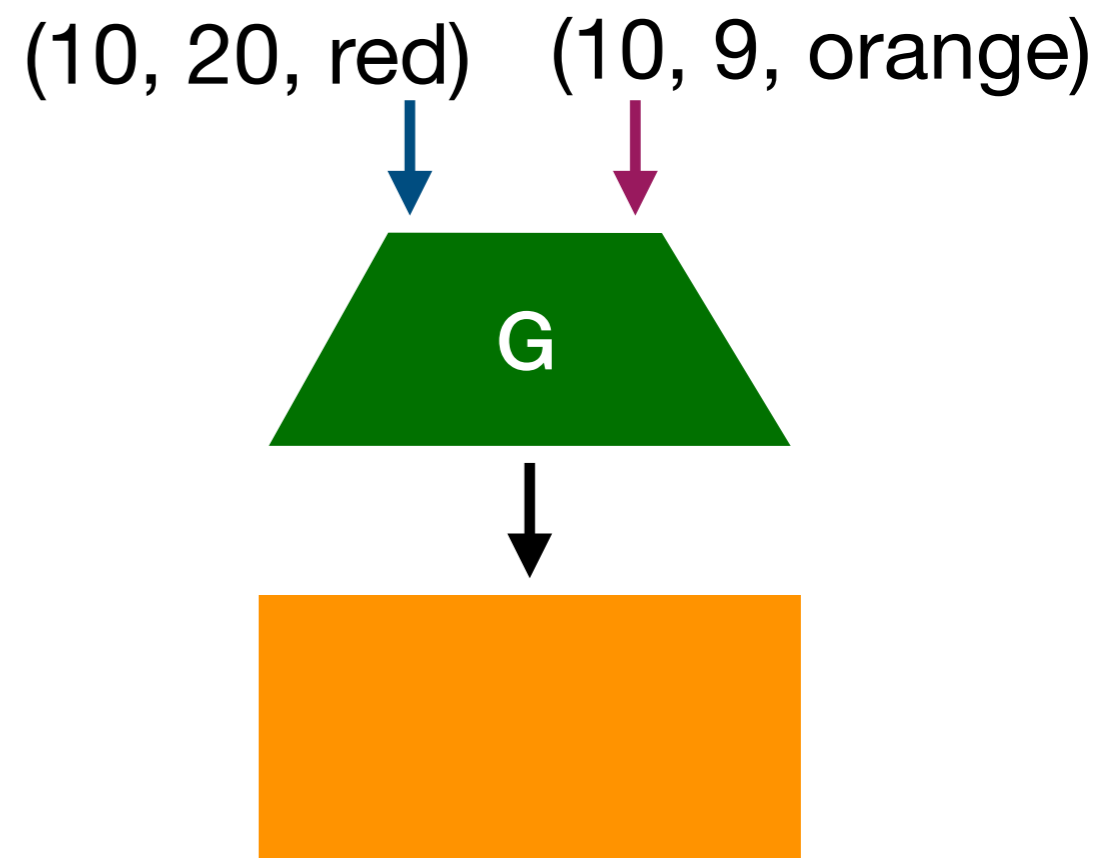
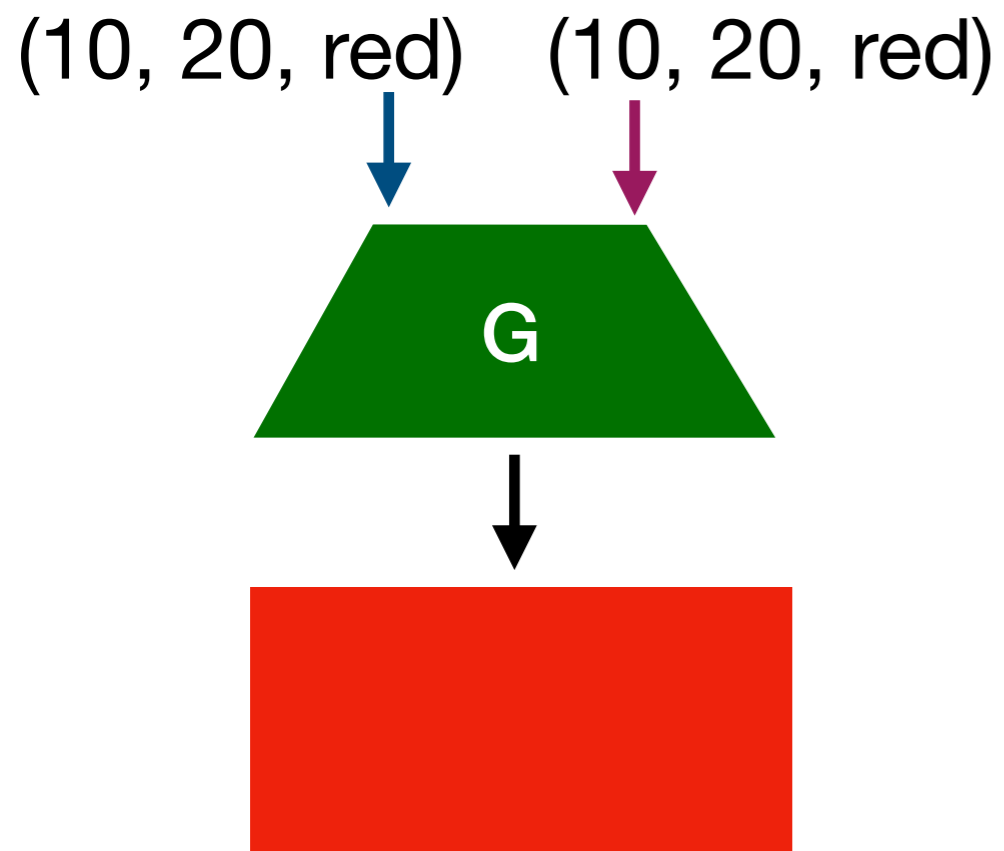
We usually want to separate the latent space, such that each part control one or more factors.



Generative model

The disentanglement is with respect to the generative model.

For example:



