

**CS597 REPORT 4**  
**Yash Pradeep Gupte**  
**MS Computer Science**  
**A20472798**

**1] CODE – StyleAttnGAN**

Repository link - <https://github.com/sidward14/Style-AttnGAN>

**(A) Datasets**

- Bird – CUB-200-2011 Number of categories: 200, Number of images: 11,788
- Coco – MS COCO - 80 object categories, 5 captions per image

**(B) Training**

- (1) Pre training DAMSM model/ encoder
- (2) Train Style- AttnGAN models
- (3) Train original AttnGAN models

**(C) Sampling**

Transformer encoder : 'gpt2'

Configuration file : eval\_bird\_style.yml

Image Transform:[Resize(size=304, interpolation=bilinear, max\_size=None, antialias=None), RandomCrop(size=(256, 256), padding=None), RandomHorizontalFlip(p=0.5)]

TextDataset:

Bbox , captions, classids, filenames, imsize (64,128,256), IndextoWords, N\_words(5450), WordtoIndex.

Dataloader :

Batchsize(8), Textdataset

Sampling function which generates images from pretrained models:

Split\_directory = 'valid'

Generator model = TRAIN.NET\_G

Text\_encoder, net\_G =

**MODEL – Architecture**

G\_NET\_STYLED(

(ca\_net): CA\_NET(

(fc): Linear(in\_features=256, out\_features=400, bias=True)

```

    (relu): GLU()
)
(map_net): StyleConditionedMappingNetwork(
  (fc_mapping_model): Sequential(
    (pixelnorm): NormalizeLayer(
      (norm): PixelNorm2d()
    )
    (fc_0): LinearEx(
      (linear): Linear(in_features=200, out_features=100, bias=True)
    )
    (nl_0): LeakyReLU(negative_slope=0.2)
    (fc_1): LinearEx(
      (linear): Linear(in_features=100, out_features=100, bias=True)
    )
    (nl_1): LeakyReLU(negative_slope=0.2)
    (fc_2): LinearEx(
      (linear): Linear(in_features=100, out_features=100, bias=True)
    )
    (nl_2): LeakyReLU(negative_slope=0.2)
    (fc_3): LinearEx(
      (linear): Linear(in_features=100, out_features=100, bias=True)
    )
    (nl_3): LeakyReLU(negative_slope=0.2)
    (fc_4): LinearEx(
      (linear): Linear(in_features=100, out_features=100, bias=True)
    )
    (nl_4): LeakyReLU(negative_slope=0.2)
    (fc_5): LinearEx(
      (linear): Linear(in_features=100, out_features=100, bias=True)
    )
    (nl_5): LeakyReLU(negative_slope=0.2)
    (fc_6): LinearEx(
      (linear): Linear(in_features=100, out_features=100, bias=True)
    )
    (nl_6): LeakyReLU(negative_slope=0.2)
    (fc_7): LinearEx(
      (linear): Linear(in_features=100, out_features=100, bias=True)
    )
    (nl_7): LeakyReLU(negative_slope=0.2)
  )
)
(h_net1): INIT_STAGE_G_STYLED(
  (gen_layers): ModuleList(
    (0): ModuleList(

```

```

(0): None
(1): StyleAddNoise()
(2): Sequential(
  (0): Conv2dBias()
  (1): LeakyReLU(negative_slope=0.2)
  (2): NormalizeLayer(
    (norm): InstanceNorm2d(128, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
  )
)
(3): LinearEx(
  (linear): Linear(in_features=100, out_features=256, bias=True)
)
(1): ModuleList(
  (0): Conv2dEx(
    (conv2d): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  )
  (1): StyleAddNoise()
  (2): Sequential(
    (0): Conv2dBias()
    (1): LeakyReLU(negative_slope=0.2)
    (2): NormalizeLayer(
      (norm): InstanceNorm2d(None, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
    )
  )
  (3): LinearEx(
    (linear): Linear(in_features=100, out_features=256, bias=True)
  )
)
(2): ModuleList(
  (0): Sequential(
    (0): Upsample(scale_factor=2.0, mode=nearest)
    (1): Conv2dEx(
      (conv2d): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    )
  )
  (2): Lambda()
)
(1): StyleAddNoise()
(2): Sequential(
  (0): Conv2dBias()
  (1): LeakyReLU(negative_slope=0.2)
  (2): NormalizeLayer(

