# MM 811 Programming Assignment: Autoregressive Model

Mahbubur Rahman (mahbubu1@ualberta.ca)

Student ID: 1870909

In this assignment, four questions have been answered. All the source codes and running results are available online: `https://github.com/Durlov1603025/MM-811-Programming-Assignment.git`

## 1 Question 1

The goal of this task is to implement an image classification model using a Convolutional Neural Network (CNN) on the MNIST dataset. The model should achieve an accuracy greater than 90% while maintaining a compact architecture. The code was divided into three main parts: `model.py`, `train_eval.py`, and `main.py`, which are described below.

The file **model.py** defines the architecture used for MNIST digit classification. The network consists of three convolutional layers followed by two fully connected layers. Each convolutional layer is followed by Batch Normalization and ReLU activation to improve learning stability.

The core architecture was implemented as follows:

```python
import torch
import torch.nn as nn

class ConvNet(nn.Module):
    def __init__(self, in_channels=1, num_classes=10):
        """
        Architecture for MNIST classification.
        :param in_channels: number of input channels (1 for grayscale)
        :param num_classes: number of output classes (10 digits)
        """

        super().__init__()


        # Feature extraction layers
        self.features = nn.Sequential(
            nn.Conv2d(in_channels, 32, kernel_size=3, padding=1),
            nn.BatchNorm2d(32),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(2),

            nn.Conv2d(32, 64, kernel_size=3, padding=1),
            nn.BatchNorm2d(64),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(2),

            nn.Conv2d(64, 128, kernel_size=3, padding=1),
            nn.BatchNorm2d(128),
            nn.ReLU(inplace=True),
            nn.AdaptiveAvgPool2d((4, 4))
        )

        # Classification layers
        self.classifier = nn.Sequential(
            nn.Flatten(),
            nn.Dropout(0.4),
```

```
37            nn.Linear(128 * 4 * 4, 256),
38            nn.ReLU(inplace=True),
39            nn.Dropout(0.3),
40            nn.Linear(256, num_classes)
41        )
42
43    def forward(self, x):
44        x = self.features(x)
45        x = self.classifier(x)
46        return x
```

The file `train_eval.py` includes helper functions for model training and validation. The `train_epoch()` function iterates through all training batches, computes loss using CrossEntropyLoss, and updates model weights with backpropagation. The evaluation function `evaluate()` is similar but runs in torch.no_grad() mode to disable gradient computation. This function measures accuracy on the test dataset.

The source code of `train_eval.py` is given below:

```
1  import torch
2  import torch.nn as nn
3
4  # --------------------------------------------------------
5  # Training and evaluation utilities
6  # --------------------------------------------------------
7
8
9  def train_epoch(model, device, dataloader, optimizer, criterion):
10     """
11     Train the model for one epoch on the training dataset.
12
13     """
14     model.train()
15     total_loss, correct, total = 0.0, 0, 0
16
17     for images, targets in dataloader:
18         images, targets = images.to(device), targets.to(device)
19         optimizer.zero_grad()
20         outputs = model(images)
21         loss = criterion(outputs, targets)
22         loss.backward()
23         optimizer.step()
24
25         total_loss += loss.item() * images.size(0)
26         _, pred = outputs.max(1)
27         correct += pred.eq(targets).sum().item()
28         total += images.size(0)
29
30     avg_loss = total_loss / total
31     acc = correct / total * 100.0
32     return avg_loss, acc
33
34
35  def evaluate(model, device, dataloader, criterion):
36     """
37     Evaluate model on validation or test dataset.
38
39     """
40     model.eval()
41     total_loss, correct, total = 0.0, 0, 0
42
43     with torch.no_grad():
44         for images, targets in dataloader:
45             images, targets = images.to(device), targets.to(device)
46             outputs = model(images)
47             loss = criterion(outputs, targets)
48
49             total_loss += loss.item() * images.size(0)
50             _, pred = outputs.max(1)
51             correct += pred.eq(targets).sum().item()
52             total += images.size(0)
53
```

```
54    avg_loss = total_loss / total
55    acc = correct / total * 100.0
56    return avg_loss, acc
```

The model accuracy is computed using the equation:

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Samples}} \times 100 \tag{1}$$

