



Xi'an Jiaotong-Liverpool University

西交利物浦大學

INT305 Assessment 2

Ruoyu Xu (2142804)

December 11, 2024

Contents

1	Q1: Introduction to CNNs	1
1.1	Convolution Kernel	1
1.2	Loss Function	2
2	Q2: Constructing CNN for MNIST	3
3	Q3: Model comparison	6
3.1	MobileNet	6
3.2	Vision Transformer(ViT)	7

1 Q1: Introduction to CNNs

Convolutional Neural Networks (CNNs), deep learning models specialized in processing grid-like data such as images, perform exceptionally well on computer vision tasks such as image classification and object detection.

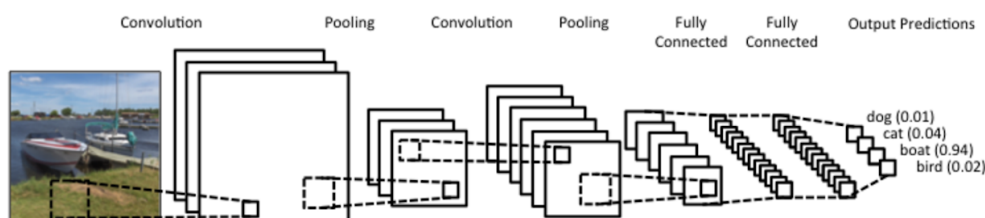


Figure 1: Process of CNN

Simply put, CNNs extract local features through convolution. Pooling reduces the dimensionality of the data while retaining key information. This process is repeated to extract more abstract and comprehensive features. Finally, the classification or prediction results are produced after fully connected layers process the extracted features [3].

1.1 Convolution Kernel

The convolution kernel is a core concept in convolution operations and even CNNs. It is a matrix of weights designed to extract different features.

The convolution operation can be written as:

$$O(i, j) = \sum_{p=0}^{m-1} \sum_{q=0}^{n-1} K(p, q) \cdot I(i + p, j + q)$$

- I represent the input matrix (e.g., image pixels).
- K represent the convolutional kernel (filter).
- O represent the output matrix (feature map).
- m, n be the dimensions of the kernel.
- i, j be the indices of the output feature map.

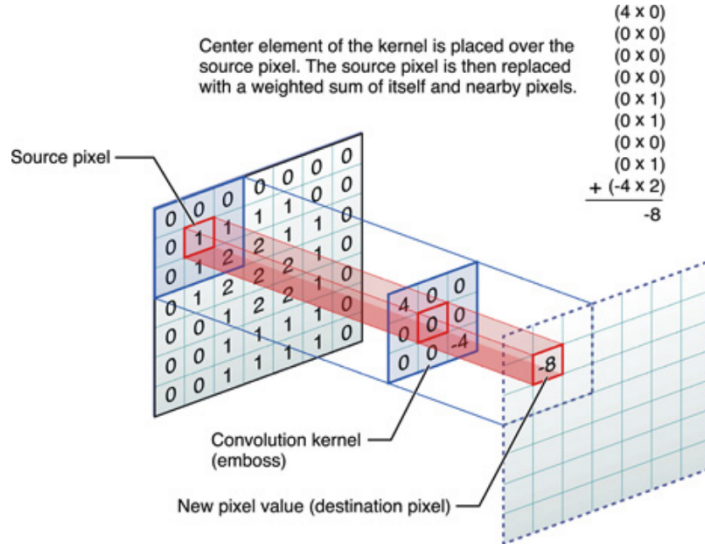


Figure 2: Operation of convolutional kernel

During runtime, the convolution kernel extracts local features from input data by sliding over it and performing element-wise multiplication and summation, producing a feature map [1]. A convolutional layer typically contains 32 or 64 convolutional kernels, each extracting distinct features, resulting in multiple feature maps. The weights of the kernels are learned from the data using backpropagation algorithms and are automatically adapted to suit the task requirements [4].

