

Deep Learning Homework 6 Report

Cheng-Liang Chi

May 6, 2025

Contents

1	Introduction	2
2	Implementation Details	2
2.1	Model Architecture	2
2.2	Training Process	3
2.3	Evaluation Process	6
2.4	Hyperparameters	9
2.5	Dataset	9
3	Discussion	11
3.1	Default Configuration	11
3.2	Same model with different time steps	11
3.2.1	500 Time Steps	11
3.2.2	100 Time Steps	12
3.3	Discussion of Results	13
3.4	Larger Model	13

1 Introduction

In this lab we build a conditional Denoising Diffusion Probabilistic Model that turns a set of color-shape phrases (e.g., “red sphere”, “yellow cube”) into a 64×64 scene containing exactly the requested objects. We will use a pretrained ResNet-18 model to evaluate the accuracy of the generated images.

We will use the CLEVR dataset, which contains 3D shapes and their corresponding color-shape phrases. The dataset is available [here](#).

2 Implementation Details

For this homework, we implemented a conditional Denoising Diffusion Probabilistic Model (DDPM) using [PyTorch](#). The model is designed to generate images based on color-shape phrases, such as "red sphere" or "yellow cube". The implementation consists of several key components, including the model architecture, training process, and evaluation metrics.

I modularized the code into several files to improve readability and maintainability. The main components of the implementation are as follows:

- `model.py`: This file contains the implementation of the conditional DDPM model. Most of the model architecture is using Hugging Face Diffusers library [1], which provides a high-level interface for building and training diffusion models.
- `dataset.py`: This file defines the dataset class for loading and preprocessing the training data, including color-shape phrases and corresponding images.
- `train.py`: This file contains the training loop, including the optimization process and logging of training metrics.
- `test.py`: This file implements the evaluation process, including generating images from the trained model and evaluating the generated image using the pretrained ResNet-18 classifier provided by the TA.

I will provide a brief overview of each component in the following sections.

2.1 Model Architecture

The model architecture is based on a U-Net structure, which is commonly used in image generation tasks. The U-Net consists of an encoder-decoder structure with skip connections, allowing the model to capture both local and global features in the image.

Since the color-shape phrase (label) is a limited set of words, I used one-hot encoding to represent the phrases. The one-hot encoded vector is then using a linear layer to project it into a higher-dimensional space, which is then provided as an additional input to the U-Net model.

When dealing with higher-resolution inputs it is reasonable to use more down and up-blocks, and keep the attention layers only at the lowest resolution (bottom) layers to reduce memory usage. However, in this homework, I used a smaller model with only 6 down and up-blocks to keep the training time reasonable and also get good results.

The implementation of the U-Net model is based on the Hugging Face Diffusers library, which provides a high-level interface for building and training diffusion models. The implementation of the U-Net model is as follows:

```

5 class ConditionalDDPM(nn.Module):
6     def __init__(self, num_classes=24, dim=512):
7         super().__init__()
8         channel = dim // 4
9         self.UNet = UNet2DModel(
10             sample_size=64,
11             in_channels=3,
12             out_channels=3,
13             layers_per_block=2,
14             block_out_channels=[
15                 channel,
16                 channel,
17                 channel * 2,
18                 channel * 2,
19             ], # type: ignore
20             down_block_types=[
21                 "DownBlock2D",
22                 "DownBlock2D",
23                 "DownBlock2D",
24                 "AttnDownBlock2D",
25             ], # type: ignore
26             up_block_types=[
27                 "AttnUpBlock2D",
28                 "UpBlock2D",
29                 "UpBlock2D",
30                 "UpBlock2D",
31             ], # type: ignore
32             class_embed_type="identity",
33         )
34         self.label_embedding = nn.Linear(num_classes, dim)
35
36     def forward(self, x, t, label):
37         embedding_label = self.label_embedding(label)
38         return self.UNet(x, t, embedding_label).sample

```

2.2 Training Process

The U-Net model is trained using a denoising diffusion process, where the model learns to predict the noise added to the image at each time step. The training process involves optimizing the model parameters using a combination of mean squared error (MSE) loss.

