

# Exercise: Text to Image Synthesis using DCGAN

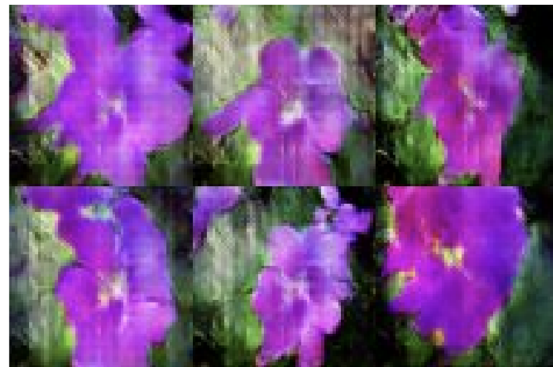
Quoc-Thai Nguyen và Quang-Vinh Dinh

*PR-Team: Đăng-Nhã Nguyễn, Minh-Châu Phạm và Hoàng-Nguyên Vũ*

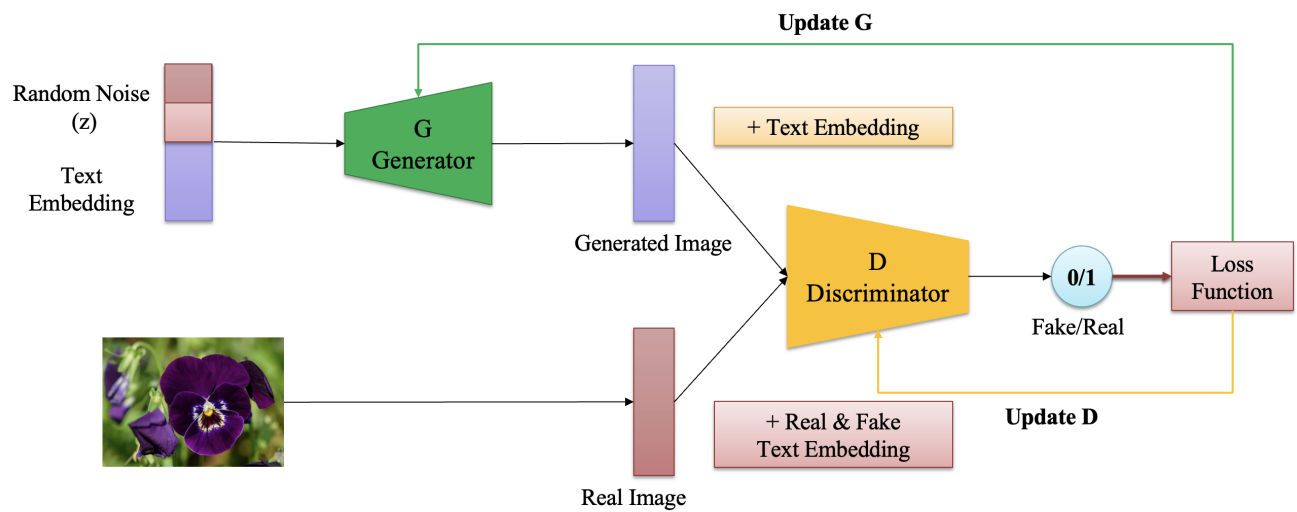
Ngày 17 tháng 3 năm 2024

## Phần I. Giới thiệu

The flower has petals that  
are bright pinkish purple  
with white stigma



Hình 1: Ví dụ sinh văn bản thành hình ảnh sử dụng mô hình DCGAN.



Hình 2: Mô hình DCGAN sinh văn bản thành hình ảnh.

**Sinh hình ảnh từ văn bản hay tổng hợp văn bản thành hình ảnh (Text to Image)** là bài toán ngày càng được ứng dụng mạnh mẽ trong lĩnh vực trí tuệ nhân tạo. Các mô hình được huấn luyện với đầu vào là đoạn văn bản và đầu ra là hình ảnh mô tả hoặc chứa các đối tượng được mô tả trong đoạn văn bản. Ví dụ về tổng hợp hình ảnh từ văn bản được mô tả trong Hình 1.

Hiện nay, có nhiều mô hình có thể xây dựng và giải quyết bài toán này có thể kể đến như GANs, Diffusion Models,... Ở trong phần này chúng ta sẽ sử dụng mô hình DCGAN (Deep Convolution Generative Adversarial Networks) để huấn luyện mô hình giải quyết bài toán này.

Mô hình DCGAN gồm 2 mạng là Generator và Discriminator. Trong đó:

1. Generator: Bao gồm các lớp mạng CNN nhận đầu vào là Noise ( $z$ ) và Text Embedding ( $e$  - Biểu diễn của đoạn văn bản đầu vào). Đầu ra của khối Generator là ma trận ảnh được sinh ra  $G(z, e)$ .
2. Discriminator: Bao gồm các lớp mạng CNN nhận đầu vào là biểu diễn của ảnh (ảnh thật hoặc ảnh fake được sinh ra từ khối Generator) sau đó kết hợp với Text Embedding để dự đoán ảnh thật hay ảnh fake.

