Vision & Perception 2019/20
# Project Presentation

Francesco Starna

*Sapienza University of Rome*

# Human Face Translation with GAN

Project
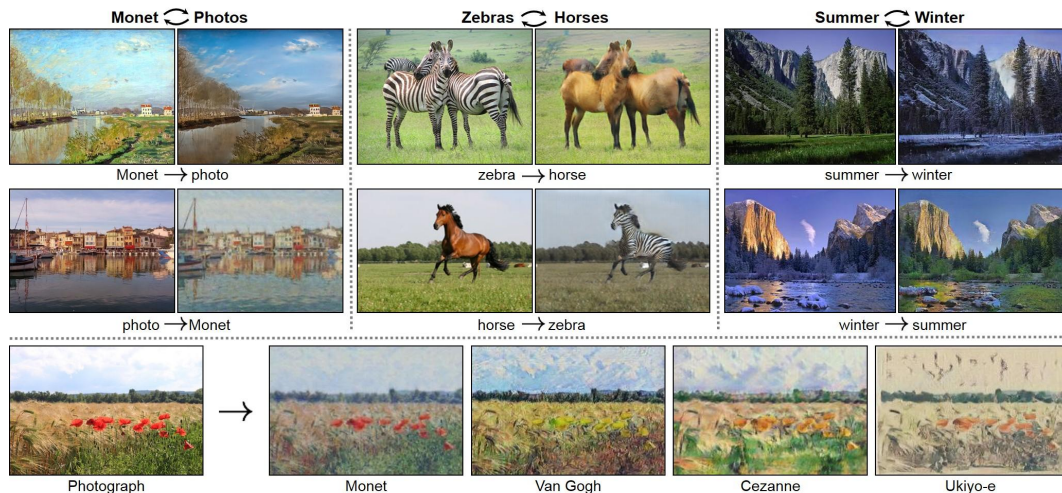
SAPIENZA

# Project Goal

**Image Translation** using Generative Adversarial Networks

## CycleGAN

(Jun-Yan Zhu et al.)
Unpaired Image-to-Image Translation using
Cycle-Consistent Adversarial Networks

# Generative Adversarial Networks

- Learns to generate new data with the same **statistic distribution** of the training set

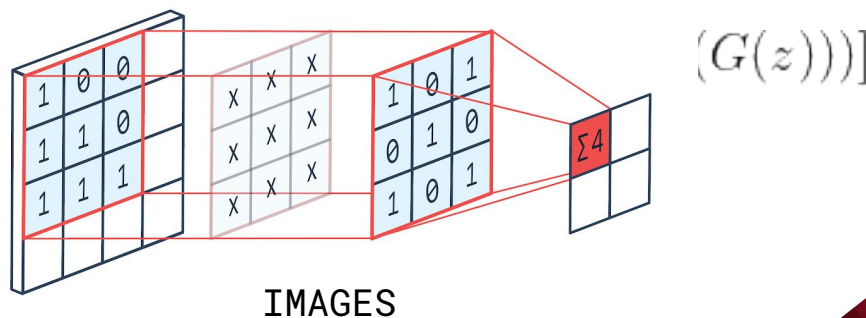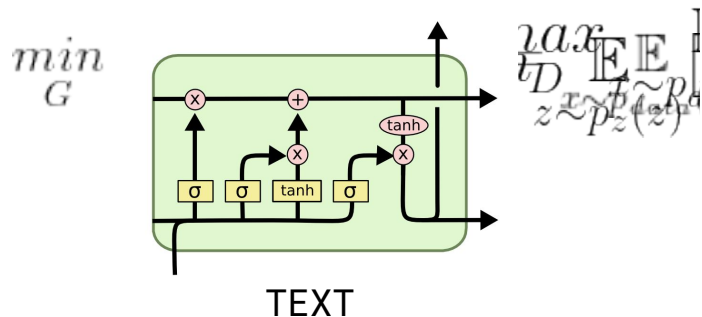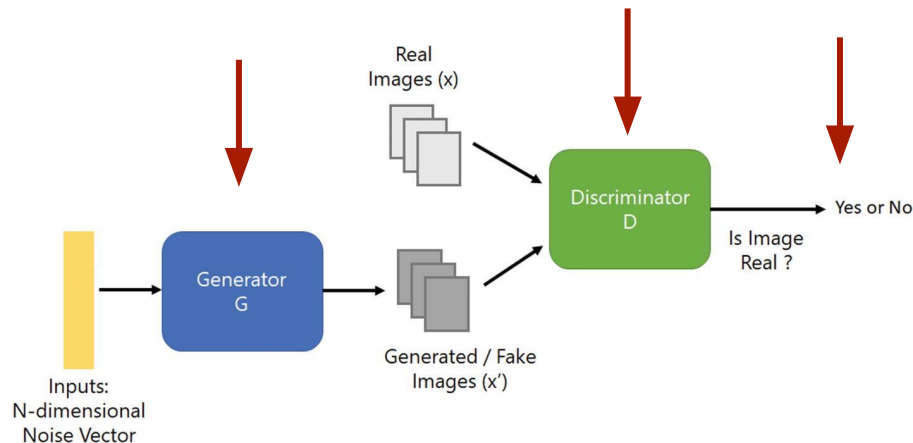- Supervised and Unsupervised learning methods

