



Learning Music

STYLE TRANSFORMATION

H O N G B O Z H U



STYLE A



eg: Country music

Violin

Man voice

STYLE B

eg: classic music

Piano

Women voice



Cycle generative adversary network

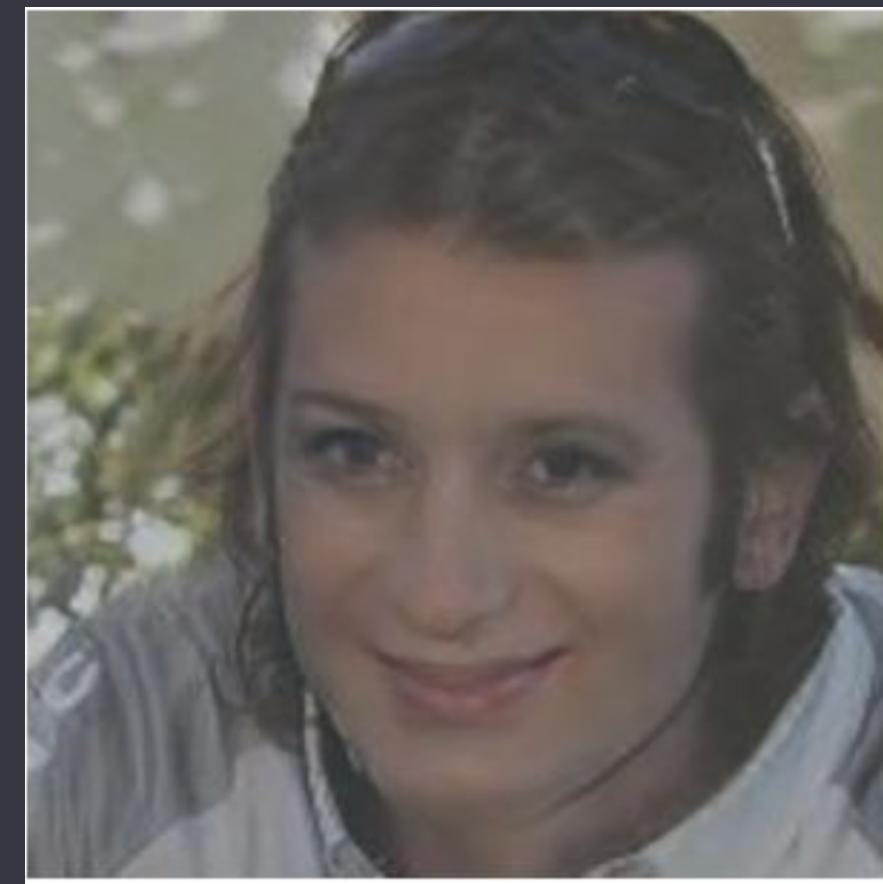


MAN



G

WOMEN



D_A



D_B





HOW TO DEAL WITH MUSIC DATA?