```

```

        (norm): InstanceNorm2d(None, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
    )
)
(3): LinearEx(
  (linear): Linear(in_features=100, out_features=256, bias=True)
)
)
(3): ModuleList(
  (0): Sequential(
    (0): Conv2dEx(
      (conv2d): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    )
  )
  (1): StyleAddNoise()
  (2): Sequential(
    (0): Conv2dBias()
    (1): LeakyReLU(negative_slope=0.2)
    (2): NormalizeLayer(
      (norm): InstanceNorm2d(None, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
    )
  )
)
(3): LinearEx(
  (linear): Linear(in_features=100, out_features=256, bias=True)
)
)
(4): ModuleList(
  (0): Sequential(
    (0): Upsample(scale_factor=2.0, mode=nearest)
    (1): Conv2dEx(
      (conv2d): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    )
    (2): Lambda()
  )
  (1): StyleAddNoise()
  (2): Sequential(
    (0): Conv2dBias()
    (1): LeakyReLU(negative_slope=0.2)
    (2): NormalizeLayer(
      (norm): InstanceNorm2d(None, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
    )
  )
)
)

```

```

(3): LinearEx(
  (linear): Linear(in_features=100, out_features=256, bias=True)
)
)
(5): ModuleList(
  (0): Sequential(
    (0): Conv2dEx(
      (conv2d): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    )
  )
  (1): StyleAddNoise()
  (2): Sequential(
    (0): Conv2dBias()
    (1): LeakyReLU(negative_slope=0.2)
    (2): NormalizeLayer(
      (norm): InstanceNorm2d(None, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
    )
  )
  (3): LinearEx(
    (linear): Linear(in_features=100, out_features=256, bias=True)
  )
)
(6): ModuleList(
  (0): Sequential(
    (0): Upsample(scale_factor=2.0, mode=nearest)
    (1): Conv2dEx(
      (conv2d): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    )
    (2): Lambda()
  )
  (1): StyleAddNoise()
  (2): Sequential(
    (0): Conv2dBias()
    (1): LeakyReLU(negative_slope=0.2)
    (2): NormalizeLayer(
      (norm): InstanceNorm2d(None, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
    )
  )
  (3): LinearEx(
    (linear): Linear(in_features=100, out_features=256, bias=True)
  )
)

```

```

(7): ModuleList(
  (0): Sequential(
    (0): Conv2dEx(
      (conv2d): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    )
  )
  (1): StyleAddNoise()
  (2): Sequential(
    (0): Conv2dBias()
    (1): LeakyReLU(negative_slope=0.2)
    (2): NormalizeLayer(
      (norm): InstanceNorm2d(None, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
    )
  )
  (3): LinearEx(
    (linear): Linear(in_features=100, out_features=256, bias=True)
  )
)
(8): ModuleList(
  (0): Sequential(
    (0): Upsample(scale_factor=2.0, mode=nearest)
    (1): Conv2dEx(
      (conv2d): Conv2d(128, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    )
    (2): Lambda()
  )
  (1): StyleAddNoise()
  (2): Sequential(
    (0): Conv2dBias()
    (1): LeakyReLU(negative_slope=0.2)
    (2): NormalizeLayer(
      (norm): InstanceNorm2d(None, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
    )
  )
  (3): LinearEx(
    (linear): Linear(in_features=100, out_features=128, bias=True)
  )
)
(9): ModuleList(
  (0): Sequential(
    (0): Conv2dEx(
      (conv2d): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)