- Large domains application:
  -Image Generation
  -Image Translation
  -Super Resolution
  -Style Transfer
  -Text to Image
  … and many more

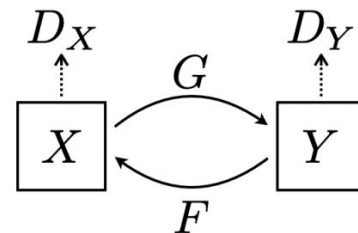this bird is red with white and has a very short beak

horse → zebra

SAPIENZA

# GAN Architecture
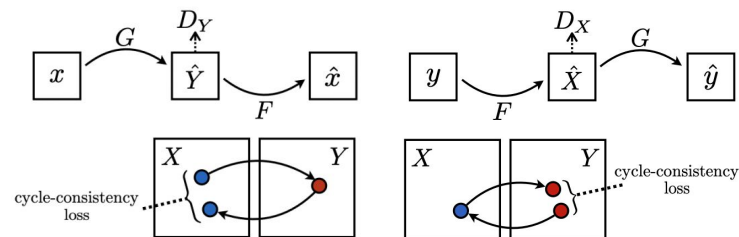
- Generator
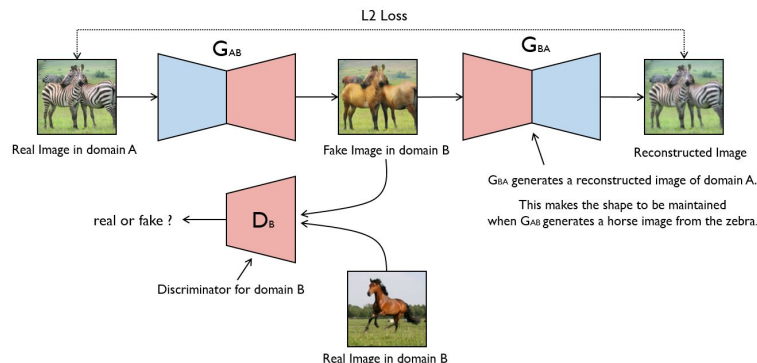
- Discriminator

- Loss Function

- Model

TEXT

IMAGES

# CycleGAN

- Unsupervised Learning of mappings
  **G: X → Y** and inverse **F: Y → X**



- **Cycle** consistency loss



- **Architecture**
  2 Generator
  2 Discriminator



L2 Loss

G_AB

G_BA

Real Image in domain A

Fake Image in domain B

Reconstructed Image

G_BA generates a reconstructed image of domain A.
This makes the shape to be maintained
when G_AB generates a horse image from the zebra.

real or fake ?