The training loop consists of the following steps:

- Load a batch of images and corresponding color-shape phrases from the dataset.
- For each image, add Gaussian noise to the image at a random time step.
- Pass the noisy image and the one-hot encoded phrase through the U-Net model to predict the noise.
- Compute the loss between the predicted noise and the actual noise added to the image.

- Backpropagate the loss and update the model parameters using an optimizer (AdamW).

The training process is implemented in the `train.py` file, which contains the training loop and logging of training metrics. The implementation of the training loop is as follows:

```

16 class Trainer:
17     def __init__(self, args):
18         self.train_loader = DataLoader(
19             IclevrDataset(args.dataset, "train"),
20             batch_size=args.batch_size,
21             shuffle=True,
22             num_workers=args.num_workers,
23         )
24         self.val_loader = DataLoader(
25             IclevrDataset(args.dataset, "test"),
26             batch_size=1,
27             shuffle=False,
28             num_workers=args.num_workers,
29         )
30
31         self.criteria = nn.MSELoss()
32         self.model = ConditionalDDPM().to(args.device)
33         self.noise_scheduler =
34             ↳ DDPM Scheduler(num_train_timesteps=args.time_steps)
35         self.time_steps = args.time_steps
36
37         self.optimizer = torch.optim.Adam(self.model.parameters(),
38             ↳ lr=args.lr)
39         self.epochs = args.epochs
40         self.save_dir = args.save_dir
41         self.save_ckpt_period = args.save_ckpt_period
42         self.save_img_period = args.save_img_period
43         self.device = args.device
44         self.writer = SummaryWriter()
45
46         self.best_loss: float = float("inf")
47         self.epoch = 0
48
49         os.makedirs(self.save_dir, exist_ok=True)
50         os.makedirs(os.path.join(self.save_dir, "imgs"), exist_ok=True)
51
52     def save_checkpoint(self, path):
53         torch.save(
54             {
55                 "epoch": self.epoch,
56                 "best_loss": self.best_loss,
57                 "model": self.model.state_dict(),
58                 "optimizer": self.optimizer.state_dict(),
59                 "time_steps": self.time_steps,
60             },
61             path,

```

```

60     )
61     tqdm.write(f"Model saved at {path}")
62
63     def train(self):
64         for epoch in trange(self.epochs, desc="Training",
65                               ↪ dynamic_ncols=True):
66             self.epoch = epoch
67             train_loss = self.train_one_epoch()
68             self.writer.add_scalar("Loss/train", train_loss, epoch)
69             tqdm.write(f"Epoch {epoch}, Train Loss: {train_loss}")
70
71             if (epoch + 1) % self.save_ckpt_period == 0:
72                 self.save_checkpoint(
73                     os.path.join(self.save_dir, f"model_epoch_{epoch}.pth")
74                 )
75
76             if (epoch + 1) % self.save_img_period == 0:
77                 self.save_images()
78
79     def train_one_epoch(self):
80         self.model.train()
81         train_loss = []
82         for i, (img, label) in enumerate(
83             tqdm(
84                 self.train_loader,
85                 desc=f"Epoch: {self.epoch}",
86                 dynamic_ncols=True,
87             )
88         ):
89             batch_size = img.shape[0]
90             img, label = img.to(self.device), label.to(self.device)
91             noise = torch.randn_like(img)
92
93             timesteps = (
94                 torch.randint(0, self.time_steps,
95                               ↪ (batch_size,)).long().to(self.device)
96             )
97             noisy_x = self.noise_scheduler.add_noise(img, noise, timesteps)
98             ↪ # type: ignore
99             output = self.model(noisy_x, timesteps, label)
100
101             loss = self.criterias(output, noise)
102
103             self.optimizer.zero_grad()
104             loss.backward()
105             self.optimizer.step()
106
107             train_loss.append(loss.item())
108
109             self.writer.add_scalar(