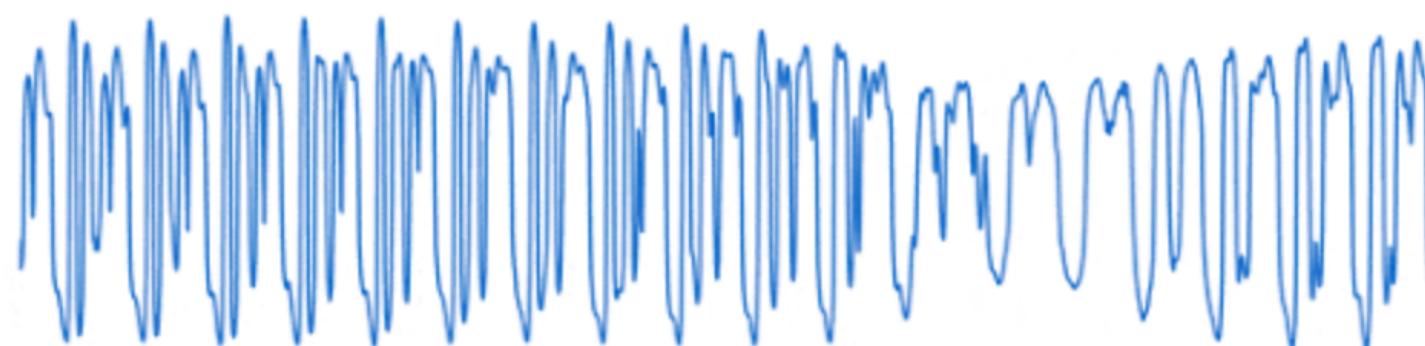
Music data representation



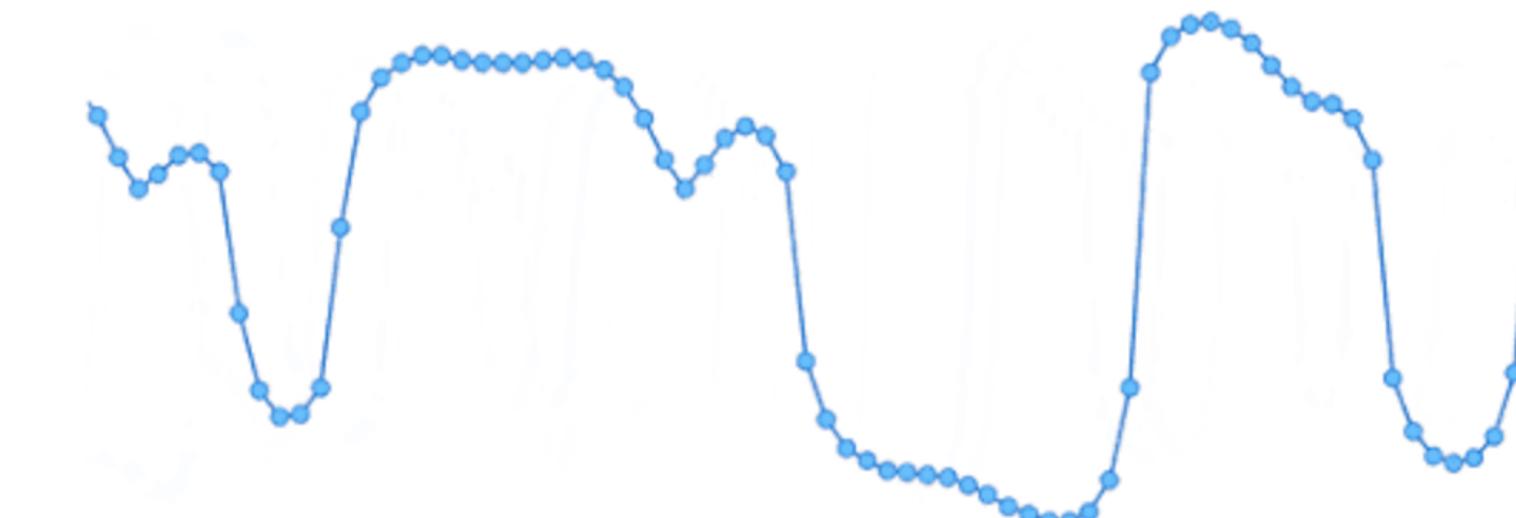
1 Second



100 milliseconds



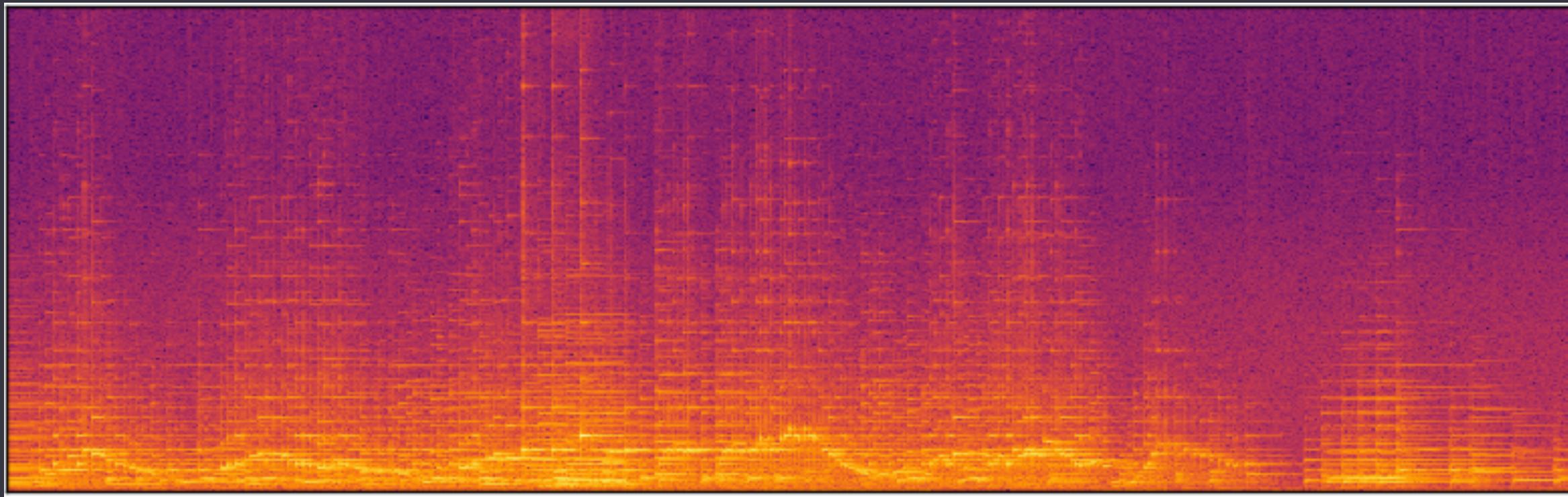
