# Deep Learning
# Lab7: Let's Play DDPM

**Name:** 許子駿

**Student ID:** 311551166

Institute of Computer Science and Engineering

2023 年 5 月 30 日

# Table of Contents

# Chapter 1

# Introduction

This task is to generate synthetic images according to multi-label conditions with conditional Denoising Diffusion Probabilistic Models (DDPM) (See Figures 1.1 and 1.2).
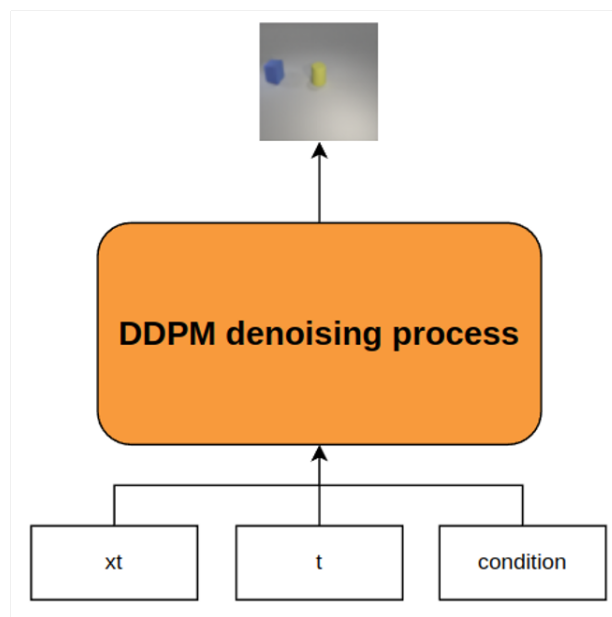


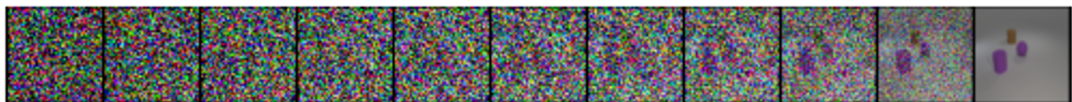Figure 1.1: Illustration of DDPM process to synthesize images.



Figure 1.2: Example of DDPM denoising process.

# Chapter 2

# Experiment setups

Experiment setups are listed in this chapter including 3 parts, DDPM, data loader, and hyperparameters.

## 2.1 DDPM

The diffusion model is Denoising Diffusion Probabilistic Models[1]. The conditioning follows the method of Classifier-Free Diffusion Guidance[2].
The model infuses timestep embeddings $t_e$ and context embeddings $c_e$ with the U-Net activations at layer $a_L$ via

$$a_{L+1} = c_e a_L + t_e$$

During training, $c_e$ is randomly set to 0 with probability 0.1, and this makes model learn to do unconditional generation, i.e., $\psi(z_t)$ for noise $z_t$ at timestep $t$, and also conditional generation, i.e., $\psi(z_t, c)$ for context $c$.
Additionally, a weight $w \geq 0$ is necessary, to guide the model to generate samples with the following equation

$$\hat{\epsilon}_t = (1 + w)\psi(z_t, c) - w\psi(z_t)$$

Increasing $w$ produces images that are more typical but less diverse.
The DDPM details are listed in this section including 3 parts, UNet, noise schedule, and loss function.