The `main.py` script acts as the entry point of the project. It initializes the dataset, model, loss function, and optimizer, and coordinates the training and evaluation process.
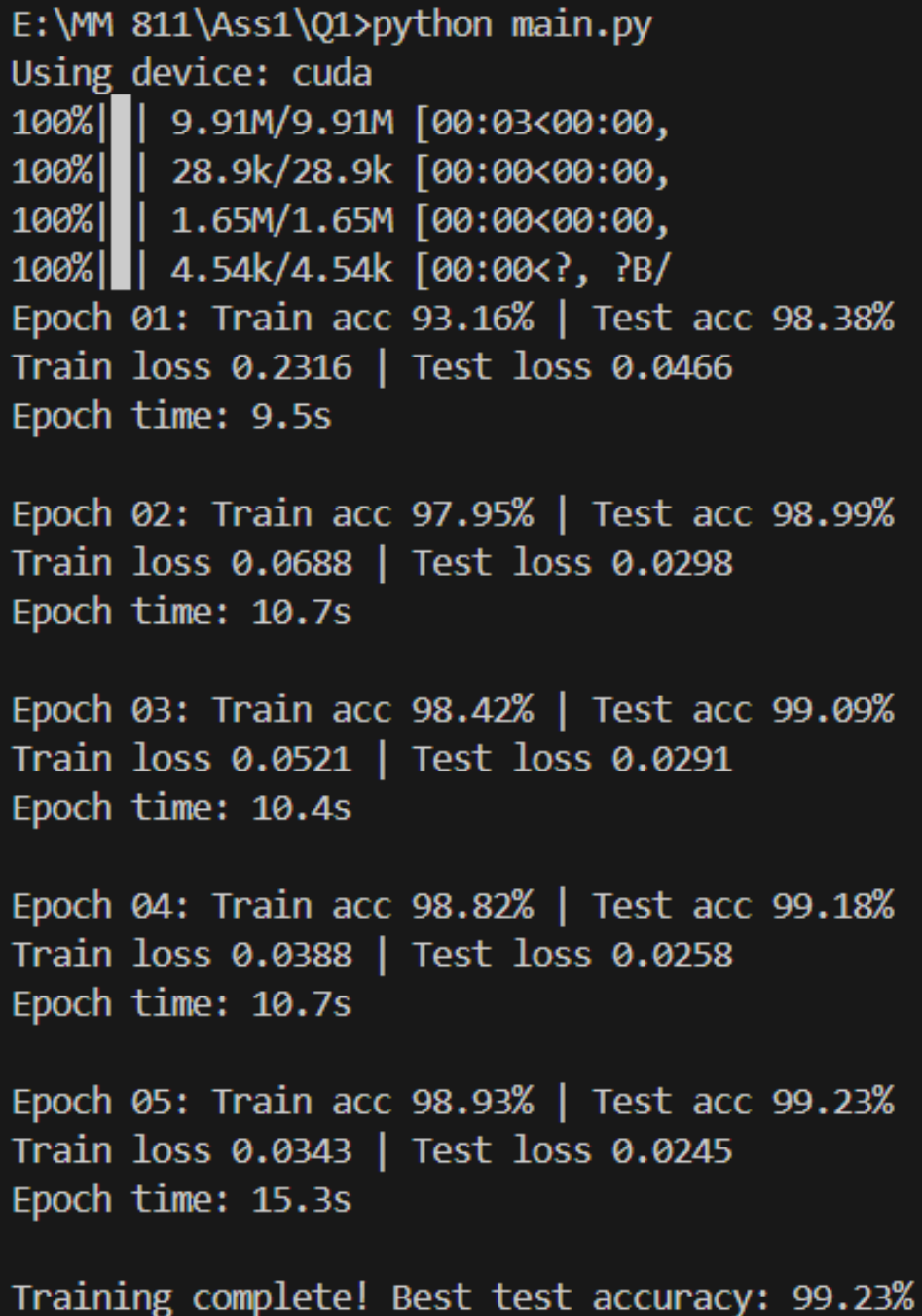
The **`main.py`** is implemented as following:

```python
1  import argparse
2  import torch
3  import torch.optim as optim
4  from torch.utils.data import DataLoader, Subset
5  from torchvision import datasets, transforms
6  from datetime import datetime
7
8  from model import ConvNet
9  from train_eval import train_epoch, evaluate
10
11
12 def main(args):
13     device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
14     print("Using device:", device)
15
16     # Data preprocessing
17     transform = transforms.Compose([
18         transforms.Resize((32, 32)),
19         transforms.ToTensor(),
20         transforms.Normalize((0.1307,), (0.3081,))
21     ])
22
23     # Load MNIST dataset
24     train_ds = datasets.MNIST(root="./data", train=True, download=True, transform=transform)
25     test_ds = datasets.MNIST(root="./data", train=False, download=True, transform=transform)
26
27     if args.subset is not None and args.subset > 0:
28         train_ds = Subset(train_ds, list(range(args.subset)))
29         print(f"Using subset of {len(train_ds)} images for training")
30
31     train_loader = DataLoader(train_ds, batch_size=args.batch_size, shuffle=True)
32     test_loader = DataLoader(test_ds, batch_size=args.batch_size, shuffle=False)
33
34     model = ConvNet(in_channels=1, num_classes=10).to(device)
35     criterion = torch.nn.CrossEntropyLoss()
36     optimizer = optim.SGD(model.parameters(), lr=args.lr, momentum=0.9)
37     scheduler = optim.lr_scheduler.StepLR(optimizer, step_size=3, gamma=0.5)
38
39     best_acc = 0.0
40     for epoch in range(1, args.epochs + 1):
41         t0 = datetime.now()
42         train_loss, train_acc = train_epoch(model, device, train_loader, optimizer, criterion)
43         test_loss, test_acc = evaluate(model, device, test_loader, criterion)
44         scheduler.step()
45
46         print(f"Epoch {epoch:02d}: Train acc {train_acc:.2f}% | Test acc {test_acc:.2f}%")
47         print(f"Train loss {train_loss:.4f} | Test loss {test_loss:.4f}")
48
49         if test_acc > best_acc:
50             best_acc = test_acc
51             torch.save(model.state_dict(), "best_model.pth")
52
53         print(f"Epoch time: {(datetime.now() - t0).total_seconds():.1f}s\n")
54
55     print(f"Training complete! Best test accuracy: {best_acc:.2f}%")
56
57
58 if __name__ == "__main__":
```

```
59    parser = argparse.ArgumentParser()
60    parser.add_argument("--epochs", type=int, default=5)
61    parser.add_argument("--batch-size", type=int, default=128)
62    parser.add_argument("--lr", type=float, default=0.01)
63    parser.add_argument("--subset", type=int, default=None)
64    args = parser.parse_args()
65
66    main(args)
```

A screenshot of program output is demonstrated in Fig. 1.



```
E:\MM 811\Ass1\Q1>python main.py
Using device: cuda
100%|     | 9.91M/9.91M [00:03<00:00,
100%|     | 28.9k/28.9k [00:00<00:00,
100%|     | 1.65M/1.65M [00:00<00:00,
100%|     | 4.54k/4.54k [00:00<?, ?B/
Epoch 01: Train acc 93.16% | Test acc 98.38%
Train loss 0.2316 | Test loss 0.0466
Epoch time: 9.5s

Epoch 02: Train acc 97.95% | Test acc 98.99%
Train loss 0.0688 | Test loss 0.0298
Epoch time: 10.7s

Epoch 03: Train acc 98.42% | Test acc 99.09%
Train loss 0.0521 | Test loss 0.0291
Epoch time: 10.4s

Epoch 04: Train acc 98.82% | Test acc 99.18%
Train loss 0.0388 | Test loss 0.0258
Epoch time: 10.7s

Epoch 05: Train acc 98.93% | Test acc 99.23%
Train loss 0.0343 | Test loss 0.0245
Epoch time: 15.3s

Training complete! Best test accuracy: 99.23%
```