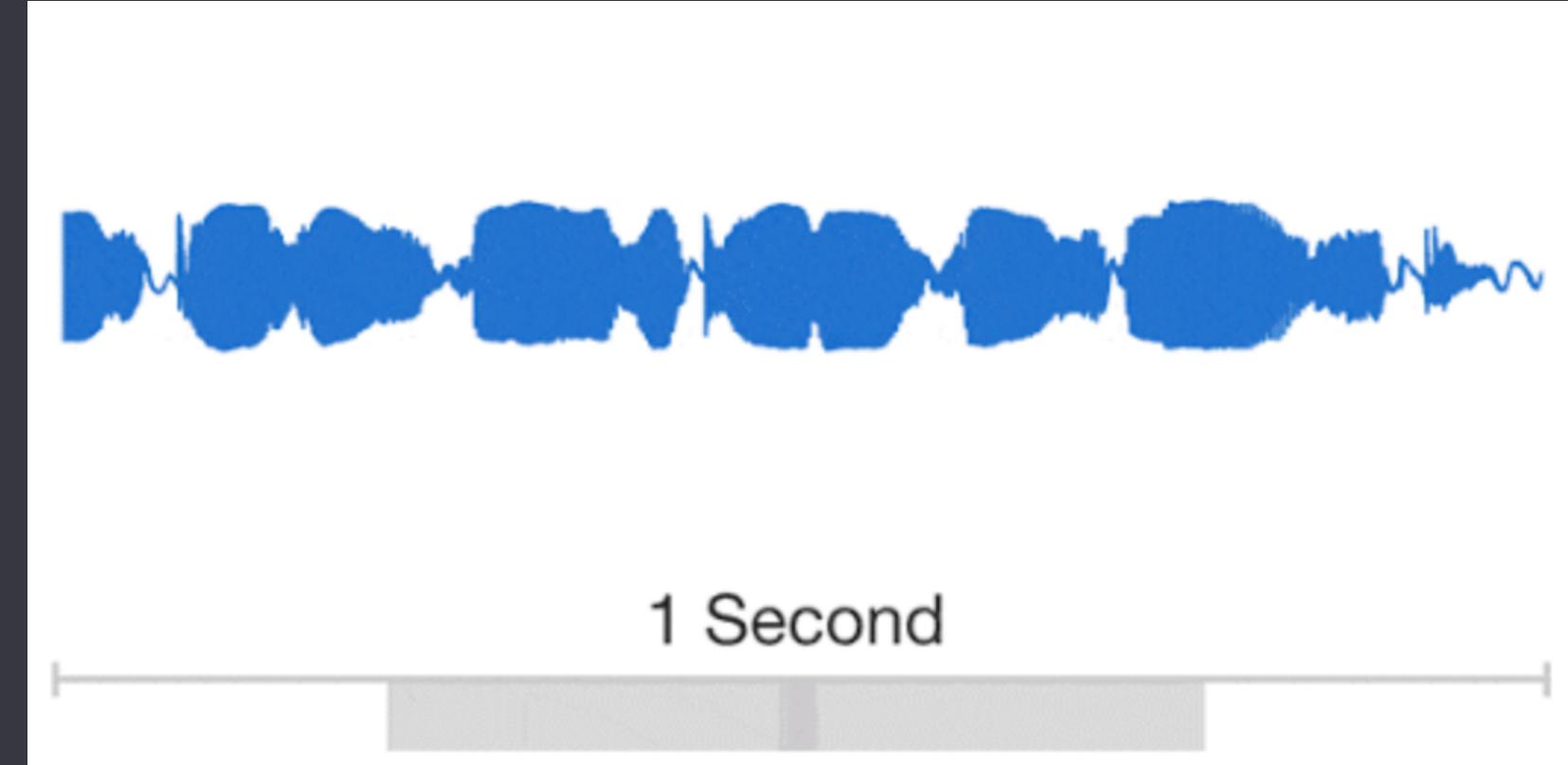
10 milliseconds



1 millisecond

Spectrogram-STFT

- Duration 30s
- Frame-size 30ms
- Frame-step 10ms
- FFT-size 512
- Audio sampling rate of 20.5kHz.

