# Deep Learning Lab7 (Let's Play DDPM)

311553060

周睦鈞

## 1. Introduction

這次的 Lab 要實作一個 conditional 的 DDPM , 透過指定顏色和形狀 的 One hot label 當做 Condition 最後生成我們要的結果圖 類似 Lab5 使用的 cVAE ,只是這次的生成模型改成使用 Diffusion model

## 2. Implement detail

2-1. Describe how you implement your model, including your choice of DDPM, UNet architectures, noise schedule, and loss functions

我所使用的 DDPM 是屬於 Pixel domain 的,UNet 架構則是使用 huggingface 所提供的 UNet2DModel ,我基於他們提供的 Model 做了一些修改加上了 condition 的資訊 (e.g. self.model.class\_embedding 預設是 None, 這邊我將它改成 nn.Linear 的形式 來對 one\_hot vector 做 embedding)

## class ConditionedUNet(nn.Module):

**UNet model architecture** 

def \_\_init\_\_(self, args, num\_class = 24, embed\_size = 512): super().\_\_init\_\_() self.model = UNet2DModel(  $sample\_size = 64$ ,  $in\_channels = 3,$ out\_channels = 3, layers\_per\_block = 2, class\_embed\_type = None,  $block\_out\_channels = (128, 256, 256, 512, 512, 1024),$ down\_block\_types=( # a regular ResNet downsampling block "DownBlock2D", "DownBlock2D", "DownBlock2D", "DownBlock2D", "AttnDownBlock2D", # a ResNet downsampling block with spatial self—attention "AttnDownBlock2D",

```
up_block_types=(
        "AttnUpBlock2D",
                              # a ResNet upsampling block with spatial self-attention
        "AttnUpBlock2D",
        "UpBlock2D",
                              # a regular ResNet upsampling block
        "UpBlock2D",
        "UpBlock2D",
        "UpBlock2D",
self.model.class_embedding = nn.Linear(num_class, embed_size)
```

```
def forward(self, x, t, y):
        return self.model(x, t, class_labels = y).sample
Loss funtion
    # Loss functions
   criterion = nn.MSELoss()
    optimizer = optim.Adam(net.parameters(), lr = args.lr)
```

## Noise

noise = torch.randn\_like(img) Dataloader

將助教提供的 json 檔透過 getData function 把對應的 Label data 轉成 one\_hot vector 並且透過 Dataloader 在 training 的時候輸出 image 和 label

filenames = list(data\_json.keys())

labels = np.array(one\_hot\_vector\_list)

return filenames, labels

# 在 testing 的時候只輸出 label

使用 Normal Distribution 來當作我的 noise

class iclevrDataset(Dataset): def \_\_init\_\_(self, mode, root='iclevr', data\_root = 'data', test\_file='test.json'): self.filenames = None self.labels = None self.mode = mode self.root = root

self.filenames, self.labels = getData(mode, data\_root, test\_file=test\_file) self.transform = data\_transform(mode) def \_\_len\_\_(self): return len(self.labels) def \_\_getitem\_\_(self, idx): if self.mode == 'train': image = Image.open(os.path.join(self.root, self.filenames[idx])).convert('RGB') image = self.transform(image) label = torch.Tensor(self.labels[idx]) return image, label else: label = torch.Tensor(self.labels[idx]) return label def getData(mode, root, test\_file='test.json'): label\_map = json.load(open(os.path.join(root, 'objects.json'))) filenames = None if mode == 'train': data\_json = json.load(open(os.path.join(root, 'train.json')))

labels\_list = list(data\_json.values()) one\_hot\_vector\_list = [] for i in range(len(labels\_list)): one\_hot\_vector = np.zeros(24).astype(int) for j in labels\_list[i]: one\_hot\_vector[label\_map[j]] = 1 one\_hot\_vector\_list.append(one\_hot\_vector) labels = np.array(one\_hot\_vector\_list) data\_json = json.load(open(os.path.join(root, test\_file))) labels\_list = data\_json one\_hot\_vector\_list = [] for i in range(len(labels\_list)): one\_hot\_vector = np.zeros(24).astype(int) for j in labels\_list[i]: one\_hot\_vector[label\_map[j]] = 1 one\_hot\_vector\_list.append(one\_hot\_vector)

2-2. Specify the hyperparameters (learning rate, epochs, etc.) def parse\_arg(): parser = argparse.ArgumentParser() parser.add\_argument('--epochs', default=60, help='Number of training epochs') parser.add\_argument('--batch\_size', default=28, help='Size of batches') parser.add\_argument('--lr', default=0.0002, help='Learning rate') parser.add\_argument('--c\_dim', default=4, help='Condition dimension') parser.add\_argument('--test\_only', action='store\_true') parser.add\_argument('--model\_path', default='ckpt', help='Path to save model checkpoint') parser.add\_argument('--log', default='log/origin', help='Path to save log') parser.add\_argument('--timesteps', default=1000, help='Time step for diffusion') parser.add\_argument('--test\_file', default='test.json', help='Test file') parser.add\_argument('--test\_batch\_size', default=32, help='Test batch size') parser.add\_argument('--figure\_file', default='figure/origin', help='Figure file') parser.add\_argument('--resume', default=False, help='Continue for training') parser.add\_argument('--ckpt', default='net.pth', help='Checkpoint for network')