Supervision

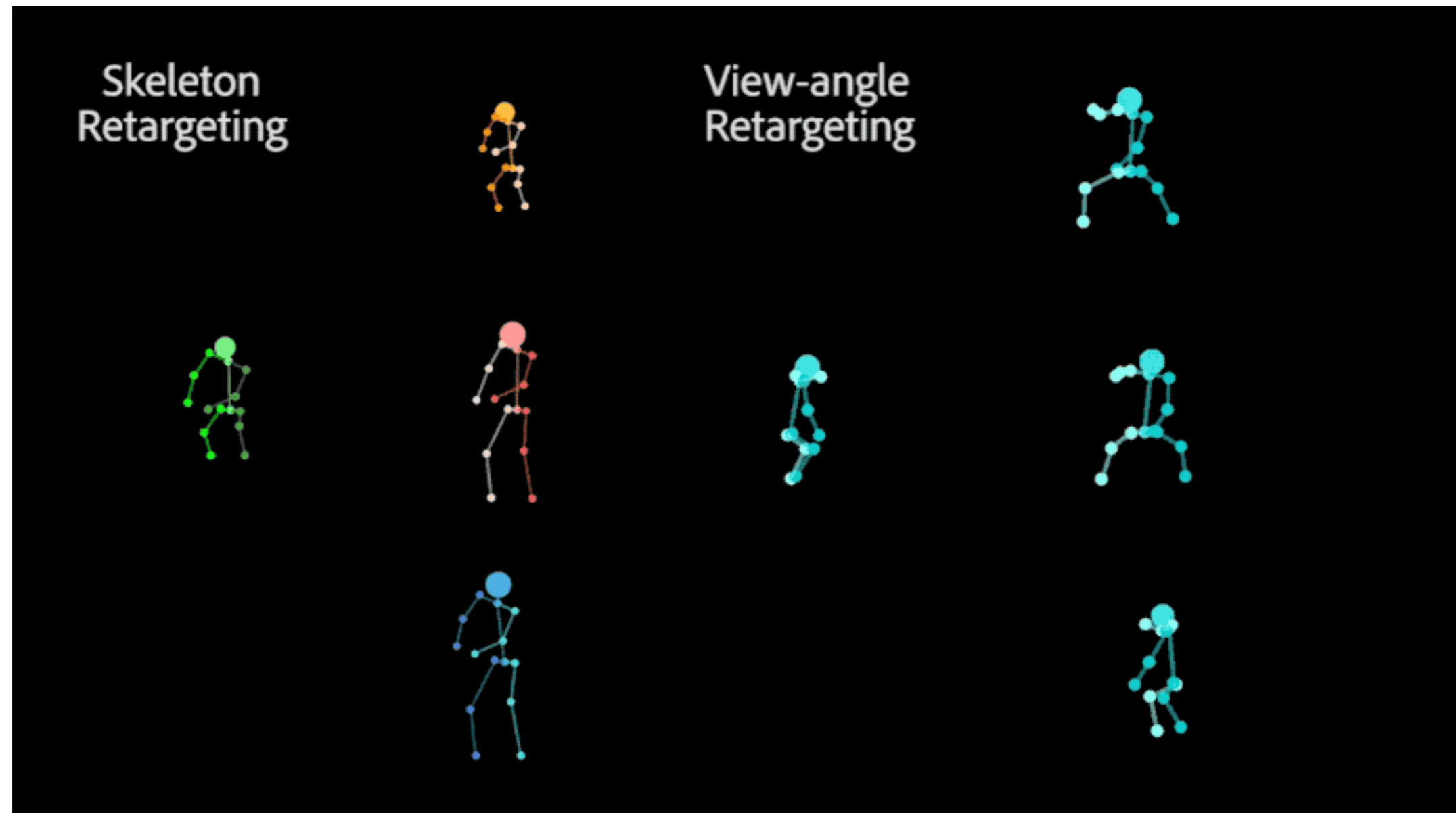


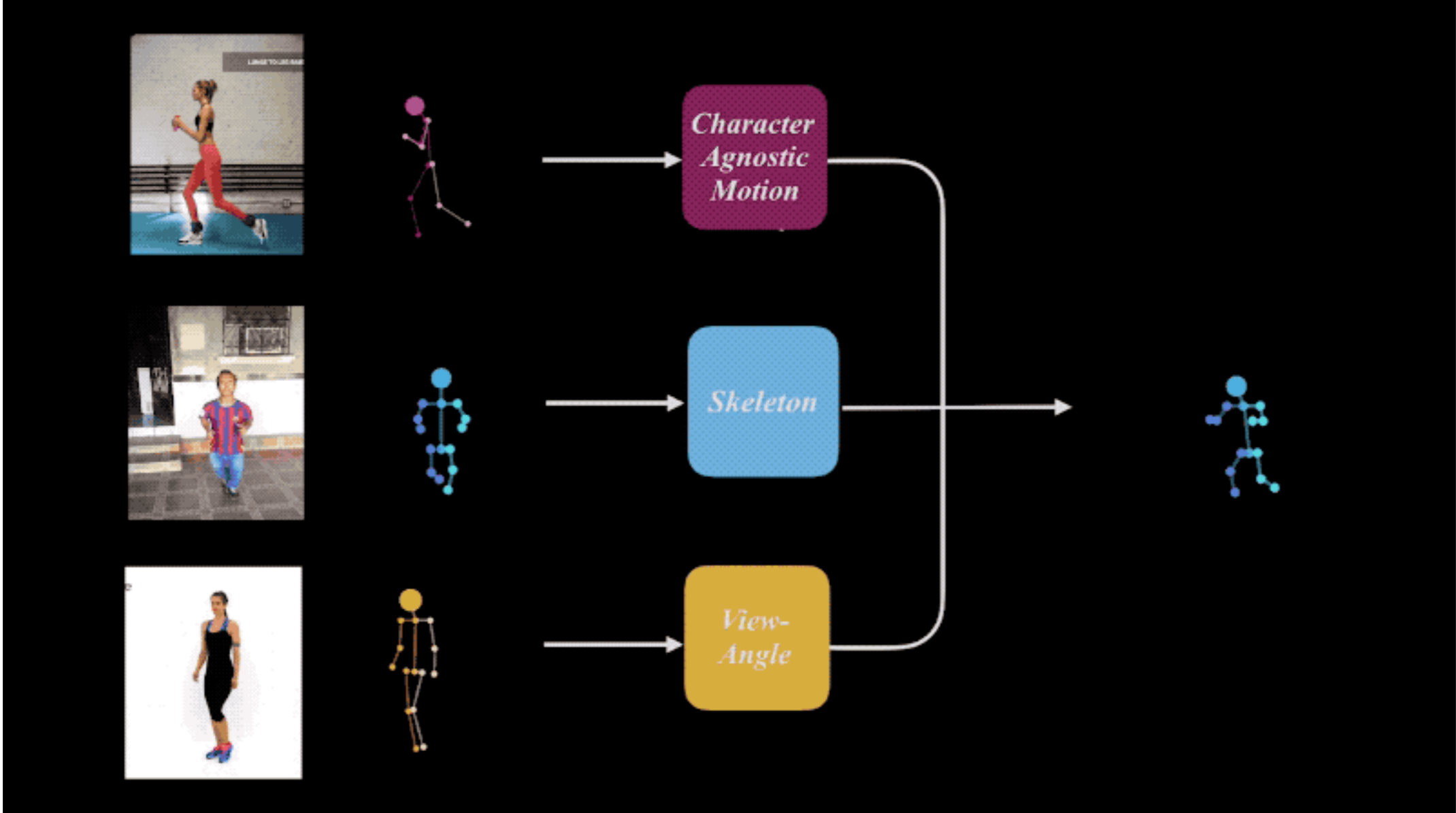
Unsupervised

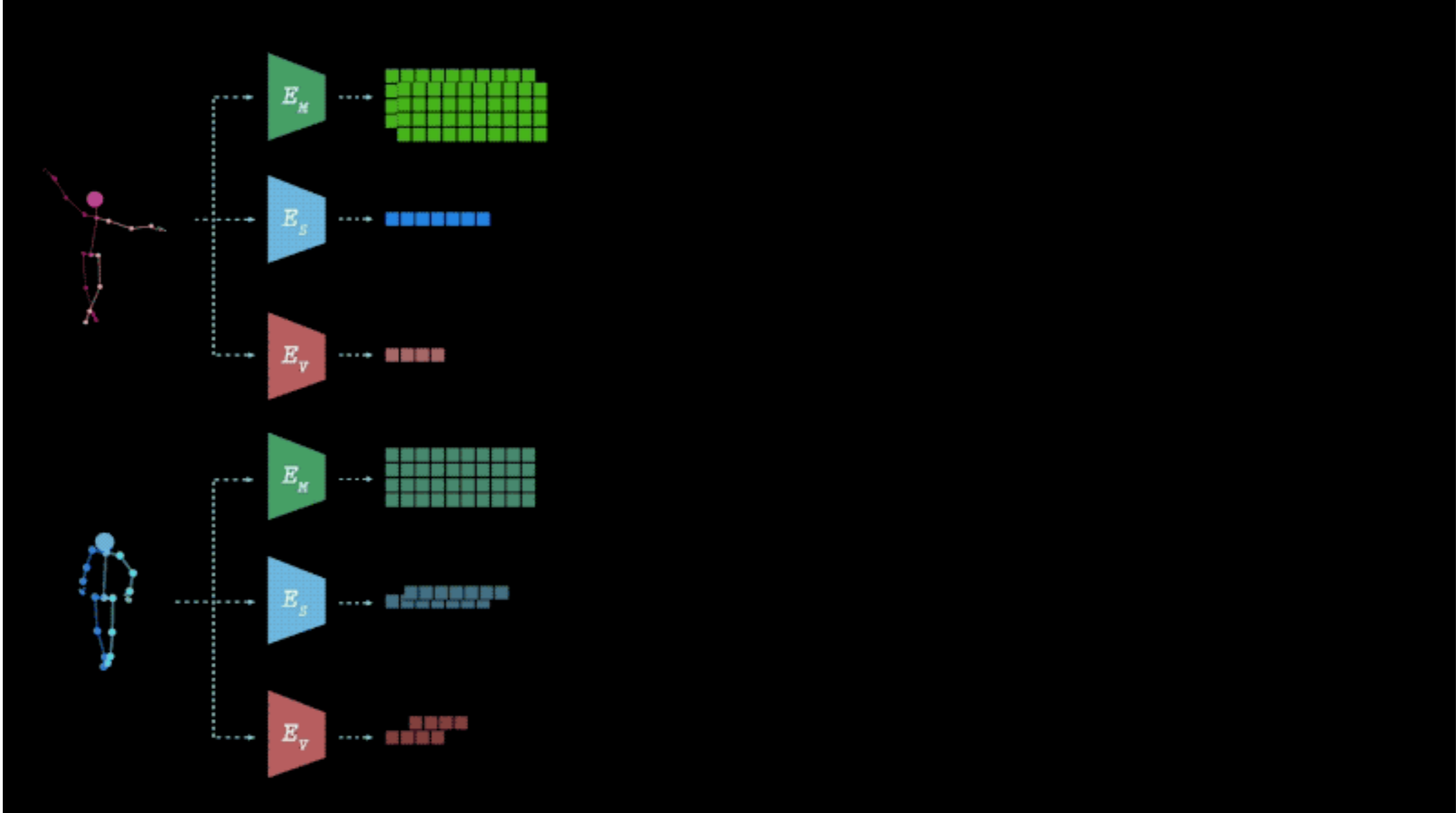
**Fully-
Supervised**

Learning Character-Agnostic Motion for Motion Retargeting in 2D

Kfir Aberman, Rundi Wu, Dani Lischinski, Baoquan Chen, Daniel Cohen-Or

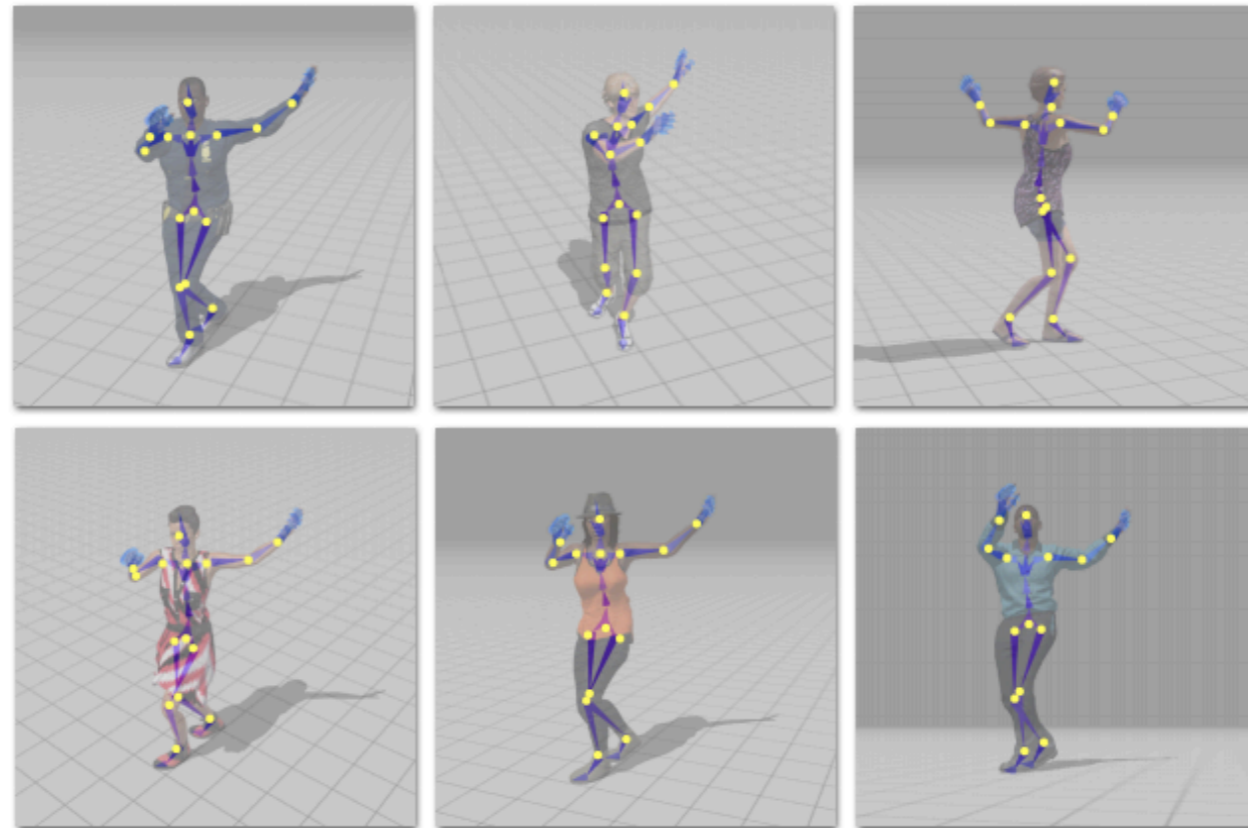






Synthetic dataset

Multiple samples of each motion, as performed by the different characters, and these motions can be projected to 2D, from arbitrary view angles.



Fully-Supervised

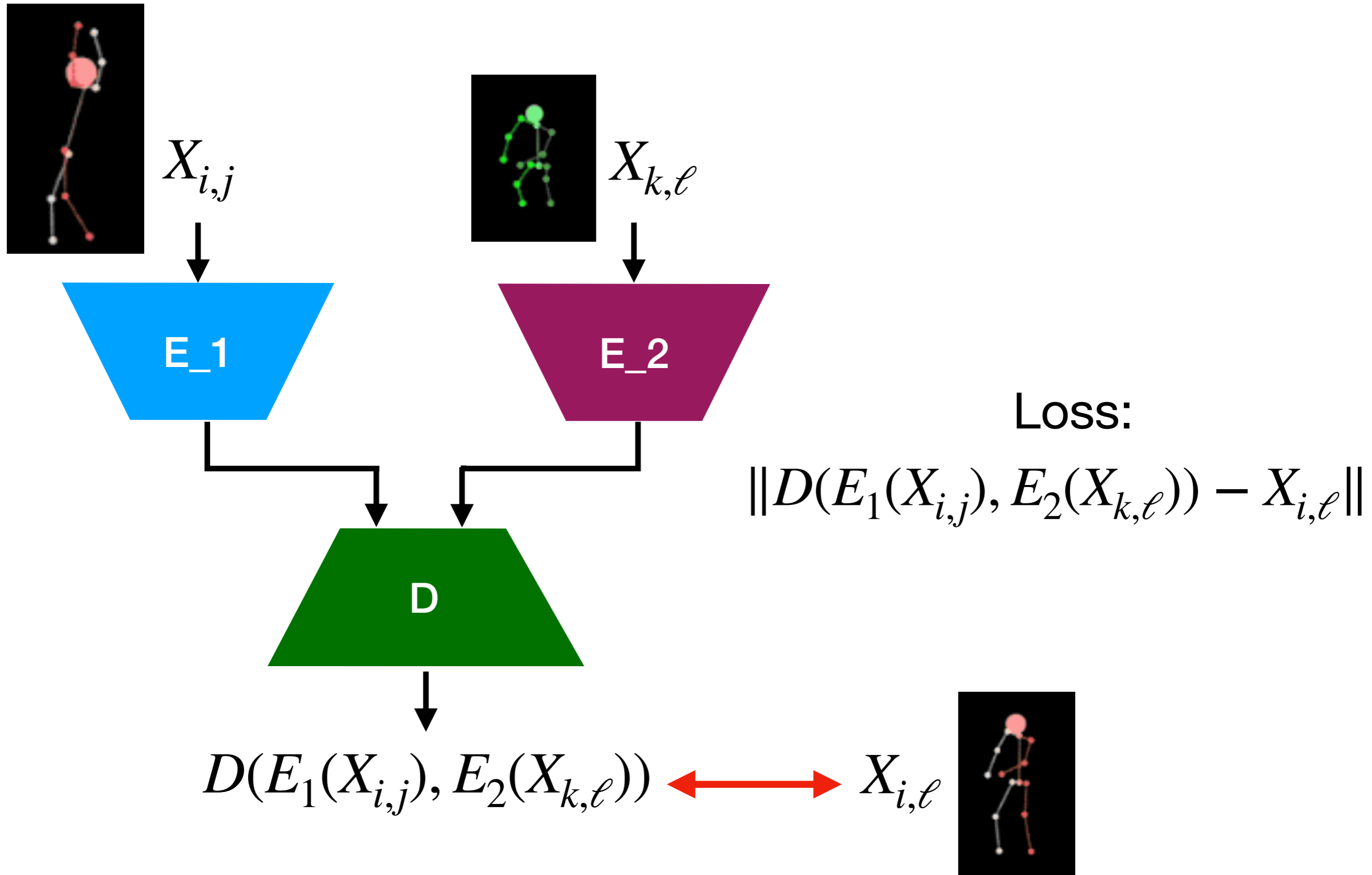
Example: $\{X_{a,b}\}_{a \in A, b \in B}$

Let A and B be two independent factors of variation, then the dataset contains triplets:

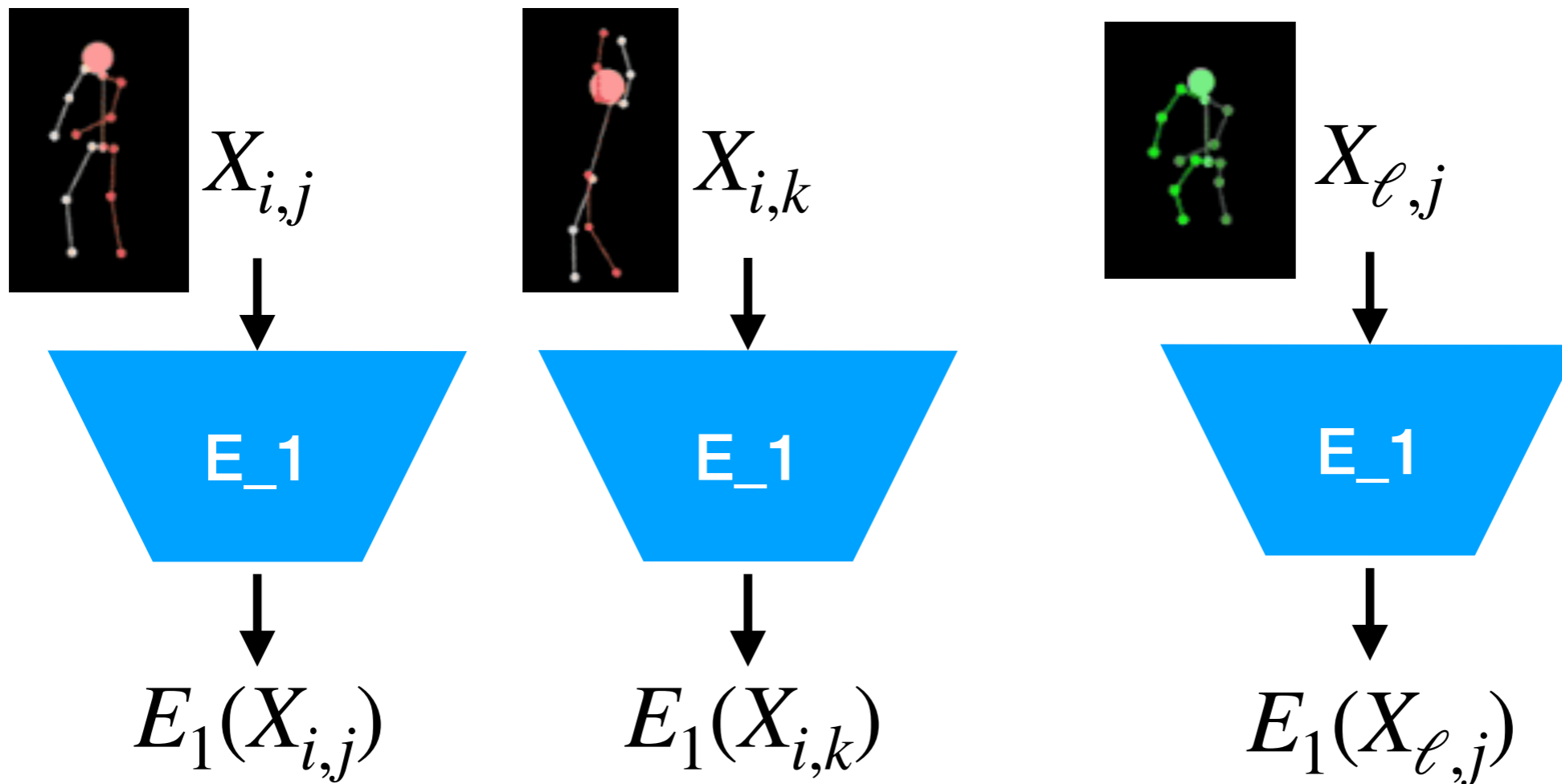
$$(X_{i,j}, X_{i,k}, X_{\ell,j})$$



Cross Loss



Triplet Loss



Loss:

$$\max (\|E_1(X_{i,j}) - E_1(X_{i,k})\| - \|E_1(X_{i,j}) - E_1(X_{\ell,j})\| + \varepsilon, 0)$$

Supervision



Unsupervised



Limited use



?