1.2 Loss Function

The loss function is a crucial component of model training, guiding the optimization algorithm to adjust model parameters by quantifying the difference between model predictions and true labels. The ultimate goal of the model is to minimize the loss function, thereby improving predictive accuracy.

The choice of the loss function influences learning dynamics. A suitable loss function can accelerate learning and penalize specific errors more effectively [6].

1. Regression Loss Functions

- **Mean Squared Error (MSE) Loss:** One of the most widely used loss functions for regression problems, MSE calculates the average squared differences between predicted and actual values. It works well with gradient-based optimization but is highly sensitive to outliers.

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- **Mean Absolute Error (MAE) Loss:** MAE is less sensitive to outliers than MSE, making it an alternative for robust regression tasks.

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Here, n represents the number of data point, y_i is the actual label, and \hat{y}_i is the predicted label.

2. Classification Loss Functions

- **Binary Cross-Entropy Loss:** Suitable for binary classification tasks, this function outputs predicted probabilities but can be sensitive to imbalanced datasets.

$$\text{Binary Cross-Entropy Loss} = -\frac{1}{n} \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

n represents the number of data point, y_i is the actual label(0 or 1), and \hat{y}_i is the predicted label.

- **Categorical Cross-Entropy Loss:** Designed for multi-class classification tasks, it measures the performance of a model outputting probability distributions across multiple classes.

$$\text{Categorical Cross-Entropy Loss} = -\sum_{i=1}^n \sum_{j=1}^k y_{ij} \log(\hat{y}_{ij})$$

n represents the number of data point, k is the number of classes, y_{ij} is the binary indicator if the class label j is the correct classification for data point i , and \hat{y}_{ij} is the predicted possibility for class j .

2 Q2: Constructing CNN for MNIST

To construct a basic CNN model for training on the MNIST dataset, we first performed data pre-processing, transforming the data into tensors and normalizing it. A simple CNN architecture

was then defined, starting with a convolutional layer to extract image features, followed by a pooling layer to down sample the feature map and reduce its dimensionality. The extracted features were subsequently mapped to 10 output classes through two fully connected layers.

For training, **CrossEntropyLoss** was selected as the loss function to measure the discrepancy between model predictions and true labels. **Stochastic Gradient Descent (SGD)** was employed as the optimizer, with a learning rate of 0.01 to adjust model parameters and enhance classification performance.

Table 1: Training and Testing Results for Basic CNN model

Epoch	Loss	Train Accuracy	Test Accuracy
1	0.6602	82.69%	91.58%
2	0.2304	93.30%	95.16%
3	0.1610	95.30%	94.91%

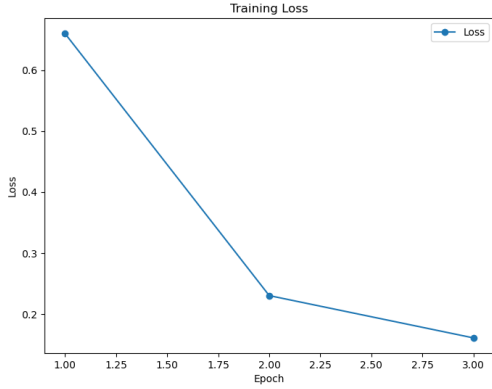


Figure 3: Loss function for CNN basic model

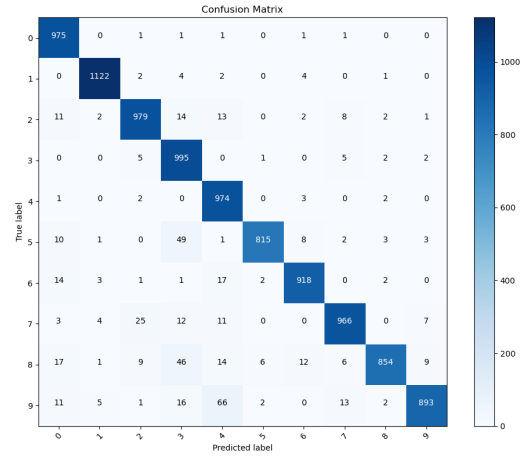


Figure 4: Confusion matrix for CNN basic model

This model achieved rapid convergence within three training cycles. The training loss decreased from 0.6602 to 0.1610, and both training and test accuracies improved significantly. Despite achieving a high test accuracy of 94.91%, the model exhibited certain limitations. Firstly, the lack of data augmentation meant only raw data was utilized for training, potentially compromising the model’s robustness. Secondly, the model’s architecture was simplistic, consisting of only one convolutional layer and one pooling layer, which limited its feature extraction capacity. Additionally, the basic optimizer and the small number of hidden units in the fully connected layers restricted the model’s learning capacity.