```

```

    )
    )
    (1): StyleAddNoise()
    (2): Sequential(
      (0): Conv2dBias()
      (1): LeakyReLU(negative_slope=0.2)
      (2): NormalizeLayer(
        (norm): InstanceNorm2d(None, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
      )
    )
    (3): LinearEx(
      (linear): Linear(in_features=100, out_features=128, bias=True)
    )
  )
  (upsampler): Upsample(scale_factor=2.0, mode=nearest)
  (nl): LeakyReLU(negative_slope=0.2)
  (norm): NormalizeLayer(
    (norm): InstanceNorm2d(None, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
  )
)
(img_net1): GET_IMAGE_G_STYLED(
  (torgb): Sequential(
    (0): Conv2dEx(
      (conv2d): Conv2d(64, 3, kernel_size=(1, 1), stride=(1, 1))
    )
    (1): Tanh()
  )
)
(h_net2): NEXT_STAGE_G_STYLED(
  (att): GlobalAttentionGeneral(
    (conv_context): Conv2d(256, 64, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (sm): Softmax(dim=1)
  )
  (residual): ResBlock(
    (block): Sequential(
      (0): Conv2dEx(
        (conv2d): Conv2d(128, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      )
      (1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
      (2): GLU()
    )
  )
)

```

```

(3): Conv2dEx(
  (conv2d): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
)
(4): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
)
)
(gen_layers): ModuleList(
  (0): ModuleList(
    (0): Sequential(
      (0): Upsample(scale_factor=2.0, mode=nearest)
      (1): Conv2dEx(
        (conv2d): Conv2d(128, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      )
      (2): Lambda()
    )
    (1): StyleAddNoise()
    (2): Sequential(
      (0): Conv2dBias()
      (1): LeakyReLU(negative_slope=0.2)
      (2): NormalizeLayer(
        (norm): InstanceNorm2d(None, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
      )
    )
    (3): LinearEx(
      (linear): Linear(in_features=100, out_features=128, bias=True)
    )
  )
  (1): ModuleList(
    (0): Sequential(
      (0): Conv2dEx(
        (conv2d): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      )
    )
    (1): StyleAddNoise()
    (2): Sequential(
      (0): Conv2dBias()
      (1): LeakyReLU(negative_slope=0.2)
      (2): NormalizeLayer(
        (norm): InstanceNorm2d(None, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
      )
    )
  )
)

```



```

        (3): LinearEx(
          (linear): Linear(in_features=100, out_features=128, bias=True)
        )
      )
    )
  (upsampler): Upsample(scale_factor=2.0, mode=nearest)
  (nl): LeakyReLU(negative_slope=0.2)
  (norm): NormalizeLayer(
    (norm): InstanceNorm2d(None, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
  )
)
(img_net2): GET_IMAGE_G_STYLED(
  (torgb): Sequential(
    (0): Conv2dEx(
      (conv2d): Conv2d(64, 3, kernel_size=(1, 1), stride=(1, 1))
    )
    (1): Tanh()
  )
)
(h_net3): NEXT_STAGE_G_STYLED(
  (att): GlobalAttentionGeneral(
    (conv_context): Conv2d(256, 64, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (sm): Softmax(dim=1)
  )
  (residual): ResBlock(
    (block): Sequential(
      (0): Conv2dEx(
        (conv2d): Conv2d(128, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      )
      (1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
      (2): GLU()
      (3): Conv2dEx(
        (conv2d): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      )
      (4): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    )
  )
)
(gen_layers): ModuleList(
  (0): ModuleList(
    (0): Sequential(
      (0): Upsample(scale_factor=2.0, mode=nearest)

```

```

    (1): Conv2dEx(
      (conv2d): Conv2d(128, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    )
    (2): Lambda()
  )
  (1): StyleAddNoise()
  (2): Sequential(
    (0): Conv2dBias()
    (1): LeakyReLU(negative_slope=0.2)
    (2): NormalizeLayer(
      (norm): InstanceNorm2d(None, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
    )
  )
  (3): LinearEx(
    (linear): Linear(in_features=100, out_features=128, bias=True)
  )
  (1): ModuleList(
    (0): Sequential(
      (0): Conv2dEx(
        (conv2d): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      )
    )
    (1): StyleAddNoise()
    (2): Sequential(
      (0): Conv2dBias()
      (1): LeakyReLU(negative_slope=0.2)
      (2): NormalizeLayer(
        (norm): InstanceNorm2d(None, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
      )
    )
    (3): LinearEx(
      (linear): Linear(in_features=100, out_features=128, bias=True)
    )
  )
  (upsampler): Upsample(scale_factor=2.0, mode=nearest)
  (nl): LeakyReLU(negative_slope=0.2)
  (norm): NormalizeLayer(
    (norm): InstanceNorm2d(None, eps=1e-08, momentum=0.1, affine=False,
track_running_stats=False)
  )