Fully-Supervised



Only synthetic datasets

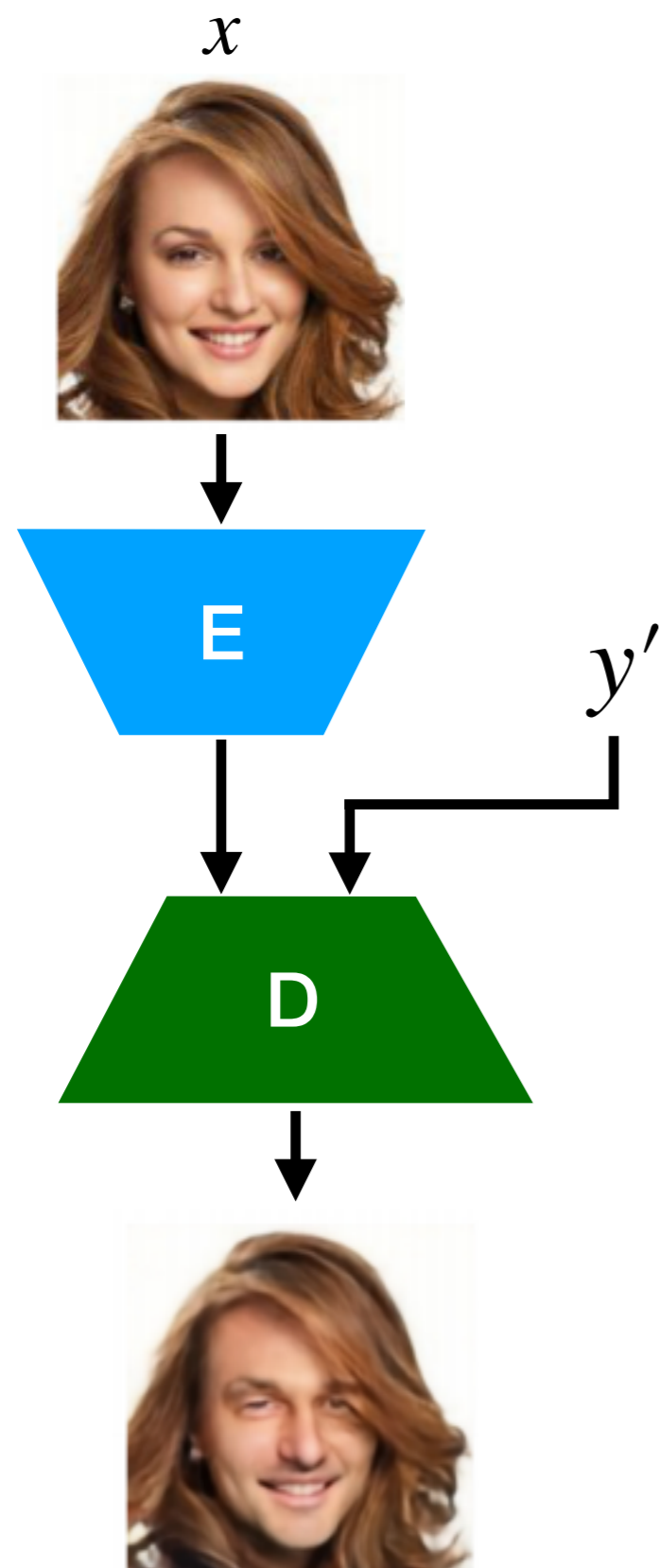
Fader Networks: Manipulating Images by Sliding Attributes

Guillaume Lample, Neil Zeghidour, Nicolas Usunier, Antoine Bordes, Ludovic Denoyer, Marc'Aurelio Ranzato

Set of images and several attributes for each image.

$$\{(x, y)\}_{x \in X, y \in \{0, 1\}^k}$$



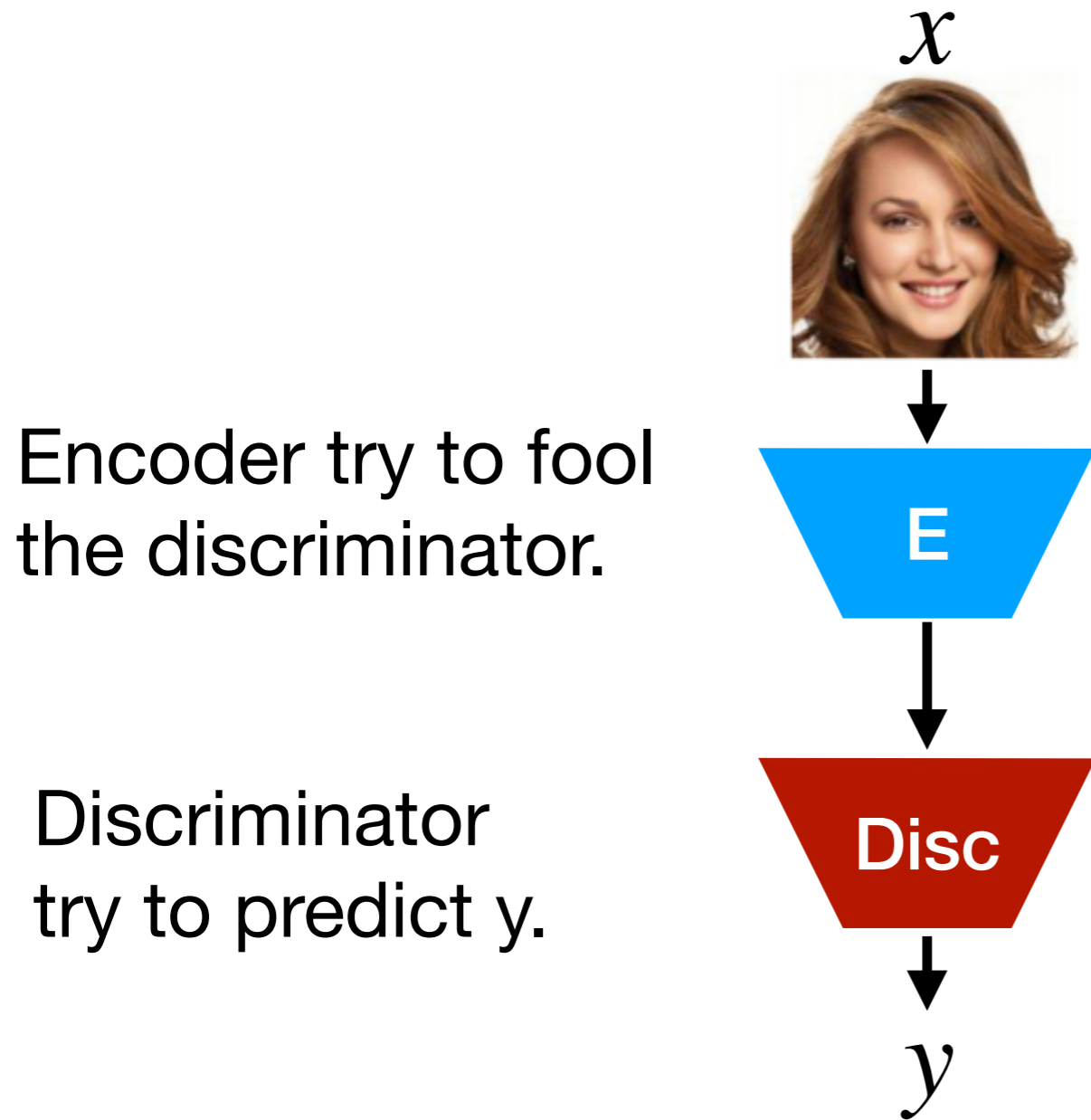


We would like to find representation of the images which is invariant to the image attributes.

Adversarial Loss

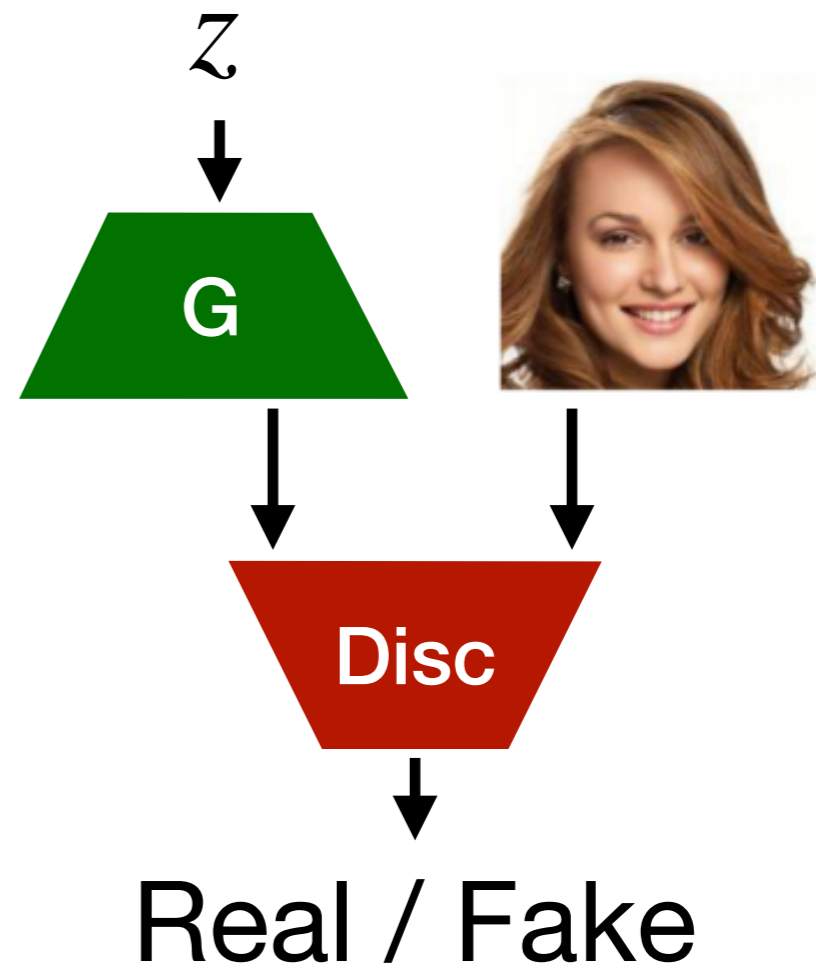
FaderNet:

GAN:



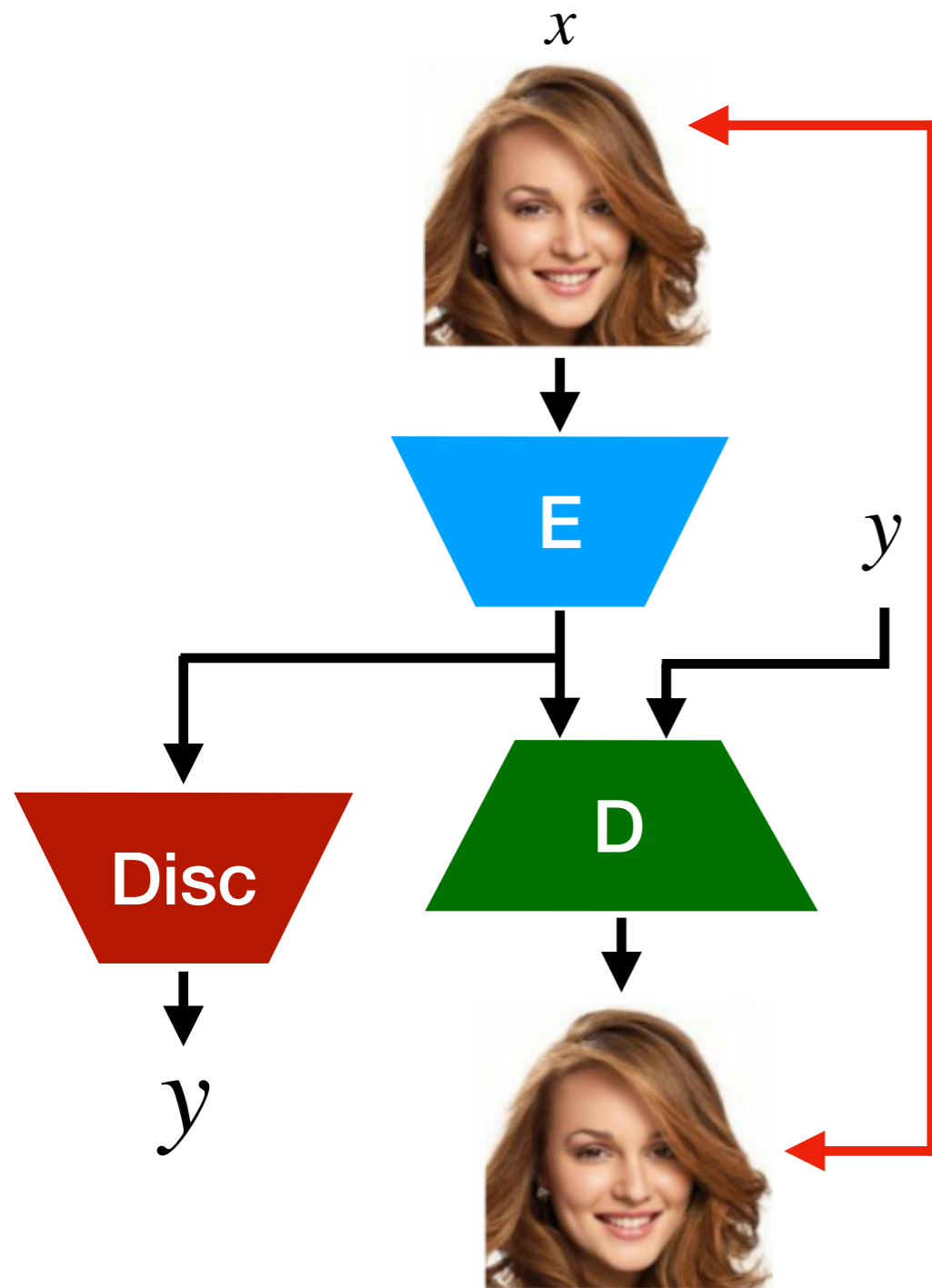
Encoder try to fool the discriminator.

Discriminator try to predict y .



Real / Fake

Loss Function



1. Adversarial loss:

Discriminator try to predict y .
Encoder try to fool the discriminator.

2. Reconstruction loss:

$$\mathcal{L}_R = \|D(E(x), y) - x\|$$

Sliding



Young → Old



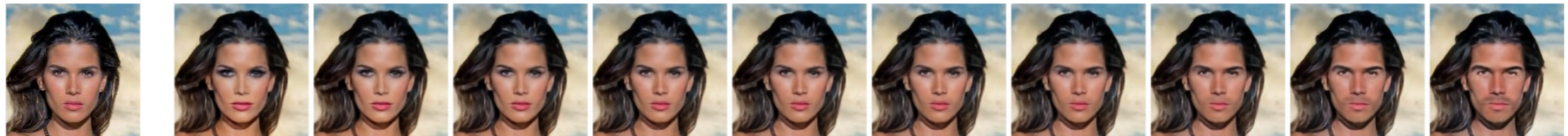
Old → Young



Male → Female



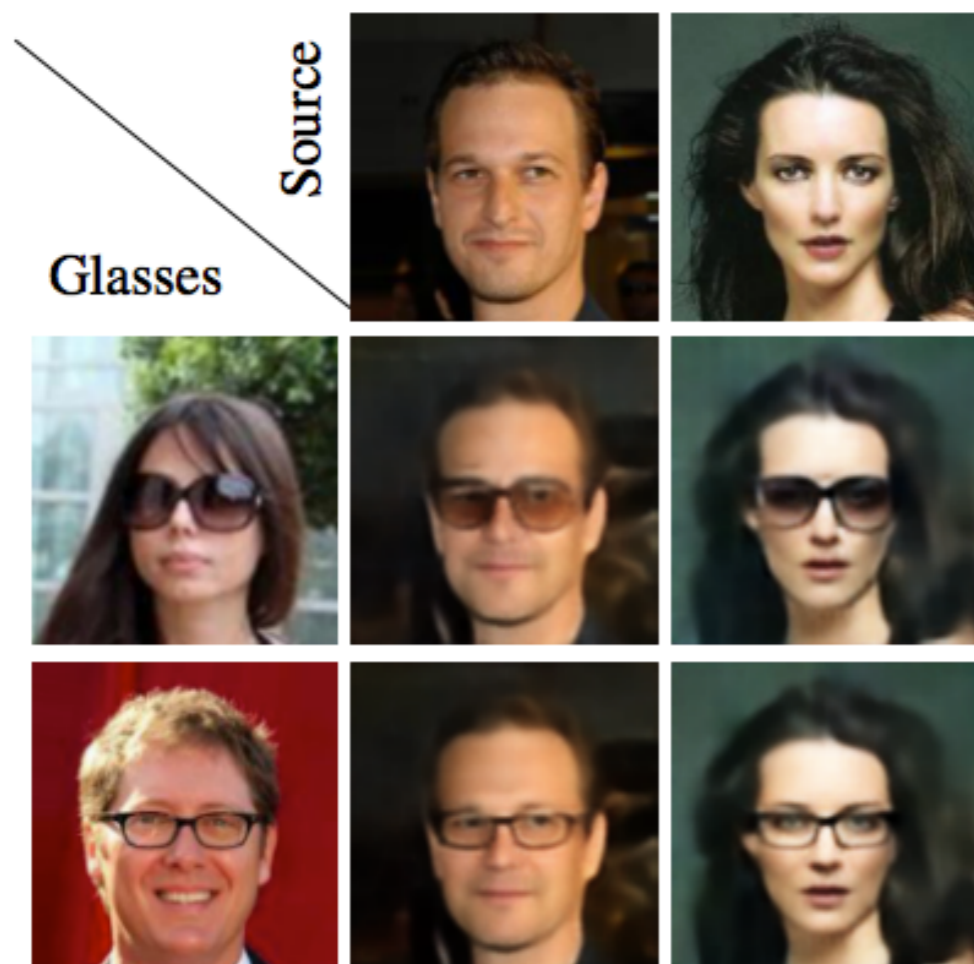
Female → Male



EMERGING DISENTANGLEMENT IN AUTO-ENCODER BASED UNSUPERVISED IMAGE CONTENT TRANSFER

Ori Press, Tomer Galanti, Sagie Benaim, Lior Wolf

Content transfer:



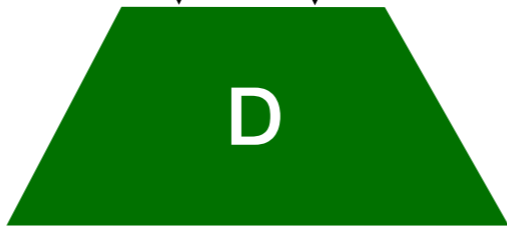
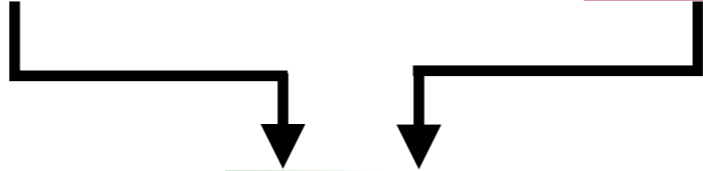
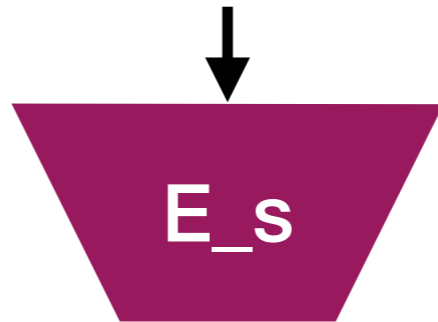
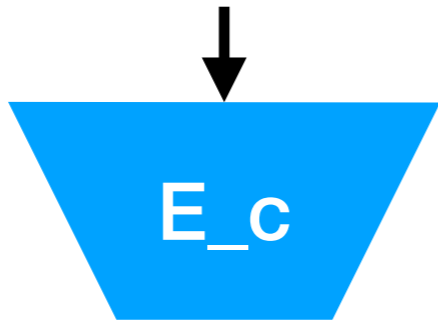
Unsupervised Content Transfer

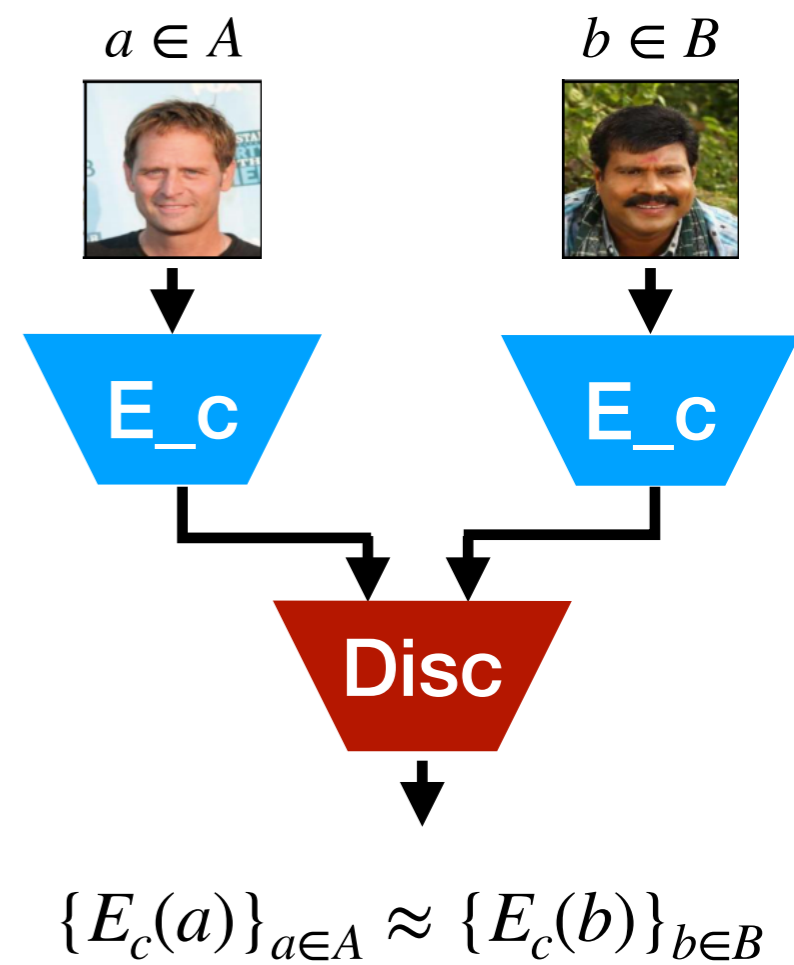
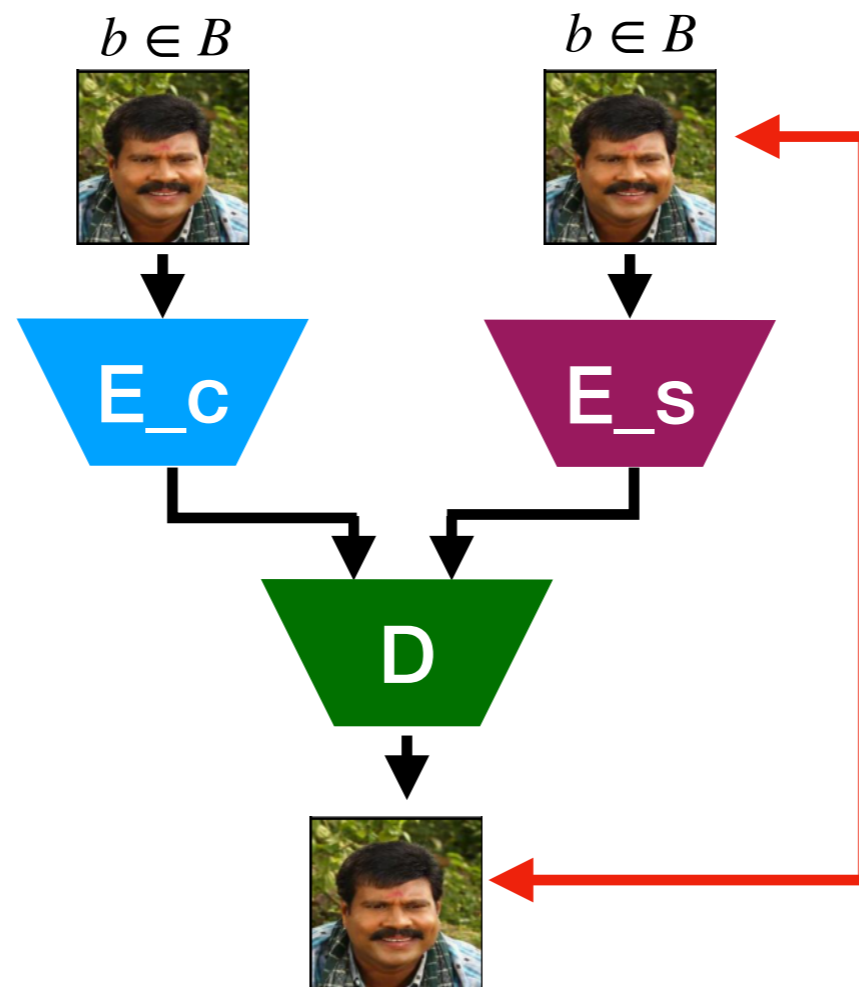
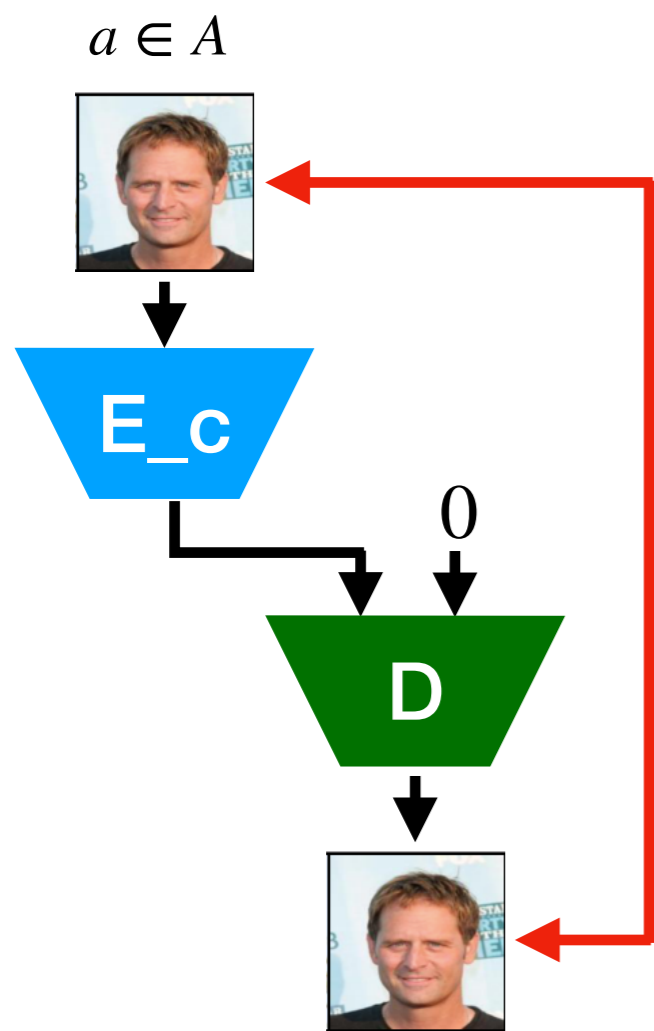


Domain A contains common content



Domain B contains common and separate content





Domain Confusion

$a \in A$

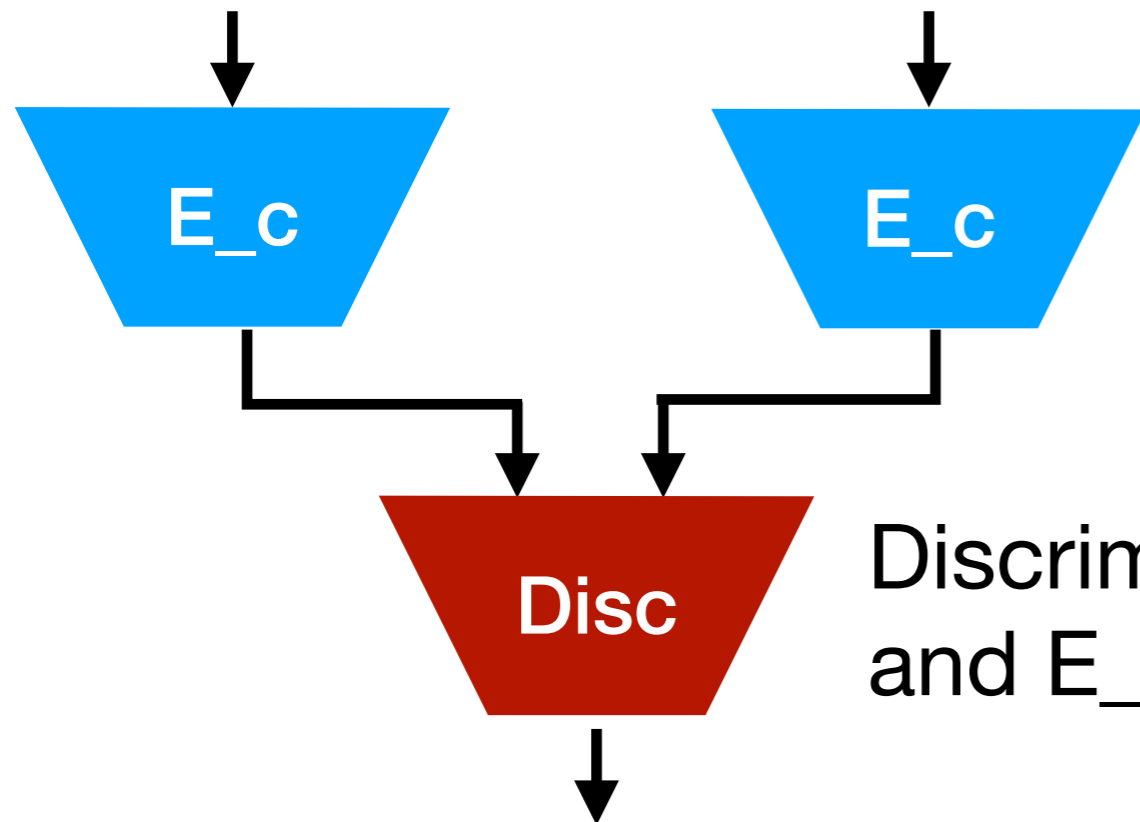


$b \in B$



Goal:

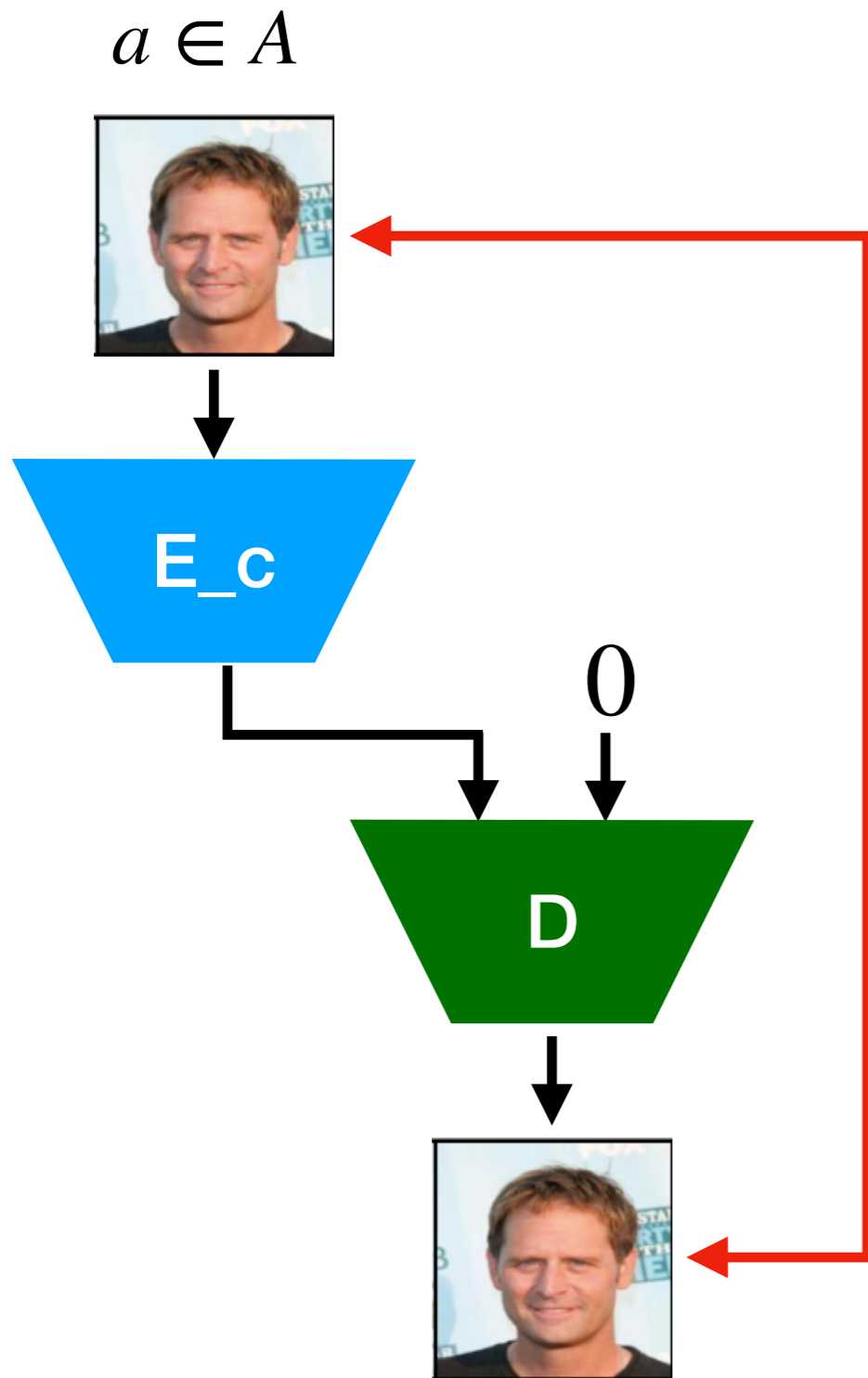
$$\{E_c(a)\}_{a \in A} \approx \{E_c(b)\}_{b \in B}$$



Encoder try to fool the discriminator

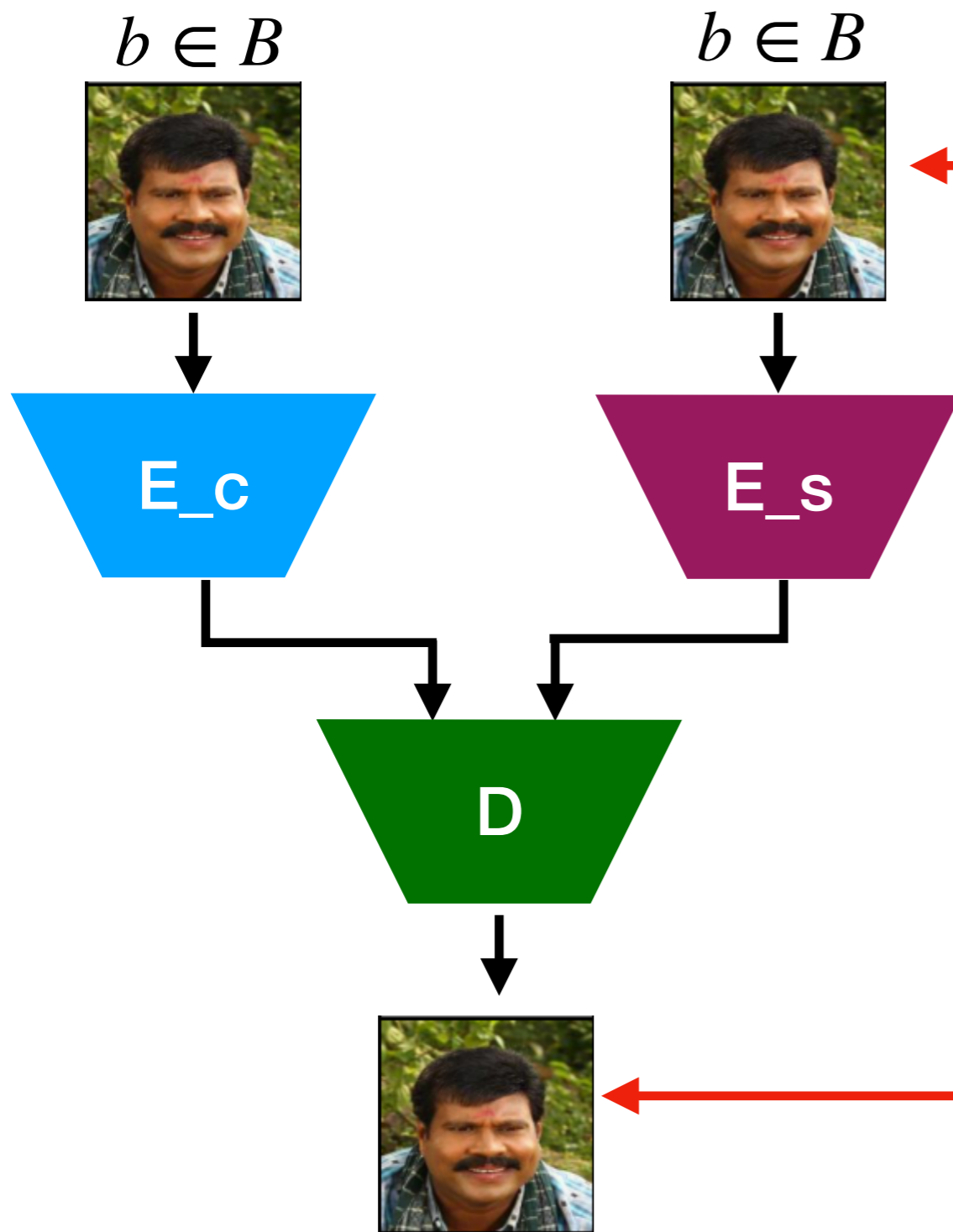
Discriminator try to distinguish $E_c(a)$ and $E_c(b)$

Reconstruction



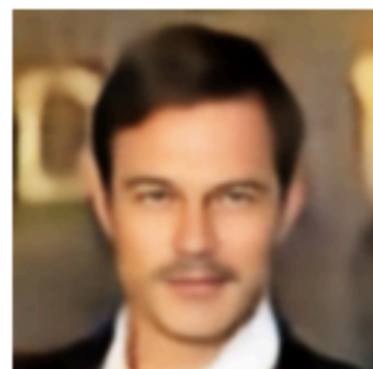
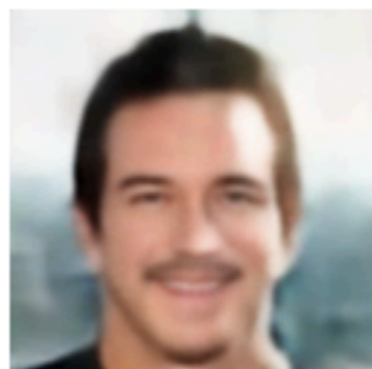
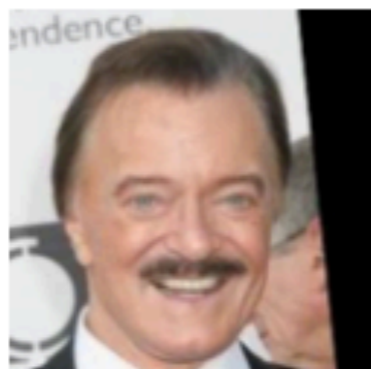
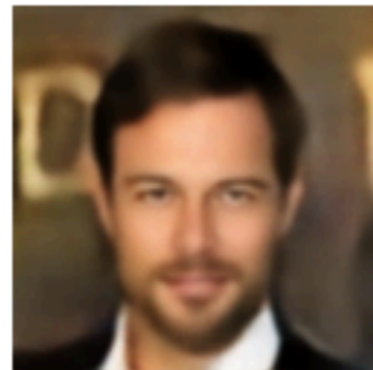
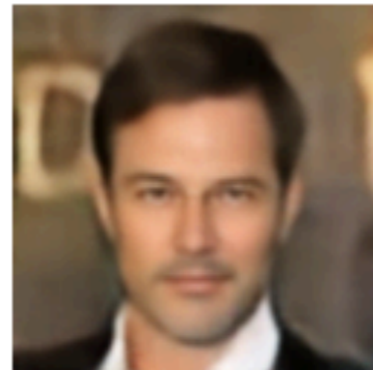
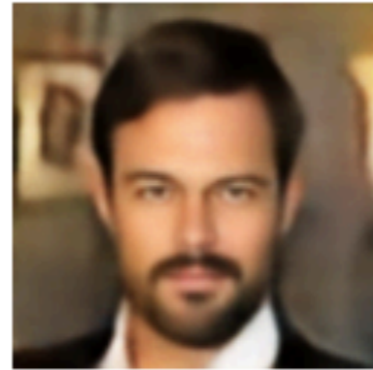
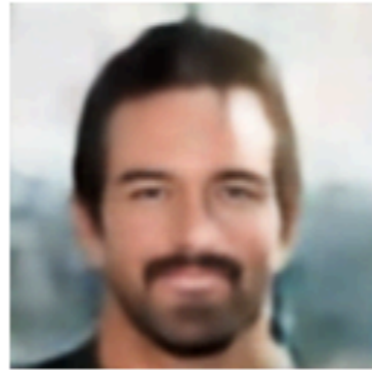
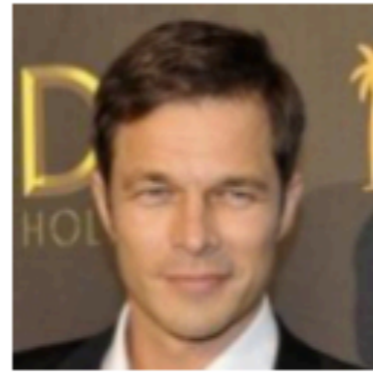
$$\mathcal{L}_R^A = \|D(E_c(a), 0) - a\|$$

Reconstruction 2



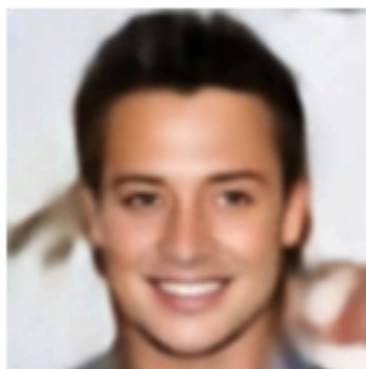
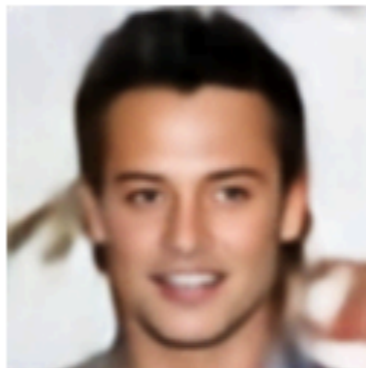
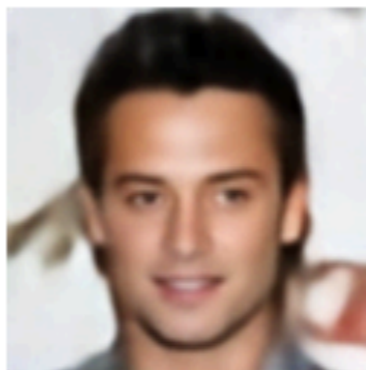
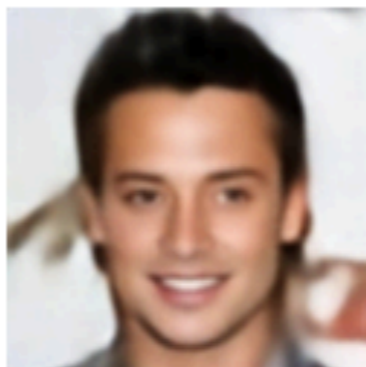
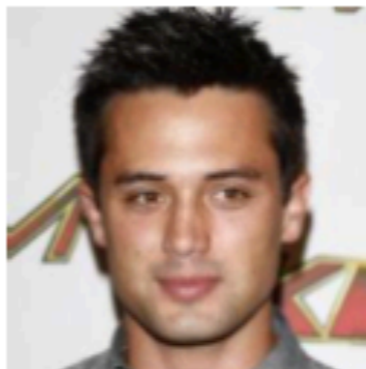
$$\mathcal{L}_R^A = \|D(E_c(b), E_c(b)) - b\|$$

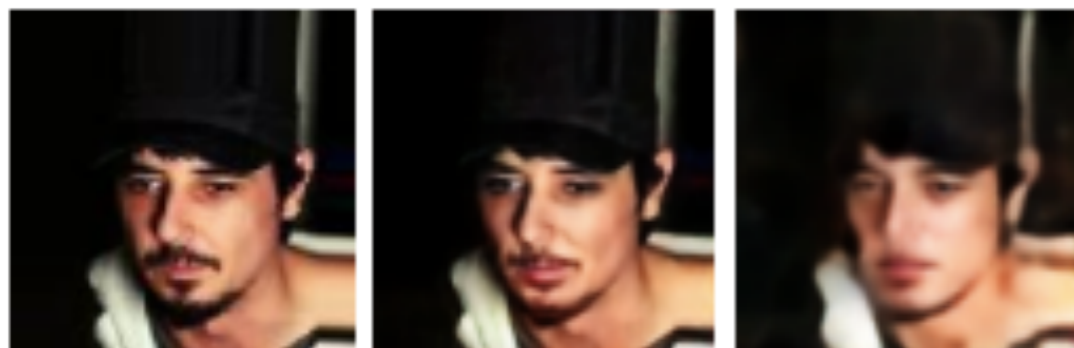
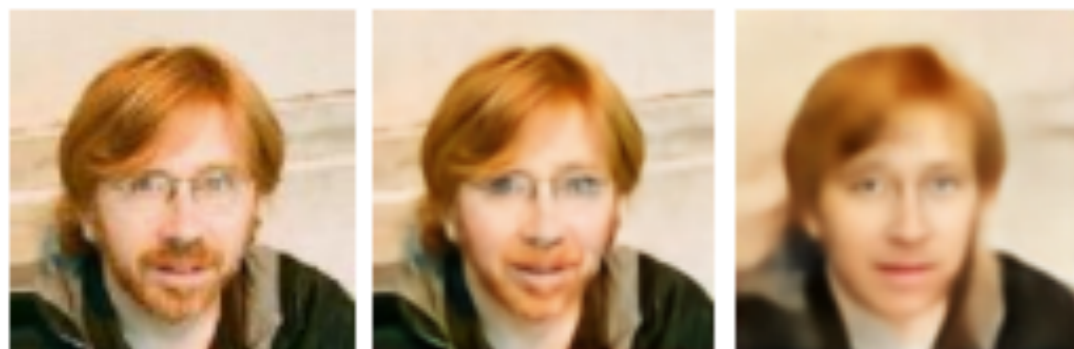
Source
Beard