```

```

107         "Loss/train-step", loss.item(), self.epoch *
108         ↪ len(self.train_loader) + i
109     )
110     os.makedirs(
111         os.path.join(self.save_dir, "imgs", f"ep{self.epoch}"),
112         ↪ exist_ok=True
113     )
114     self.writer.add_scalar("Loss/epoch", np.mean(train_loss),
115         ↪ self.epoch)
116     return np.mean(train_loss)
117
118 def save_images(self):
119     self.model.eval()
120
121     for idx, (y, label) in enumerate(
122         tqdm(self.val_loader, desc=f"Epoch: {self.epoch}",
123             ↪ dynamic_ncols=True)
124     ):
125         y = y.to(self.device)
126         x = torch.randn(1, 3, 64, 64).to(self.device)
127         denoising_result = []
128         for i, t in enumerate(self.noise_scheduler.timesteps):
129             with torch.no_grad():
130                 residual = self.model(x, t, y)
131
132                 x = self.noise_scheduler.step(residual, t, x).prev_sample
133                 ↪ # type: ignore
134                 if i % (len(self.noise_scheduler.timesteps) // 10) == 0:
135                     denoising_result.append(x.squeeze(0))
136
137             denoising_result.append(x.squeeze(0))
138             denoising_result = torch.stack(denoising_result)
139             row_image = make_grid(
140                 (denoising_result + 1) / 2, nrow=denoising_result.shape[0],
141                 ↪ pad_value=0
142             )
143             save_image(

```

2.3 Evaluation Process

The evaluation process involves generating images from the trained model and evaluating the generated images using a pretrained ResNet-18 classifier. The evaluation process consists of the following steps:

- Load the trained model and the pretrained ResNet-18 classifier.
- Generate images from the trained model using a set of color-shape phrases.

- Evaluate the generated images using the pretrained ResNet-18 classifier to compute the accuracy of the generated images.
- Save the generated images and their corresponding accuracy scores.

The evaluation process is implemented in the `test.py` file, which contains the evaluation loop and logging of evaluation metrics. The implementation of the evaluation loop is as follows:

```

15 class Tester:
16     def __init__(self, args):
17         self.test_loader = DataLoader(
18             IclevrDataset(args.dataset, "test"),
19             batch_size=1,
20             shuffle=False,
21             num_workers=args.num_workers,
22         )
23         self.new_test_loader = DataLoader(
24             IclevrDataset(args.dataset, "new_test"),
25             batch_size=1,
26             shuffle=False,
27             num_workers=args.num_workers,
28         )
29         self.manual_test_loader = DataLoader(
30             IclevrDataset(args.dataset, "manual_test"),
31             batch_size=1,
32             shuffle=False,
33             num_workers=args.num_workers,
34         )
35
36         self.device = args.device
37         self.model = ConditionalDDPM().to(self.device)
38         self.model.load_state_dict(torch.load(args.ckpt)["model"])
39         self.model.eval()
40
41         self.eval_model = evaluation_model()
42
43         self.time_steps = torch.load(args.ckpt)["time_steps"]
44         self.noise_scheduler =
45             ↳ DDPMScheduler(num_train_timesteps=self.time_steps)
46
47         self.save_dir = args.save_dir
48         os.makedirs(self.save_dir, exist_ok=True)
49         os.makedirs(os.path.join(self.save_dir, "test"), exist_ok=True)
50         os.makedirs(os.path.join(self.save_dir, "new_test"), exist_ok=True)
51         os.makedirs(os.path.join(self.save_dir, "manual_test"),
52             ↳ exist_ok=True)
53
54     def test(self):
55         test_acc = self.inference(
56             self.test_loader,