# Phần II. Text To Image Synthesis using DCGAN

Trong phần này chúng ta sẽ xây dựng và huấn luyện mô hình DCGAN trên bộ dữ liệu flowers. Bộ dữ liệu flowers bao gồm:

- (a) Thư mục chứa hình ảnh: [link](#)
- (b) Thư mục chứa các đoạn văn bản mô tả hình ảnh: [link](#)
- (c) Thư mục chứa các vector biểu diễn cho các đoạn văn bản mô tả hình ảnh: [link](#)

## 1. Data Preparing

```
1 # install libs
2 !pip install -q torchfile
3
4 # import libs
5 import os
6 import h5py
7 import torchfile
8 from glob import glob
9
10 # read dataset and save h5py file
11 images_path = './data/'
12 embedding_path = './data/flowers_icml/'
13 text_path = './data/cvpr2016_flowers/text_c10/'
14 save_dataset_path = './data/flowers.hdf5'
15
16 train_path = './data/flowers_icml/trainclasses.txt'
17 val_path = './data/flowers_icml/valclasses.txt'
18 test_path = './data/flowers_icml/testclasses.txt'
19
20 train_classes = open(train_path).read().splitlines()
21 val_classes = open(val_path).read().splitlines()
22 test_classes = open(test_path).read().splitlines()
23
24 f = h5py.File(save_dataset_path, 'w')
25 train = f.create_group('train')
26 val = f.create_group('val')
27 test = f.create_group('test')
28
29 for _class in sorted(os.listdir(embedding_path)):
30     split = ''
31     if _class in train_classes:
32         split = 'train'
33     elif _class in val_classes:
34         split = 'val'
35     elif _class in test_classes:
36         split = 'test'
37
38     data_path = os.path.join(embedding_path, _class)
39     txt_path = os.path.join(text_path, _class)
40     for example, txt_file in zip(sorted(glob(data_path + "/*.t7")), sorted(glob(
41         txt_path + "/*.txt"))):
42         example_data = torchfile.load(example)
43         img_path = example_data[b'img'].decode("utf-8")
```

```

43     embeddings = example_data[b'txt']
44     example_name = img_path.split('/')[-1][: -4]
45
46     f = open(txt_file, "r")
47     txt = f.readlines()
48     f.close()
49
50     img_path = os.path.join(images_path, img_path)
51     img = open(img_path, 'rb').read()
52
53     print(example_name, txt)
54     txt_choice = np.random.choice(range(10), 5)
55
56     embeddings = embeddings[txt_choice]
57     txt = np.array(txt)
58     txt = txt[txt_choice]
59     dt = h5py.special_dtype(vlen=str)
60
61     for c, e in enumerate(embeddings):
62         ex = split.create_group(example_name + '_' + str(c))
63         ex.create_dataset('name', data=example_name)
64         ex.create_dataset('img', data=np.void(img))
65         ex.create_dataset('embeddings', data=e)
66         ex.create_dataset('class', data=_class)
67         ex.create_dataset('txt', data=txt[c].astype(object), dtype=dt)
68
69     print(example_name, txt[1], _class)

```

## 2. Dataset

```

1 import io
2 import torch
3 import numpy as np
4 from PIL import Image
5 from torch.utils.data import Dataset
6
7 class Text2ImageDataset(Dataset):
8
9     def __init__(self, dataset_path, transform=None, split=0):
10         self.dataset_path = dataset_path
11         self.transform = transform
12         self.dataset = h5py.File(self.dataset_path, mode='r')
13         self.split = 'train' if split == 0 else 'val' if split == 1 else 'test'
14         self.dataset_keys = [str(k) for k in self.dataset[self.split].keys()]
15         self.h5py2int = lambda x: int(np.array(x))
16
17     def __len__(self):
18         length = len(self.dataset[self.split])
19         return length
20
21     def __getitem__(self, idx):
22         example_name = self.dataset_keys[idx]
23         example = self.dataset[self.split][example_name]
24
25         right_image = bytes(np.array(example['img']))
26         right_embed = np.array(example['embeddings'], dtype=float)
27         wrong_image = bytes(np.array(self.find_wrong_image(example['class'])))
28
29         right_image = Image.open(io.BytesIO(right_image)).resize((64, 64))
30         wrong_image = Image.open(io.BytesIO(wrong_image)).resize((64, 64))

```

```

31
32     right_image = self.validate_image(right_image)
33     wrong_image = self.validate_image(wrong_image)
34
35     try:
36         txt = np.array(example['txt']).astype(str)
37     except:
38
39         txt = np.array([example['txt'][(0)].decode('utf-8', errors='replace')])
40         txt = np.char.replace(txt, ' í ', ' ').astype(str)
41
42     sample = {
43         'right_images': torch.FloatTensor(right_image),
44         'right_embed': torch.FloatTensor(right_embed),
45         'wrong_images': torch.FloatTensor(wrong_image),
46         'txt': str(txt)
47     }
48
49     sample['right_images'] = sample['right_images'].sub_(127.5).div_(127.5)
50     sample['wrong_images'] = sample['wrong_images'].sub_(127.5).div_(127.5)
51
52     return sample
53
54     def find_wrong_image(self, category):
55         idx = np.random.randint(len(self.dataset_keys))
56         example_name = self.dataset_keys[idx]
57         example = self.dataset[self.split][example_name]
58         _category = example['class']
59
60         if _category != category:
61             return example['img']
62
63         return self.find_wrong_image(category)
64
65
66     def validate_image(self, img):
67         img = np.array(img, dtype=float)
68         return img.transpose(2, 0, 1)
69
70 train_dataset = Text2ImageDataset(save_dataset_path, split=0)
71 val_dataset = Text2ImageDataset(save_dataset_path, split=1)
72 test_dataset = Text2ImageDataset(save_dataset_path, split=2)
73
74 # dataloader
75 from torch.utils.data import DataLoader
76
77 batch_size = 2048
78 train_loader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True)
79 val_loader = DataLoader(val_dataset, batch_size=batch_size)
80 test_loader = DataLoader(test_dataset, batch_size=batch_size)