Mouth

Source





Original

Fader

Our



Original

Fader

Our

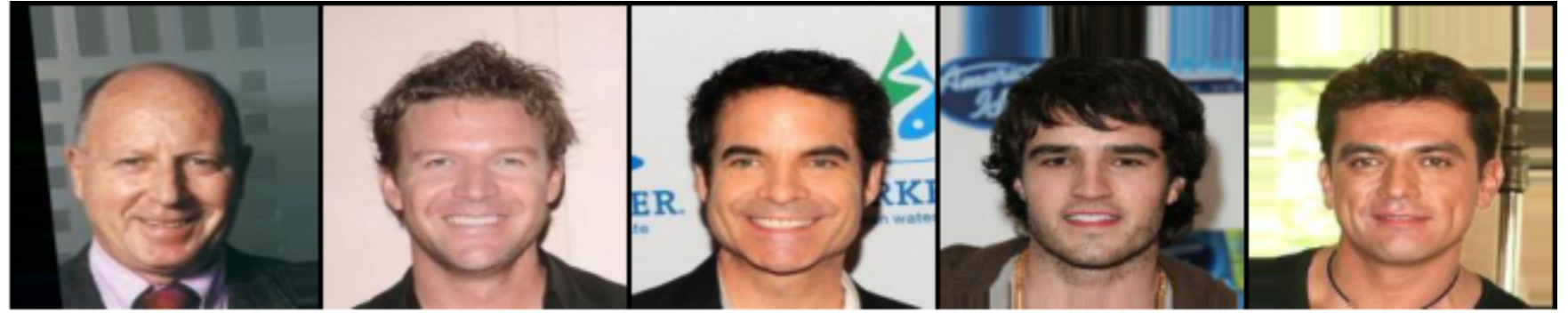
MASK BASED UNSUPERVISED CONTENT TRANSFER

Ron Mokady, Sagie Benaim, Lior Wolf, Amit Bermano

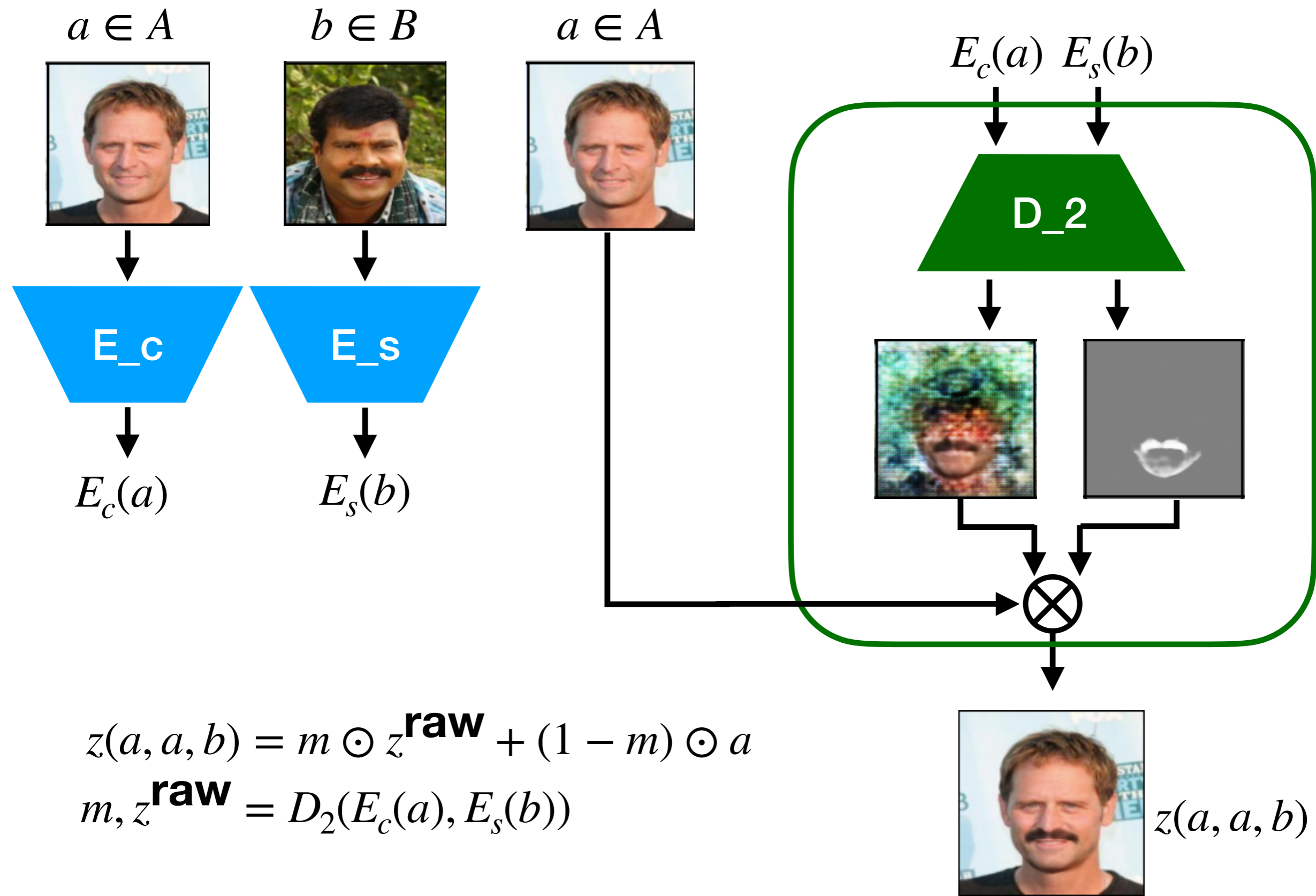
- Transfer content using a mask for state-of-the-art quality.
- Use the mask as a weakly-supervised semantic segmentation.

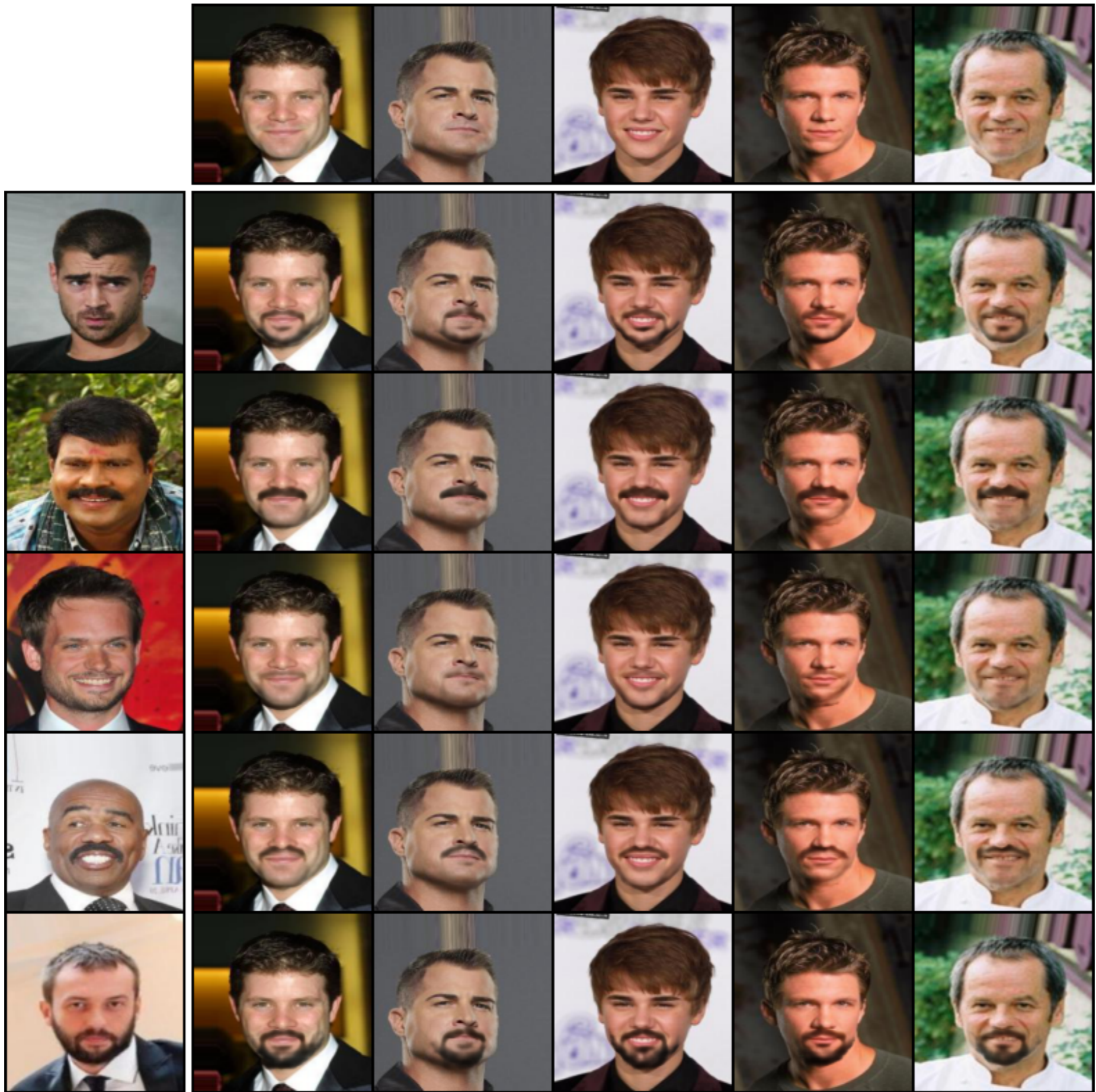
Comparison





Using a mask

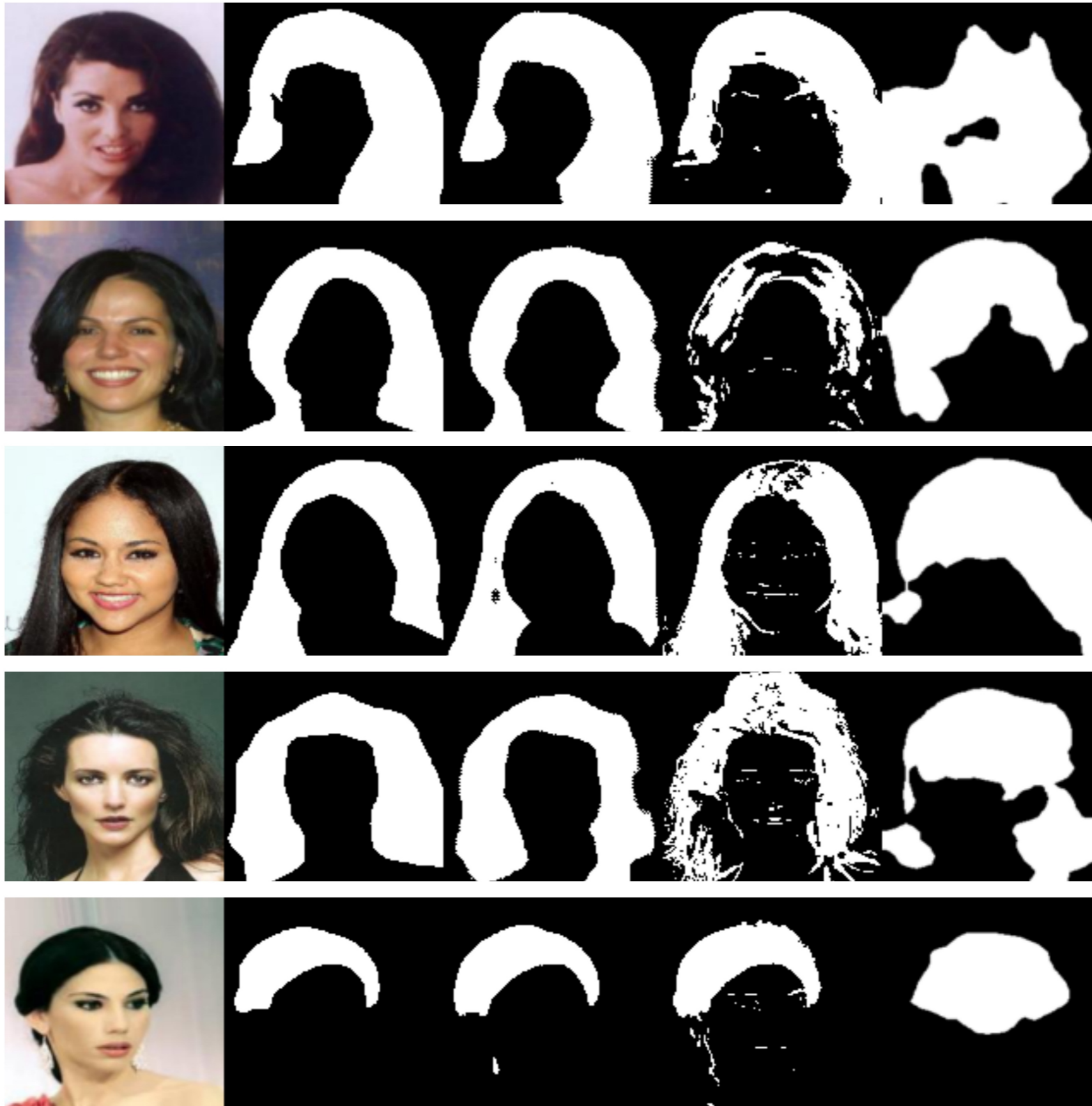












GT

Ours

Press et al.

Ahn et al.