---

[1]`https://arxiv.org/abs/2006.11239`

[2]`https://arxiv.org/abs/2207.12598`

### 2.1.1 UNet

The model uses asymmetric UNet, composed of 2 down-sampling layers and 3 up-sampling layers, and the final output is output of third up-sampling layer with input image concatenated. Timestep and context are embedded with 2 linear layers with one `GELU` activation layer in between. (See Listings 2.1, A.1, and A.2 for details).

```python
class ContextUnet(nn.Module):
    def forward(self,
            x: torch.FloatTensor, c: torch.FloatTensor, t: torch.FloatTensor,
            context_mask: torch.FloatTensor) -> torch.FloatTensor:
        # x: image, c: context, t: timestep.
        x = self.init_conv(x)
        down1 = self.down1(x)
        down2 = self.down2(down1)
        hiddenvec = self.to_vec(down2)

        context_mask = (-1 * (1 - context_mask)) # flip 0 <-> 1
        c = c * context_mask

        # embed context, time step
        cemb1 = self.contextembed1(c).view(-1, self.n_feat * 2, 1, 1)
        temb1 = self.timeembed1(t).view(-1, self.n_feat * 2, 1, 1)
        cemb2 = self.contextembed2(c).view(-1, self.n_feat, 1, 1)
        temb2 = self.timeembed2(t).view(-1, self.n_feat, 1, 1)

        up1 = self.up0(hiddenvec)
        up2 = self.up1(cemb1 * up1 + temb1, down2)
        up3 = self.up2(cemb2 * up2 + temb2, down1)
        out = self.out(torch.cat((up3, x), 1))
        return out
```

Listing 2.1: Python code of **ContextUnet** of DDPM (some code is omitted).

### 2.1.2 Noise schedule

Noise schedule follows equation

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$$

Listing 2.2 shows the noise schedule function calculates each part of above equation. During sampling, the pre-calculated values are used to sample images.

4

```python
def ddpm_schedules(beta1, beta2, T):
    assert beta1 < beta2 < 1.0, 'beta1 and beta2 must be in (0, 1)'

    beta_t = (beta2 - beta1) * torch.arange(0, T + 1, dtype=torch.float32) / T
        + beta1
    sqrt_beta_t = torch.sqrt(beta_t)
    alpha_t = 1 - beta_t
    log_alpha_t = torch.log(alpha_t)
    alphabar_t = torch.cumsum(log_alpha_t, dim=0).exp()

    sqrtab = torch.sqrt(alphabar_t)
    oneover_sqrta = 1 / torch.sqrt(alpha_t)

    sqrtmab = torch.sqrt(1 - alphabar_t)
    mab_over_sqrtmab_inv = (1 - alpha_t) / sqrtmab

    return {
        'alpha_t': alpha_t,  # \alpha_t
        'oneover_sqrta': oneover_sqrta,  # 1/\sqrt{\alpha_t}
        'sqrt_beta_t': sqrt_beta_t,  # \sqrt{\beta_t}
        'alphabar_t': alphabar_t,  # \bar{\alpha_t}
        'sqrtab': sqrtab,  # \sqrt{\bar{\alpha_t}}
        'sqrtmab': sqrtmab,  # \sqrt{1-\bar{\alpha_t}}
        'mab_over_sqrtmab': mab_over_sqrtmab_inv,  # (1-\alpha_t)/\sqrt{1-\bar
    {\alpha_t}}
    }
```

Listing 2.2: Python code of **ddpm_schedules** of DDPM.

### 2.1.3   Loss function

MSE loss is used to calculate loss between input image and generated image.

## 2.2   Data loader

Images are pre-processed with `Resize` to $64 \times 64$ and then `Normalize` with mean $(0.5, 0.5, 0.5)$ and standard deviation $(0.5, 0.5, 0.5)$. And labels are converted to one-hot vectors with total 24 classes (See Listings 2.3 and A.3 for details).

```python
class _ClevrTrainingDataset(Dataset):
    def __init__(self, dir_dataset: str, dir_root='iclevr'):
        self.transforms = transforms.Compose([
            transforms.Resize((64, 64)),
            transforms.ToTensor(),
```

```python
 6            transforms.Normalize((0.5, 0.5, 0.5),(0.5, 0.5, 0.5))])
 7
 8        def __getitem__(self, index):
 9            img = Image.open(
10                Path(self.dir_dataset, self.dir_root, self.img_names[index])).
   convert('RGB')
11            img = self.transforms(img)
12            cond = self.int2one_hot(self.img_conds[index])
13
14            return img, cond
15
16        def int2one_hot(self, int_list):
17            one_hot = torch.zeros(self.num_classes)
18            for i in int_list:
19                one_hot[i] = 1.
20
21            return one_hot
```

Listing 2.3: Python code of data pre-process and loader.