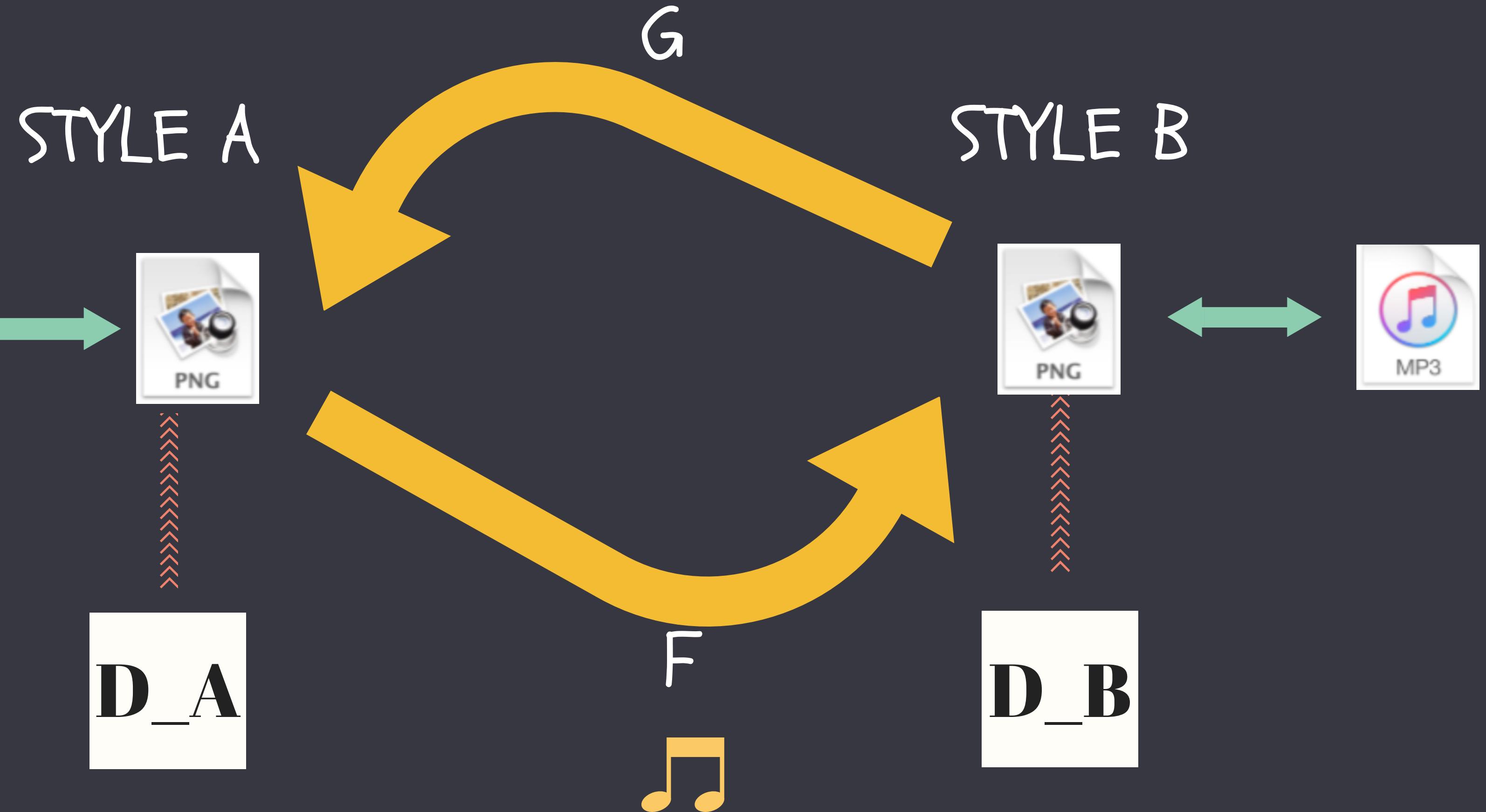


$\log()$

