



# TEXT TO ARTISTIC IMAGE GENERATION USING GANS

YUXING CHEN<sup>1</sup>, ZHEFAN WANG<sup>2</sup>

SYMBOLIC SYSTEMS PROGRAM, STANFORD UNIVERSITY<sup>1</sup>; DEPARTMENT OF ELECTRICAL ENGINEERING, STANFORD UNIVERSITY<sup>2</sup>



## ABSTRACT

Generating images from texts has been a trending research topic in computer vision. Style transferring between photos and artwork is also a popular subfield. We build an application that combines these two, which allows the user to not only generate ordinary photo-like images from sentences, but also get the certain artistic style of images specified by the user.

## DATASET

### Microsoft COCO dataset 2014:

328,000 RGB images, size of  $256 \times 256$ , 5 captions per image, 91 object categories, 80K/40K train/val split



5 Captions for Current Image:  
(1) a giraffe standing next to a forest filled with trees.  
(2) a giraffe eating food from the top of the tree.  
(3) two giraffes standing in a tree filled area.  
(4) a giraffe mother with its baby in the forest.  
(5) a giraffe standing up nearby a tree.

Filename: COCO\_train2014\_000000000025.jpg

### Text to Image Generation

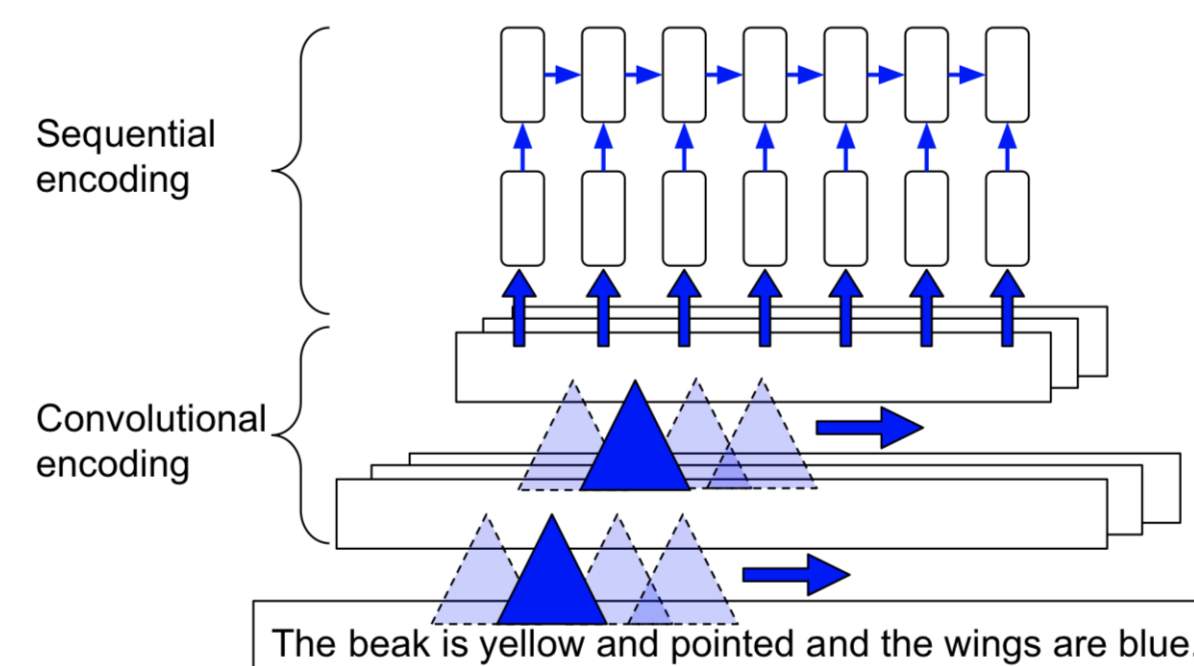
40,000 training images, 2,000 validation images