Figure 1: Output of Question 1

# 2 Question 2

The binary autoencoder was implemented in the file `autoencoder_model.py`. The architecture consists of three main parts: an encoder, a binarizer, and a decoder. The encoder compresses the input image of size $1 \times 32 \times 32$ into a smaller latent feature map of size $3 \times 4 \times 4$ using a series of convolutional, batch normalization, and ReLU layers. Immediately after the encoder, the latent representation is passed through a binarization function that converts the continuous values into discrete $\{-1, +1\}$ activations using a Straight-Through Estimator (STE) [1]. The decoder reconstructs the original image from these binary latent values using transposed convolutional layers. A final `tanh` activation ensures that the reconstructed image values lie within the normalized range $[-1, 1]$. The source code (***autoencoder_model.py***) is shown as follows:

```python
import torch
import torch.nn as nn

# BinarizeFunction
class BinarizeFunction(torch.autograd.Function):
    @staticmethod
    def forward(ctx, x):
        return x.sign()  # outputs -1 or +1

    @staticmethod
    def backward(ctx, grad_output):
        return grad_output  # straight-through estimator


def binarize(x: torch.Tensor) -> torch.Tensor:
    return BinarizeFunction.apply(x)


class AutoEncoder(nn.Module):
    def __init__(self, use_binary: bool = True):
        super().__init__()
        self.use_binary = use_binary

        # Encoder: 1x32x32 -> 3x4x4
        self.encoder = nn.Sequential(
            nn.Conv2d(1, 32, 3, 1, 1),  # 1x32x32 -> 32x32x32
            nn.BatchNorm2d(32),
            nn.ReLU(inplace=True),
            nn.Conv2d(32, 64, 4, 2, 1),  # -> 64x16x16
            nn.BatchNorm2d(64),
            nn.ReLU(inplace=True),
            nn.Conv2d(64, 128, 4, 2, 1),  # -> 128x8x8
            nn.BatchNorm2d(128),
            nn.ReLU(inplace=True),
            nn.Conv2d(128, 256, 4, 2, 1),  # -> 256x4x4
            nn.BatchNorm2d(256),
            nn.ReLU(inplace=True),
            nn.Conv2d(256, 3, 3, 1, 1),  # -> 3x4x4 (latent logits)
        )

        # Decoder: 3x4x4 -> 1x32x32
        self.decoder = nn.Sequential(
            nn.ConvTranspose2d(3, 128, 4, 2, 1),  # -> 128x8x8
            nn.BatchNorm2d(128),
            nn.ReLU(inplace=True),
            nn.ConvTranspose2d(128, 64, 4, 2, 1),  # -> 64x16x16
            nn.BatchNorm2d(64),
            nn.ReLU(inplace=True),
            nn.ConvTranspose2d(64, 32, 4, 2, 1),  # -> 32x32x32
            nn.BatchNorm2d(32),
            nn.ReLU(inplace=True),
            nn.Conv2d(32, 1, 3, 1, 1),  # -> 1x32x32
            nn.Tanh(),  # output values in {-1, +1}
        )

    def forward(self, x):
        latent_logits = self.encoder(x)  # (N, 3, 4, 4), real-valued
```

```
58         if self.use_binary:
59             latent = binarize(latent_logits)  # strictly +1 or -1
60         else:
61             latent = latent_logits  # warm-up (no binarization)
62         recon = self.decoder(latent)
63         # Return both for training utilities
64         return recon, latent, latent_logits
```

The training and evaluation procedures are implemented in the file train_autoencoder.py. Two key functions are defined: train_epoch() and evaluate(). The train_epoch() function performs one full epoch of training using either L1 or MSE reconstruction loss.

The evaluate() function computes the average PSNR on the dataset using the functions compute_psnr_torch() and compute_psnr_sigs(). The source code of train_autoencoder.py is given below:

```
1  import torch
2  import torch.nn.functional as F
3
4
5  def compute_psnr_torch(imgs, refs, eps=1e-8):
6      mse = F.mse_loss(imgs, refs, reduction='none').mean(dim=[1, 2, 3])
7      psnr = 20 * torch.log10(255.0 / torch.sqrt(mse + eps))
8      return psnr
9
10
11 def compute_psnr_sigs(imgs, refs):
12     imgs_norm = imgs + 1.0
13     refs_norm = refs + 1.0
14     imgs_norm *= 0.5
15     refs_norm *= 0.5
16     imgs_norm *= 255.0
17     refs_norm *= 255.0
18     psnr_vals = compute_psnr_torch(imgs_norm, refs_norm)
19     return psnr_vals
20
21
22 def train_epoch(model, dataloader, device, optimizer, *,
23                 use_l1=True, reg_lambda=1e-3, push_margin=1.0):
24     model.train()
25     total_loss = 0.0
26     total_psnr = 0.0
27
28     for imgs, _ in dataloader:
29         imgs = imgs.to(device)
30         imgs = imgs * 2.0 - 1.0
31
32         optimizer = torch.optim.Adam(model.parameters(), lr=1e-3)
33         optimizer.zero_grad(set_to_none=True)
34         recon, binary_latent, latent_logits = model(imgs)
35
36         # Reconstruction loss
37         if use_l1:
38             rec_loss = F.l1_loss(recon, imgs)
39         else:
40             rec_loss = F.mse_loss(recon, imgs)
41
42         commit = F.relu(push_margin - latent_logits.abs()).mean()
43         loss = rec_loss + reg_lambda * commit
44
45         loss.backward()
46         optimizer.step()
47
48         total_loss += rec_loss.item()
49         total_psnr += compute_psnr_sigs(recon.detach(), imgs).mean().item()
50
51     avg_loss = total_loss / len(dataloader)
52     avg_psnr = total_psnr / len(dataloader)
53     return avg_loss, avg_psnr
54
55
56 def evaluate(model, dataloader, device):
```

```
57     model.eval()
58     total_psnr = 0.0
59     with torch.no_grad():
60         for imgs, _ in dataloader:
61             imgs = imgs.to(device)
62             imgs = imgs * 2.0 - 1.0
63             recon, _, _ = model(imgs)
64             total_psnr += compute_psnr_sigs(recon, imgs).mean().item()
65
66     avg_psnr = total_psnr / len(dataloader)
67     return avg_psnr
```