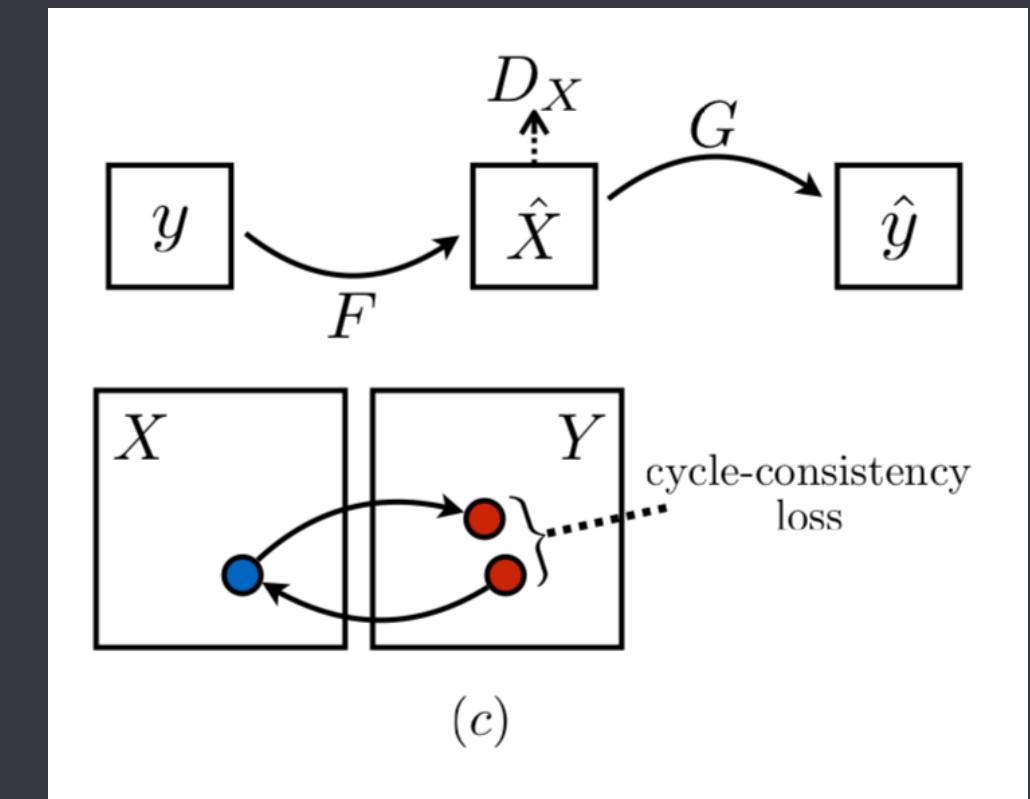
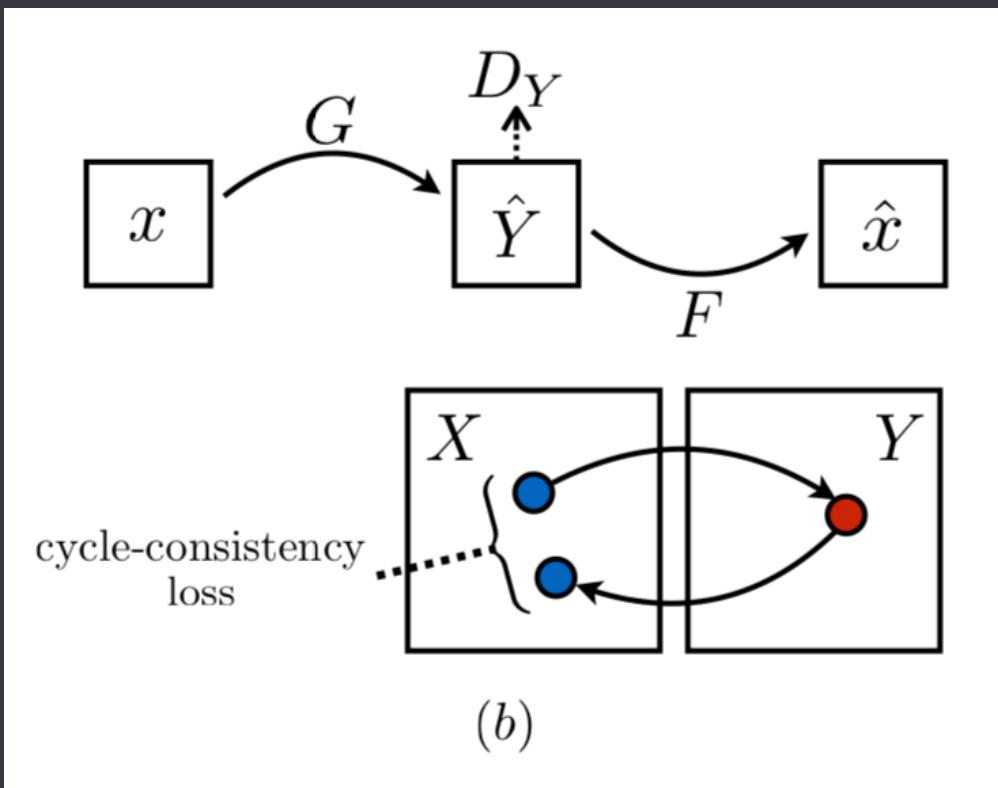