### Style Transfer

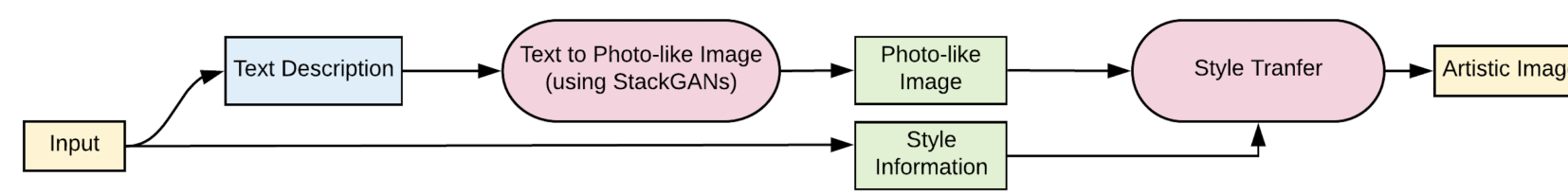
8,000 training images, 3,000 validation images  
No annotation needed

## TEXT ENCODING

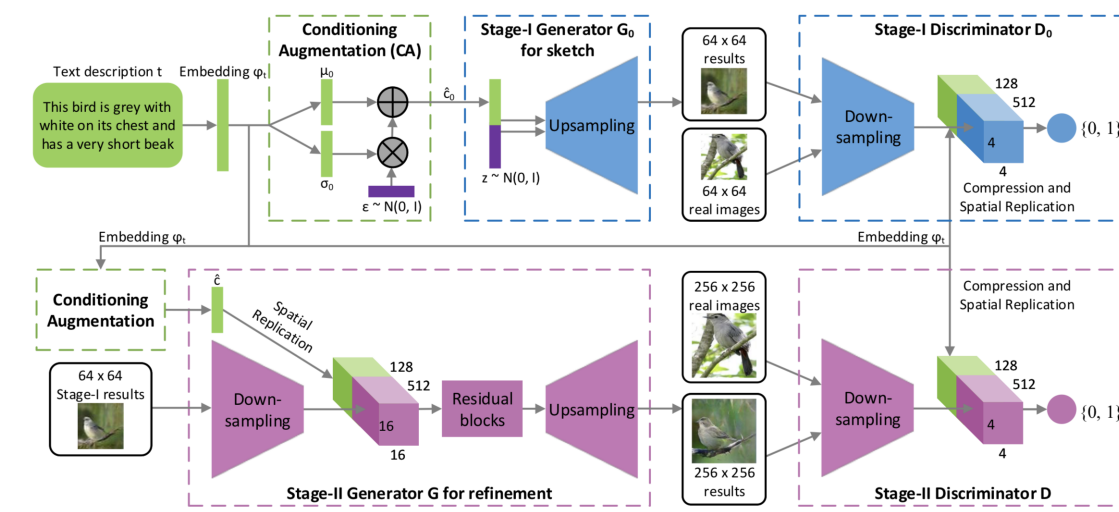
A widely-applied way of encoding text descriptions is to use deep convolution and recurrent text encoder (i.e. char-CNN-RNN model) which learn the correspondence function with images. The idea of this approach is that a recurrent neural network is stacked on top of a temporal convolutional neural network hidden layer.