```

```

55         os.path.join(self.save_dir, "test"),
56     )
57     new_test_acc = self.inference(
58         self.new_test_loader,
59         os.path.join(self.save_dir, "new_test"),
60     )
61     manual_test_acc = self.inference(
62         self.manual_test_loader,
63         os.path.join(self.save_dir, "manual_test"),
64     )
65     print(f"Test accuracy: {test_acc:.4f}")
66     print(f"New test accuracy: {new_test_acc:.4f}")
67     print(f"Manual test accuracy: {manual_test_acc:.4f}")
68
69     def inference(self, loader, save_dir):
70         all_results = []
71         accs = []
72         for idx, (y, label) in enumerate(pbar := tqdm(loader)):
73             y = y.to(self.device)
74             x = torch.randn(1, 3, 64, 64).to(self.device)
75             denoising_results = []
76             for i, t in enumerate(self.noise_scheduler.timesteps):
77                 with torch.no_grad():
78                     residual = self.model(x, t, y)
79
80                     x = self.noise_scheduler.step(residual, t, x).prev_sample
81                     ↪ # type: ignore
82                     if i % (self.time_steps // 10) == 0:
83                         denoising_results.append(x.squeeze(0))
84             acc = self.eval_model.eval(x, y)
85             tqdm.write(f"image: {idx}, label: {label}, accuracy:
86                 ↪ {acc:.4f}")
87             accs.append(acc)
88             pbar.set_postfix_str(f"image: {idx}, accuracy:
89                 ↪ {np.mean(accs):.4f}")
90             pbar.refresh()
91
92             denoising_results.append(x.squeeze(0))
93             denoising_results = torch.stack(denoising_results)
94             row_image = make_grid(
95                 (denoising_results + 1) / 2,
96                 nrow=denoising_results.shape[0],
97                 pad_value=0,
98             )
99             save_image(
100                 row_image,
101                 os.path.join(save_dir, f"{idx}.png"),
102             )
103             all_results.append(denoising_results[-1])
104         all_results = torch.stack(all_results)

```



```

102     save_image(
103         make_grid(
104             (all_results + 1) / 2,
105             nrow=8,
106             pad_value=0,
107         ),
108         os.path.join(save_dir, "all_results.png"),
109     )
110     return np.mean(accs)

```

2.4 Hyperparameters

The hyperparameters used in the training process are as follows:

- batch_size: 32
- learning_rate: 1e-4
- epochs: 10
- num_steps: 1000

2.5 Dataset

The dataset used for training the model is the CLEVR dataset, which contains 3D shapes and their corresponding color-shape phrases. I used PyTorch's `DataLoader` class to load the dataset in batches and apply necessary transformations, such as resizing and normalization. The dataset class is implemented in the `dataset.py` file, which defines the dataset class for loading and preprocessing the training data. The implementation of the dataset class is as follows:

```

9  def transform_img(img):
10      transform = transforms.Compose(
11          [
12              transforms.Resize((64, 64)),
13              transforms.ToTensor(),
14              transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5)),
15          ]
16      )
17      return transform(img)
18
19
20  class IcleivrDataset(Dataset):
21      def __init__(self, root: str, mode="train"):
22          super().__init__()
23          assert mode in [
24              "train",
25              "test",
26              "new_test",
27              "manual_test",