The file `main.py` manages the data loading, model initialization, and training loop. During the first few epochs (warm-up phase), binarization is disabled to allow the decoder to learn from continuous latent features. After the warm-up period, binarization is activated so that the latent space becomes strictly $\{-1, +1\}$. The model is trained and evaluated on the same subset of the MNIST dataset to test its memorization capability. The best-performing model is automatically saved based on PSNR improvement. The source code of `main.py` is given below:

```
1  import argparse
2  import torch
3  import torch.optim as optim
4  from torch.utils.data import DataLoader
5  from torchvision import datasets, transforms
6
7  from autoencoder_model import AutoEncoder
8  from train_autoencoder import train_epoch, evaluate
9
10 def get_dataloaders(batch_size: int, num_workers: int = 2):
11     tfm = transforms.Compose([
12         transforms.Resize(32),
13         transforms.ToTensor(),
14
15     ])
16     train_ds = datasets.MNIST(root='./data', train=True, download=True, transform=tfm)
17     test_ds = train_ds
18
19     train_loader = DataLoader(
20         train_ds, batch_size=batch_size, shuffle=True,
21         num_workers=num_workers, pin_memory=True
22     )
23     val_loader = DataLoader(
24         test_ds, batch_size=batch_size, shuffle=False,
25         num_workers=num_workers, pin_memory=True
26     )
27     return train_loader, val_loader
28
29
30 def main():
31     parser = argparse.ArgumentParser()
32     parser.add_argument('--epochs', type=int, default=200)
33     parser.add_argument('--warmup_epochs', type=int, default=10, help='Train without binarization for first N
        epochs')
34     parser.add_argument('--batch_size', type=int, default=128)
35     parser.add_argument('--lr', type=float, default=1e-3)
36     parser.add_argument('--reg_lambda', type=float, default=1e-3, help='Weight for latent push-away regularizer
        ')
37     parser.add_argument('--use_l1', action='store_true', help='Use L1 loss (default off if flag absent)')
38     parser.add_argument('--no_l1', action='store_true', help='Force MSE even if --use_l1 set elsewhere')
39     args = parser.parse_args()
40
41     # Resolve L1 flag
42     use_l1 = True
43     if args.no_l1:
44         use_l1 = False
45     elif args.use_l1:
46         use_l1 = True
47
48     device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
49     print('Using device:', device)
50
```

```
51     train_loader, val_loader = get_dataloaders(args.batch_size)
52
53     model = AutoEncoder(use_binary=False)   # start with warm-up by default
54     model.to(device)
55
56     optimizer = optim.Adam(model.parameters(), lr=args.lr)
57     scheduler = optim.lr_scheduler.ReduceLROnPlateau(
58         optimizer, mode='max', factor=0.5, patience=10
59 )
60
61
62     best_psnr = 0.0
63
64     for epoch in range(1, args.epochs + 1):
65         # Toggle binarization after warmup period
66         if epoch == args.warmup_epochs + 1:
67             model.use_binary = True
68
69
70         train_loss, train_psnr = train_epoch(
71             model, train_loader, device, optimizer,
72             use_l1=use_l1, reg_lambda=args.reg_lambda
73         )
74         val_psnr = evaluate(model, val_loader, device)
75         scheduler.step(val_psnr)
76
77         print(f"Epoch {epoch:03d} | TrainLoss: {train_loss:.4f} | TrainPSNR: {train_psnr:.2f} | ValPSNR: {
    val_psnr:.2f}")
78
79         # save best
80         if val_psnr > best_psnr:
81             best_psnr = val_psnr
82             torch.save(model.state_dict(), 'autoencoder_best.pth')
83             print(f" -> Saved new best (Val PSNR: {best_psnr:.2f})")
84
85         # Early exit if goal reached
86         if val_psnr > 25.0:
87             print(f"\nGoal achieved! PSNR  25.0 (Current: {val_psnr:.2f}) ***")
88             break
89
90     print(f"\nTraining complete. Best Validation PSNR: {best_psnr:.2f}")
91
92
93
94 if __name__ == '__main__':
95     main()
```

Table 1: AutoEncoder Network Architecture

| Part | Layer Type | Kernel / Filter Size | Stride | Padding | Output Size |
|---|---|---|---|---|---|
| **Encoder** | Conv2d(1, 32) | $3 \times 3$ | 1 | 1 | $32 \times 32 \times 32$ |
| | Conv2d(32, 64) | $4 \times 4$ | 2 | 1 | $64 \times 16 \times 16$ |
| | Conv2d(64, 128) | $4 \times 4$ | 2 | 1 | $128 \times 8 \times 8$ |
| | Conv2d(128, 256) | $4 \times 4$ | 2 | 1 | $256 \times 4 \times 4$ |
| | Conv2d(256, 3) | $3 \times 3$ | 1 | 1 | $3 \times 4 \times 4$ |
| **Decoder** | ConvTranspose2d(3, 256) | $3 \times 3$ | 1 | 1 | $256 \times 4 \times 4$ |
| | ConvTranspose2d(256, 128) | $4 \times 4$ | 2 | 1 | $128 \times 8 \times 8$ |
| | ConvTranspose2d(128, 64) | $4 \times 4$ | 2 | 1 | $64 \times 16 \times 16$ |
| | ConvTranspose2d(64, 32) | $4 \times 4$ | 2 | 1 | $32 \times 32 \times 32$ |
| | Conv2d(32, 1) | $3 \times 3$ | 1 | 1 | $1 \times 32 \times 32$ |

An overview of AutoEncoder Network Architecture can be found from Table 1. The output is shown in Fig. 2

Figure 2: Output of question 2

## 3 Question 3

### 3.1 Sequence to Token

The goal of this experiment is to train an autoregressive model that can memorize binary latent sequences extracted from the trained AutoEncoder (Question 2) and generate new MNIST-like digit images token by token. The success criterion is achieving more than 80% memorization rate, defined as the proportion of generated sequences that exactly match the training data.

**Dataset Preparation**

- Uses the **binary latent codes** produced by the trained autoencoder (Question 2).

- Each latent tensor $(3 \times 4 \times 4)$ elements is flattened into a 48-token sequence.

- Continuous values in $[-1, 1]$ are mapped to binary values (either +1 or -1) using the transformation.

- Each training sample is split into prefix–target pairs

```
1    class LatentSequenceDataset(Dataset):
2    def __init__(self, autoencoder, images, device):
3        self.device = device
4        autoencoder.eval()
5
6        # Extract binary latent codes
7        with torch.no_grad():
8            imgs_normalized = images.to(device) * 2 - 1
9            _, latents, _ = autoencoder(imgs_normalized)
10
11        self.latents = latents.view(latents.size(0), -1)  # (N, 48), values {-1, +1}
12        self.tokens = ((self.latents + 1) / 2).long()
13
14        # Create autoregressive training pairs
15        seq_len = self.tokens.size(1)  # 48
16        self.inputs, self.labels, self.lengths = [], [], []
17
18        for seq in self.tokens:
19            for i in range(1, seq_len):
```

```
20              self.inputs.append(seq[:i])
21              self.labels.append(seq[i])
22              self.lengths.append(i)
23
24      def __len__(self):
25          return len(self.inputs)
26
27      def __getitem__(self, idx):
28          x = self.inputs[idx].float()
29          y = self.labels[idx].float()
30          length = self.lengths[idx]
31          return x, y, length
32
33
34  def collate_fn(batch):
35      sequences, labels, lengths = zip(*batch)
36      # Pad sequences to max length in this batch
37      sequences_padded = pad_sequence(sequences, batch_first=True, padding_value=0)
38      labels = torch.stack(labels)
39      lengths = torch.tensor(lengths)
40
41      return sequences_padded, labels, lengths
```

## Model Architecture

The model, named **Seq2TokenLSTM**, is designed for next-token prediction. The model can be defined as:

- **Embedding Layer:** Maps binary tokens into 128-dimensional continuous embeddings.