```

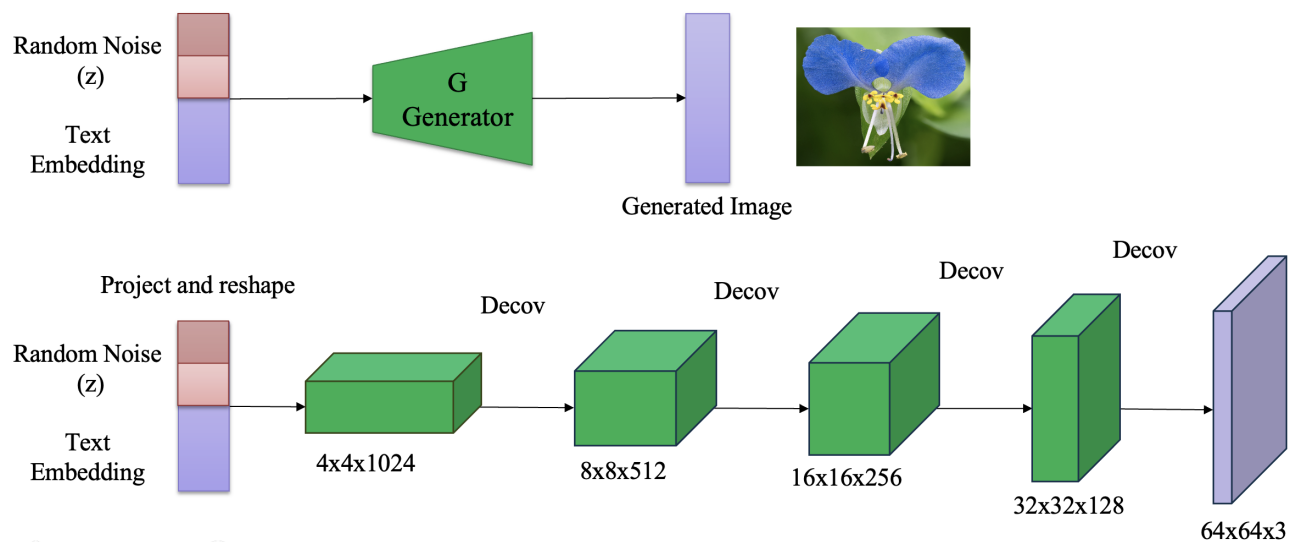
### 3. Model

Trong phần này chúng ta xây dựng mô hình DCGAN bao gồm 2 mạng Generator và Discriminator:

#### 3.1. Generator

Khối Generator sinh ra ảnh từ văn bản đầu vào:

- (a) Input: Nhận đầu vào là vector: Random Noise ( $z$ ) có kích thước  $R$ , được nối với Vector Embedding (Vector  $e$  - biểu diễn cho cả đoạn văn bản đầu vào) có kích thước là  $D$ . Vì vậy, vector đầu vào mạng Generator là  $R + D$ .
- (b) Output: Sau khi học mối quan hệ để sinh ảnh, giá trị đầu ra của Generator sẽ là biểu diễn các điểm ảnh dự đoán trong không gian 3 chiều có kích thước là Channel x Width x Height (Ví dụ,  $C \times W \times H$  -  $3 \times 64 \times 64$ )



Hình 3: Khối Generator trong DCGAN.

```

1 import torch
2 import torch.nn as nn
3
4 # The Generator model
5 class Generator(nn.Module):
6     def __init__(self, channels, noise_dim=100, embed_dim=1024, embed_out_dim=128):
7         super(Generator, self).__init__()
8         self.channels = channels
9         self.noise_dim = noise_dim
10        self.embed_dim = embed_dim
11        self.embed_out_dim = embed_out_dim
12
13        # Text embedding layers
14        self.text_embedding = nn.Sequential(
15            nn.Linear(self.embed_dim, self.embed_out_dim),
16            nn.BatchNorm1d(self.embed_out_dim),
17            nn.LeakyReLU(0.2, inplace=True)
18        )
19
20        # Generator architecture
21        model = []