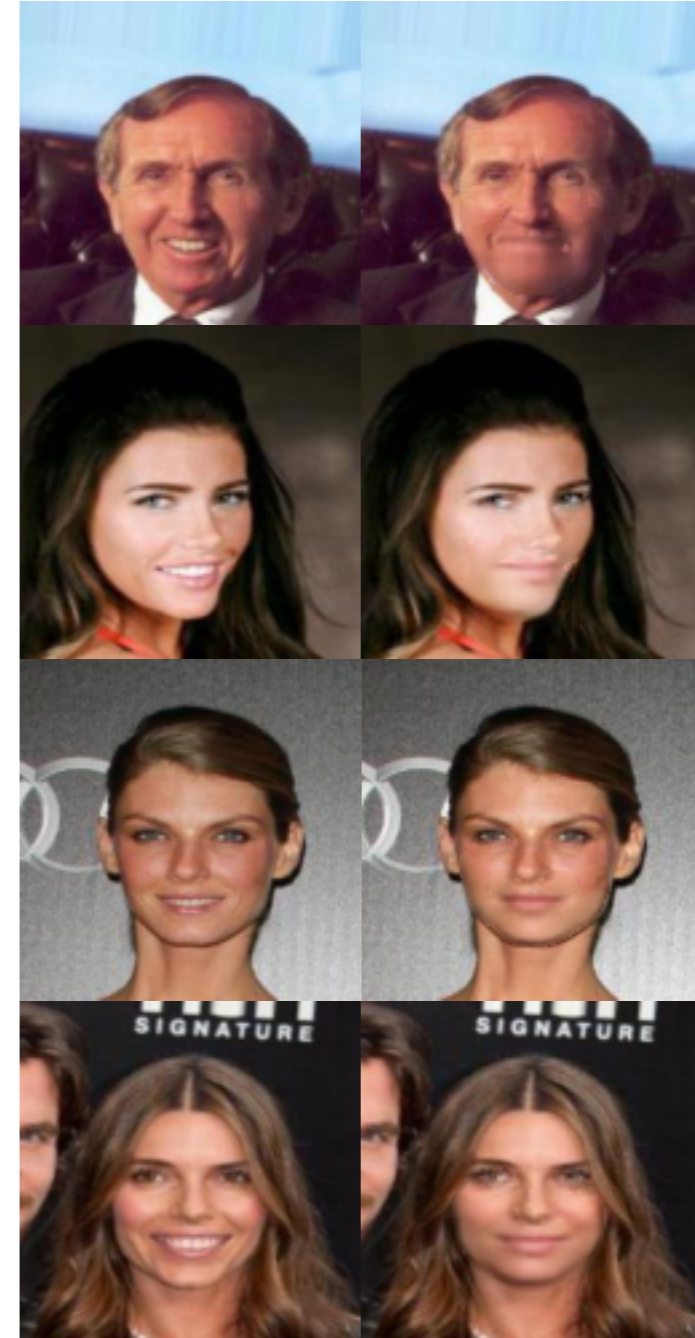


GT

Ours

Press et al. Ahn et al.

Removal



Generate Unseen Images



Domain A: glasses without mustache.



Domain B: mustache without glasses.

Can we generate face with both mustache and glasses? Or one without mustache and glasses?

Domain Intersection and Domain Difference

Sagie Benaim, Michael Khaitov , Tomer Galanti, Lior Wolf

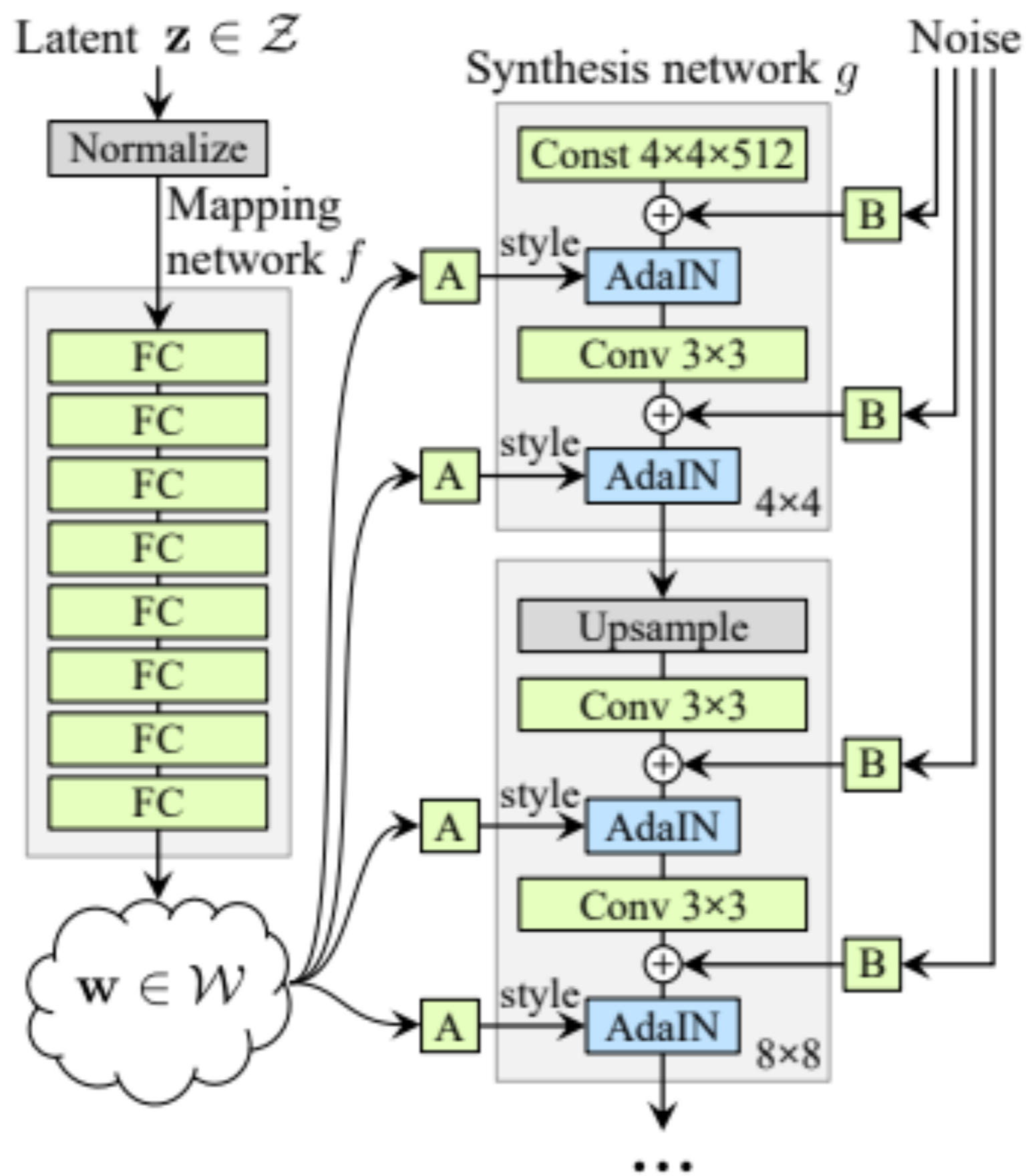
Three encoders, one for the common and two for the attributes.

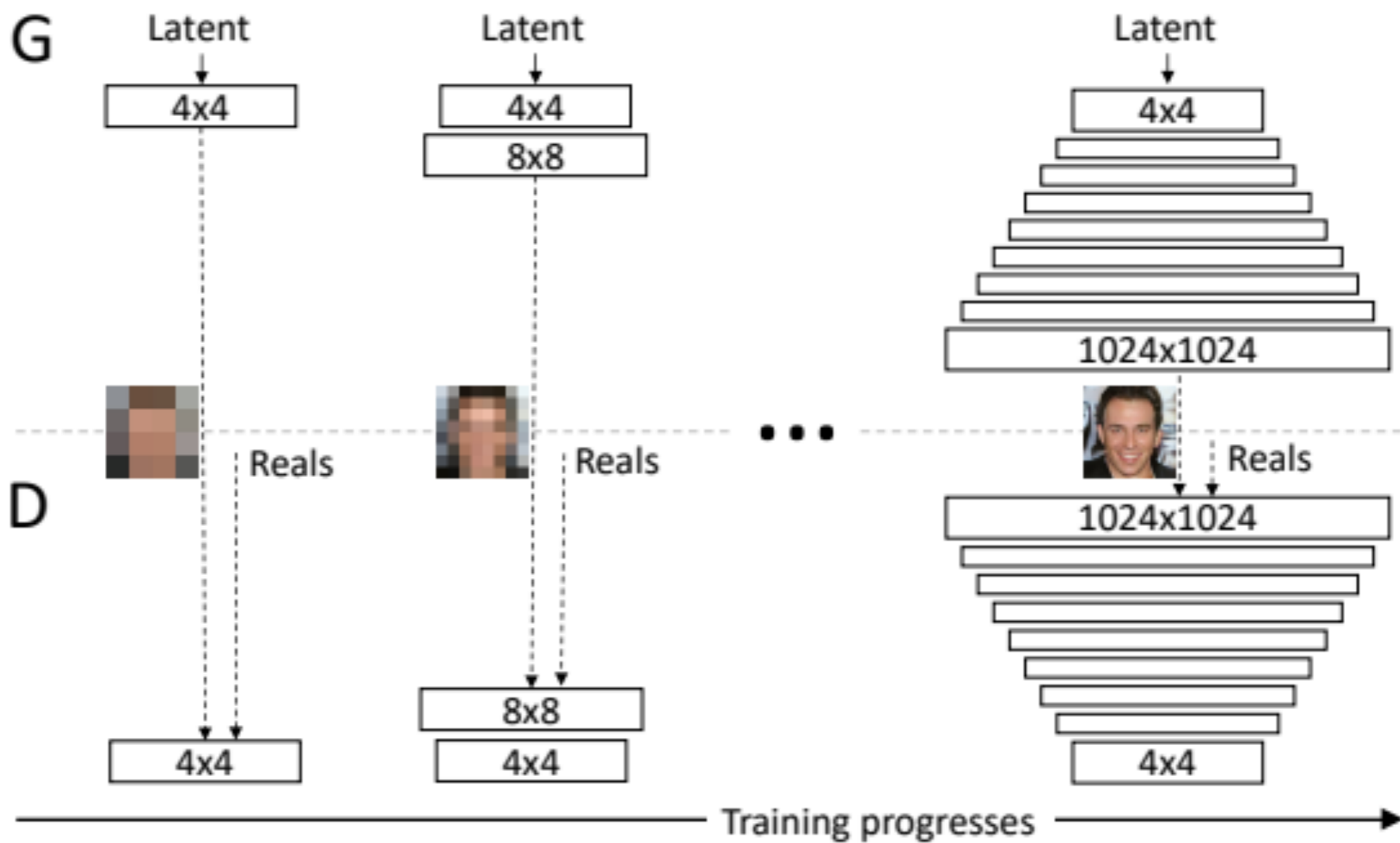
Same Reconstruction and Domain-Confusion loss as Press et al. with Additional zero loss.

A Style-Based Generator Architecture for Generative Adversarial Networks

Tero Karras, Samuli Laine, Timo Aila







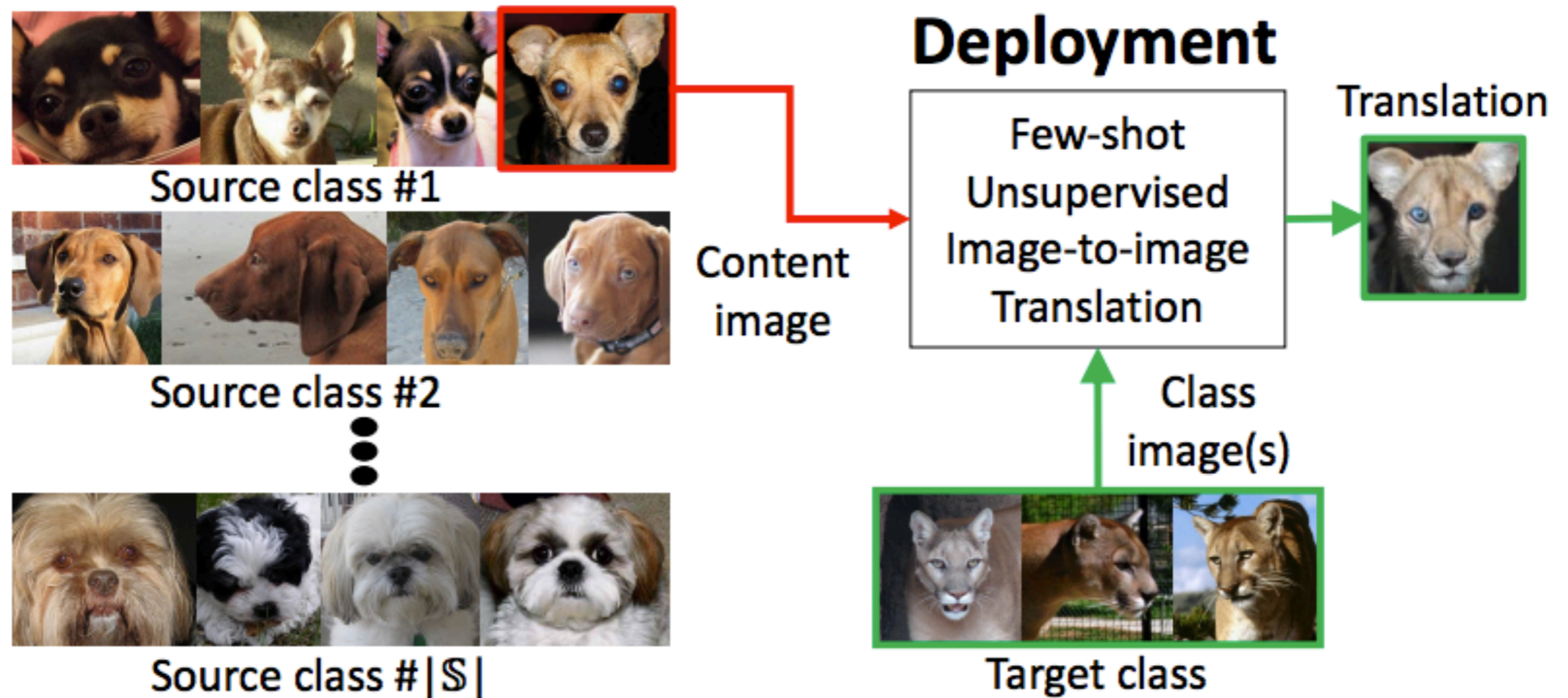
A Style-Based Generator Architecture for Generative Adversarial Networks