- **LSTM [2] Layers:** Two stacked layers with hidden size 128, capturing temporal dependencies between tokens.

- **Fully Connected Layers:** Linear(128$\rightarrow$64) $\rightarrow$ ReLU $\rightarrow$ Linear(64$\rightarrow$1) $\rightarrow$ Sigmoid structure producing the probability for the next token.

Table 2: Seq2Token Network Architecture

| Layer Type | Input Size | Output Size | Activation | Notes |
|---|---|---|---|---|
| Embedding(2, 128) | binary token | 128-dim vector | – | Binary token $\in \{-1, +1\}$ embedded into 128-dim vector |
| LSTM Layer 1 | 128 | 128 | – | Hidden size = 128, batch_first=True |
| LSTM Layer 2 | 128 | 128 | – | num_layers = 2, no dropout |
| Linear(128, 64) | 128 | 64 | ReLU | First fully connected layer |
| Linear(64, 1) | 64 | 1 | Sigmoid | Predicts next binary token probability |

An overview of the Sequence to Token Network Architecture is shown in the Table 2. The model can be defined as follows:

```
1   class Seq2TokenLSTM(nn.Module):
2   def __init__(self, input_size=1, hidden_size=128, num_layers=2, dropout=0.0):
3       super().__init__()
4
5       self.hidden_size = hidden_size
6       self.num_layers = num_layers
7
8       # Embedding layer
9       self.embedding = nn.Embedding(num_embeddings=2, embedding_dim=hidden_size)
10
11      # LSTM layers
12      self.lstm = nn.LSTM(
13          input_size=hidden_size,
14          hidden_size=hidden_size,
15          num_layers=num_layers,
16          batch_first=True,
17          dropout=0  # Removed dropout to encourage memorization
18      )
```

```
19
20      # Output layers
21      self.fc = nn.Sequential(
22          nn.Linear(hidden_size, hidden_size // 2),
23          nn.ReLU(),
24          nn.Linear(hidden_size // 2, 1),
25          nn.Sigmoid()
26      )
27
28  def forward(self, x, lengths=None):
29      batch_size = x.size(0)
30
31      # Embed tokens
32      embedded = self.embedding(x.long())
33      if lengths is not None:
34          packed = pack_padded_sequence(
35              embedded, lengths.cpu(),
36              batch_first=True,
37              enforce_sorted=False
38          )
39          lstm_out, (hidden, cell) = self.lstm(packed)
40          # Use last hidden state
41          output = hidden[-1]
42      else:
43          # Process full sequence
44          lstm_out, (hidden, cell) = self.lstm(embedded)
45          output = hidden[-1]
46
47      # Predict next token
48      pred = self.fc(output).squeeze(-1)
49
50      return pred
```

## Training Procedure:

- Training uses Binary Cross-Entropy (BCE) loss [3]:

$$L = -[y \log(p) + (1-y) \log(1-p)] \tag{2}$$

  where $y$ is the true token and $p$ is the predicted probability.

- Adam optimizer with learning rate $1 \times 10^{-1}$ is used for stable convergence.

The training process is done by following source code:

```
1   def train_seq2token(model, loader, device, epochs=50, lr=1e-3):
2       optimizer = torch.optim.Adam(model.parameters(), lr=lr)
3       loss_fn = nn.BCELoss()
4
5       history = {'loss': [], 'acc': []}
6
7       for epoch in range(1, epochs + 1):
8           model.train()
9           total_loss = 0.0
10          total_correct = 0
11          total_samples = 0
12
13          pbar = tqdm(loader, desc=f"Epoch {epoch}/{epochs}", leave=False)
14          for x, y, lengths in pbar:
15              x, y = x.to(device), y.to(device)
16              lengths = lengths.to(device)
17
18              # Forward pass
19              pred = model(x, lengths)
20              loss = loss_fn(pred, y)
21
22              # Backward pass
23              optimizer.zero_grad()
24              loss.backward()
25              torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)
```

```
26              optimizer.step()
27
28              # Track metrics
29              total_loss += loss.item()
30              predicted_class = (pred > 0.5).float()
31              total_correct += (predicted_class == y).sum().item()
32              total_samples += y.size(0)
33
34              pbar.set_postfix({
35                  'loss': f'{loss.item():.4f}',
36                  'acc': f'{total_correct/total_samples:.3f}'
37              })
38
39          avg_loss = total_loss / len(loader)
40          avg_acc = total_correct / total_samples
41          history['loss'].append(avg_loss)
42          history['acc'].append(avg_acc)
43
44          print(f"Epoch {epoch:03d} | Loss: {avg_loss:.4f} | Acc: {avg_acc:.4f}")
45
46      return history
```

## Token Generation

The trained model generates binary sequences autoregressively. Starting from an empty input, it predicts each token in order:

$$t_i \sim \text{Bernoulli}(P(t_i = 1)) \tag{3}$$

Each new token depends on all previous tokens, ensuring autoregressive dependency. After generating 48 tokens, each sequence is converted back to $\{-1, +1\}$, reshaped into $(3, 4, 4)$, and decoded through the AutoEncoder's decoder to obtain MNIST-like images.