## 2.3  Hyperparameters

Hyperparameters of this experiment are listed below:

- Embedding size: 128.

- Betas of noise schedule: $(10^{-4}, 0.02)$.

- Dropout of context: 0.1.

- Sampling timestep: 400.

- Batch size: 128.

- Number of training epochs: 300.

- Optimizer: Adam.

- Learning rate: $10^{-4}$.

- Seed: 42.

# Chapter 3

# Experimental results

Testing labels are evaluated at once, i.e., generate 32 images with total size $[32, 3, 64, 64]$, and evaluate by pre-trained ResNet. Two different guidance $w$ 0.5 and 2.0 are used.

Scores are listed in Table 3.1.

| Labels $w$ | test.json | new_test.json |
|---|---|---|
| 0.5 | 0.85 (Figure 3.1) | 0.83 (Figure 3.3) |
| 2.0 | 0.82 (Figure 3.2) | 0.81 (Figure 3.4) |

Table 3.1: Scores with guidance $w$ 0.5 and 2.0.



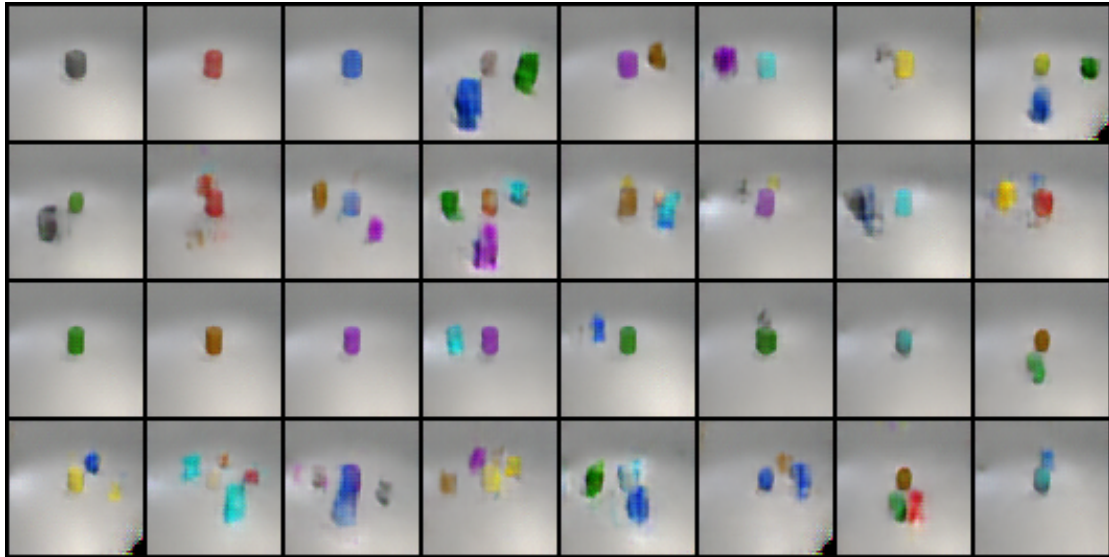Figure 3.1: Results of test.json with guidance $w$ 0.5.

Figure 3.2: Results of test.json with guidance $w$ 2.0.



Figure 3.3: Results of new_test.json with guidance $w$ 0.5.

Figure 3.4: Results of new_test.json with guidance $w$ 2.0.

# Chapter 4

# Discussion

As model architecture listed in Chapter 2, it's sufficient for model to learn and generate objects with specific multi-labels.
If activation of embedding layers, i.e., `GELU` activation layer between two fully connected layers, is removed, the score of generated images with testing labels would drop to about **0.7**. This might be caused by the non-linearity is removed.
Also, mask is predicted during image generation, and it's more intuitive and simple to optimize through MSE loss between input images and generated images.

