# GANs N' Roses

- Link to paper: https://arxiv.org/pdf/2106.06561.pdf

- Abstract

    ⇒ Inputs
        ① Content code derived from human image
        ② Style code chosen randomly
    ⇒ Output = anime image taking the pose of the human in the human image.

    ⇒ New adversarial loss based on new definition of content and style

- Intro

    ⇒ The paper achieves better human-to-anime pose transfer by defining new losses based on new definition of "content" and "style"
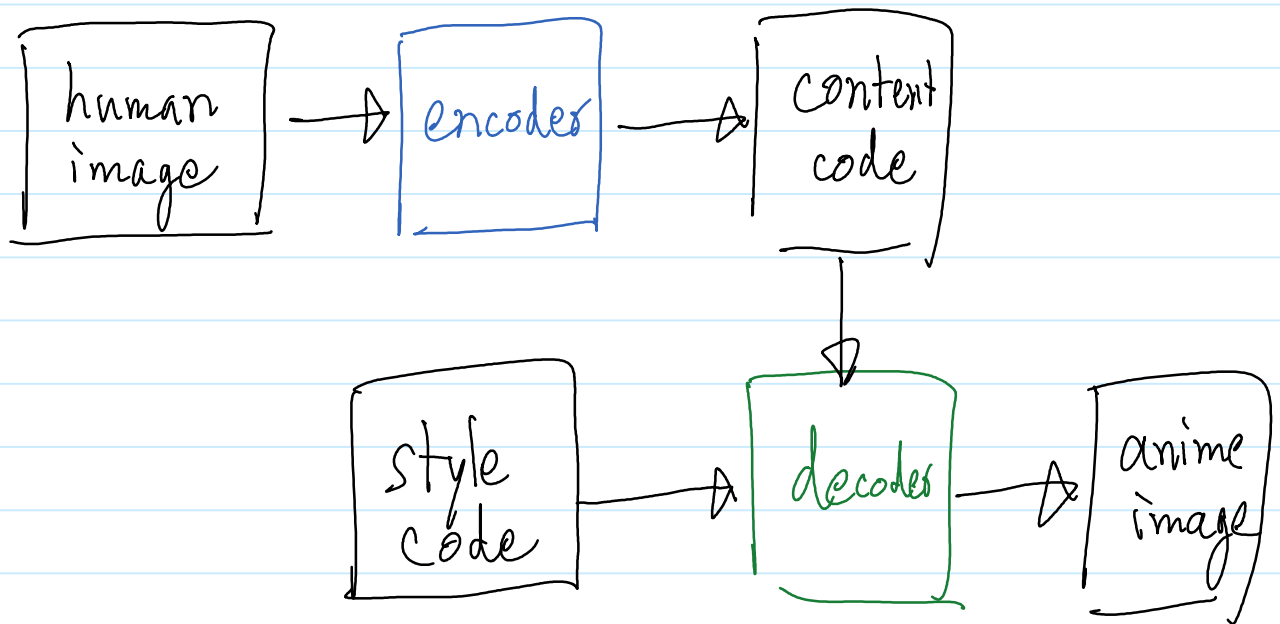
        - Should it just use the word "pose" for "content" and "appearance" for "style"?

    ⇒ Idea

        - "Content" is what changes if face image undergoes a family of data augmentation transformations.

          - Style is what does not change

— Augmentations = scaling, rotating, cropping.

— How this works



— To make sure that the anime image has the same pose as the human image, we must ensure that the anime image has the same content code as the input human image

— However, we cannot use cycle consistency loss because we do not want a 1-1 mapping between human and anime face.

— The paper proposes a way to ensure the content code is the same without using cycle consistency loss

- The Related Works session has two paper related to anime.
  - CountcilGAN [LINK]
    - Generate diverse anime face from a single human face using multiple generators working in parallel
    - Problems
      - ☐ Cannot capture diverse anime styles.
      - ☐ Mode collapses.
  - AniGAN [LINK]
    - New normalizations that
      - ☐ Transfer color and texture styles
      - ☐ Maintaining global structure
    - Problem: Style not diverse enough.