```

```

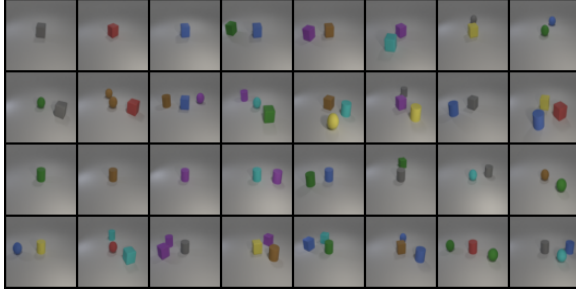
28     ], f"mode {mode} not in [train, test, new_test, manual_test]"
29     assert os.path.exists(root), f"{root} does not exist"
30
31     self.root = root
32     self.mode = mode
33
34     json_path = os.path.join(root, f"{mode}.json")
35     with open(json_path, "r") as json_file:
36         self.json_data = json.load(json_file)
37         match mode:
38             case "train":
39                 self.img_paths = list(self.json_data.keys())
40                 self.labels = list(self.json_data.values())
41             case _:
42                 self.labels = list(self.json_data)
43
44     with open(os.path.join(root, "objects.json"), "r") as json_file:
45         self.objects_dict = json.load(json_file)
46     self.labels_one_hot = torch.zeros(len(self.labels),
47     ↪ len(self.objects_dict))
48
49     for i, label in enumerate(self.labels):
50         self.labels_one_hot[i][[self.objects_dict[j] for j in label]] =
51         ↪ 1
52
53     def __len__(self):
54         return len(self.labels)
55
56     def __getitem__(self, index) -> tuple[torch.Tensor, torch.Tensor |
57     ↪ list]:
58         match self.mode:
59             case "train":
60                 img_path = os.path.join(self.root, "iclevr",
61                 ↪ self.img_paths[index])
62                 img = Image.open(img_path).convert("RGB")
63                 img = transform_img(img)
64                 label_one_hot = self.labels_one_hot[index]
65                 return img, label_one_hot
66             case _:
67                 label_one_hot = self.labels_one_hot[index]
68                 semantic_label = self.labels[index]
69                 return label_one_hot, semantic_label

```

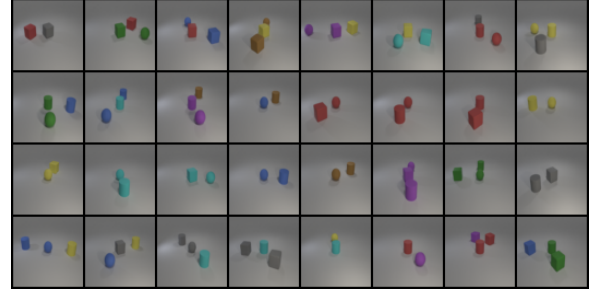
3 Discussion

3.1 Default Configuration

The results of the default configuration are shown in Figure 1. The default configuration is using four down and up-blocks, and keep the attention layers only at the lowest resolution (bottom) layers to reduce memory usage. And the model is trained for 300 epochs with a batch size of 32 and time steps of 1000.



(a) Default Configuration Results on Test Set with accuracy of 0.9010



(b) Default Configuration Results on New Test Set with accuracy of 0.9271

Figure 1: Default Configuration Results and Accuracy

And the results of the manual test with “red sphere”, “cyan cylinder”, and “cyan cube” are shown in Figure 2.

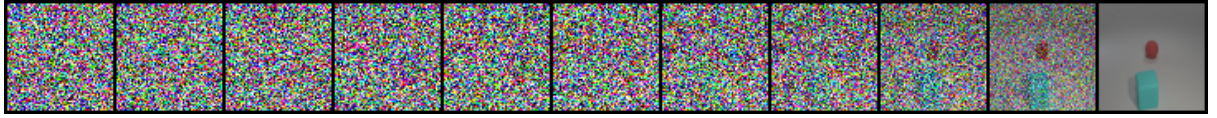


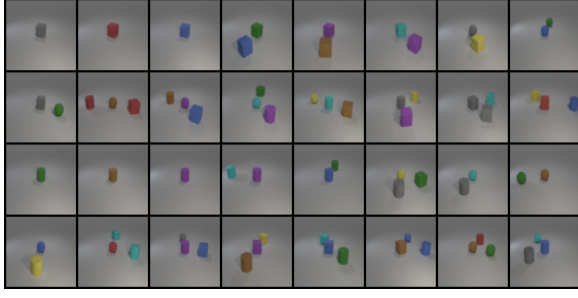
Figure 2: Default Configuration Results on Manual Test Set with accuracy of 1.0000

And the loss curve of the default configuration is shown in Figure 7.

