

# Synthesizing Novel Pairs of Image and Text: Repurposing Generative Adversarial Networks as Semantic Autoencoders



Jason Xie, Tingwen Bao

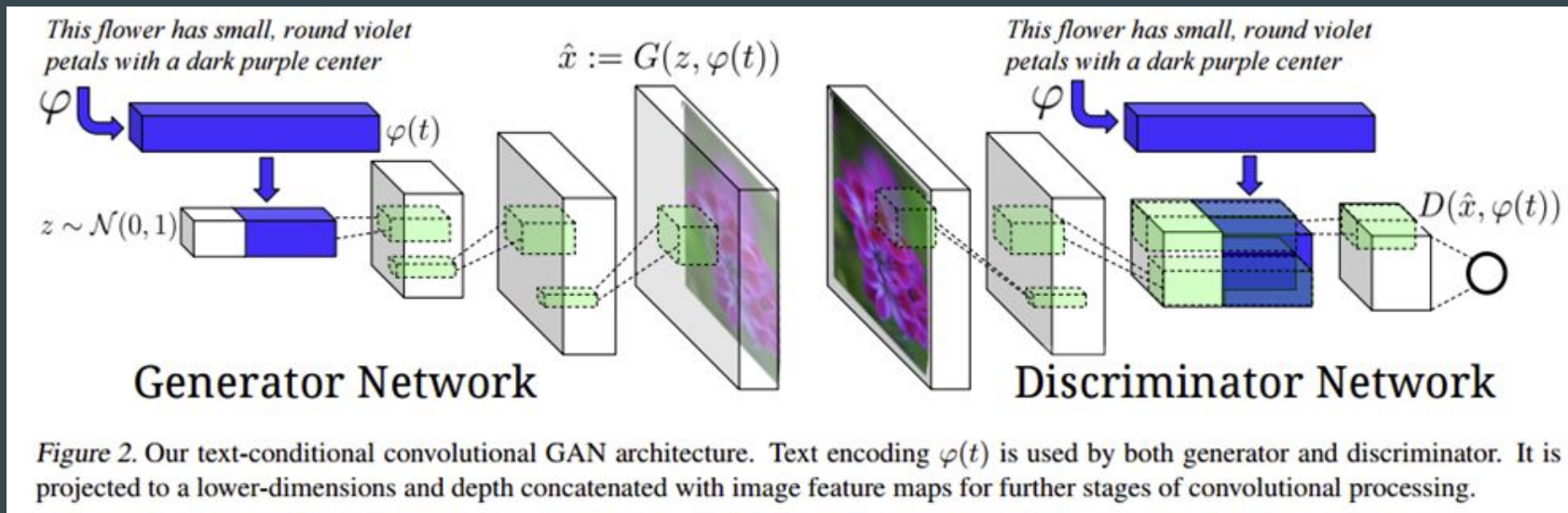
# Overview

In this project, we present strategies for generating coherent image and caption pairs based on an existing captioning dataset. Our model takes advantage of recent advances in generative adversarial networks and sequence-to-sequence modeling. We were able to generate novel paired samples from multiple domains and it can further be repurposed as semantic autoencoders.



# Background

- Generative Adversarial Networks (GAN)
- Text to Image Synthesis
- Image to Text Caption



# Our Model

$\text{Image}_{\text{new}}, \text{Text}_{\text{new}} \sim S(\text{Image}, \text{Text})$ , where  $\{\text{Image}_{\text{new}}, \text{Text}_{\text{new}}\} \notin \{\text{Image}_i, \text{Text}_i\}_N$



## 1. Source Domain Generation:

$\text{Image}_{\text{new}} \sim I(\text{Embedding}\{\text{Image}_i\})$  &  $\text{Text}_{\text{new}} \sim T(\text{Embedding}\{\text{Text}_i\})$

- Prototype Based for Text:  $\lambda * \text{Embedding}\{i_1\} + (1-\lambda) * \text{Embedding}\{i_2\}$
- Density Based for Image: Fit GMM with Image Embedding and sample from Gaussian Distribution

## 2. Target Domain Generation:

$\text{Text}_{\text{new}} \sim F(\text{Image}_{\text{new}})$  &  $\text{Image}_{\text{new}} \sim G(\text{Text}_{\text{new}})$

- Text to Image: GAN-CLS, loss function:

$$L_D < -\log(D(x, h)) + (\log(1 - D(x, h_{\text{fake}})) + \log(1 - D(x_{\text{fake}}, h))) / 2$$

- Image to Text: LSTM on the last convolution layer of DCGAN, maximize by the likelihood of correct caption:

$$\theta^* = \operatorname{argmax}_{\theta} \sum_{(\text{Image}_i, \text{Text}_i)} \log p(\text{Text}_i | \text{Image}_i; \theta)$$

# Our Model--Continued

- **Formulating a Circle: Image  $\rightarrow$  Text  $\rightarrow$  Image  $\rightarrow$  Text  $\rightarrow$ .....**

$$\text{Image}_{\text{new}} \sim G(\text{Embedding}\{F(\text{Embedding}\{\text{Image}_{\text{new}}\})\})$$

$$\text{Text}_{\text{new}} \sim F(\text{Embedding}\{G(\text{Embedding}\{\text{Text}_{\text{new}}\})\})$$