```

```

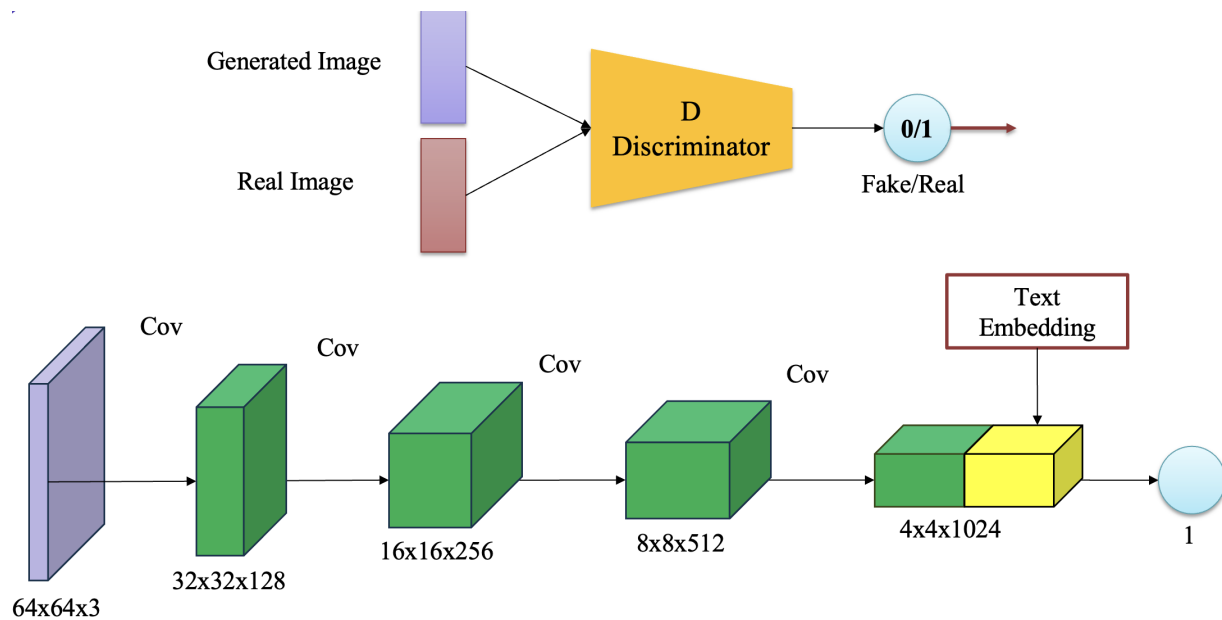
22     model += self._create_layer(self.noise_dim + self.embed_out_dim, 512, 4,
stride=1, padding=0)
23     model += self._create_layer(512, 256, 4, stride=2, padding=1)
24     model += self._create_layer(256, 128, 4, stride=2, padding=1)
25     model += self._create_layer(128, 64, 4, stride=2, padding=1)
26     model += self._create_layer(64, self.channels, 4, stride=2, padding=1, output=
True)
27
28     self.model = nn.Sequential(*model)
29
30     def _create_layer(self, size_in, size_out, kernel_size=4, stride=2, padding=1,
output=False):
31         layers = [nn.ConvTranspose2d(size_in, size_out, kernel_size, stride=stride,
padding=padding, bias=False)]
32         if output:
33             layers.append(nn.Tanh()) # Tanh activation for the output layer
34         else:
35             layers += [nn.BatchNorm2d(size_out), nn.ReLU(True)] # Batch normalization
and ReLU for other layers
36         return layers
37
38     def forward(self, noise, text):
39         # Apply text embedding to the input text
40         text = self.text_embedding(text)
41         text = text.view(text.shape[0], text.shape[1], 1, 1) # Reshape to match the
generator input size
42         z = torch.cat([text, noise], 1) # Concatenate text embedding with noise
43         return self.model(z)
44
45     def weights_init(m):
46         classname = m.__class__.__name__
47         if classname.find('Conv') != -1:
48             m.weight.data.normal_(0.0, 0.02)
49         elif classname.find('BatchNorm') != -1:
50             m.weight.data.normal_(1.0, 0.02)
51             m.bias.data.fill_(0)
52
53 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
54 embed_dim = 1024
55 noise_dim = 100
56 embed_out_dim = 64
57 generator = Generator(
58     channels=3, embed_dim=embed_dim, noise_dim=noise_dim, embed_out_dim=embed_out_dim
59 ).to(device)
60 generator.apply(weights_init)

```

### 3.2. Discriminator

Khối Discriminator dự đoán hình ảnh Fake/Real, bao gồm các lớp CNN với:

- Input: Nhận đầu vào là ma trận điểm ảnh và Text Embedding. Ma trận điểm ảnh (CxHxW) sau khi qua các lớp CNN để học các đặc trưng của ảnh đầu vào sẽ được nối với vector embedding.
- Output: Vector sau khi được nối sẽ được sử dụng để dự đoán 0(Fake) hoặc 1(Real).



Hình 4: Khối Discriminator trong DCGAN.

- Đầu tiên chúng ta sẽ định nghĩa lớp Embedding để nối các đặc trưng của ảnh và Text Embedding.

```

1 # The Embedding model
2 class Embedding(nn.Module):
3     def __init__(self, size_in, size_out):
4         super(Embedding, self).__init__()
5         self.text_embedding = nn.Sequential(
6             nn.Linear(size_in, size_out),
7             nn.BatchNorm1d(size_out),
8             nn.LeakyReLU(0.2, inplace=True)
9         )
10
11     def forward(self, x, text):
12         embed_out = self.text_embedding(text)
13         embed_out_resize = embed_out.repeat(4, 4, 1, 1).permute(2, 3, 0, 1) # Resize
14         # to match the discriminator input size
15         out = torch.cat([x, embed_out_resize], 1) # Concatenate text embedding with
16         # the input feature map
17         return out

```

- Sau đó, chúng ta định nghĩa kiến trúc mạng Discriminator:

```

1 # The Discriminator model
2 class Discriminator(nn.Module):
3     def __init__(self, channels, embed_dim=1024, embed_out_dim=128):
4         super(Discriminator, self).__init__()
5         self.channels = channels
6         self.embed_dim = embed_dim
7         self.embed_out_dim = embed_out_dim
8
9         # Discriminator architecture
10        self.model = nn.Sequential(
11            *self._create_layer(self.channels, 64, 4, 2, 1, normalize=False),
12            *self._create_layer(64, 128, 4, 2, 1),
13            *self._create_layer(128, 256, 4, 2, 1),
14            *self._create_layer(256, 512, 4, 2, 1)
15        )

```



```

16         self.text_embedding = Embedding(self.embed_dim, self.embed_out_dim) # Text
embedding module
17         self.output = nn.Sequential(
18             nn.Conv2d(512 + self.embed_out_dim, 1, 4, 1, 0, bias=False), nn.Sigmoid()
19         )
20
21     def _create_layer(self, size_in, size_out, kernel_size=4, stride=2, padding=1,
normalize=True):
22         layers = [nn.Conv2d(size_in, size_out, kernel_size=kernel_size, stride=stride,
padding=padding)]
23         if normalize:
24             layers.append(nn.BatchNorm2d(size_out))
25             layers.append(nn.LeakyReLU(0.2, inplace=True))
26         return layers
27
28     def forward(self, x, text):
29         x_out = self.model(x) # Extract features from the input using the
discriminator architecture
30         out = self.text_embedding(x_out, text) # Apply text embedding and concatenate
with the input features
31         out = self.output(out) # Final discriminator output
32         return out.squeeze(), x_out
33
34 discriminator = Discriminator(
35     channels=3, embed_dim=embed_dim, embed_out_dim=embed_out_dim
36 ).to(device)
37 discriminator.apply(weights_init)