```

```

)
(img_net3): GET_IMAGE_G_STYLED(
  (torgb): Sequential(
    (0): Conv2dEx(
      (conv2d): Conv2d(64, 3, kernel_size=(1, 1), stride=(1, 1))
    )
    (1): Tanh()
  )
)
)
)
)

```

```

Net_G.eval()
  Text_encoder : RNN_ENCODER(
    (encoder): Embedding(5450, 300)
    (drop): Dropout(p=0.5, inplace=False)
    (rnn): LSTM(300, 128, batch_first=True, dropout=0.5,
bidirectional=True)
  )

```

**Batch\_size** = 8

**Noise** : tensor dimension = (8,100) -> (batch\_size, nz)

**Dataloader** (for 1<sup>st</sup> batch)

```

Keys - ['036.Northern_Flicker/Northern_Flicker_0072_28678',
'004.Groove_billed_Ani/Groove_Billed_Ani_0031_1588',
'112.Great_Grey_Shrike/Great_Grey_Shrike_0086_106533',
'166.Golden_winged_Warbler/Golden_Winged_Warbler_0087_794810',
'098.Scott_Oriole/Scott_Oriole_0084_795860',
'180.Wilson_Warbler/Wilson_Warbler_0050_175573',
'023.Brandt_Cormorant/Brandt_Cormorant_0073_23259',
'029.American_Crow/American_Crow_0130_25163']

```

Imgs –

```

0 – (8,3,64,64)
1 – (8,3,128,128)
2 – (8,3,256,256)

```

Class\_ids : [ 36 4 112 166 98 180 23 29]

Captions : Tensor(8,18)

Cap\_lens = 8 captions in one batch – Max length of words in captions = 18

Word\_embeddings : size(8,256,18)

Sentence embeddings : size(8,256)

#Generate Fake images

fake\_imgs, \_, \_ = netG(noise, sent\_emb, words\_embs, mask)

Pass the noise, sentence embeddings, word embeddings and mask

#### **(D) Losses ->**

*Discriminator loss – BCE Binary cross entropy*

--- Losses.py file ---

if netD.UNCOND\_DNET is not None:

errD = ((real\_errD + cond\_real\_errD) / 2. + (fake\_errD + cond\_fake\_errD +  
cond\_wrong\_errD) / 3.)

else:

errD = cond\_real\_errD + (cond\_fake\_errD + cond\_wrong\_errD) / 2.

*Generator Loss – BCE*

--- Losses.py file LINE :187 ---

errG = nn.BCELoss()(logits, real\_labels)

g\_loss = errG + cond\_errG

*Rankingloss – word features and sentence code*

w\_loss0, w\_loss1, \_ = words\_loss(region\_features, words\_embs,  
match\_labels, cap\_lens,  
class\_ids, batch\_size)

w\_loss = (w\_loss0 + w\_loss1) \* cfg.TRAIN.SMOOTH.LAMBDA

s\_loss0, s\_loss1 = sent\_loss(cnn\_code, sent\_emb, match\_labels, class\_ids,  
batch\_size)

s\_loss = (s\_loss0 + s\_loss1) \* cfg.TRAIN.SMOOTH.LAMBDA

#### **(E) Training – config file : bird\_attn2\_style.yml**

Command : python main.py --cfg cfg/bird\_attn2\_style.yml --gpu 0

<class 'model.G\_NET\_STYLED'>

<class 'model.D\_NET\_STYLED64'>

<class 'model.D\_NET\_STYLED128'>

<class 'model.D\_NET\_STYLED256'>

# of netsD 3

Optimizers – Adam optimizer with LR = 0.002