The following function `generate_sequences` autoregressively builds sequences:

```
1   def sample_sequences(model, device, num_samples=64, seq_len=48, temperature=1.0):
2       model.eval()
3       generated_tokens = torch.zeros(num_samples, seq_len, device=device)
4
5       with torch.no_grad():
6           for i in range(seq_len):
7               if i == 0:
8                   # First token: start with empty sequence
9                   dummy = torch.zeros(num_samples, 1, device=device)
10                  probs = torch.ones(num_samples, device=device) * 0.5
11              else:
12                  # Use generated tokens so far
13                  current_seq = generated_tokens[:, :i]
14                  lengths = torch.full((num_samples,), i, device=device)
15                  probs = model(current_seq, lengths)
16
17              # Apply temperature for diversity
18              if temperature != 1.0:
19                  probs = torch.sigmoid(torch.logit(probs) / temperature)
20
21              token = torch.bernoulli(probs)
22              generated_tokens[:, i] = token
23
24      return generated_tokens
```

The generated binary sequences are later decoded by the autoencoder's decoder to reconstruct MNIST-like images as follows:

```
1   def decode_latents(autoencoder, token_samples, device):
2       autoencoder.eval()
3
4       # Convert tokens
5       latent_samples = token_samples * 2 - 1
6
7       # Reshape to (N, 3, 4, 4) for decoder
```

```
8      latent_samples = latent_samples.view(-1, 3, 4, 4)
9
10     with torch.no_grad():
11         imgs = autoencoder.decoder(latent_samples.to(device))
12
13     return imgs
```

## Memorization Evaluation

The generated sequences are compared with training sequences to measure memorization. Exact matches are identified using:

```
1      def check_memorization(generated_tokens, train_tokens):
2      diffs = (generated_tokens.unsqueeze(1) - train_tokens.unsqueeze(0)).abs().sum(dim=-1)
3
4      # A perfect match has distance 0
5      matches = (diffs == 0).any(dim=1).float()
6
7      return matches.mean().item()
```

In this experiment, 64 sequences were generated, 57 of which matched training data, giving an **89.06% memorization rate**

The output of this code is shown in Fig. 3

```
Epoch 078 | Loss: 0.1071 | Acc: 0.9445
Epoch 079 | Loss: 0.1057 | Acc: 0.9435
Epoch 080 | Loss: 0.1068 | Acc: 0.9405
Epoch 081 | Loss: 0.1062 | Acc: 0.9418
Epoch 082 | Loss: 0.1059 | Acc: 0.9407
Epoch 083 | Loss: 0.1062 | Acc: 0.9435
Epoch 084 | Loss: 0.1067 | Acc: 0.9446
Epoch 085 | Loss: 0.1074 | Acc: 0.9438
Epoch 086 | Loss: 0.1048 | Acc: 0.9435
Epoch 087 | Loss: 0.1059 | Acc: 0.9433
Epoch 088 | Loss: 0.1041 | Acc: 0.9450
Epoch 089 | Loss: 0.1055 | Acc: 0.9430
Epoch 090 | Loss: 0.1059 | Acc: 0.9423
Epoch 091 | Loss: 0.1134 | Acc: 0.9377
Epoch 092 | Loss: 0.1655 | Acc: 0.9234
Epoch 093 | Loss: 0.1809 | Acc: 0.9166
Epoch 094 | Loss: 0.1757 | Acc: 0.9209
Epoch 095 | Loss: 0.1427 | Acc: 0.9270
Epoch 096 | Loss: 0.1417 | Acc: 0.9325
Epoch 097 | Loss: 0.1263 | Acc: 0.9337
Epoch 098 | Loss: 0.1126 | Acc: 0.9400
Epoch 099 | Loss: 0.1071 | Acc: 0.9415
Epoch 100 | Loss: 0.1036 | Acc: 0.9433

OK - Saved model -> seq2token_lstm.pth

Memorization Ratio: 89.06%
(57/64 generated sequences match training data)


========================================
EXPERIMENT SUMMARY Sequence to Token
========================================
Training samples: 128
Training epochs: 100
Final loss: 0.1036
Final accuracy: 0.9433
Generated samples: 64
Memorization rate: 89.06%
========================================
```

Figure 3: Output of Sequence to Token

## 3.2   Sequence to Sequence

The objective of this part is designing a direct **sequence-to-sequence (Seq2Seq)** model.

## Model Architecture

The model uses a stack of fully-connected residual layers. This model is implemented as:

```python
class Seq2Seq(nn.Module):
def __init__(self, seq_len=48, hidden_dim=256, num_blocks=12):
    super().__init__()

    self.seq_len = seq_len

    # Input projection
    self.input_proj = nn.Sequential(
        nn.Linear(seq_len, hidden_dim),
        nn.LayerNorm(hidden_dim),
        nn.GELU()
    )

    # Residual blocks
    self.blocks = nn.ModuleList([
        ResidualFC(hidden_dim, hidden_scale=4)
        for _ in range(num_blocks)
    ])

    # Output projection
    self.output_proj = nn.Sequential(
        nn.Linear(hidden_dim, hidden_dim // 2),
        nn.GELU(),
        nn.Linear(hidden_dim // 2, seq_len),
        nn.Sigmoid()
    )

def forward(self, x):
    # Project to hidden dim
    h = self.input_proj(x)

    # Process through residual blocks
    for block in self.blocks:
        h = block(h)

    # Project to output
    out = self.output_proj(h)

    return out
```

It consists of three major components:

- **Input Projection:** Projects each 48-dimensional binary input into a 256-dimensional hidden representation using a linear layer, layer normalization, and GELU activation.

- **Residual Blocks:** Twelve (or more) `ResidualFC` blocks, each containing two linear layers with a skip connection.

- **Output Projection:** Maps the hidden representation back to a 48-dimensional output that produces probabilities for each token.

The ResidualFC blocks are defined as follows:

```python
class ResidualFC(nn.Module):
    def __init__(self, dim, hidden_scale=4):
        super().__init__()
        self.block = nn.Sequential(
            nn.Linear(dim, dim * hidden_scale),
            nn.LayerNorm(dim * hidden_scale),
            nn.GELU(),
            nn.Linear(dim * hidden_scale, dim),
        )

    def forward(self, x):
        return x + self.block(x)
```

An overview of Sequence to Sequence Network Architecture is shown in Table 3.