D_B

Discriminator for domain B

Real Image in domain B

SAPIENZA

# Generator

## ENCODER

- **Convolutional Layer**
- Highlights Extraction
- Downsample

## TRANSFORMER

- **Residual Connection**
- Join Features
- Same dimension

## DECODER

- **Transposed Convolutional Layer**
- Image Construction from low-level
- Upsample

| | |
|---|---|
| 64 Filters, 7×7, s=1 | k |
| 128 Filters, 3×3, s=2 | k/2 |
| 256 Filters, 3×3, s=2 | k/4 |
| 256 Filters, 3×3, s=2 | k/4 |
| 256 Filters, 3×3, s=2 | k/4 |
| 256 Filters, 3×3, s=2 | k/4 |
| 256 Filters, 3×3, s=2 | k/4 |
| 256 Filters, 3×3, s=2 | k/4 |
| 256 Filters, 3×3, s=2 | k/4 |
| 128 Filters, 3×3, s=2 | k/4 |
| 64 Filters, 3×3, s=2 | k/2 |
| 3 Filters, 7×7, s=1 | k |

# Discriminator

- ## PatchGAN

    (Isola et al.)
    Image to-image translation
    with conditional adversarial
    networks.

- ## Layers

    - Convolutional Layer
    - Instance Normalization
    - Leaky ReLU (0.01x if x < 0 )

- ## Mapping

    - 256x256 to NxN array
    - Average to classify Real or Fake



64 Filters, 4×4, s=2

128 Filters, 4×4, s=2

256 Filters, 4×4, s=2

512 Filters, 4×4, s=2

1 Filters, 4×4, s=1

Convolutional Layer

Instance Normalization, ReLU

Sigmoid

SAPIENZA

# Objective

- 2 Generators **G** and **F**

- 2 Discriminator **Dx** and **Dy**

- 1st Adversarial Loss

$$Loss_{advers}\left(G, D_y, X, Y\right) = \frac{1}{m} \sum \left(1 - D_y\left(G\left(x\right)\right)\right)^2$$

- 2nd Adversarial Loss

$$Loss_{advers}\left(F, D_x, Y, X\right) = \frac{1}{m} \sum \left(1 - D_x\left(F\left(y\right)\right)\right)^2$$

- **Cycle Consistency Loss**

$$Loss_{cyc}\left(G, F, X, Y\right) = \frac{1}{m} \left[\left(F\left(G\left(x_i\right)\right) - x_i\right) + \left(G\left(F\left(y_i\right)\right) - y_i\right)\right]$$

# Training

**generator loss**      =  total cycle loss + identity loss

**discriminator loss** =  disc real loss + disc fake loss

# Dataset

- ## Domain **X**
  FLICKR FACE

- ## Domain **Y**
  SIMPSON FACE
  ANIMAL FACE
  BITMOJI FACE

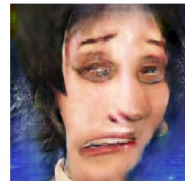**1000** Images Train
**100** Images Test
Each

# Preprocessing

- ## Normalization
  [0,255] to [-1,1]