To address these drawbacks, the following optimizations were implemented:

1. **Data Augmentation:** Techniques such as random flipping and rotation were applied to augment data with an ultimate goal to enhance the model’s robustness.
2. **Model Enhancement:** Two convolutional layers were introduced to extract both low-level and high-level features, with pooling layers following each convolutional layer to reduce feature map size and mitigate overfitting.

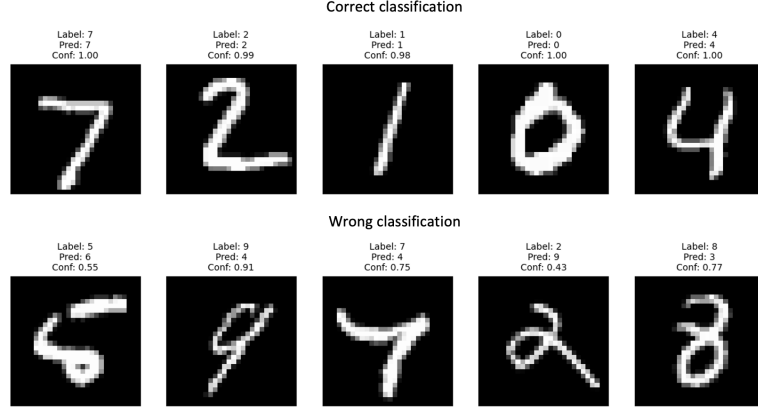


Figure 5: Samples of correct classification(CNN Basic)

3. **Optimizer Improvement:** The Adam optimizer was used to accelerate convergence by automatically adjusting the learning rate for each parameter, leveraging momentum and gradient squares. Additionally, the learning rate was reduced to 0.001 to enable fine-grained parameter adjustments, promoting stable convergence.

Table 2: Training and Testing Results for Optimized CNN model

Epoch	Loss	Train Accuracy	Test Accuracy
1	0.1593	95.12%	98.46%
2	0.0449	98.60%	98.65%
3	0.0311	99.05%	98.81%
4	0.0216	99.34%	98.55%
5	0.0167	99.47%	98.95%

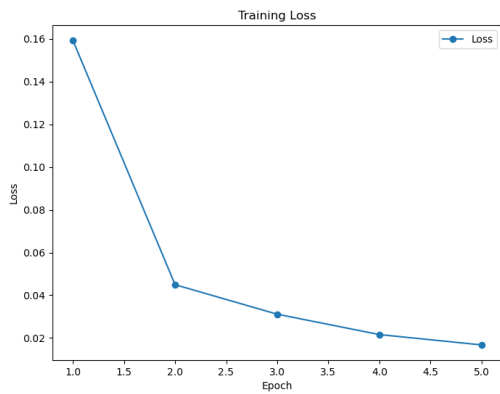


Figure 6: Loss function for CNN optimized model

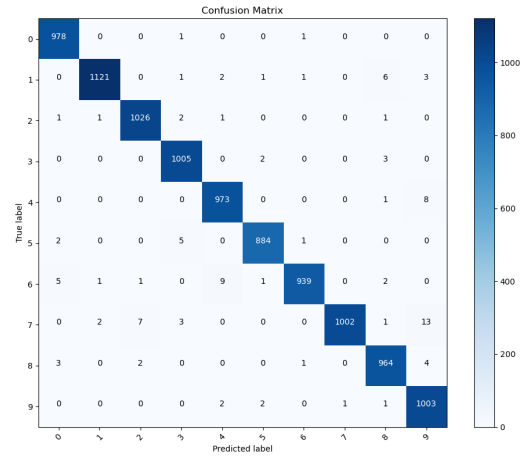


Figure 7: Confusion matrix for CNN optimized model

The improved model demonstrated superior performance, achieving a final test accuracy of 98.95%. The training loss decreased consistently from 0.1593 to 0.0167 in the final epoch, indicating effective learning and convergence. Both training and test accuracies remained closely