```

#### 4. Training

```

1 # setting up Adam optimizer for Generator and Discriminator
2 learning_rate = 0.0002
3 optimizer_G = torch.optim.Adam(
4     generator.parameters(), lr=learning_rate, betas=(0.5, 0.999)
5 )
6 optimizer_D = torch.optim.Adam(
7     discriminator.parameters(), lr=learning_rate, betas=(0.5, 0.999)
8 )
9
10 # loss functions
11 criterion = nn.BCELoss()
12 l2_loss = nn.MSELoss()
13 l1_loss = nn.L1Loss()
14
15 num_epochs = 20
16 real_label = 1.
17 fake_label = 0.
18 l1_coef = 50
19 l2_coef = 100
20
21 D_losses = []
22 G_losses = []
23
24 for epoch in range(num_epochs):
25     epoch_D_loss = []
26     epoch_G_loss = []
27     batch_time = time.time()
28
29     for idx, batch in enumerate(train_loader):
30

```

```

31     images = batch['right_images'].to(device)
32     wrong_images = batch['wrong_images'].to(device)
33     embeddings = batch['right_embed'].to(device)
34     batch_size = images.size(0)
35
36     # ===== #
37     #                               Train the discriminator                               #
38     # ===== #
39
40     # Clear gradients for the discriminator
41     optimizer_D.zero_grad()
42
43     # Generate random noise
44     noise = torch.randn(batch_size, noise_dim, 1, 1, device=device)
45
46     # Generate fake image batch with the generator
47     fake_images = generator(noise, embeddings)
48
49     # Forward pass real batch and calculate loss
50     real_out, real_act = discriminator(images, embeddings)
51     d_loss_real = criterion(real_out, torch.full_like(real_out, real_label, device
=device))
52
53     # Forward pass wrong batch and calculate loss
54     wrong_out, wrong_act = discriminator(wrong_images, embeddings)
55     d_loss_wrong = criterion(wrong_out, torch.full_like(wrong_out, fake_label,
=device))
56
57     # Forward pass fake batch and calculate loss
58     fake_out, fake_act = discriminator(fake_images.detach(), embeddings)
59     d_loss_fake = criterion(fake_out, torch.full_like(fake_out, fake_label, device
=device))
60
61     # Compute total discriminator loss
62     d_loss = d_loss_real + d_loss_wrong + d_loss_fake
63
64     # Backpropagate the gradients
65     d_loss.backward()
66
67     # Update the discriminator
68     optimizer_D.step()
69
70     # ===== #
71     #                               Train the generator                               #
72     # ===== #
73
74     # Clear gradients for the generator
75     optimizer_G.zero_grad()
76
77     # Generate new random noise
78     noise = torch.randn(batch_size, noise_dim, 1, 1, device=device)
79
80     # Generate new fake images using Generator
81     fake_images = generator(noise, embeddings)
82
83     # Get discriminator output for the new fake images
84     out_fake, act_fake = discriminator(fake_images, embeddings)
85
86     # Get discriminator output for the real images
87     out_real, act_real = discriminator(images, embeddings)

```

```

88
89     # Calculate losses
90     g_bce = criterion(out_fake, torch.full_like(out_fake, real_label, device=
device))
91     g_l1 = l1_coef * l1_loss(fake_images, images)
92     g_l2 = l2_coef * l2_loss(torch.mean(act_fake, 0), torch.mean(act_real, 0).
detach())
93
94     # Compute total generator loss
95     g_loss = g_bce + g_l1 + g_l2
96
97     # Backpropagate the gradients
98     g_loss.backward()
99
100    # Update the generator
101    optimizer_G.step()
102
103    epoch_D_loss.append(d_loss.item())
104    epoch_G_loss.append(g_loss.item())
105
106    print('Epoch {} [{} / {}] loss_D: {:.4f} loss_G: {:.4f} time: {:.2f}'.format(
107        epoch+1, idx+1, len(train_loader),
108        d_loss.mean().item(),
109        g_loss.mean().item(),
110        time.time() - batch_time)
111    )
112    D_losses.append(sum(epoch_D_loss)/len(epoch_D_loss))
113    G_losses.append(sum(epoch_G_loss)/len(epoch_G_loss))
114
115    # save model
116    model_save_path = './save_model'
117    torch.save(generator.state_dict(), os.path.join(model_save_path, 'generator.pth'))
118    torch.save(discriminator.state_dict(), os.path.join(model_save_path, 'discriminator.pth
'))

```

## 5. Prediction

```

1 import matplotlib.pyplot as plt
2
3 example = next(iter(test_dataset))
4
5 # show real image
6 plt.imshow(example['right_images'].permute(1, 2, 0))
7 plt.show()
8
9 # prediction
10 embeddings = example['right_embed'].to(device)
11 noise = torch.randn(1, noise_dim, 1, 1, device=device)
12 pred = generator(noise, embeddings.unsqueeze(0))
13
14 # show the prediction image
15 plt.imshow(pred[0].cpu().detach().permute(1, 2, 0))
16 plt.show()

```

# Phần III. Text To Image Synthesis using DCGAN and BERTs

Trong phần này chúng ta sẽ cải tiến mô hình tạo hình ảnh từ văn bản bằng cách sử dụng các mô hình ngôn ngữ giàu ngữ cảnh hơn để biểu diễn văn bản thành các vector như BERTs,...