Table 3: Seq2Seq Network Architecture

| Layer Type | Input Size | Output Size | Activation | Notes |
|---|---|---|---|---|
| Linear(48, 256) | 48 | 256 | GELU | Input projection + LayerNorm |
| ResidualFC Block (×12) | 256 | 256 | GELU | Each block: Linear(256, 1024) → LN → GELU → Linear(1024, 256) |
| Linear(256, 128) | 256 | 128 | GELU | Output projection (part 1) |
| Linear(128, 48) | 128 | 48 | Sigmoid | Output projection (part 2); reconstructs full sequence |

**Training Procedure**

Training uses the same logic as Part 1 with some key differences: Instead of next-token prediction, the model learns to denoise corrupted binary sequences. The denoising autoencoder framework [4] is used during training. Tokens are randomly flipped with probability noise_prob (default 0.05). This makes the model robust while still memorizing. Cosine Annealing LR scheduler has been used for smooth learning-rate decay.

```python
def train_seq2token(model, loader, device, epochs=50, lr=1e-3):
    optimizer = torch.optim.Adam(model.parameters(), lr=lr)
    loss_fn = nn.BCELoss()

    history = {'loss': [], 'acc': []}

    for epoch in range(1, epochs + 1):
        model.train()
        total_loss = 0.0
        total_correct = 0
        total_samples = 0

        pbar = tqdm(loader, desc=f"Epoch {epoch}/{epochs}", leave=False)
        for x, y, lengths in pbar:
            x, y = x.to(device), y.to(device)
            lengths = lengths.to(device)

            # Forward pass
            pred = model(x, lengths)
            loss = loss_fn(pred, y)

            # Backward pass
            optimizer.zero_grad()
            loss.backward()
            torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)
            optimizer.step()

            # Track metrics
            total_loss += loss.item()
            predicted_class = (pred > 0.5).float()
            total_correct += (predicted_class == y).sum().item()
            total_samples += y.size(0)

            pbar.set_postfix({
                'loss': f'{loss.item():.4f}',
                'acc': f'{total_correct/total_samples:.3f}'
            })

        avg_loss = total_loss / len(loader)
        avg_acc = total_correct / total_samples
        history['loss'].append(avg_loss)
        history['acc'].append(avg_acc)

        print(f"Epoch {epoch:03d} | Loss: {avg_loss:.4f} | Acc: {avg_acc:.4f}")

    return history
```

### Sequence Generation

Unlike the other method, this model generates all tokens at once. Starting from random binary sequences, the network iteratively refines its predictions for several passes until convergence.

```python
def generate_sequences_direct(model, device, num_samples=64, seq_len=48):
    model.eval()

    with torch.no_grad():
        sequences = torch.rand(num_samples, seq_len, device=device)
        sequences = (sequences > 0.5).float()

        # Let model refine (multiple passes for convergence)
        for _ in range(10):
            sequences = model(sequences)
            sequences = (sequences > 0.5).float()

    return sequences
```

Each generated sequence is then mapped back to $-1, +1$ and decoded by the AutoEncoder's decoder to reconstruct MNIST-like images.

```python
def decode_sequences(autoencoder, token_sequences, device):
    autoencoder.eval()
    latent_sequences = token_sequences * 2 - 1
    latent_sequences = latent_sequences.view(-1, 3, 4, 4)

    with torch.no_grad():
        imgs = autoencoder.decoder(latent_sequences.to(device))

    return imgs
```

### Evaluation and Memorization

The memorization evaluation is almost identical to Part 1. In this experiment, 64 sequences were generated, 60 of which matched training data, giving an **93.75% memorization rate**

The output of Sequence to Sequence is shown in Fig. 4

# Question 4

The goal of this experiment is to reproduce the memorization task using a **non-image dataset**. Instead of MNIST images, we employ a **character-level text dataset** and train an autoregressive model to memorize and regenerate the same sentences under extremely limited data conditions using character-level modeling [5].

### Dataset Preparation

A very small text dataset is used, consisting of two or more sentences. Each sentence is converted into a sequence of characters.

- Each character is mapped to an index (character-level encoding).

- Autoregressive training pairs are created as prefix–next-token pairs.

### Model Architecture

The network is a character-level LSTM with a small number of parameters, trained to predict the next token given a sequence. An overview of the architecture with the source code is given in the Table 4:

```
PS E:\MM 811\Ass1\Q2-3> python q3_seq2seq.py
Using device: cuda

OK - Loaded autoencoder_best.pth

Dataset: 64 sequences of length 48
Epoch 010 | Loss: 0.3136 | Acc: 0.8613 | LR: 0.000994
Epoch 020 | Loss: 0.1704 | Acc: 0.9281 | LR: 0.000976
Epoch 030 | Loss: 0.1111 | Acc: 0.9596 | LR: 0.000946
Epoch 040 | Loss: 0.0786 | Acc: 0.9727 | LR: 0.000905
Epoch 050 | Loss: 0.0430 | Acc: 0.9854 | LR: 0.000854
Epoch 060 | Loss: 0.0405 | Acc: 0.9847 | LR: 0.000794
Epoch 070 | Loss: 0.0494 | Acc: 0.9847 | LR: 0.000727
Epoch 080 | Loss: 0.0362 | Acc: 0.9863 | LR: 0.000655
Epoch 090 | Loss: 0.0405 | Acc: 0.9876 | LR: 0.000578
Epoch 100 | Loss: 0.0211 | Acc: 0.9932 | LR: 0.000500
Epoch 110 | Loss: 0.0108 | Acc: 0.9971 | LR: 0.000422
Epoch 120 | Loss: 0.0256 | Acc: 0.9948 | LR: 0.000345
Epoch 130 | Loss: 0.0076 | Acc: 0.9971 | LR: 0.000273
Epoch 140 | Loss: 0.0078 | Acc: 0.9971 | LR: 0.000206
Epoch 150 | Loss: 0.0023 | Acc: 1.0000 | LR: 0.000146
Epoch 160 | Loss: 0.0060 | Acc: 0.9971 | LR: 0.000095
Epoch 170 | Loss: 0.0025 | Acc: 0.9993 | LR: 0.000054
Epoch 180 | Loss: 0.0020 | Acc: 0.9997 | LR: 0.000024
Epoch 190 | Loss: 0.0057 | Acc: 0.9980 | LR: 0.000006
Epoch 200 | Loss: 0.0021 | Acc: 0.9993 | LR: 0.000000
Memorization Ratio: 93.75%
(60/64 sequences match training)

Saved visualization -> q3_seq2seq_fc_results.png
=================================================
SUMMARY (Sequence to Sequence)
=================================================
Training samples: 64
Epochs: 200
Final loss: 0.0021
Final accuracy: 0.9993
Memorization: 93.75%
=================================================
```