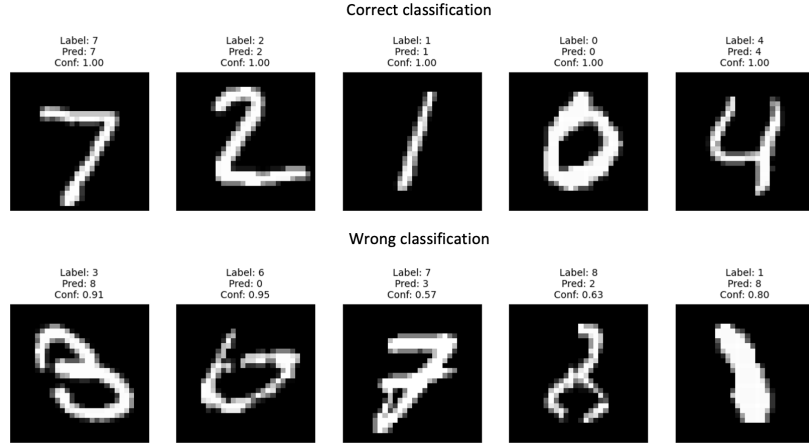


Figure 8: Samples of correct classification(CNN Optimized)

aligned throughout the training process, confirming the model’s robustness. These results underscore the effectiveness of the optimized CNN for the MNIST classification task.

3 Q3: Model comparison

Although traditional CNNs perform well on computer vision tasks, they do have notable limitations. To address these, various CNN variants have been developed, including MobileNet, which focuses on lightweight efficiency.

3.1 MobileNet

MobileNet employs a unique technique: depthwise separable convolution [2]. This technique splits the traditional convolution operation into two steps:

1. **Depthwise Convolution** Convolution is applied independently to each input channel, extracting spatial features within each channel. This step does not alter inter-channel relationships therefore reducing computational complexity.
2. **Pointwise Convolution** A 1×1 convolution kernel is used to fuse features across channels, capturing inter-channel relationships and producing a new channel representation.

By separating these two steps, MobileNet achieves feature extraction with significantly reduced computational cost compared to standard convolutions. Each layer in MobileNet consists of paired depthwise and pointwise convolutions, followed by a global pooling layer and a fully connected layer for classification.

For the MNIST dataset, MobileNet achieved a test accuracy of 98.88%, comparable to CNN’s 98.95%. MobileNet displayed slightly lower loss values during training, indicating better optimization. However, the computational efficiency of MobileNet is not evident in this

Table 3: Training and Testing Results for MobileNet

Epoch	Loss	Train Accuracy	Test Accuracy
1	0.1980	94.89%	97.84%
2	0.0413	98.77%	98.66%
3	0.0297	99.05%	98.41%
4	0.0241	99.25%	98.96%
5	0.0202	99.38%	98.88%