- Framework

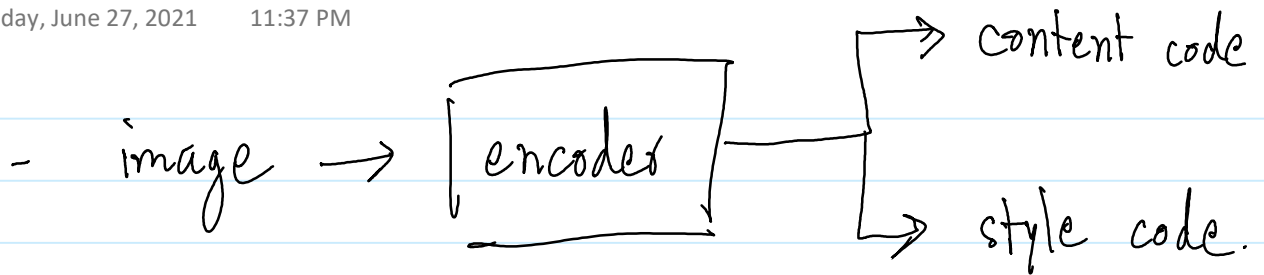$\Rightarrow$ Two domains $X$ (human face)

$Y$ (anime face)

$\Rightarrow$ Goal: Given $x \in X$, generate a subset $\hat{Y} \subseteq Y$ such that each $y \in \hat{Y}$ contains similar semantic content as $x$.

$\Rightarrow$ While goal is only in direction $X \to Y$, we need a mechanism to do $Y \to X$ as well. So, 2 directions in total.

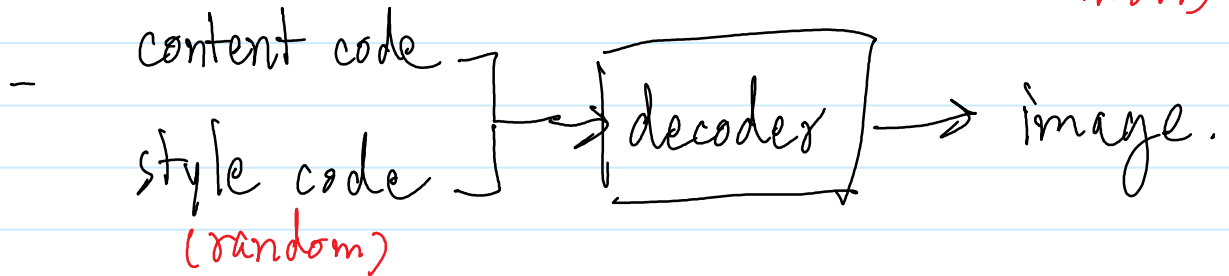$\Rightarrow$ For each direction, we need

① an encoder, denoted by $E_{X \to Y}$, $E_{Y \to X}$

② a decoder, denoted by $F_{X \to Y}$, $F_{Y \to X}$

- image $\longrightarrow$ [ encoder ] $\longrightarrow$ content code

  $\longrightarrow$ style code.
  (random)

- content code
  style code $\longrightarrow$ [ decoder ] $\longrightarrow$ image.
  (random)

- Generator = encoder + decoder

- Main idea : Choose a relevant collection of data augmentation. Under it:

  Content is what changes

  Style is what does not change

- Ensuring style diversity

  $\Rightarrow$ Existing strategy.

    - Generate from random style code

      $\Rightarrow$ The decoder might ignore style code.

    - Make sure that the style code can be recovered from the generated image

      $\Rightarrow$ The decoder might hide it in a few pixels

- Force outputs from different style codes to be different,

$\Rightarrow$ No gaurantee that this is the right diversity.

$\Rightarrow$ Let $P(X)$ denote probability distribution of $X$.

$T(\cdot) =$ a function that applies a random augmentation that changes content and preserve style.

$P(C) =$ distribution of content codes

$P(Y) =$ distribution of $Y$.

$P(\hat{Y}) =$ distribution of $F_{X \to Y}(c(x), s_Z)$
where $x \sim P(X)$, $s_Z \sim N(0,1)$

<span style="color:red">content code of $x$</span>

<span style="color:red">normal distribution</span>

$\Rightarrow$ Note that $c(x_i) \sim P(C)$ if $x_i \sim P(x)$

$\Rightarrow$ Requirement on $T$: $c(T(x_i)) \sim P(C)$

(i.e. while the exact content code would change, the overall distribution does not)

⇒ The paper note that the last requirement
is reasonable. Otherwise, augmentation used
when training classifiers would not work.

⇒ IMHO, I doubt this. Rotation changes the
pose significantly.

    ⇒ If you have a dataset with mainly heads
    in upright position, rotating the image
    can change the range of head angles.

⇒ After we generate $\hat{y} = F_{X \to Y}(C(x_i), s_z)$

we pass it to a discriminator $D$ to

judge whether the generated image is real

or fake.

⇒ The paper proposing generating a fake batch
in the following way:

    (a)  choose a single $X \in X$.