## MODEL

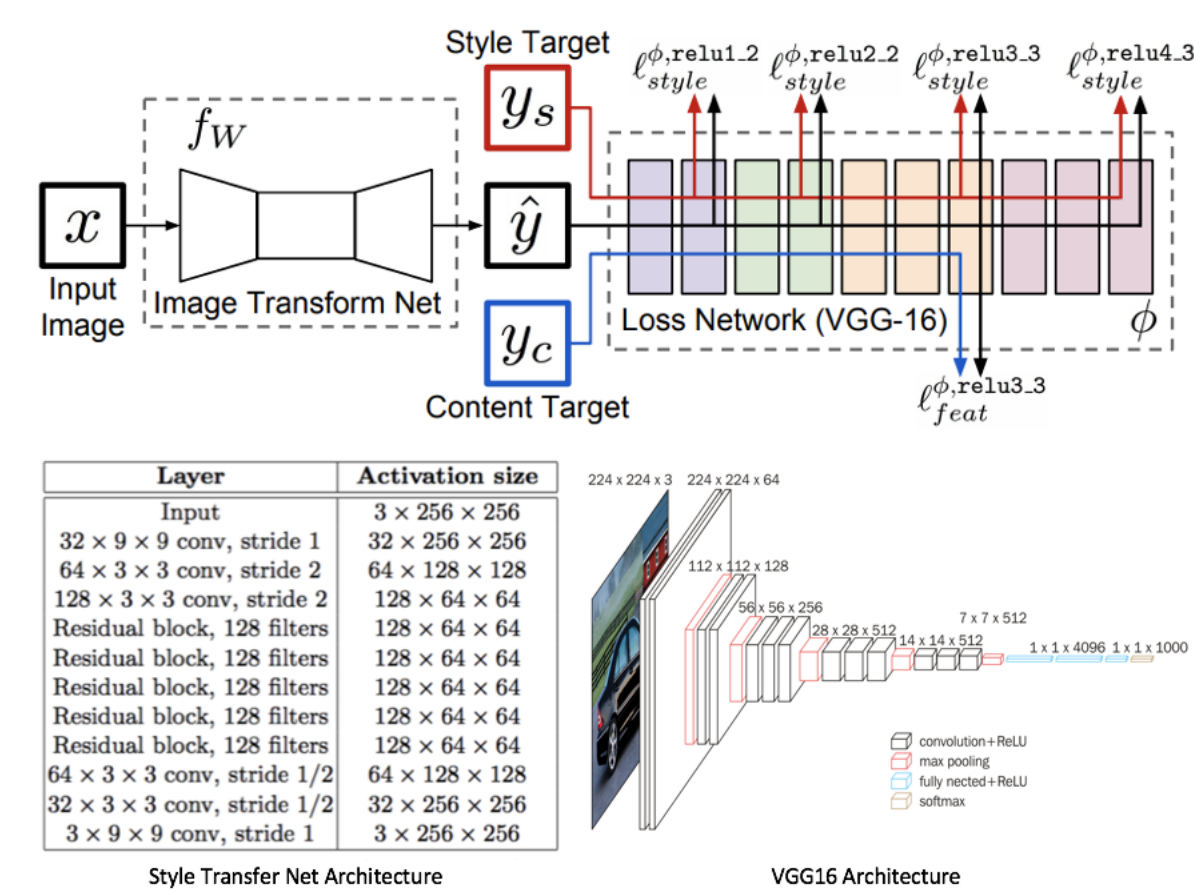


### Text to Image Generation



Each up-sampling block contains the nearest-neighbor up-sampling followed by a  $3 \times 3$  stride 1 conv2d layer, Batch normalization and ReLU activation. Residual block consists of a  $3 \times 3$  stride 1 conv2d layer followed by Batch normalization, ReLU activation, another  $3 \times 3$  stride 1 conv2d layer, Batch normalization and ReLU activation. Each down-sampling block has a  $4 \times 4$  stride 2 conv2d layer, followed by Batch normalization and LeakyReLU, except that the first down-sampling block does not contain Batch normalization.