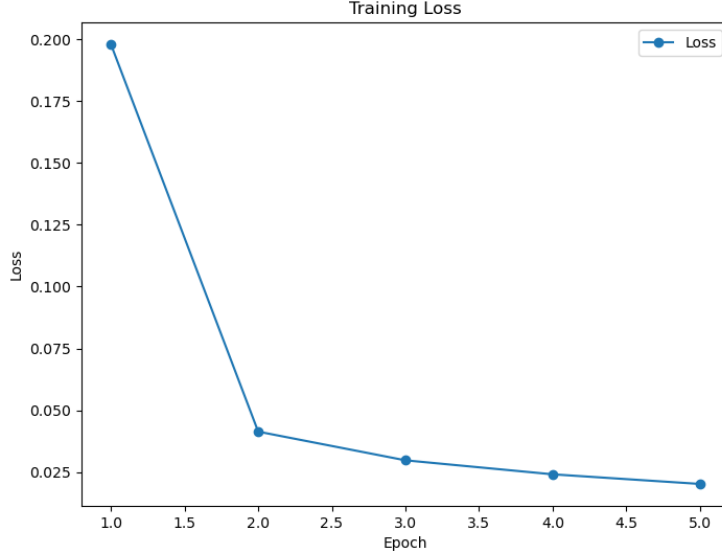


Figure 9: Loss function for MobileNet

simple task involving low-resolution grayscale images. While CNN is more suitable for MNIST due to its faster training time and marginally better performance, MobileNet’s architectural innovations are likely to excel in complex tasks or larger datasets [5].

3.2 Vision Transformer(ViT)

In 2020, the Vision Transformer (ViT) introduced a novel approach by combining attention mechanisms with convolutional architectures. ViT segments images into fixed-size patches, flattens them, and maps them into a fixed-dimensional vector space via linear projections [7]. These vectors, along with positional encodings, are fed into a Transformer architecture. Using self-attention mechanisms, ViT captures global image features by computing relationships between patches.

For MNIST, ViT achieved a test accuracy of 98.90%, comparable to CNN and MobileNet. However, it required more training iterations (10 epochs) to converge and lacked the spatial inductive bias of CNNs, resulting in weaker initial performance on small datasets. ViT is more suitable for high-dimensional, complex datasets where its global feature extraction capabilities can be fully leveraged.

Table 4: Training and Testing Results for ViT

Epoch	Loss	Train Accuracy	Test Accuracy
1	0.6988	76.61%	95.02%
2	0.2802	91.17%	96.82%
3	0.2161	93.16%	97.65%
4	0.1346	95.70%	98.60%
5	0.1121	96.52%	98.77%
6	0.1058	96.67%	98.74%
7	0.1014	96.91%	98.86%
8	0.0993	96.83%	98.91%
9	0.0961	97.04%	98.92%
10	0.0959	96.95%	98.90%

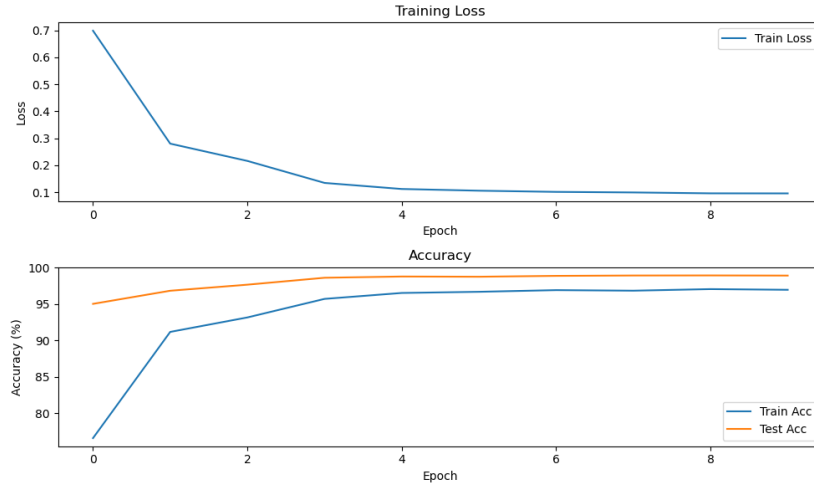


Figure 10: Loss function and accuracy for ViT

References

- [1] David Haussler et al. Convolution kernels on discrete structures. Technical report, Citeseer, 1999.
- [2] Andrew G. Howard et al. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv*, Apr. 2017. Cornell University.
- [3] Tim Kattenborn, Jens Leitloff, Frank Schiefer, and Stefan Hinz. Review on convolutional neural networks (cnn) in vegetation remote sensing. *ISPRS Journal of Photogrammetry and Remote Sensing*, 173:24–49, Mar. 2021.
- [4] David W Romero, Anna Kuzina, Erik J Bekkers, Jakub M Tomczak, and Mark Hoogenboom. Ckconv: Continuous kernel convolution for sequential data. *arXiv preprint arXiv:2102.02611*, 2021.
- [5] Anumol C S. Advancements in cnn architectures for computer vision: A comprehensive review. In *2023 Annual International Conference on Emerging Research Areas: International Conference on Intelligent Systems (AICERA/ICIS)*, pages 1–7, 2023.