Tero Karras, Samuli Laine, Timo Aila



Few-Shot Unsupervised Image-to-Image Translation

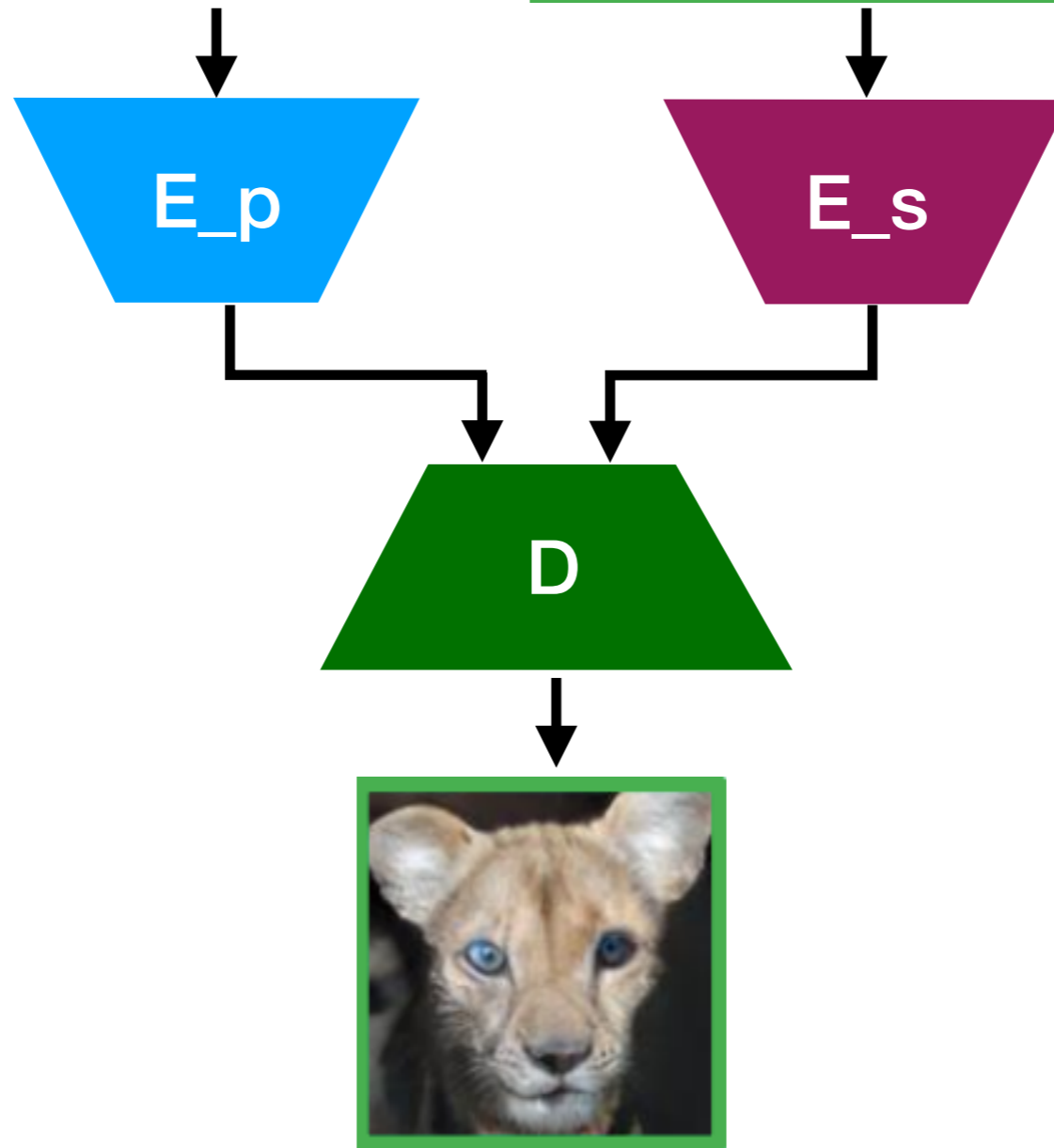
Ming-Yu Liu, Xun Huang, Arun Mallya, Tero Karras, Timo Aila, Jaakko Lehtinen, Jan Kautz



Source



Target



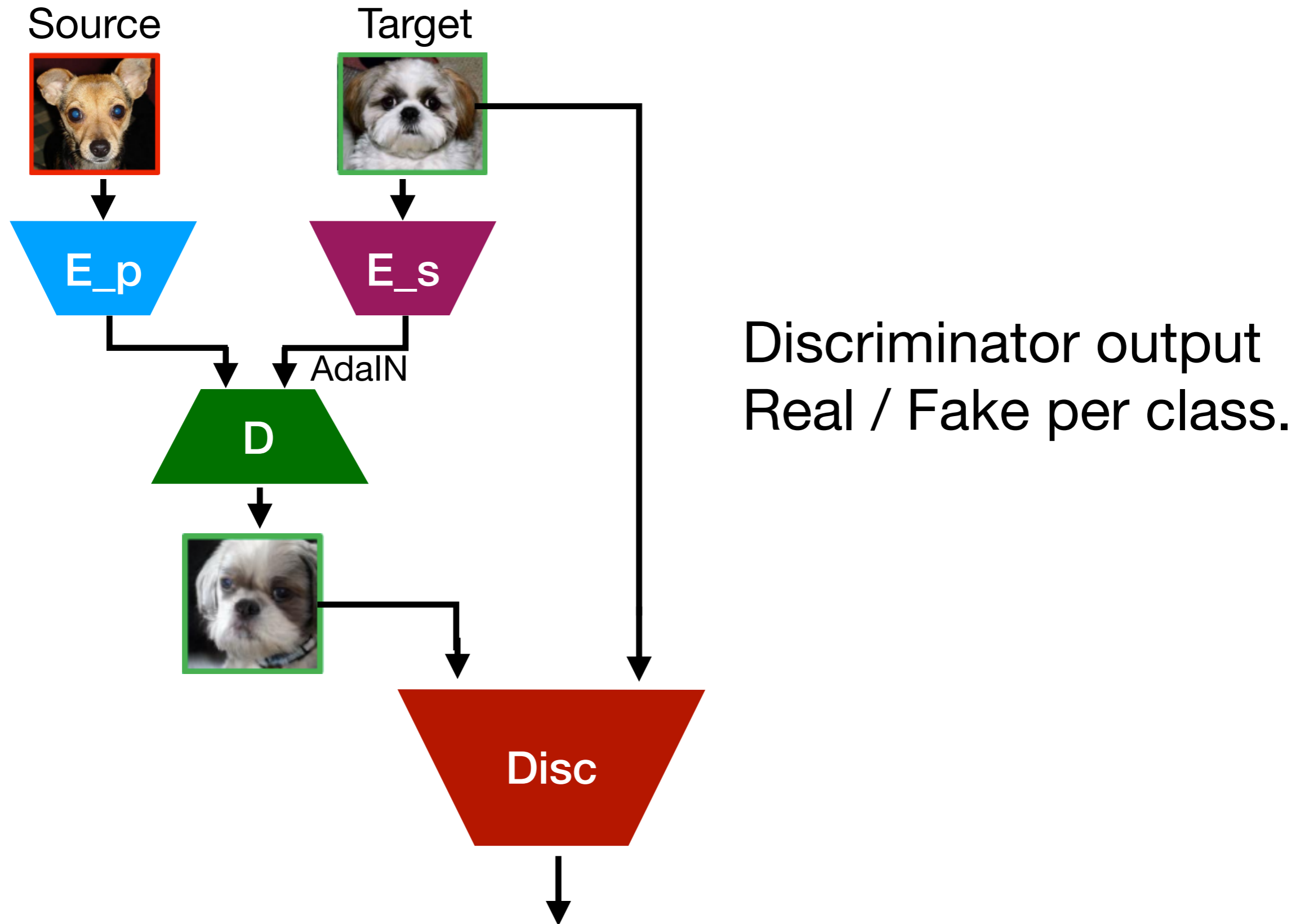
Style Encoder

Map each image to latent space then computes the mean.

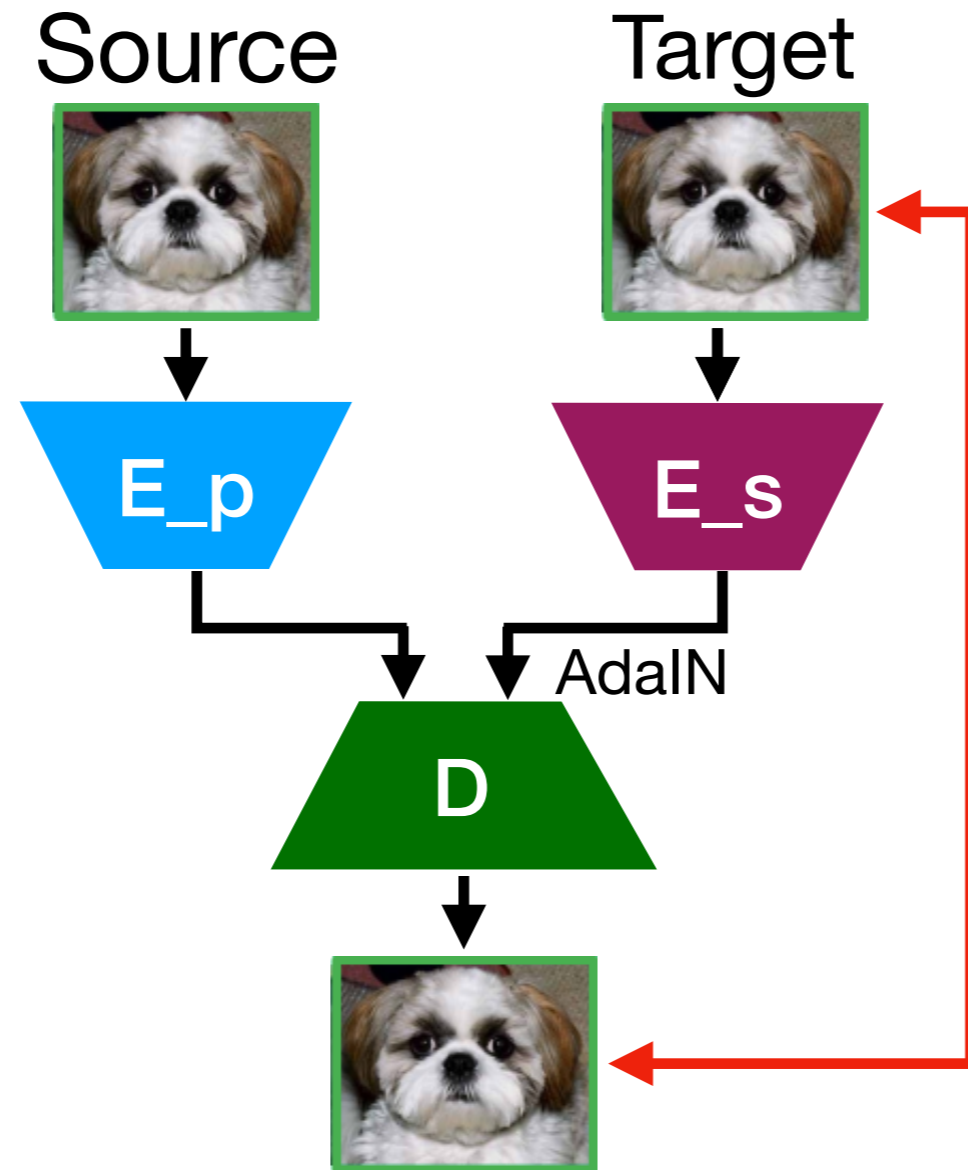
Decoder

Since target domain used only for style, we feed the target encoding (style encoding) to the decoder via the AdaIN layers. We let the target images control the global look (e.g., object appearance), while the content image determines the local structure.

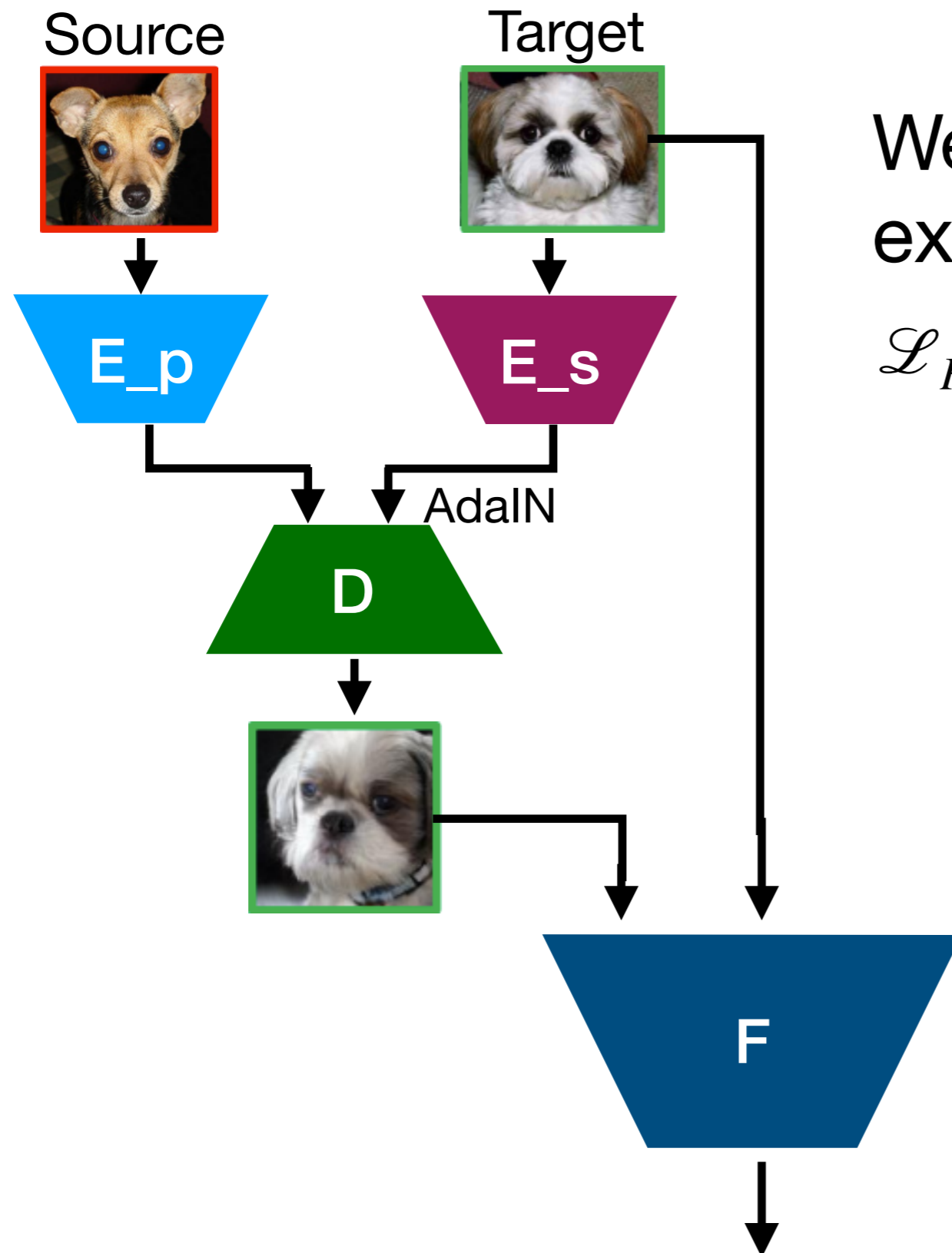
GAN Loss



Reconstruction



Feature Matching Loss



We use feature extractor F .

$$\mathcal{L}_F = \|F(D(E_c(x), E_s(y))) - F(y)\|$$

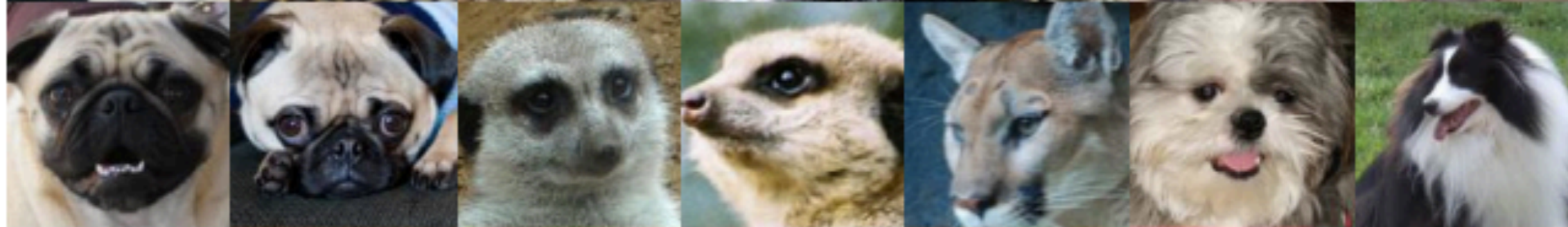
Disentanglement?

Network architecture. Style encoding can't change the local structure.

y_1



y_2



x



\bar{x}



Conclusions

- Transfer some factors from image to image (content, style, rotation, etc.).
- Easier to work over latent space (e.g. adversarial loss).
- Enables to generate out-of-domain images.

Questions?