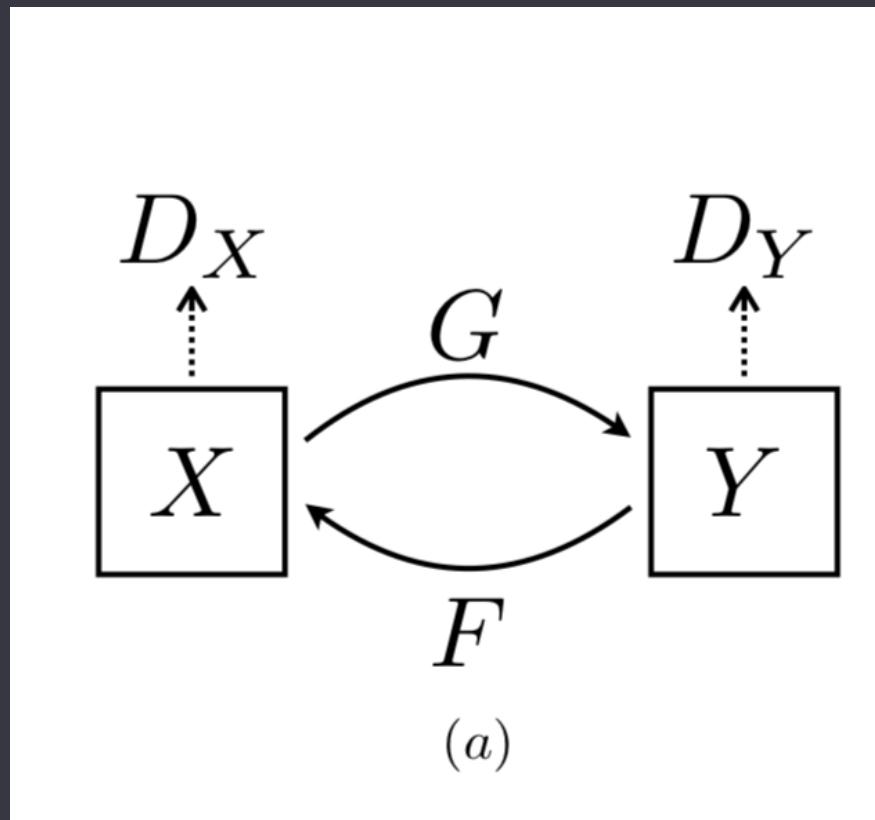