- [6] Qi Wang, Yue Ma, Kun Zhao, and Yingjie Tian. A comprehensive survey of loss functions in machine learning. *Annals of Data Science*, pages 1–26, 2020.
- [7] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 558–567, 2021.

Appendix

Code for model construction

Optimized CNN

```
1
2  # 2. Define Optimized CNN Model
3  class SimpleCNN(nn.Module):
4      def __init__(self):
5          super(SimpleCNN, self).__init__()
6          self.conv1 = nn.Conv2d(1, 32, kernel_size=3, stride=1, padding=1)
7          self.pool = nn.MaxPool2d(kernel_size=2, stride=2)
8          self.conv2 = nn.Conv2d(32, 64, kernel_size=3, stride=1, padding=1)
9          self.fc1 = nn.Linear(64 * 7 * 7, 128)
10         self.fc2 = nn.Linear(128, 10)
11
12     def forward(self, x):
13         x = self.pool(torch.relu(self.conv1(x)))
14         x = self.pool(torch.relu(self.conv2(x)))
15         x = x.view(-1, 64 * 7 * 7)
16         x = torch.relu(self.fc1(x))
17         x = self.fc2(x)
18         return x
19
20 model = SimpleCNN().to(device)
21
22 # 3. Define Loss Function and Optimizer
23 criterion = nn.CrossEntropyLoss()
24 optimizer = optim.Adam(model.parameters(), lr=0.001)
```

MobileNet

```
1
2 # Define Depth-wise separable convolution
3 class DepthwiseSeparableConv(nn.Module):
4     def __init__(self, in_channels, out_channels, stride=1):
5         super(DepthwiseSeparableConv, self).__init__()
6         self.depthwise = nn.Conv2d(in_channels, in_channels,
7                                     ↪ kernel_size=3, stride=stride, padding=1, groups=in_channels)
8         self.pointwise = nn.Conv2d(in_channels, out_channels,
9                                     ↪ kernel_size=1)
10
11     def forward(self, x):
12         x = self.depthwise(x)
13         x = self.pointwise(x)
14         return x
15
16 # Define MobileNet Model
17 class MobileNet(nn.Module):
18     def __init__(self, num_classes=10):
19         super(MobileNet, self).__init__()
20         self.features = nn.Sequential(
21             nn.Conv2d(1, 32, kernel_size=3, stride=1, padding=1), #
22             nn.BatchNorm2d(32),
23             nn.ReLU(inplace=True),
24
25             DepthwiseSeparableConv(32, 64),
26             nn.BatchNorm2d(64),
27             nn.ReLU(inplace=True),
28
29             DepthwiseSeparableConv(64, 128, stride=2),
30             nn.BatchNorm2d(128),
31             nn.ReLU(inplace=True),
32
33             DepthwiseSeparableConv(128, 128),
34             nn.BatchNorm2d(128),
35             nn.ReLU(inplace=True),
36
37             DepthwiseSeparableConv(128, 256, stride=2),
38             nn.BatchNorm2d(256),
39             nn.ReLU(inplace=True),
40
41             DepthwiseSeparableConv(256, 256),
42             nn.BatchNorm2d(256),
43             nn.ReLU(inplace=True),
44
45             nn.AdaptiveAvgPool2d(1) #
```

```

44         )
45         self.classifier = nn.Linear(256, num_classes)
46
47     def forward(self, x):
48         x = self.features(x)
49         x = x.view(x.size(0), -1)  #
50         x = self.classifier(x)
51         return x
52 model = MobileNet().to(device)
53
54 # Define Loss Function and Optimizer
55 criterion = nn.CrossEntropyLoss()
56 optimizer = optim.Adam(model.parameters(), lr=0.001)