### st\_only /home/mcchou/anaconda3/envs/Lab/lib/python3.10/site-packages/torchvision/models/ \_utils.py:208: UserWarning: The parameter 'pretrained' is deprecated since 0.13

return parser.parse\_args()

3. Results and discussion

/home/mcchou/anaconda3/envs/Lab/lib/python3.10/site-packages/torchvision/models/ \_utils.py:223: UserWarning: Arguments other than a weight enum or `None` for 'we ights' are deprecated since 0.13 and will be removed in 0.15. The current behavi

and will be removed in 0.15, please use 'weights' instead.

3-1. Show your results based on the testing data (including images)

or is equivalent to passing `weights=None`. warnings.warn(msg)

warnings.warn(

warnings.warn(msg)

1000it [01:51, 8.94it/s]

"AttnDownBlock2D",

"DownBlock2D",

"UpBlock2D",

"UpBlock2D",

"UpBlock2D",

"UpBlock2D",

"UpBlock2D",

def forward(self, x, t, y):

bs, c, w, h = x.shape

"AttnUpBlock2D",

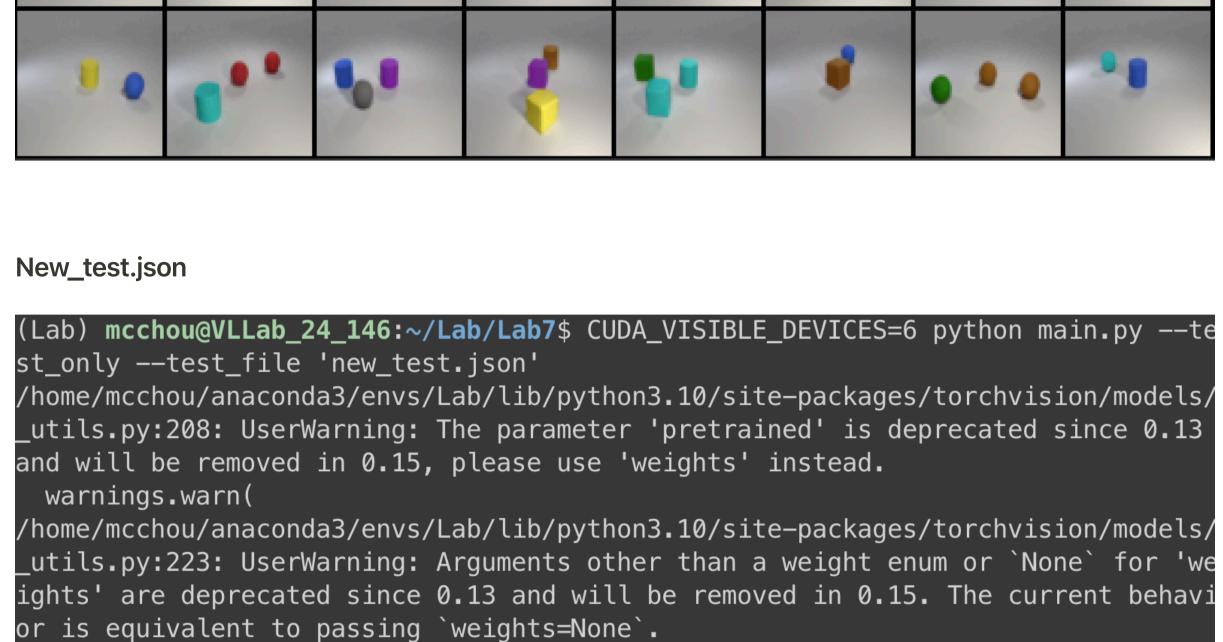
up\_block\_types=(

acc: 0.8095238095238095

Test.json

1000it [01:52, 8.91it/s] acc: 0.8333333333333334

(Lab) mcchou@VLLab\_24\_146:~/Lab/Lab7\$ CUDA\_VISIBLE\_DEVICES=3 python main.py --te



## 3-2. Discuss the results of different model architectures 我總共測試了四種架構,但是都是基於 huggingface 的 UNet2DModel ,差別在於網路的 channel 數目以及 網路 down / up block type 使用 attention 的多寡,我發現網路的 channel 數目如果越以及 UNet 的深度越深效果會越好,且 attention 使用越 多,對於記憶體的要求就會越大 a. One\_hot: 不對 One hot label 做 embedding 直接當作 condition block\_out\_channels = (128, 128, 256, 256, 512, 512), down\_block\_types=( "DownBlock2D", # a regular ResNet downsampling block "DownBlock2D", "DownBlock2D", "DownBlock2D",