## 1. Data Preparing

```

1 images_path = './data/'
2 text_path = './data/cvpr2016_flowers/text_c10/'
3 save_dataset_path = './data/flowers.hdf5' # from section 2
4
5 # load BERTs
6 import torch
7 from transformers import AutoTokenizer, DistilBertModel
8
9 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
10 tokenizer = AutoTokenizer.from_pretrained("distilbert-base-uncased")
11 model = DistilBertModel.from_pretrained("distilbert-base-uncased")
12 model.to(device)
13 model.eval()
14
15 # Max Pooling - Take the max value over time for every dimension.
16 def max_pooling(model_output, attention_mask):
17     token_embeddings = model_output[0] #First element of model_output contains all
    token_embeddings
18     input_mask_expanded = attention_mask.unsqueeze(-1).expand(token_embeddings.size())
    .float()
19     token_embeddings[input_mask_expanded == 0] = -1e9 # Set padding tokens to large
    negative value
20     return torch.max(token_embeddings, 1)[0]
21
22 def convert_text_to_feature(sentences, max_length=50):
23     inputs = tokenizer.batch_encode_plus(
24         sentences, padding='max_length', max_length=max_length, truncation=True,
    return_tensors='pt'
25     )
26     input_ids = inputs['input_ids'].to(device)
27     attention_mask = inputs['attention_mask'].to(device)
28     outputs = model(input_ids=input_ids, attention_mask=attention_mask)
29     sentence_embeddings = max_pooling(outputs, attention_mask)
30     return sentence_embeddings

```

## 2. Dataset

```

1 import io
2 import torch
3 import numpy as np
4 from PIL import Image
5 from torch.utils.data import Dataset
6
7 class Text2ImageDataset(Dataset):
8
9     def __init__(self, dataset_path, transform=None, split=0):
10         self.dataset_path = dataset_path
11         self.transform = transform
12         self.dataset = h5py.File(self.dataset_path, mode='r')

```

```

13     self.split = 'train' if split == 0 else 'val' if split == 1 else 'test'
14     self.dataset_keys = [str(k) for k in self.dataset[self.split].keys()]
15     self.h5py2int = lambda x: int(np.array(x))
16
17     def __len__(self):
18         length = len(self.dataset[self.split])
19         return length
20
21     def __getitem__(self, idx):
22         example_name = self.dataset_keys[idx]
23         example = self.dataset[self.split][example_name]
24
25         right_image = bytes(np.array(example['img']))
26         right_embed = np.array(example['embeddings'], dtype=float)
27         wrong_image = bytes(np.array(self.find_wrong_image(example['class'])))
28
29         right_image = Image.open(io.BytesIO(right_image)).resize((64, 64))
30         wrong_image = Image.open(io.BytesIO(wrong_image)).resize((64, 64))
31
32         right_image = self.validate_image(right_image)
33         wrong_image = self.validate_image(wrong_image)
34
35         try:
36             txt = np.array(example['txt']).astype(str)
37         except:
38
39             txt = np.array([example['txt'][(0)].decode('utf-8', errors='replace')])
40             txt = np.char.replace(txt, ' í ', ' ').astype(str)
41
42         embeddings = convert_text_to_feature([str(txt)]).detach().cpu()
43
44         sample = {
45             'right_images': torch.FloatTensor(right_image),
46             'right_embed': embeddings[0],
47             'wrong_images': torch.FloatTensor(wrong_image),
48             'txt': str(txt)
49         }
50
51         sample['right_images'] = sample['right_images'].sub_(127.5).div_(127.5)
52         sample['wrong_images'] = sample['wrong_images'].sub_(127.5).div_(127.5)
53
54         return sample
55
56     def find_wrong_image(self, category):
57         idx = np.random.randint(len(self.dataset_keys))
58         example_name = self.dataset_keys[idx]
59         example = self.dataset[self.split][example_name]
60         _category = example['class']
61
62         if _category != category:
63             return example['img']
64
65         return self.find_wrong_image(category)
66
67     def validate_image(self, img):
68         img = np.array(img, dtype=float)
69         return img.transpose(2, 0, 1)
70
71 # dataset
72 train_dataset = Text2ImageDataset(save_dataset_path, split=0)

```

```

73 val_dataset = Text2ImageDataset(save_dataset_path, split=1)
74 test_dataset = Text2ImageDataset(save_dataset_path, split=2)
75
76 # dataloader
77 from torch.utils.data import DataLoader
78
79 batch_size = 2048
80 train_loader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True)
81 val_loader = DataLoader(val_dataset, batch_size=batch_size)
82 test_loader = DataLoader(test_dataset, batch_size=batch_size)

```

### 3. Model

```

1 import torch
2 import torch.nn as nn
3
4 # The Generator model
5 class Generator(nn.Module):
6     def __init__(self, channels, noise_dim=100, embed_dim=1024, embed_out_dim=128):
7         super(Generator, self).__init__()
8         self.channels = channels
9         self.noise_dim = noise_dim
10        self.embed_dim = embed_dim
11        self.embed_out_dim = embed_out_dim
12
13        # Text embedding layers
14        self.text_embedding = nn.Sequential(
15            nn.Linear(self.embed_dim, self.embed_out_dim),
16            nn.BatchNorm1d(self.embed_out_dim),
17            nn.LeakyReLU(0.2, inplace=True)
18        )
19
20        # Generator architecture
21        model = []
22        model += self._create_layer(self.noise_dim + self.embed_out_dim, 512, 4,
stride=1, padding=0)
23        model += self._create_layer(512, 256, 4, stride=2, padding=1)
24        model += self._create_layer(256, 128, 4, stride=2, padding=1)
25        model += self._create_layer(128, 64, 4, stride=2, padding=1)
26        model += self._create_layer(64, self.channels, 4, stride=2, padding=1, output=
True)
27
28        self.model = nn.Sequential(*model)
29
30        def _create_layer(self, size_in, size_out, kernel_size=4, stride=2, padding=1,
output=False):
31            layers = [nn.ConvTranspose2d(size_in, size_out, kernel_size, stride=stride,
padding=padding, bias=False)]
32            if output:
33                layers.append(nn.Tanh()) # Tanh activation for the output layer
34            else:
35                layers += [nn.BatchNorm2d(size_out), nn.ReLU(True)] # Batch normalization
and ReLU for other layers
36            return layers
37
38        def forward(self, noise, text):
39            # Apply text embedding to the input text
40            text = self.text_embedding(text)
41            text = text.view(text.shape[0], text.shape[1], 1, 1) # Reshape to match the
generator input size