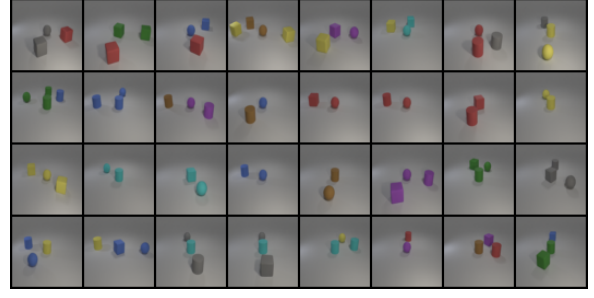
3.2 Same model with different time steps

3.2.1 500 Time Steps

The results of the same model with different time steps are shown in Figure 3. The model is trained for 300 epochs with a batch size of 32 and time steps of 500.



(a) 500 Time Steps Results on Test Set with accuracy of 0.9740



(b) 500 Time Steps Results on New Test Set with accuracy of 0.8906

Figure 3: 500 Time Steps Results and Accuracy

And the results of the manual test with “red sphere”, “cyan cylinder”, and “cyan cube” are shown in Figure 4.

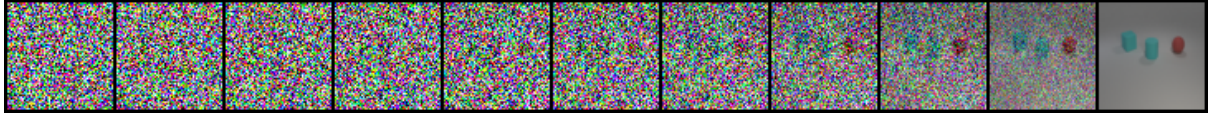
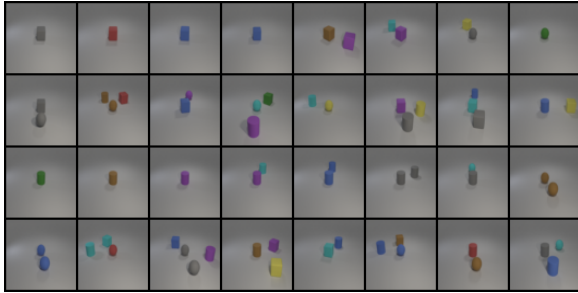


Figure 4: 500 Time Steps Results on Manual Test Set with accuracy of 1.0000

And the loss curve of the 500 time steps is shown in Figure 7.

3.2.2 100 Time Steps

The results of the same model with different time steps are shown in Figure 5. The model is trained for 300 epochs with a batch size of 32 and time steps of 100.



(a) 100 Time Steps Results on Test Set with accuracy of 0.8021



(b) 100 Time Steps Results on New Test Set with accuracy of 0.7865

Figure 5: 100 Time Steps Results and Accuracy

And the results of the manual test with “red sphere”, “cyan cylinder”, and “cyan cube” are shown in Figure 6.

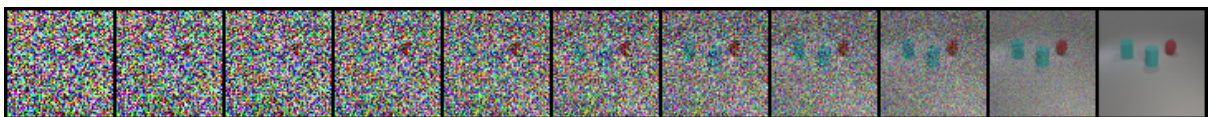


Figure 6: 100 Time Steps Results on Manual Test Set with accuracy of 1.0000

And the loss curve of the 100 time steps is shown in Figure 7.

3.3 Discussion of Results

The results of the default configuration and the same model with different time steps are shown in Figure 7. The model with 1000 time steps has the lowest loss, and the model with 500 or 100 time steps has a higher loss. We can see that the more time steps we have, the lower the loss is.