```

ViT

```

1
2 # Define transformer
3 class TransformerEncoderBlock(nn.Module):
4     def __init__(self, dim, num_heads, mlp_ratio=4.0, dropout=0.1):
5         super(TransformerEncoderBlock, self).__init__()
6         self.norm1 = nn.LayerNorm(dim)
7         self.attn = nn.MultiheadAttention(dim, num_heads, dropout=dropout,
8             ↪ batch_first=True)
9         self.norm2 = nn.LayerNorm(dim)
10
11         mlp_hidden_dim = int(dim * mlp_ratio)
12         self.mlp = nn.Sequential(
13             nn.Linear(dim, mlp_hidden_dim),
14             nn.GELU(),
15             nn.Dropout(dropout),
16             nn.Linear(mlp_hidden_dim, dim),
17             nn.Dropout(dropout)
18         )
19
20     def forward(self, x):
21         # Self-attention
22         x = x + self.attn(self.norm1(x), self.norm1(x), self.norm1(x))[0]
23         # MLP
24         x = x + self.mlp(self.norm2(x))
25         return x
26
27 # Define the ViT Model
28 class VisionTransformer(nn.Module):

```

```

28 def __init__(self, image_size=28, patch_size=7, dim=128, depth=4,
    ↪ num_heads=4, mlp_ratio=4.0, num_classes=10, dropout=0.1):
29     super(VisionTransformer, self).__init__()
30     assert image_size % patch_size == 0, "image_size must be divisible
    ↪ by patch_size"
31     self.num_patches = (image_size // patch_size) ** 2
32     self.patch_size = patch_size
33     patch_dim = patch_size * patch_size # MNIST1
34     self.dim = dim
35
36     # patchdim
37     self.patch_embed = nn.Linear(patch_dim, dim)
38
39     # CLS token
40     self.cls_token = nn.Parameter(torch.randn(1, 1, dim))
41     #
42     self.pos_embed = nn.Parameter(torch.randn(1, self.num_patches + 1,
    ↪ dim))
43     self.dropout = nn.Dropout(dropout)
44
45     # Transformer Encoder Block
46     self.blocks = nn.ModuleList([
47         TransformerEncoderBlock(dim, num_heads, mlp_ratio,
    ↪ dropout=dropout) for _ in range(depth)
48     ])
49
50     self.norm = nn.LayerNorm(dim)
51     self.head = nn.Linear(dim, num_classes)
52
53 def forward(self, x):
54     # x: [B,1,28,28]
55     B, C, H, W = x.shape
56     # patch
57     patches = x.unfold(2, self.patch_size, self.patch_size).unfold(3,
    ↪ self.patch_size, self.patch_size)
58     # patches: [B, 1, H/ps, W/ps, ps, ps]
59     patches = patches.contiguous().view(B, -1,
    ↪ self.patch_size*self.patch_size)
60     x = self.patch_embed(patches) # [B, num_patches, dim]
61
62     # CLS token
63     cls_tokens = self.cls_token.expand(B, -1, -1) # [B,1,dim]
64     x = torch.cat((cls_tokens, x), dim=1) # [B, num_patches+1, dim]
65
66     #
67     x = x + self.pos_embed[:, :x.size(1), :]
68     x = self.dropout(x)
69

```

```

70     # Transformer Encoder
71     for blk in self.blocks:
72         x = blk(x)
73
74     x = self.norm(x)
75     # CLS token
76     cls_out = x[:,0]
77     logits = self.head(cls_out)
78     return logits
79
80 net = VisionTransformer().to(device)
81
82 criterion = nn.CrossEntropyLoss()
83 optimizer = optim.AdamW(net.parameters(), lr=0.001, weight_decay=1e-4)
84 scheduler = optim.lr_scheduler.StepLR(optimizer, step_size=3, gamma=0.1)

```