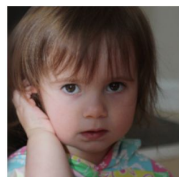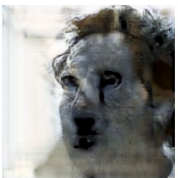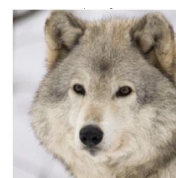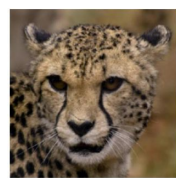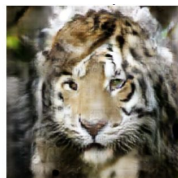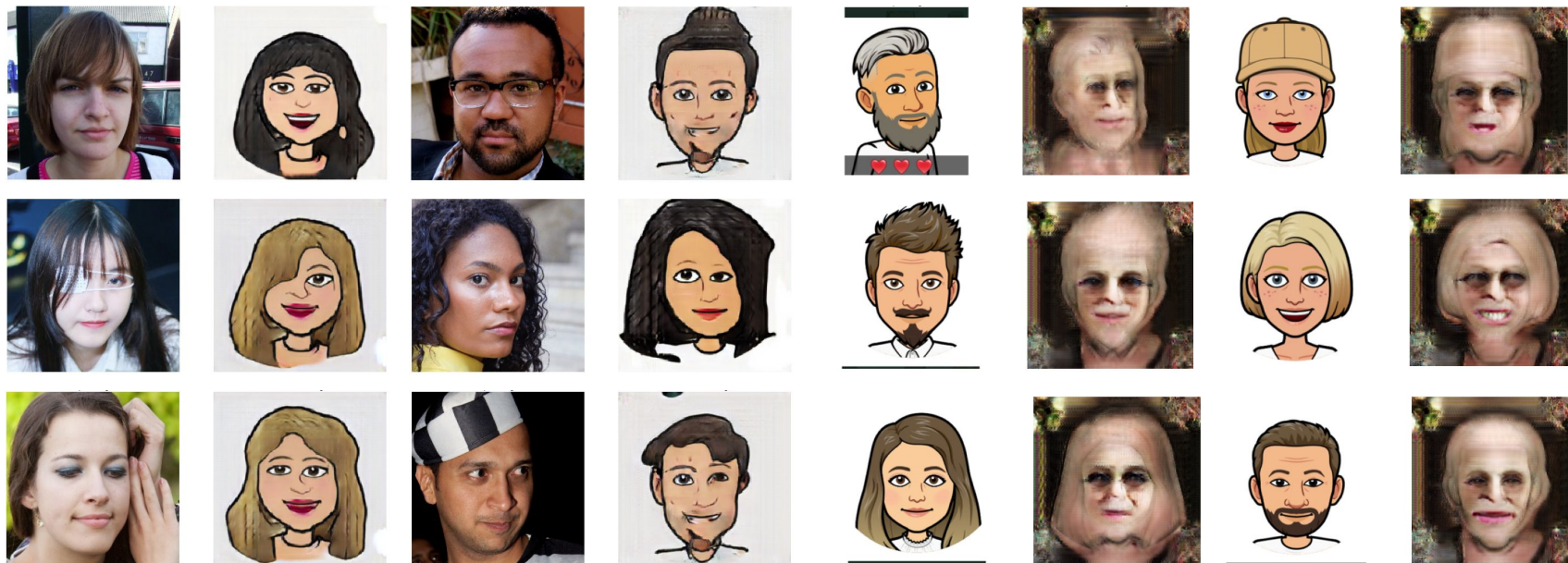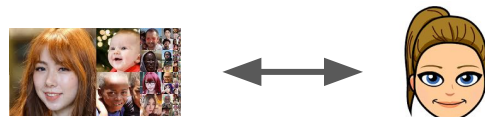- ## Data Augmentation
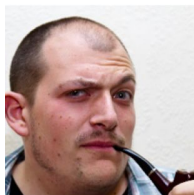  Random Jittering (resize, crop, flip)
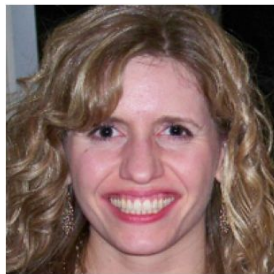
# Results (Human & Simpson)

# Results (Human & Animal)

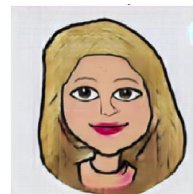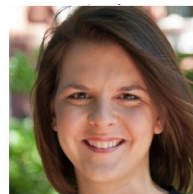# Results (Human & Bitmoji)

# Results

# Comments

- <u>Data Augmentation</u> > more training samples

- Large <u>Geometric Shifts</u> are not Successful



**Generator G result after 200 epochs training**

- <u>Visual Inspection</u> is better than more epochs
  Simpson    ~    100 epochs
  Animal      ~    150 epochs
  Bitmoji     ~    120 epochs

- CycleGAN can be <u>improved</u>
  Reduce oscillation by feeding the discriminator with
  a history of n generated images rather than last ones

SAPIENZA

# Thanks