# a ResNet downsampling block with spatial self—attention

# a ResNet upsampling block with spatial self-attention

# a regular ResNet upsampling block

```
class_cond = y.view(bs, y.shape[1], 1, 1).expand(bs, y.shape[1], w, h)
        net_input = torch.cat((x, class_cond), 1)
        return self.model(net_input, t).sample
Embedding:
Down / UP type 同上 , 只是使用 embedding
           block_out_channels = (128, 128, 256, 256, 512, 512),
            down_block_types=(
                "DownBlock2D",
                                     # a regular ResNet downsampling block
               "DownBlock2D",
                "DownBlock2D",
               "DownBlock2D",
               "AttnDownBlock2D", # a ResNet downsampling block with spatial self—attention
               "DownBlock2D",
            up_block_types=(
               "UpBlock2D",
               "AttnUpBlock2D",
                                     # a ResNet upsampling block with spatial self-attention
               "UpBlock2D",
               "UpBlock2D",
                                     # a regular ResNet upsampling block
               "UpBlock2D",
               "UpBlock2D",
        self.model.class_embedding = nn.Linear(num_class, embed_size)
```

# a regular ResNet downsampling block

# a regular ResNet upsampling block

# a regular ResNet downsampling block

# a ResNet downsampling block with spatial self—attention

# a ResNet downsampling block with spatial self—attention

# a ResNet upsampling block with spatial self-attention

### 改變 Down / Up block type , 多使用一層 attention block\_out\_channels = (128, 128, 256, 256, 512, 512), down\_block\_types=(

Embedding\_DDDDAA

def forward(self, x, t, y):

"DownBlock2D",

"DownBlock2D",

"DownBlock2D",

"DownBlock2D",

up\_block\_types=(

"AttnDownBlock2D",

"AttnDownBlock2D",

"AttnUpBlock2D",

"AttnUpBlock2D",

"UpBlock2D",

"UpBlock2D",

"UpBlock2D",

"UpBlock2D",

def forward(self, x, t, y):

down\_block\_types=(

up\_block\_types=(

"DownBlock2D",

"DownBlock2D",

"DownBlock2D", "DownBlock2D",

"AttnDownBlock2D",

"AttnDownBlock2D",

return  $self.model(x, t, class_labels = y).sample$ 

return  $self.model(x, t, class_labels = y).sample$ **Big Block** 基本上架構同於 Embedding\_DDDDAA 差別在於 channel 數目加大

self.model.class\_embedding = nn.Linear(num\_class, embed\_size)

block\_out\_channels = (128, 256, 256, 512, 512, 1024),

```
"AttnUpBlock2D",
               "AttnUpBlock2D",
                                   # a ResNet upsampling block with spatial self—attention
               "UpBlock2D",
                                   # a regular ResNet upsampling block
               "UpBlock2D",
               "UpBlock2D",
               "UpBlock2D",
       self.model.class_embedding = nn.Linear(num_class, embed_size)
   def forward(self, x, t, y):
       return self.model(x, t, class_labels = y).sample
Comparison:
從結果圖來比較,可以發現單使用 one_hot vector 當作 condition 效果是最差的,另外對於網路多使用 attention block 且
channel 數目加大會有較好的效果,因此我最後使用的 model 就是採用 Big Block 的架構
                                               Testing acc
                                   - Big Block - Embedding-DDDDAA - One_hot - Embedding
```

4. Experimental results (test.json and new\_test.json) Test.json (Lab) mcchou@VLLab\_24\_146:~/Lab/Lab7\$ CUDA\_VISIBLE\_DEVICES=3 python main.py --te st\_only /home/mcchou/anaconda3/envs/Lab/lib/python3.10/site-packages/torchvision/models/ \_utils.py:208: UserWarning: The parameter 'pretrained' is deprecated since 0.13 and will be removed in 0.15, please use 'weights' instead.

/home/mcchou/anaconda3/envs/Lab/lib/python3.10/site-packages/torchvision/models/

\_utils.py:223: UserWarning: Arguments other than a weight enum or `None` for 'we

ights' are deprecated since 0.13 and will be removed in 0.15. The current behavi

or is equivalent to passing `weights=None`.

st\_only --test\_file 'new\_test.json'

warnings.warn(

warnings.warn(

warnings.warn(msg)

1000it [01:52, 8.91it/s]

acc: 0.8333333333333334

New\_test.json

(Lab) mcchou@VLLab\_24\_146:~/Lab/Lab7\$ CUDA\_VISIBLE\_DEVICES=6 python main.py --te

/home/mcchou/anaconda3/envs/Lab/lib/python3.10/site-packages/torchvision/models/

\_utils.py:208: UserWarning: The parameter 'pretrained' is deprecated since 0.13

and will be removed in 0.15, please use 'weights' instead.

## /home/mcchou/anaconda3/envs/Lab/lib/python3.10/site-packages/torchvision/models/ \_utils.py:223: UserWarning: Arguments other than a weight enum or `None` for 'we ights' are deprecated since 0.13 and will be removed in 0.15. The current behavi or is equivalent to passing `weights=None`. warnings.warn(msg) 1000it [01:51, 8.94it/s] acc: 0.8095238095238095