- **Semantic Autoencoder**

Text to Image: G is learning the true distribution of the data  $p_{\text{data}}(\text{image}|\text{text})$

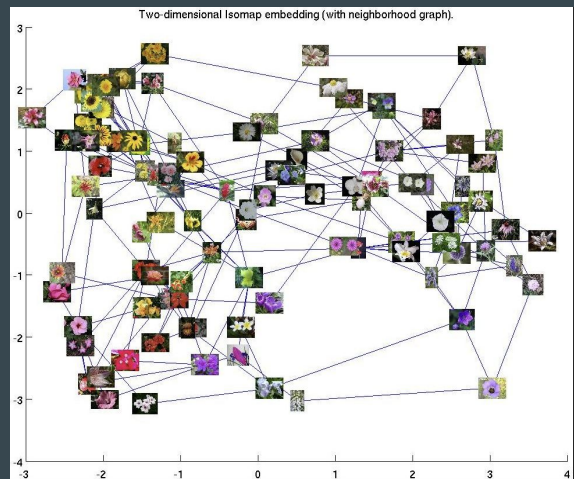
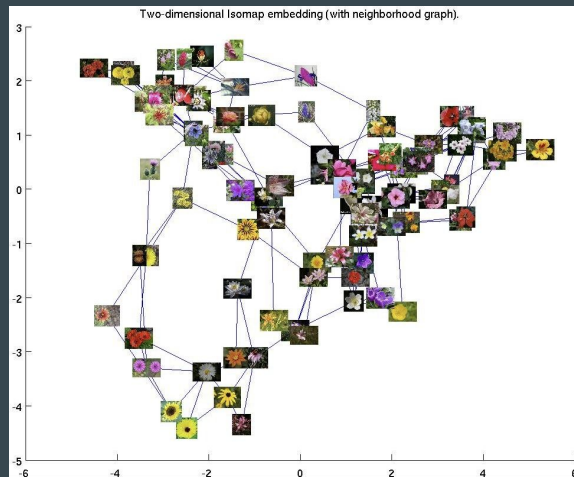
Image to Text: we are maximizing  $p_{\text{data}}(\text{text}|\text{image})$  directly

Note that  $p_{\text{data}}(\text{image}|\text{text})$  and  $p_{\text{data}}(\text{text}|\text{image})$  only preserving information that are shared between image and text. Information not shared are lossed.

# Dataset

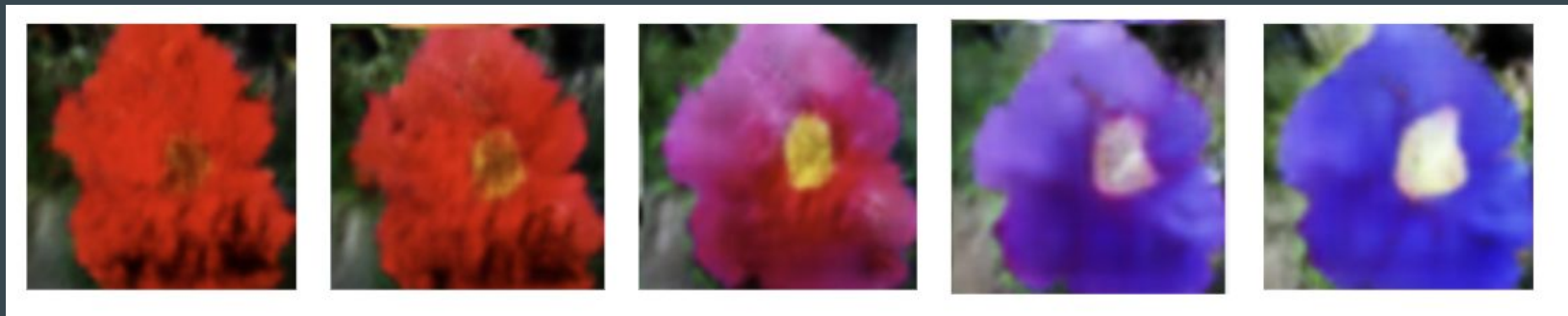
Oxford 102 flowers with caption:

- 102 Classes of flower
- 8000 images
- 10 captions per image



# Results

$$\text{Image}_{\text{new}} \sim F(\text{Text}_{\text{new}})$$



$$\lambda * \text{Embedding}\{\text{"The flower is red."}\} + (1-\lambda) * \text{Embedding}\{\text{"The flower is blue."}\}$$



$$\text{Text}_{\text{new}} \sim F(\text{Image}_{\text{new}})$$



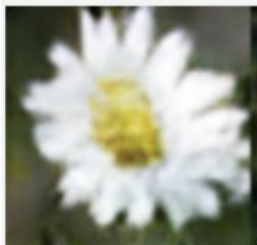
+



=

This flower has  
petals that are  
pink and has  
yellow stamen

# Cycle-- Image-> Text -> Image



→ this flower has a white circular center overlapping petals , and pale yellow stamen



→ this flower is red in color , and has petals that are fluffy and ruffled



→ the petals on this flower are yellow and fully enclosed folds



# Cycle-- Text -> Image -> Text

this flower has wing like  
petals with sharp long  
orange leaves



this flower is orange, with  
petals that are skinny and  
pointed.