# Appendices

# Appendix A

# Code

```python
class ContextUnet(nn.Module):
    def __init__(self, in_channels: int, n_feat=128, n_classes=24):
        super().__init__()

        self.in_channels = in_channels
        self.n_feat = n_feat
        self.n_classes = n_classes

        self.init_conv = ResidualConvBlock(in_channels, n_feat, is_res=True)

        self.down1 = UnetDown(n_feat, n_feat)
        self.down2 = UnetDown(n_feat, 2 * n_feat)

        self.to_vec = nn.Sequential(
            nn.AvgPool2d(7),
            nn.GELU())

        self.timeembed1 = EmbedFC(1, 2 * n_feat)
        self.timeembed2 = EmbedFC(1, 1 * n_feat)
        self.contextembed1 = EmbedFC(n_classes, 2 * n_feat)
        self.contextembed2 = EmbedFC(n_classes, 1 * n_feat)

        self.up0 = nn.Sequential(
            nn.ConvTranspose2d(2 * n_feat, 2 * n_feat, 8, 8),
            nn.GroupNorm(8, 2 * n_feat),
            nn.ReLU(),
        )

        self.up1 = UnetUp(4 * n_feat, n_feat)
        self.up2 = UnetUp(2 * n_feat, n_feat)
        self.out = nn.Sequential(
```

```python
32              nn.Conv2d(2 * n_feat, n_feat, 3, 1, 1),
33              nn.GroupNorm(8, n_feat),
34              nn.ReLU(),
35              nn.Conv2d(n_feat, self.in_channels, 3, 1, 1),
36          )
37
38      def forward(self,
39              x: torch.FloatTensor, c: torch.FloatTensor, t: torch.FloatTensor,
40              context_mask: torch.FloatTensor) -> torch.FloatTensor:
41          # x is (noisy) image, c is context label, t is timestep,
42          # context_mask says which samples to block the context on
43
44          x = self.init_conv(x)
45          down1 = self.down1(x)
46          down2 = self.down2(down1)
47          hiddenvec = self.to_vec(down2)
48
49          context_mask = (-1 * (1 - context_mask)) # flip 0 <-> 1
50          c = c * context_mask
51
52          # embed context, time step
53          cemb1 = self.contextembed1(c).view(-1, self.n_feat * 2, 1, 1)
54          temb1 = self.timeembed1(t).view(-1, self.n_feat * 2, 1, 1)
55          cemb2 = self.contextembed2(c).view(-1, self.n_feat, 1, 1)
56          temb2 = self.timeembed2(t).view(-1, self.n_feat, 1, 1)
57
58          up1 = self.up0(hiddenvec)
59          up2 = self.up1(cemb1 * up1 + temb1, down2)
60          up3 = self.up2(cemb2 * up2 + temb2, down1)
61          out = self.out(torch.cat((up3, x), 1))
62          return out
```

Listing A.1: Full Python code of **ContextUnet** of DDPM.

```python
1  class ResidualConvBlock(nn.Module):
2      def __init__(self,
3              in_channels: int, out_channels: int, is_res: bool=False) -> None:
4          super().__init__()
5
6          self.same_channels = (in_channels==out_channels)
7          self.is_res = is_res
8          self.conv1 = nn.Sequential(
9              nn.Conv2d(in_channels, out_channels, 3, 1, 1),
10             nn.BatchNorm2d(out_channels),
11             nn.GELU(),
12         )
```

```python
13          self.conv2 = nn.Sequential(
14              nn.Conv2d(out_channels, out_channels, 3, 1, 1),
15              nn.BatchNorm2d(out_channels),
16              nn.GELU(),
17          )
18
19      def forward(self, x: torch.FloatTensor) -> torch.FloatTensor:
20          if self.is_res:
21              x1 = self.conv1(x)
22              x2 = self.conv2(x1)
23
24              if self.same_channels:
25                  out = x + x2
26              else:
27                  out = x1 + x2
28
29              return out / 1.414
30          else:
31              x1 = self.conv1(x)
32              x2 = self.conv2(x1)
33
34              return x2
35
36  class UnetDown(nn.Module):
37      def __init__(self, in_channels: int, out_channels: int):
38          super().__init__()
39
40          self.layers = nn.Sequential(
41              ResidualConvBlock(in_channels, out_channels),
42              nn.MaxPool2d(2)
43          )
44
45      def forward(self, x: torch.FloatTensor) -> torch.FloatTensor:
46          return self.layers(x)
47
48  class UnetUp(nn.Module):
49      def __init__(self, in_channels: int, out_channels: int):
50          super().__init__()
51
52          self.layers = nn.Sequential(
53              nn.ConvTranspose2d(in_channels, out_channels, 2, 2),
54              ResidualConvBlock(out_channels, out_channels),
55              ResidualConvBlock(out_channels, out_channels),
56          )
57
```

```python
58    def forward(self, x: torch.FloatTensor, skip: torch.FloatTensor) -> torch.
      FloatTensor:
59        x = torch.cat((x, skip), 1)
60        x = self.layers(x)
61
62        return x
63
64 class EmbedFC(nn.Module):
65    def __init__(self, input_dim: int, emb_dim: int):
66        super().__init__()
67
68        self.input_dim = input_dim
69        self.layers = nn.Sequential(
70            nn.Linear(input_dim, emb_dim),
71            nn.GELU(),
72            nn.Linear(emb_dim, emb_dim))
73
74    def forward(self, x: torch.FloatTensor) -> torch.FloatTensor:
75        x = x.view(-1, self.input_dim)
76
77        return self.layers(x)
```