Cycle generative adversary network



Tweak on the architecture



$$\begin{aligned} \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) = & \mathbb{E}_{y \sim p_{\text{data}}(y)} [\log D_Y(y)] \\ & + \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log(1 - D_Y(G(x)))] \end{aligned} \quad (1)$$

$$\begin{aligned} \mathcal{L}_{\text{cyc}}(G, F) = & \mathbb{E}_{x \sim p_{\text{data}}(x)} [\|F(G(x)) - x\|_1] \\ & + \mathbb{E}_{y \sim p_{\text{data}}(y)} [\|G(F(y)) - y\|_1]. \end{aligned} \quad (2)$$

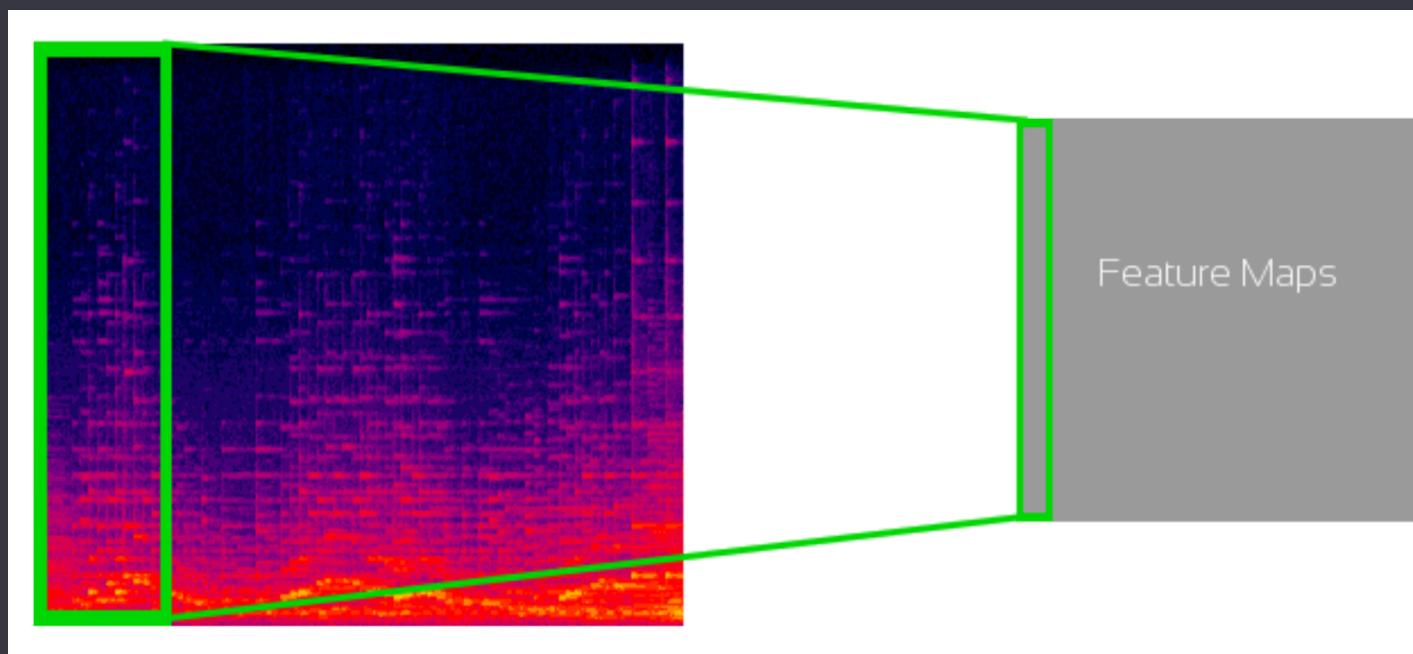
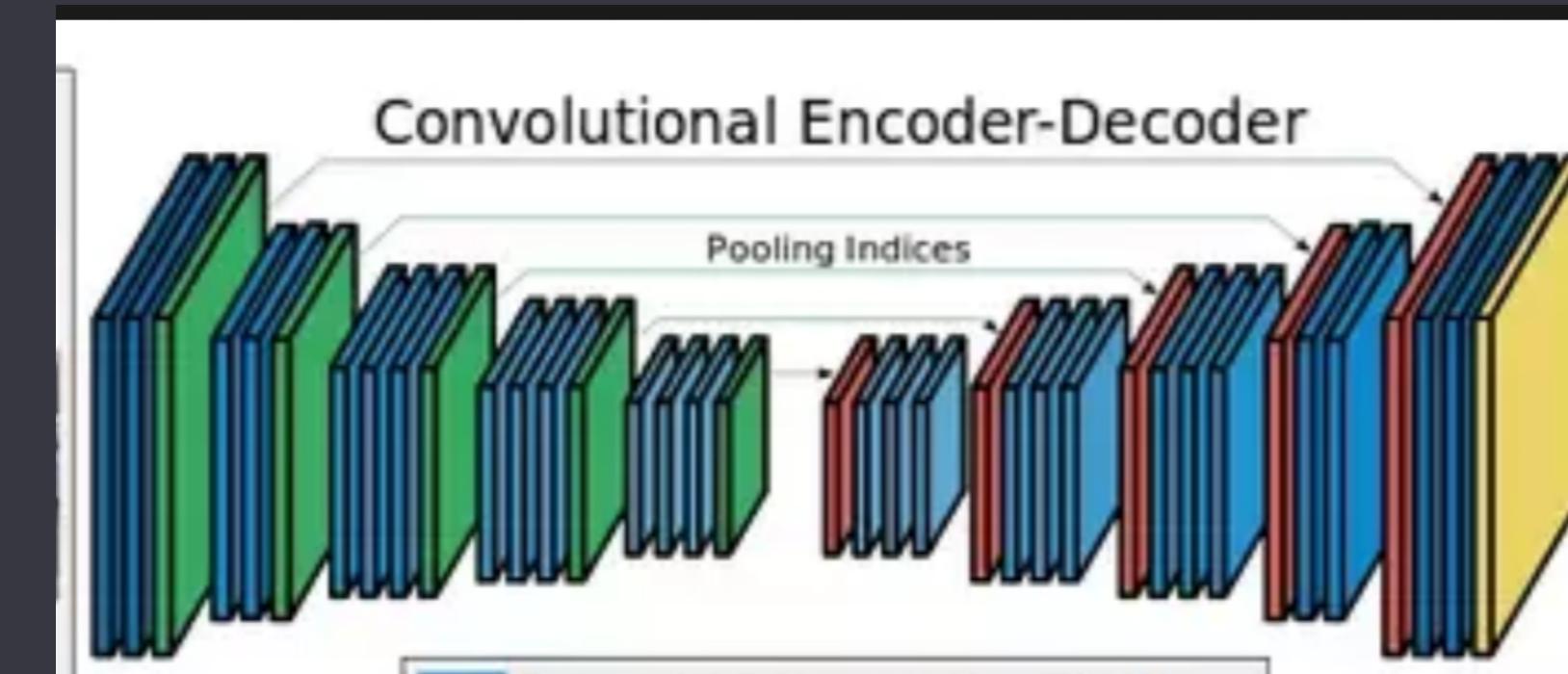
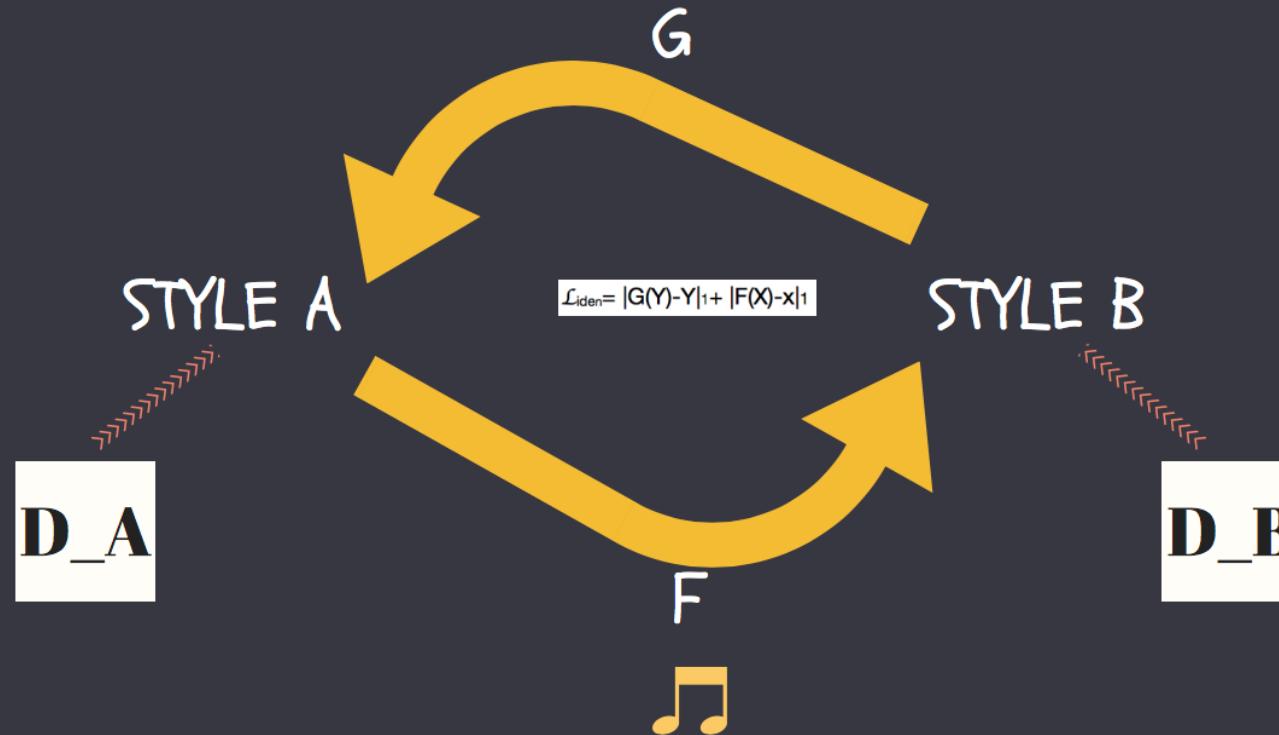
$$\begin{aligned} \mathcal{L}(G, F, D_X, D_Y) = & \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) \\ & + \mathcal{L}_{\text{GAN}}(F, D_X, Y, X) \\ & + \lambda \mathcal{L}_{\text{cyc}}(G, F), \end{aligned} \quad (3)$$

+

$\mathcal{L}_{\text{iden}} = |G(Y) - Y|_1 + |F(X) - X|_1$



Tweak on Generator G and F



- 2 convolution layers (stride=2, width=3)+
- 9 Resnet layers (stride=1, width=3)+
- 2 transpose convolution layers +
- No pooling
- Instance normal each layer
- Relu as activation



Tweak on discriminators



- 5 convolution layers (stride=[2,2,2,2,1], width=4) +
- 1 dense layer
- No pooling
- Instance normal each layer
- Relu as activation