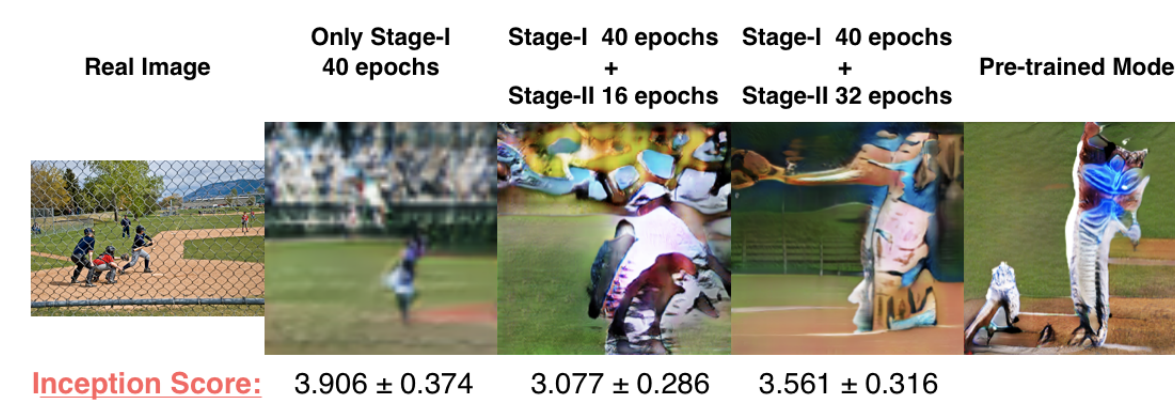
### Style Transfer



Loss = lambda\_s \* (phi\_style^phi\_relu1.2(y\_hat, y\_s) + phi\_style^phi\_relu2.2(y\_hat, y\_s) + phi\_style^phi\_relu3.3(y\_hat, y\_s) + phi\_style^phi\_relu4.3(y\_hat, y\_s)) + lambda\_c \* (phi\_feat^phi\_relu3.3(y\_hat, y\_c))