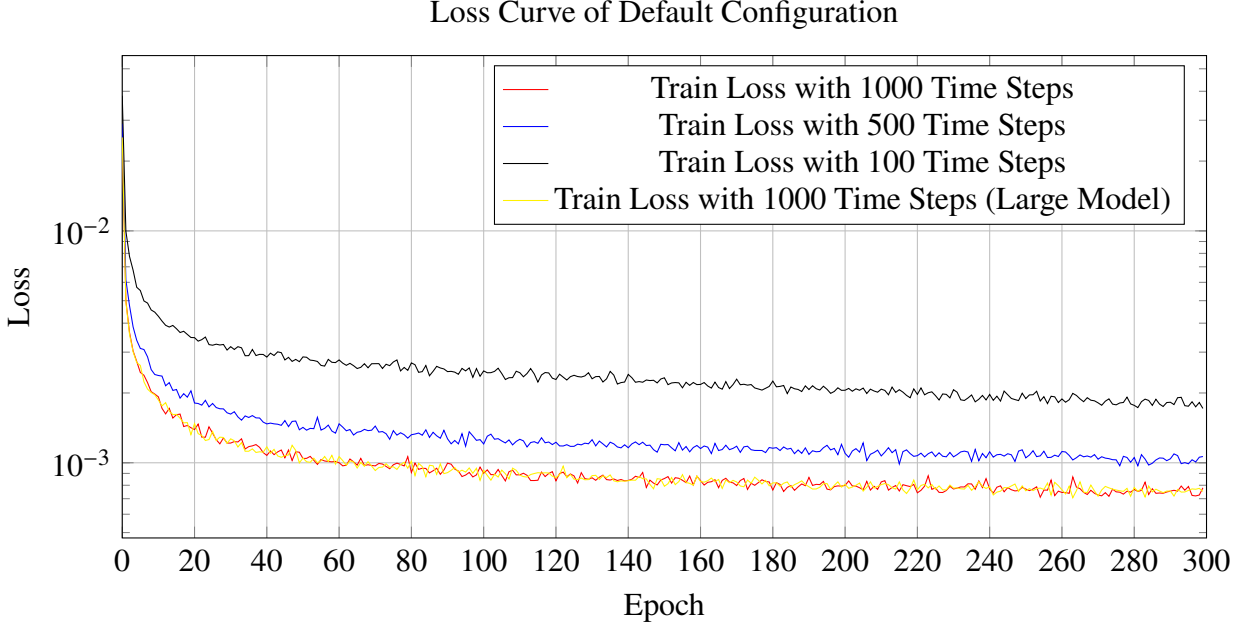


Figure 7: Loss Curve of Different Time Steps

3.4 Larger Model

I also tried to use 6 down and up-blocks, and keep the attention layers only at the lowest two resolution (bottom) layers to reduce memory usage. The model is trained for 300 epochs with a batch size of 32 and time steps of 1000.

Since the results shows similar with the previous model, I didn't include the results here. But I will include the loss curve of the larger model. The loss curve is shown in Figure ??.

References

- [1] P. Pernias, D. Rampas, M. L. Richter, C. J. Pal, and M. Aubreville, *Wuerstchen: An efficient architecture for large-scale text-to-image diffusion models*, 2023. arXiv: [2306.00637](https://arxiv.org/abs/2306.00637) [cs.CV].
- [2] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *CoRR*, vol. abs/2006.11239, 2020. arXiv: [2006.11239](https://arxiv.org/abs/2006.11239). [Online]. Available: <https://arxiv.org/abs/2006.11239>.
- [3] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *CoRR*, vol. abs/1312.6114, 2013. [Online]. Available: <https://api.semanticscholar.org/CorpusID:216078090>.

- [4] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, *et al.*, “Generative adversarial nets,” in *Neural Information Processing Systems*, 2014. [Online]. Available: <https://api.semanticscholar.org/CorpusID:261560300>.