Figure 4: Output of Sequence to Sequence

Table 4: Architecture for Text Memorization

| Layer Type | Input Size | Output Size | Activation | Notes |
| --- | --- | --- | --- | --- |
| Embedding | vocab_size | hidden_size | – | Converts each character token into a dense vector. |
| LSTM (×3 layers) | hidden_size | hidden_size | – | Learns sequential dependencies across characters. |
| Linear | hidden_size | vocab_size | ReLU / Softmax | Produces probability distribution over next character. |

```
1    class TextAutoRegressiveLSTM(nn.Module):
2    def __init__(self, vocab_size, hidden_size=256, num_layers=3):
3        super().__init__()
4
5        self.vocab_size = vocab_size
6        self.hidden_size = hidden_size
7        self.num_layers = num_layers
8
9        # Embedding layer
10        self.embedding = nn.Embedding(
11            num_embeddings=vocab_size,
12            embedding_dim=hidden_size
13        )
14
15        # LSTM layers
16        self.lstm = nn.LSTM(
17            input_size=hidden_size,
18            hidden_size=hidden_size,
19            num_layers=num_layers,
20            batch_first=True,
21            dropout=0
22        )
23
24        # Output layers - produce distribution over vocab
25        self.fc = nn.Sequential(
26            nn.Linear(hidden_size, hidden_size),
27            nn.ReLU(),
28            nn.Linear(hidden_size, vocab_size)  # Output: logits for each character
29        )
30
31    def forward(self, x, lengths=None):
32        # Embed
33        embedded = self.embedding(x.long())
34
35        # LSTM
36        if lengths is not None:
37            from torch.nn.utils.rnn import pack_padded_sequence
38            packed = pack_padded_sequence(
39                embedded, lengths.cpu(),
40                batch_first=True,
41                enforce_sorted=False
42            )
43            lstm_out, (hidden, cell) = self.lstm(packed)
44            output = hidden[-1]
45        else:
46            lstm_out, (hidden, cell) = self.lstm(embedded)
47            output = hidden[-1]
48
49        # Predict distribution over vocab
50        logits = self.fc(output)
51        probs = F.softmax(logits, dim=1)
52
53        return logits, probs
```

**Training process**

The model is trained using the `CrossEntropyLoss` to predict the next token in each sequence.

- Loss: CrossEntropy between predicted and true next tokens.

- Optimizer: Adam, LR = $1 \times 10^{-3}$.

- Scheduler: CosineAnnealingLR for gradual decay.

- Metrics: tracks token-level accuracy and overall loss.

```
1   def train_text_model(model, loader, device, epochs=200, lr=1e-3):
2   """Train with CrossEntropyLoss (for multi-class distribution)"""
3   optimizer = torch.optim.Adam(model.parameters(), lr=lr)
4   scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, T_max=epochs)
5   loss_fn = nn.CrossEntropyLoss()
6
7   history = {'loss': [], 'acc': []}
8
9   for epoch in range(1, epochs + 1):
10      model.train()
11      total_loss = 0.0
12      total_correct = 0
13      total_samples = 0
14
15      pbar = tqdm(loader, desc=f"Epoch {epoch}/{epochs}", leave=False)
16      for x, y, lengths in pbar:
17          x, y = x.to(device), y.to(device)
18          lengths = lengths.to(device)
19
20          # Forward
21          logits, probs = model(x, lengths)
22          loss = loss_fn(logits, y)
23
24          # Backward
25          optimizer.zero_grad()
26          loss.backward()
27          torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)
28          optimizer.step()
29
30          # Metrics
31          total_loss += loss.item()
32          predicted = logits.argmax(dim=1)
33          total_correct += (predicted == y).sum().item()
34          total_samples += y.size(0)
35
36          pbar.set_postfix({
37              'loss': f'{loss.item():.3f}',
38              'acc': f'{total_correct/total_samples:.3f}'
39          })
40
41      scheduler.step()
42
43      avg_loss = total_loss / len(loader)
44      avg_acc = total_correct / total_samples
45      history['loss'].append(avg_loss)
46      history['acc'].append(avg_acc)
47
48      if epoch % 20 == 0 or epoch == epochs:
49          print(f"Epoch {epoch:03d} | Loss: {avg_loss:.4f} | Acc: {avg_acc:.4f}")
50
51   return history
```

## Generation (Autoregression)

- Starts from initial token.

- Repeatedly feeds generated tokens back to the model to predict the next one.

- Uses **temperature sampling** or **greedy decoding** to control randomness.

- Stops generation when a period (".") is predicted.

```
1   def generate_text_sequences(model, dataset, device, num_samples=10,
2                               max_len=50, temperature=0.5):
3   model.eval()
4
5   generated_texts = []
6
7   with torch.no_grad():
```

```python
    for sample_idx in range(num_samples):
        generated_tokens = []

        # Use greedy decoding (argmax) for deterministic generation
        start_char_idx = sample_idx % dataset.vocab_size
        generated_tokens.append(start_char_idx)

        for i in range(1, max_len):
            current_seq_input = torch.tensor([generated_tokens], device=device).float()
            length = torch.tensor([len(generated_tokens)], device=device)

            logits, probs = model(current_seq_input, length)

            # Use temperature sampling (low temp = more deterministic)
            if temperature < 0.1:
                # Greedy: always pick most likely
                token = logits.argmax(dim=1).item()
            else:
                # Sample from distribution with temperature
                logits_scaled = logits / temperature
                probs_scaled = F.softmax(logits_scaled, dim=1)
                token = torch.multinomial(probs_scaled[0], num_samples=1).item()

            generated_tokens.append(token)

            # Stop at period
            if dataset.idx2char.get(token, '') == '.':
                break

        # Decode to text
        text = dataset.decode_sequence(generated_tokens)
        generated_texts.append(text)

    return generated_texts, None
```

## Memorization Evaluation

The generated sentences are compared to the original training sentences. Here 5/20 generated texts match training, giving **25% memorization rate**.

The output is shown in Fig. 5



Figure 5: Output of question 4

# References

[1] Matthieu Courbariaux, Itay Hubara, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks: Training deep neural networks with weights and activations constrained to+ 1 or-1. *arXiv preprint arXiv:1602.02830*, 2016.

[2] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[3] Maxim Sorokin. Binary Cross-Entropy: Mathematical Insights and Python Implementation — vergotten. `https://medium.com/@vergotten/binary-cross-entropy-mathematical-insights-and-python-implementation-31e5a4df78f3`. [Accessed 21-10-2025].

[4] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103, 2008.

[5] Andrej Karpathy, Justin Johnson, and Li Fei-Fei. Visualizing and understanding recurrent networks. *arXiv preprint arXiv:1506.02078*, 2015.