Listing A.2: Python code of components of **ContextUnet**.

```python
1 def get_train_loader(
2        batch_size: int, dir_dataset: str,
3        num_workers=8, shuffle=True) -> tuple[DataLoader, int]:
4    dataset = _ClevrTrainingDataset(dir_dataset)
5
6    return DataLoader(
7        dataset, batch_size=batch_size, num_workers=num_workers, shuffle=
      shuffle), \
8    dataset.num_classes, dataset.data_shape
9
10 def get_test_labels(dir_dataset: str) -> list:
11    with open(Path(dir_dataset, 'objects.json'), 'r') as file:
12        classes = json.load(file)
13
14    with open(Path(dir_dataset, 'test.json'), 'r') as file:
15        conds_list = json.load(file)
16
17    labels = torch.zeros(len(conds_list), len(classes))
18    for i, conds in enumerate(conds_list):
19        for cond in conds:
20            labels[i, int(classes[cond])] = 1.
21
```

```python
22        return labels
23
24  class _ClevrTrainingDataset(Dataset):
25      def __init__(self, dir_dataset: str, dir_root='iclevr'):
26          self.dir_dataset = dir_dataset
27          self.dir_root = dir_root
28
29          with open(Path(dir_dataset, 'objects.json'), 'r') as file:
30              self.classes = json.load(file)
31          self.num_classes = len(self.classes)
32
33          self.transforms = transforms.Compose([
34              transforms.Resize((64, 64)),
35              transforms.ToTensor(),
36              transforms.Normalize((0.5, 0.5, 0.5),(0.5, 0.5, 0.5))])
37
38          self.max_objects = 0
39          self.img_names = []
40          self.img_conds = []
41
42          with open(Path(dir_dataset, 'train.json'), 'r') as file:
43              dict = json.load(file)
44              for name, conds in dict.items():
45                  self.img_names.append(name)
46                  self.max_objects = max(self.max_objects, len(conds))
47                  self.img_conds.append([self.classes[cond] for cond in conds])
48
49          img, _ = self.__getitem__(0)
50          self.data_shape = img.shape
51
52      def __len__(self):
53          return len(self.img_names)
54
55      def __getitem__(self, index):
56          img = Image.open(
57              Path(self.dir_dataset, self.dir_root, self.img_names[index])).
58      convert('RGB')
58          img = self.transforms(img)
59          cond = self.int2one_hot(self.img_conds[index])
60
61          return img, cond
62
63      def int2one_hot(self,int_list):
64          one_hot = torch.zeros(self.num_classes)
65          for i in int_list:
```

```
66            one_hot[i] = 1.
67
68        return one_hot
```

Listing A.3: Python code of data pre-process and loader.