```

```

42         z = torch.cat([text, noise], 1) # Concatenate text embedding with noise
43         return self.model(z)
44
45
46 # The Embedding model
47 class Embedding(nn.Module):
48     def __init__(self, size_in, size_out):
49         super(Embedding, self).__init__()
50         self.text_embedding = nn.Sequential(
51             nn.Linear(size_in, size_out),
52             nn.BatchNorm1d(size_out),
53             nn.LeakyReLU(0.2, inplace=True)
54         )
55
56     def forward(self, x, text):
57         embed_out = self.text_embedding(text)
58         embed_out_resize = embed_out.repeat(4, 4, 1, 1).permute(2, 3, 0, 1) # Resize
to match the discriminator input size
59         out = torch.cat([x, embed_out_resize], 1) # Concatenate text embedding with
the input feature map
60         return out
61
62
63 # The Discriminator model
64 class Discriminator(nn.Module):
65     def __init__(self, channels, embed_dim=1024, embed_out_dim=128):
66         super(Discriminator, self).__init__()
67         self.channels = channels
68         self.embed_dim = embed_dim
69         self.embed_out_dim = embed_out_dim
70
71         # Discriminator architecture
72         self.model = nn.Sequential(
73             *self._create_layer(self.channels, 64, 4, 2, 1, normalize=False),
74             *self._create_layer(64, 128, 4, 2, 1),
75             *self._create_layer(128, 256, 4, 2, 1),
76             *self._create_layer(256, 512, 4, 2, 1)
77         )
78         self.text_embedding = Embedding(self.embed_dim, self.embed_out_dim) # Text
embedding module
79         self.output = nn.Sequential(
80             nn.Conv2d(512 + self.embed_out_dim, 1, 4, 1, 0, bias=False), nn.Sigmoid()
81         )
82
83     def _create_layer(self, size_in, size_out, kernel_size=4, stride=2, padding=1,
normalize=True):
84         layers = [nn.Conv2d(size_in, size_out, kernel_size=kernel_size, stride=stride,
padding=padding)]
85         if normalize:
86             layers.append(nn.BatchNorm2d(size_out))
87             layers.append(nn.LeakyReLU(0.2, inplace=True))
88         return layers
89
90     def forward(self, x, text):
91         x_out = self.model(x) # Extract features from the input using the
discriminator architecture
92         out = self.text_embedding(x_out, text) # Apply text embedding and concatenate
with the input features
93         out = self.output(out) # Final discriminator output
94         return out.squeeze(), x_out

```

```

95
96 def weights_init(m):
97     classname = m.__class__.__name__
98     if classname.find('Conv') != -1:
99         m.weight.data.normal_(0.0, 0.02)
100     elif classname.find('BatchNorm') != -1:
101         m.weight.data.normal_(1.0, 0.02)
102         m.bias.data.fill_(0)
103
104 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
105 embed_dim = 1024
106 noise_dim = 100
107 embed_out_dim = 64
108 generator = Generator(
109     channels=3, embed_dim=embed_dim, noise_dim=noise_dim, embed_out_dim=embed_out_dim
110 ).to(device)
111 generator.apply(weights_init)
112
113 discriminator = Discriminator(
114     channels=3, embed_dim=embed_dim, embed_out_dim=embed_out_dim
115 ).to(device)
116 discriminator.apply(weights_init)

```

#### 4. Training

```

1 # setting up Adam optimizer for Generator and Discriminator
2 learning_rate = 0.0002
3 optimizer_G = torch.optim.Adam(
4     generator.parameters(), lr=learning_rate, betas=(0.5, 0.999)
5 )
6 optimizer_D = torch.optim.Adam(
7     discriminator.parameters(), lr=learning_rate, betas=(0.5, 0.999)
8 )
9
10 # loss functions
11 criterion = nn.BCELoss()
12 l2_loss = nn.MSELoss()
13 l1_loss = nn.L1Loss()
14
15 num_epochs = 20
16 real_label = 1.
17 fake_label = 0.
18 l1_coef = 50
19 l2_coef = 100
20
21 D_losses = []
22 G_losses = []
23
24 for epoch in range(num_epochs):
25     epoch_D_loss = []
26     epoch_G_loss = []
27     batch_time = time.time()
28
29     for idx, batch in enumerate(train_loader):
30
31         images = batch['right_images'].to(device)
32         wrong_images = batch['wrong_images'].to(device)
33         embeddings = batch['right_embed'].to(device)
34         batch_size = images.size(0)
35