    (b)  Compute $X_1 = T(x)$, $X_2 = T(x)$, ..., $X_s = T(x)$
        (These are randomly augmented examples)

    (c)  compute $c_i = E_{X \to Y}(x_i)$

(d) Same random style code $z_1, z_2, \ldots, z_k \sim N(0,1)$

(e) Generate $\hat{y}_i = F_{x \to y}(c_i, z_i)$

$\Rightarrow$ The paper proposes that the batch $\{\hat{y}_1, \hat{y}_2, \ldots, \hat{y}_k\}$ should be indistinguishable from a batch of samples from $P(Y)$ of the same size.

This is the basis of its adversarial loss func

Note that this goal ensures that a single content code can be translated to multiple, diverse style.

- Losses

  ⇒ <u>Style consistency loss</u>. For each batch generated as above the style must be the same. So, the variance on the style code must be low:

$$\mathcal{L}_{scon} = Var\left( s(x_1), s(x_2), ..., s(x_k) \right)$$

  ⇒ <u>Cycle consistency loss</u>

  Let $\hat{y}_i = F_{X \to Y}(c_i, z_i)$

  $\hat{x}_i = F_{Y \to X}(cc(\hat{y}_i), s(x_i))$ <span style="color:red">← this should equals $x_i$</span>

  However, the style codes can be different and this can allow information about content to leak through it. The paper thus shuffles the style codes before computing $\hat{x}_i$

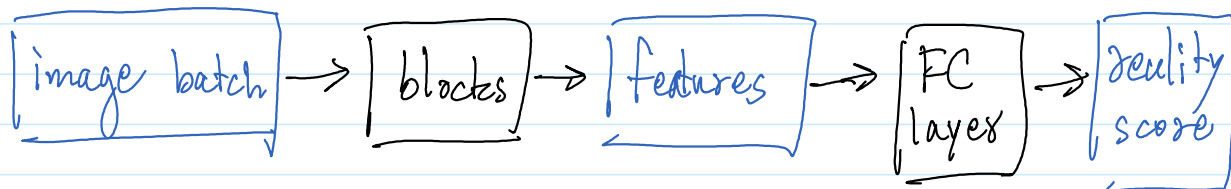So, now $\hat{x}_i = F_{Y \to X}(c(\hat{y}_i), S(x_{\pi(i)}))$ for some random permutation $\pi$,

The loss is then

https://arxiv.org/abs/1801.03924

$$\mathcal{L}_{cyc} = E\left[\|x_i - \hat{x}_i\|_2 + \lambda \cdot LPIPS(x_i, \hat{x}_i)\right]$$

$\Rightarrow$ Adversarial loss

- Normal discriminator

image batch $\to$ blocks $\to$ features $\to$ FC layers $\to$ reality score

- The paper uses the minibatch standard deviation trick used in the progressive GAN paper.

— Modified discriminator

image batch → blocks → features → FC → reality score   #1

features → σ → stddev of features → FC → reality score   #2

— Both reality scores are then use in the standard non-saturating GAN loss.

— The paper also uses the R1 regularization term from [Mescheder et. al. 2018] ( https://arxiv.org/pdf/1801.04406.pdf ) on both reality scores. This involves restricting gradients w.r.t. to inputs of the discriminator output.

$\Rightarrow$ The whole loss

$$\mathcal{L} = \lambda_{adv} \mathcal{L}_{adv} + \lambda_{scon} \mathcal{L}_{scon} + \lambda_{cyc} \mathcal{L}_{cyc}$$

(above the terms, in red: **1**, **10**, **20**)

- Implementation details

  $\Rightarrow$ Architecture = StyleGAN 2

  $\Rightarrow$ Style code $\in \mathbb{R}^8$ (Huh?)

  $\Rightarrow$ Batch size $\doteq 7$

  $\Rightarrow$ Adam with learning rate 0.002. for 300k iterations.

  $\Rightarrow$ Augmentations

  - horizontal flip
  - rotation between $(-20°, 20°)$
  - scaling between $(0.9, 1.1)$

  - translation up to $(0.1, 0.1)$
  - shearing up to $0.15$
  - Upscale to $286 \times 286$ then crop to $256 \times 256$ randomly.

⇒ Datasets ——— selfie2anime

AFHQ = animal faces

- Experiments

⇒ GANs N' Roses (GNR) produced more divese images from same human face + random style codes than DRIT++, CouncilGAN, and AniGAN.

http://vllab.ucmerced.edu/hylee/DRIT_pp/

⇒ The paper observed that the batch standard deviation trick is important to ensure diversity

⇒ The paper also outperformed others in several metrics: FID, DFID (original), LPIPS pairwise distances.