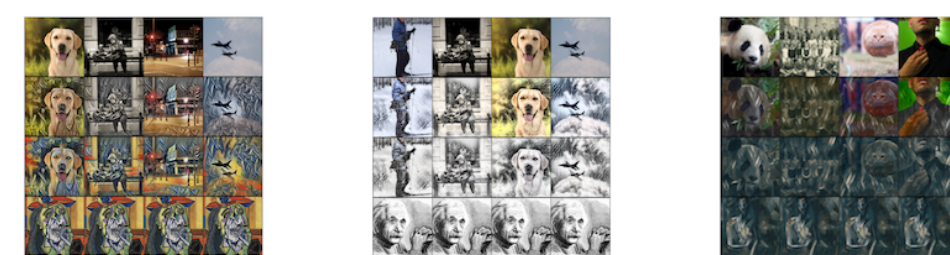
## EXPERIMENTS & RESULTS

### Text to Image Generation



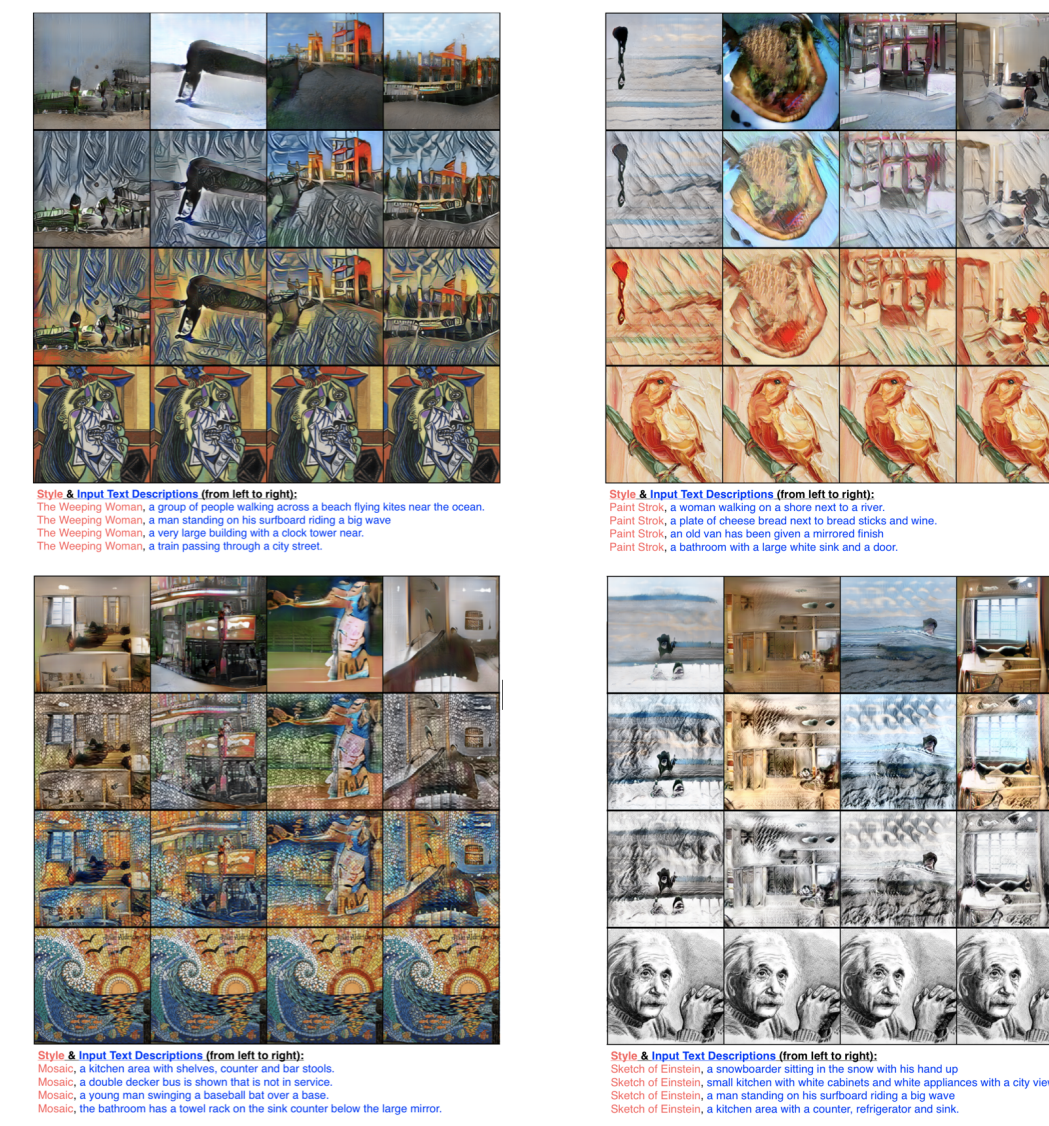
Inception Score: 3.906 ± 0.374 3.077 ± 0.286 3.561 ± 0.316  
Text Description: a young man swinging a baseball bat over a base.

### Style Transfer



Row 1: input validation images. Row 2: style transferred images with original color preserved. The third row is the style transferred images without color preservation. The last row is the style images.

### Overall Results



## EVALUATION & ANALYSIS

### Text to Image Generation

We find that our current model did well in plotting the "background" indicated in the text descriptions, while it often failed to "fill in" the details of the human-like objects, as shown in the figure below.



Captions (from left to right):  
baseball players practicing their batting skills in a filled arena  
a giraffe walking across a lush green field.  
a woman walking on a shore next to a river.  
a very cute goat standing in some very tall grass.

First row: real images. Second row: fake images.

### Style Transfer

- loss = lambda\_c \* content\_loss + lambda\_s \* style\_loss. The larger the ratio lambda\_s / lambda\_c, the stronger the influence of the style image.
- Color preservation mechanism sometimes introduces undesirable artifacts. Eg: darker results, white artifacts around strokes, etc.

## FUTURE WORK

- Compare StackGAN with other creative models (eg. AttnGANs).
- Improve StackGAN by using VisDial dialogues along with MS COCO captions.
- Find better approaches to preserving color.

## REFERENCE

[1] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaolei Huang, Xiaogang Wang, and Dimitris Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. arXiv:1612.03242, 2017.  
[2] Justin Johnson, Alexandre Alahi, and Fei-Fei Li. Perceptual losses for real-time style transfer and super-resolution. arXiv:1603.08155, 2016.