```



```

36     # ===== #
37     #                               Train the discriminator                               #
38     # ===== #
39
40     # Clear gradients for the discriminator
41     optimizer_D.zero_grad()
42
43     # Generate random noise
44     noise = torch.randn(batch_size, noise_dim, 1, 1, device=device)
45
46     # Generate fake image batch with the generator
47     fake_images = generator(noise, embeddings)
48
49     # Forward pass real batch and calculate loss
50     real_out, real_act = discriminator(images, embeddings)
51     d_loss_real = criterion(real_out, torch.full_like(real_out, real_label, device
=device))
52
53     # Forward pass wrong batch and calculate loss
54     wrong_out, wrong_act = discriminator(wrong_images, embeddings)
55     d_loss_wrong = criterion(wrong_out, torch.full_like(wrong_out, fake_label,
device=device))
56
57     # Forward pass fake batch and calculate loss
58     fake_out, fake_act = discriminator(fake_images.detach(), embeddings)
59     d_loss_fake = criterion(fake_out, torch.full_like(fake_out, fake_label, device
=device))
60
61     # Compute total discriminator loss
62     d_loss = d_loss_real + d_loss_wrong + d_loss_fake
63
64     # Backpropagate the gradients
65     d_loss.backward()
66
67     # Update the discriminator
68     optimizer_D.step()
69
70     # ===== #
71     #                               Train the generator                               #
72     # ===== #
73
74     # Clear gradients for the generator
75     optimizer_G.zero_grad()
76
77     # Generate new random noise
78     noise = torch.randn(batch_size, noise_dim, 1, 1, device=device)
79
80     # Generate new fake images using Generator
81     fake_images = generator(noise, embeddings)
82
83     # Get discriminator output for the new fake images
84     out_fake, act_fake = discriminator(fake_images, embeddings)
85
86     # Get discriminator output for the real images
87     out_real, act_real = discriminator(images, embeddings)
88
89     # Calculate losses
90     g_bce = criterion(out_fake, torch.full_like(out_fake, real_label, device=
device))
91     g_l1 = l1_coef * l1_loss(fake_images, images)

```

```

92     g_l2 = l2_coef * l2_loss(torch.mean(act_fake, 0), torch.mean(act_real, 0).
    detach())
93
94     # Compute total generator loss
95     g_loss = g_bce + g_l1 + g_l2
96
97     # Backpropagate the gradients
98     g_loss.backward()
99
100    # Update the generator
101    optimizer_G.step()
102
103    epoch_D_loss.append(d_loss.item())
104    epoch_G_loss.append(g_loss.item())
105
106    print('Epoch {} [{} / {}] loss_D: {:.4f} loss_G: {:.4f} time: {:.2f}'.format(
107        epoch+1, idx+1, len(train_loader),
108        d_loss.mean().item(),
109        g_loss.mean().item(),
110        time.time() - batch_time)
111    )
112    D_losses.append(sum(epoch_D_loss)/len(epoch_D_loss))
113    G_losses.append(sum(epoch_G_loss)/len(epoch_G_loss))
114
115    # save model
116    model_save_path = './save_model'
117    torch.save(generator.state_dict(), os.path.join(model_save_path, 'generator.pth'))
118    torch.save(discriminator.state_dict(), os.path.join(model_save_path, 'discriminator.pth'
    ))

```

## 5. Prediction

```

1 generator.eval()
2 example = next(iter(test_dataset))
3 sentence = example['txt'] # this flower consists of large yellow-green petals that are
    vertically oriented.
4 embeddings = convert_text_to_feature([str(sentence)])
5 noise = torch.randn(1, noise_dim, 1, 1, device=device)
6 pred = generator(noise, embeddings.detach())
7
8 plt.imshow(pred[0].cpu().detach().permute(1, 2, 0))
9 plt.show()

```

## Phần 4. Câu hỏi trắc nghiệm

**Câu hỏi 1** Mục đích của bài toán Text to Image là gì?

- a) Tạo hình ảnh tương ứng với thông tin mô tả từ văn bản
- b) Phân loại hình ảnh
- c) Phân loại văn bản
- d) Phát hiện cạnh của các đối tượng trong hình ảnh

**Câu hỏi 2** Mô hình nào được xây dựng để giải quyết bài toán Text-to-Image?

- a) GAN-CLS
- b) BERT
- c) ResNet
- d) VGG

**Câu hỏi 3** Khối Generator được sử dụng để làm gì?

- a) Sinh hình ảnh mới từ Noise và Text Embedding
- b) Dự đoán hình ảnh Real/Fake
- c) Sinh hình ảnh mới từ hình ảnh cũ
- d) Sinh văn bản mới từ Noise và Text Embedding

**Câu hỏi 4** Khối Discriminator được sử dụng để làm gì?

- a) Sử dụng kết hợp hình ảnh và Text Embedding dự đoán hình ảnh Fake/Real
- b) Sử dụng Text Embedding dự đoán hình ảnh Fake/Real
- c) Cả 2 đáp án trên đều đúng
- d) Cả 2 đáp án trên đều sai

**Câu hỏi 5** Kiến trúc của Generator và Discriminator không bao gồm layer nào sau đây?

- a) CNN
- b) Batch Normalization
- c) LeakyReLU Activation
- d) RNN

**Câu hỏi 6** Số chiều biểu diễn văn bản được sử dụng trong phần thực nghiệm 2 là bao nhiêu?

- a) 512
- b) 1024
- c) 768
- d) 2048

**Câu hỏi 7** Dựa vào code thực nghiệm trong phần 2, lớp Convolution cuối cùng trong khối Discriminator có tham số ‘in channels’ là bao nhiêu?

- a) 512

- b) 64
- c) 576
- d) 1024

**Câu hỏi 8** Hàm loss nào sau đây không được sử dụng trong phần thực nghiệm?

- a) BCELoss
- b) MSELoss
- c) L1Loss
- d) CTCLoss

**Câu hỏi 9** Mô hình pre-trained language model nào được sử dụng trong thực nghiệm?

- a) BERT
- b) DistilBERT
- c) RoBERTa
- d) T5

**Câu hỏi 10** Dựa vào phần code thực nghiệm trong phần 3, kích thước chiều embedding của text của mô hình pre-trained language model được sử dụng là?

- a) 512
- b) 768
- c) 576
- d) 1024

**- Hết -**